

# ML Methods

E. Le Pennec



Fall 2021

# Outline

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

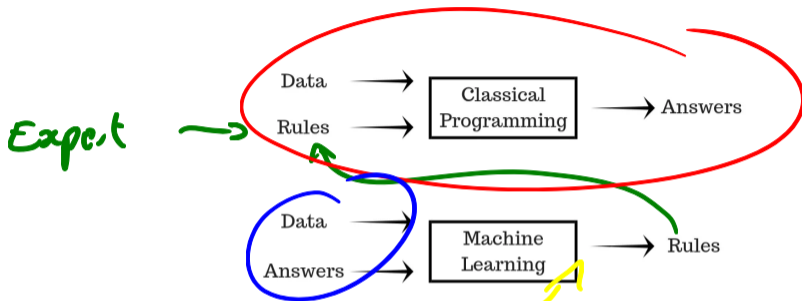
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References





A screenshot of the Google News homepage. The page features a search bar at the top, navigation tabs for "Headlines", "Local", "For You", and "U.S.", and a sidebar with various news categories. The main content area displays "Top Stories" with several news items, including one about Sarah Huckabee Sanders and another about a Time Magazine cover.





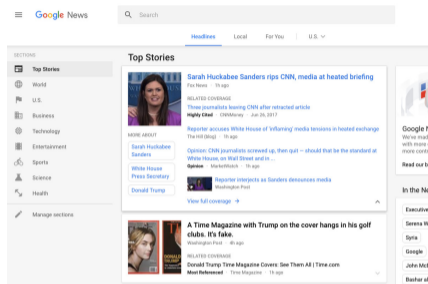
A definition by Tom Mitchell (<http://www.cs.cmu.edu/~tom/>)

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.



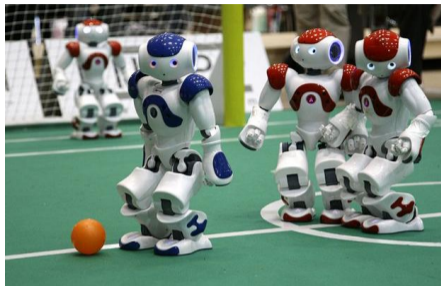
A detection algorithm:

- **Task:** say if an object is present or not in the image
- **Performance:** number of errors
- **Experience:** set of previously seen labeled images



## An article clustering algorithm:

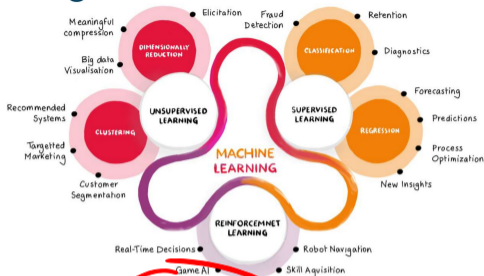
- **Task:** group articles corresponding to the same news
- **Performance:** quality of the clusters
- **Experience:** set of articles



A robot endowed with a set of sensors playing football:

- **Task:** play football
- **Performance:** score evolution
- **Experience:**
  - past games
  - current environment and action outcome,

# Three Kinds of Learning



### Unsupervised Learning

- **Task:** Clustering/DR
- **Performance:** Quality
- **Experience:** Raw dataset (No Ground Truth)

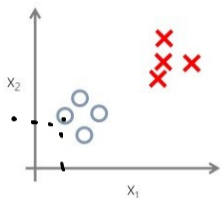
### Supervised Learning

- **Task:** Prediction/Classification
- **Performance:** Average error
- **Experience:** Good Predictions (Ground Truth)

### Reinforcement Learning

- **Task:** Action
- **Performance:** Total reward
- **Experience:** Reward from env. (Interact. with env.)

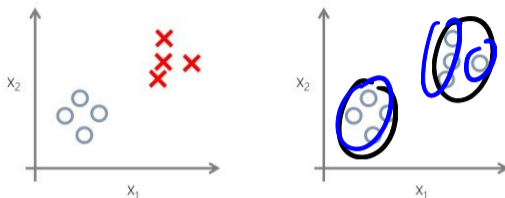
● **Timing:** Offline/Batch (learning from past data) vs Online (continuous learning)



$$f: X \mapsto \beta(x) \in Y$$

## Supervised Learning (Imitation)

- **Goal:** Learn a function  $f$  predicting a variable  $Y$  from an individual  $\underline{X}$ .
- **Data:** Learning set with labeled examples  $(\underline{X}_i, Y_i)$
- **Assumption:** Future data behaves as past data!
- **Predicting is not explaining!**



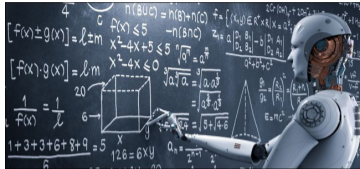
## Supervised Learning (Imitation)

- **Goal:** Learn a function  $f$  predicting a variable  $Y$  from an individual  $\underline{X}$ .
- **Data:** Learning set with labeled examples  $(\underline{X}_i, Y_i)$
- **Assumption:** Future data behaves as past data!
- **Predicting is not explaining!**

## Unsupervised Learning (Structure Discovery)

- **Goal:** Discover a structure within a set of individuals  $(\underline{X}_i)$ .
- **Data:** Learning set with unlabeled examples  $(\underline{X}_i)$
- Unsupervised learning is not a well-posed setting...





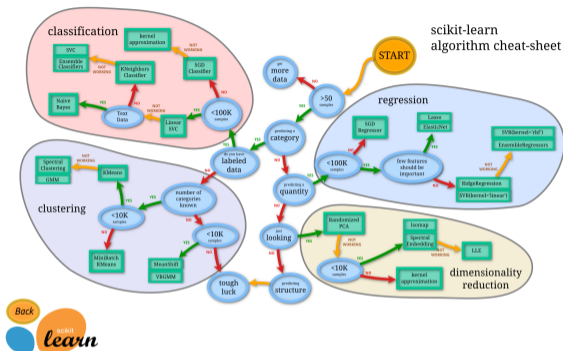
## Machine Can

- Forecast (Prediction using the past)
- Detect some changes
- Memorize/Reproduce
- Take a decision very quickly
- Learn from huge dataset
- Optimize a single task
- Replace/Help some humans

## Machine Cannot

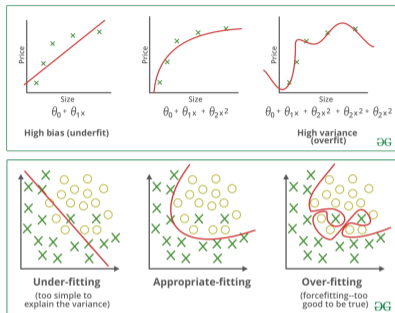
- Predict something never seen before
- Detect any new behaviour
- Create something brand new
- Understand the world
- Get smart really fast
- Go beyond their task
- Kill all humans

● Some progresses but still very far from the *singularity*...



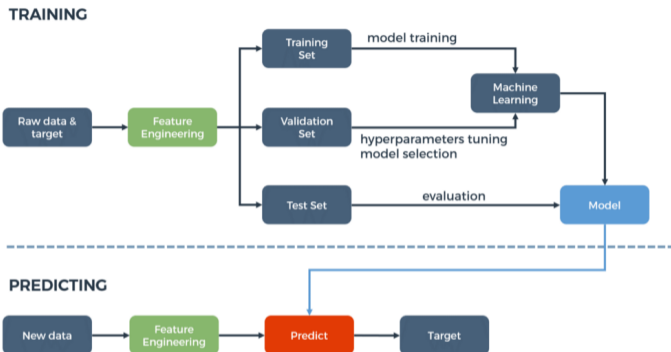
## Machine Learning Methods

- Huge catalog of methods,
- Need to define the performance,
- Numerous tricks: feature design, hyperparameter selection. . .



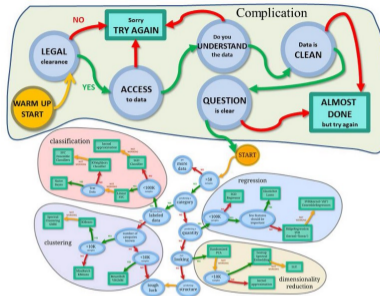
## Finding the Right Complexity

- What is best?
  - A simple model that is stable but false? (*oversimplification*)
  - A very complex model that could be correct but is unstable? (*conspiracy theory*)
- Neither of them: tradeoff that depends on the dataset.



## Learning pipeline

- Test and compare models.
- Deployment pipeline is different!



## Main DS difficulties

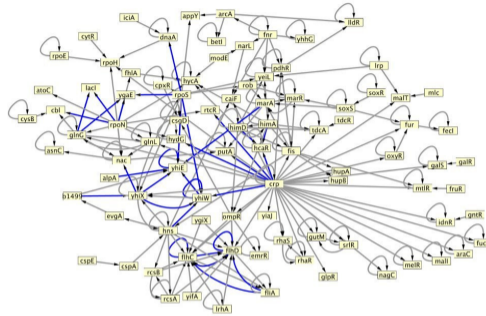
- Figuring out the problem,
- Formalizing it,
- Storing and accessing the data,
- Deploying the solution,
- Not (always) the Machine Learning part!

- 1 Introduction
  - Machine Learning
  - **Motivation**
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

0	0	0	0	0
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

## Reading a ZIP code on an envelop

- **Task:** give a number from an image.
- **Data:**  $X$  = image /  $Y$  = corresponding number.
- **Performance measure:** error rate.



## Predicting protein interaction

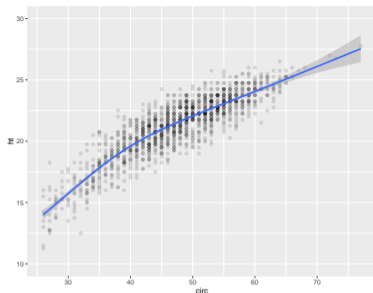
- **Task:** Predict (unknown) interactions between proteins.
- **Data:**  $X$  = pair of proteins /  $Y$  = existence or no of interaction.
- **Performance measure:** error rate.
- Numerous similar questions in bio(informatics): genomic,...





## Face detection

- **Task:** Detect the position of faces in an image
- Different setting?
- Reformulation as a supervised learning problem.
- **Goal:** Detect the presence of faces at several positions and scales.
- **Data:**  $X$  = sub image /  $Y$  = presence or no of a face...
- **Performance measure:** error rate.
- Lots of detections in an image: post processing required...
- **Performance measure:** box precision.



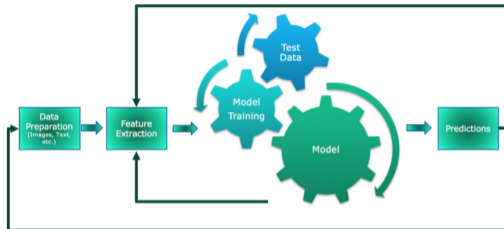
## Height estimation

- Simple (and classical) dataset.
- **Task:** predict the height from circumference.
- **Data:**  $X$  = circumference /
- $Y$  = height.
- **Performance measure:** means squared error.

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 **A Practical View**
  - **Method or Models**
  - **Interpretability**
  - **Metric Choice**
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

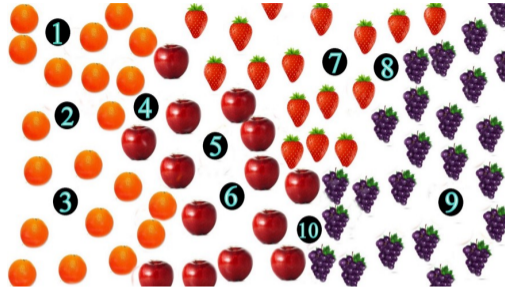
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 **A Practical View**
  - **Method or Models**
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

## A Standard Machine Learning Pipeline



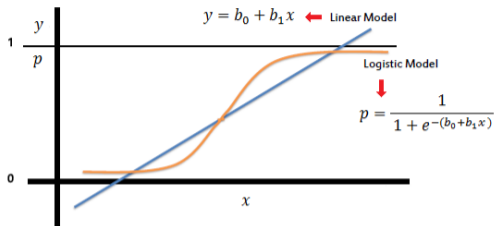
## A Learning Method

- Formula/Algorithm allowing to make predictions
- Algorithm allowing to chose this formula/algorithm
- Data preprocessing (cleansing, coding...)
- Optimization criterion for the choice!



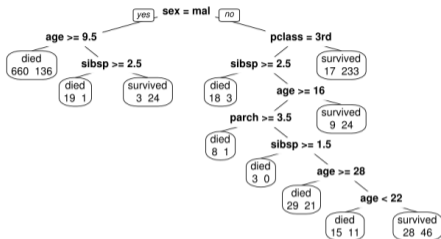
## Similarity

- Imitate the answer to give by mixing answers to similar questions (**k nearest neighbors**)
- Require to search for those similar questions for each request
- Not always very efficient but fast to build (less to use...)
- Easy to understand and rather stable



## Linear Method

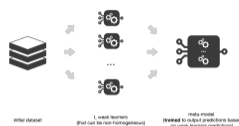
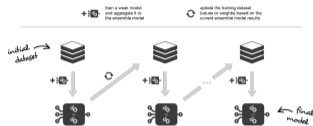
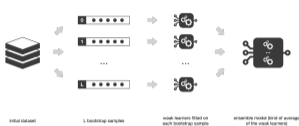
- Simple formula:  $a_0 + a_1 X^{(1)} + \dots + a_d X^{(d)}$
- Imitate the answer to give (**linear regression**) or a transformation of the conditional probability of the category (**logistic regression**)
- Numerous variations on the parameter optimization (**penalization, SVM,...**)
- Pretty efficient and fast to build
- Easy to understand and rather stable



## Tree

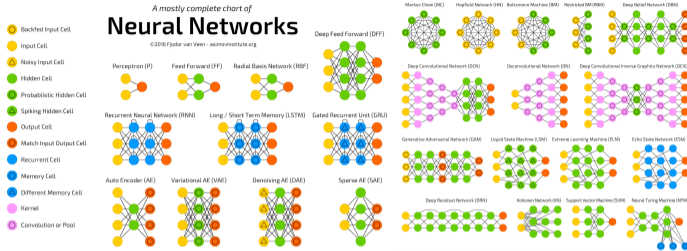
- Construction of a **decision tree**
- Impossible to really optimize but good tree can be obtained
- Not always very efficient but very quick to build
- Very easy to understand but not really stable





## Ensemble Methods

- Strategy:
  - **Bagging:** construction of variations in parallel and averaging (**random forest**)
  - **Boosting:** construction of sequential improvements (**XGBoost, Lightgbm**)
  - **Stacking:** Use of a first set of predictors as features
- Very good performance for structured data but quite slow to build
- Stable but hard to understand



## Deep Learning

- Chain of simple formulae (**Neural Network**)
- Joint optimization
- Very good performance for unstructured data but slow to build
- Mildly stable and very hard to understand

Method	Performance	Training Speed	Inf. Speed	Stability	Interpretability
Similarity	-	∅	-	+	+
Linear	+	++	++	++	+
Tree	-	++	++	-	++
Ensemble	++	-	+	++	-
Deep	++	-	-	-	-

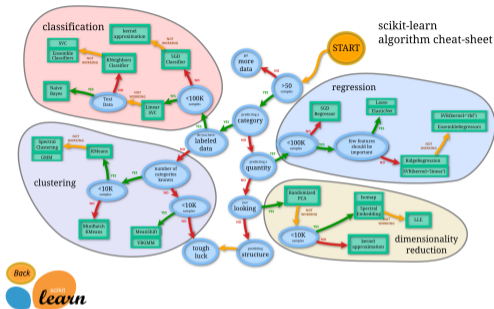
## Take Away Message

- No unanimously best solution
- Impossible to guess which method is going to be the best!
- A good practice is to always try a linear method as well as an ensemble one for structured data or deep one for unstructured data



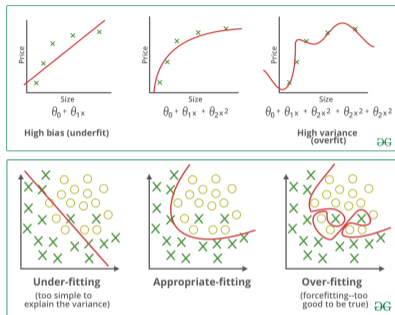
## Preprocessing

- Art of creating sophisticated representations of initial data
- Key for good performances
- Examples: individual transformation, variable combination, category (and text) coding. . .
- **Important part of the learning method**



## ML Methods

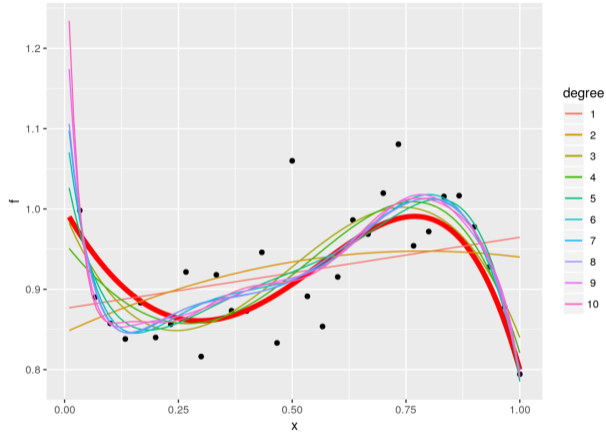
- Huge catalog of methods,
- Need to define the performance,
- Need to represent well the data
- Need to choose the **best** method yielding a good model



## Finding the Right Complexity

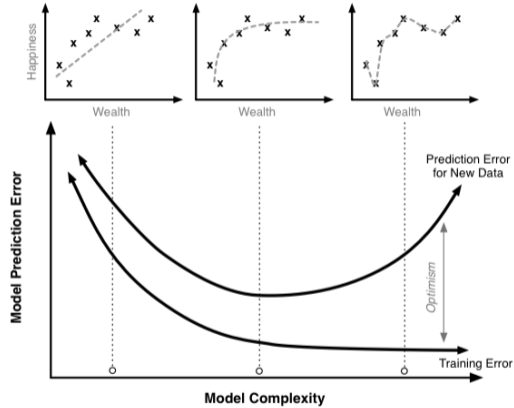
- What is best?
  - A simple model that is stable but false? (*oversimplification*)
  - A very complex model that could be correct but is unstable? (*conspiracy theory*)
- Neither of them: tradeoff that depends on the dataset.

# Which method to use?

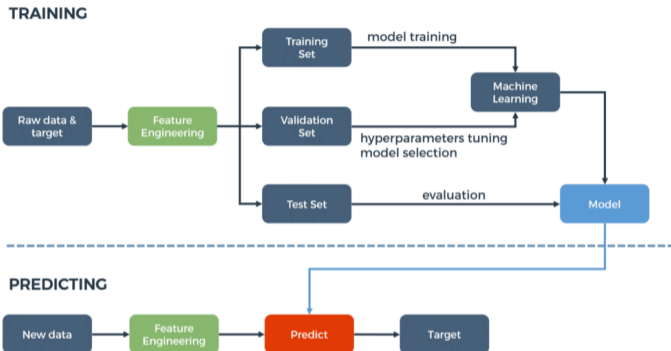


- Competition between several polynomial models.

# Over-fitting





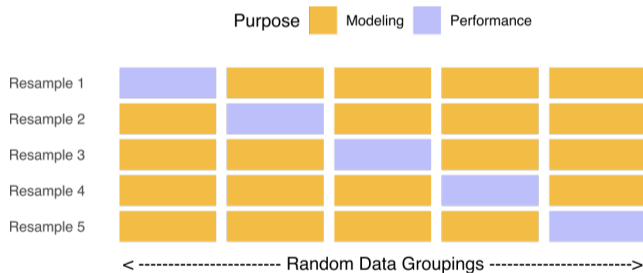


## Learning pipeline

- Test and compare models.
- Deployment pipeline is different!

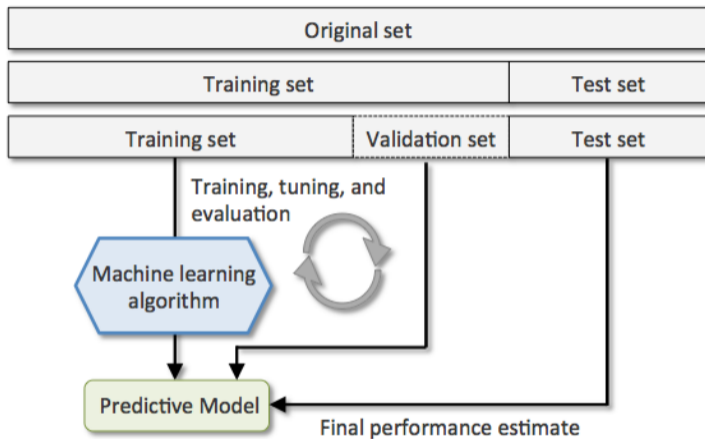


- Train a model and check its quality on different pieces of the data.



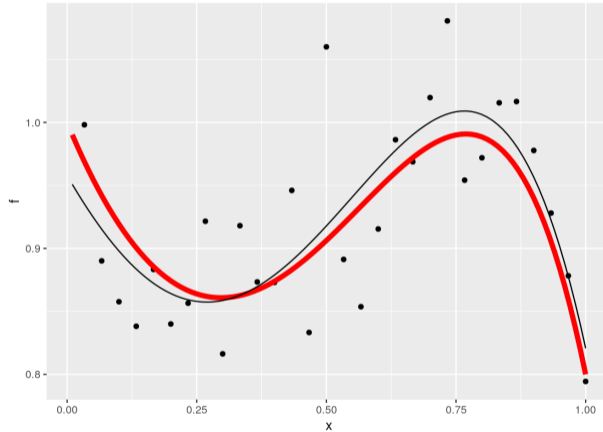
- Check the quality of a method by repeating the previous approach.
- **Beware:** a different predictor is learnt for each split.

# The Full Cross Validation Scheme



- Most important part of machine learning.
- Automatic choice of model possible by (intelligent ?) exploration...

# Best Polynomial



- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - **Interpretability**
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

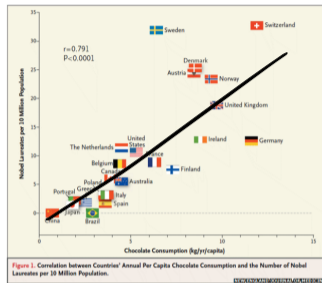


Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

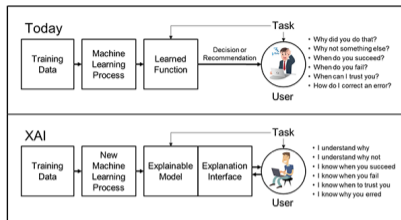
Nombre de prix Nobel par dix millions d'habitants en fonction de la consommation nationale de chocolat en kilogrammes par personne et par an.

Image : Franz H. Messerli, The New England Journal of Medicine 367(14) (2012), p. 1562-1564

- Simple formula setting:

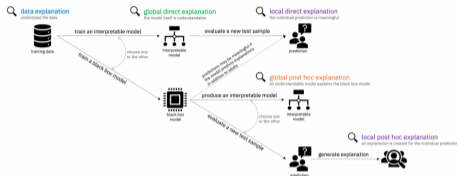
$$Y \simeq f(X) = a_0 + a_1X^{(1)} + a_2X^{(2)} + \dots + a_dX^{(d)}$$

- Beware of the interpretation!
- Everything being equal... Correlation is not causality...



## Intepretability or Explainability

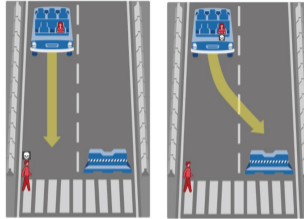
- Interpretability: possibility to give a causal aspect to the formula.
- Explainability: possibility to find the variables having an effect on the decision and their effect.
- Explainability is much easier than interpretability.
- Transparency (on the datasets, the criterion optimized and the algorithms) yields already a lot of informations.



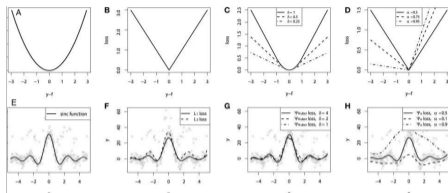
- Data Explaination.
- Use of explainable methods (linear?).
- Use of black box methods:
  - Global explanation (variable importance)
  - Local explanation (linear approximationn, alternative scenario. . .)
- Causality very hard to access without a real experimental plan with interventions!



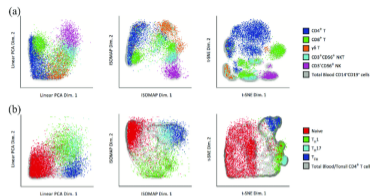
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 **A Practical View**
  - Method or Models
  - Interpretability
  - **Metric Choice**
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



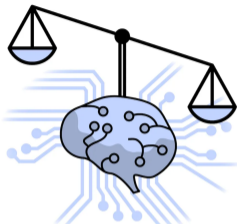
- Quality metric has a strong impact on the solution.
- Implicite encoding rather than an explicit one!
- Often simplified criterion in the optimization part.
- More involved criterion can be used in evaluation.



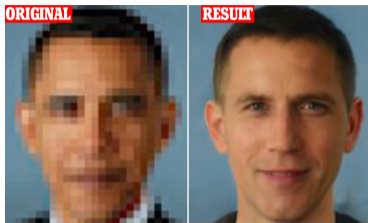
- Measure of the cost of not being perfect!
- Criterion used to *optimize* the predictor and/or *evaluate* its interest.
- Classical metrics: quadratic error, zero/one error.
- Many other possible choices, ideally encoding domain expertise (asymmetry...)
- The criterion can be different between optimization and evaluation because of computation requirements.
- Very important factor (too) often neglected.



- Measure the quality of the result!
- Dimension Reduction / Representation: reconstruction quality, relationship preservation. . .
- Clustering: measure of intra-group proximity and inter-group difference?
- Very subjective criterion!
- Hard to define the right distances especially for discrete variables.
- In practice, quality often evaluated by the a posteriori interest.



- Very hard to specify criterion.
- No consensus on its definition:
  - faithful reproduction of the reality?
  - correction of its bias?
- Current approaches through constraints in the optimization.
- A posteriori verification unavoidable!

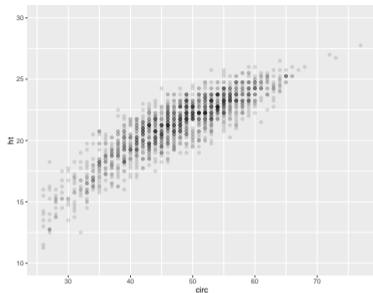


- Central assumption: representativity of the data!
- Optimization made in this setting.
- Possible training data bias:
  - selection bias in the data
  - population evolution
  - (historical) bias in the targets
- Correction possible at least up to a certain point for the 2 first cases if one is aware of the situation.

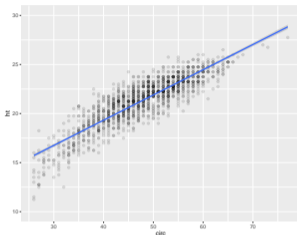
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 **A Better Point of View**
  - **The Example of Univariate Linear Regression**
  - **Supervised Learning**
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 **A Better Point of View**
  - **The Example of Univariate Linear Regression**
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References





- Simple (and classical) dataset.
- Goal: predict the height from circumference
- $X = \text{circ} = \text{circumference}$ .
- $Y = \text{ht} = \text{height}$ .



## Linear Model

- Parametric model:

$$f_{\beta}(\text{circ}) = \beta^{(1)} + \beta^{(2)}\text{circ}$$

- How to choose  $\beta = (\beta^{(1)}, \beta^{(2)})$ ?

## Methodology

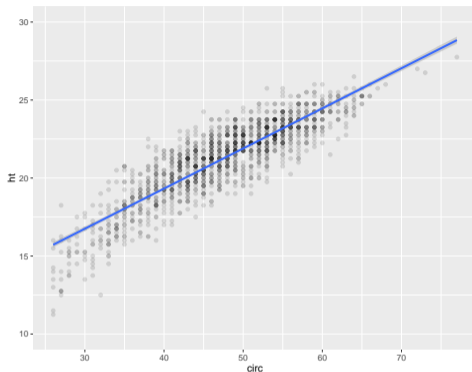
- Natural goodness criterion:

$$\begin{aligned}\sum_{i=1}^n |Y_i - f_{\beta}(X_i)|^2 &= \sum_{i=1}^n |ht_i - f_{\beta}(\text{circ}_i)|^2 \\ &= \sum_{i=1}^n |ht_i - (\beta^{(1)} + \beta^{(2)} \text{circ}_i)|^2\end{aligned}$$

- Choice of  $\beta$  that minimizes this criterion!

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^2}{\operatorname{argmin}} \sum_{i=1}^n |h_i - (\beta^{(1)} + \beta^{(2)} \text{circ}_i)|^2$$

- Easy minimization with an explicit solution!



## Prediction

- Linear prediction for the height:

$$\widehat{ht} = f_{\widehat{\beta}}(\text{circ}) = \widehat{\beta}^{(1)} + \widehat{\beta}^{(2)} \text{circ}$$

## Linear Regression

- **Statistical model:**  $(\text{circ}_i, \text{ht}_i)$  **i.i.d.** with the same law than a generic  $(\text{circ}, \text{ht})$ .

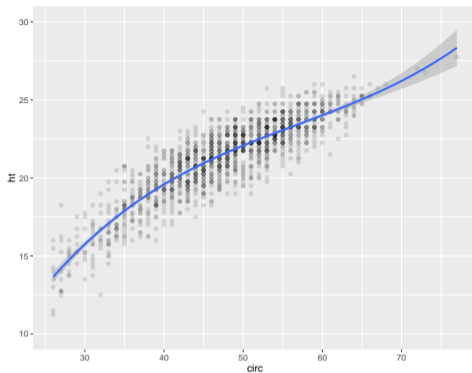
- **Performance criterion:** Look for  $f$  with a **small average error**

$$\mathbb{E} \left[ |\text{ht} - f(\text{circ})|^2 \right]$$

- **Empirical criterion:** Replace the unknown law by its **empirical** counterpart

$$\frac{1}{n} \sum_{i=1}^n |\text{ht}_i - f(\text{circ}_i)|^2$$

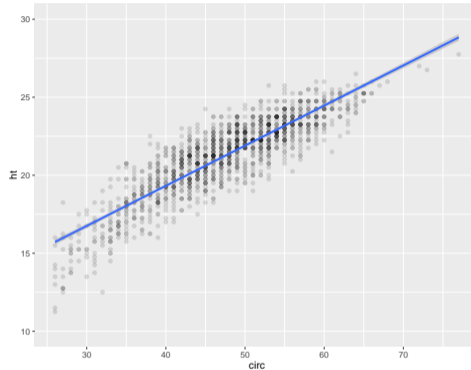
- **Predictor model:** As the minimum over all function is 0 (if all the  $\text{circ}_i$  are different), **restrict** to the linear functions  $f(\text{circ}) = \beta^{(1)} + \beta^{(2)}\text{circ}$  to avoid over-fitting.
- **Model fitting:** Explicit formula here.
- This model can be **too simple!**



## Polynomial Model

- Polynomial model:  $f_{\beta}(\text{circ}) = \sum_{l=1}^p \beta^{(l)} \text{circ}^{l-1}$
- Linear in  $\beta$ !
- Easy least squares estimation for any degree!

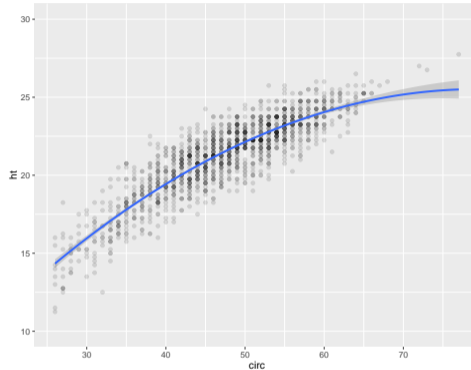
# Which Degree?



## Models

- Increasing degree = increasing complexity and better fit on the data

# Which Degree?

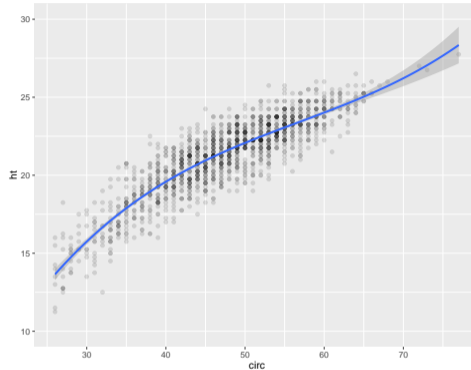


## Models

- Increasing degree = increasing complexity and better fit on the data



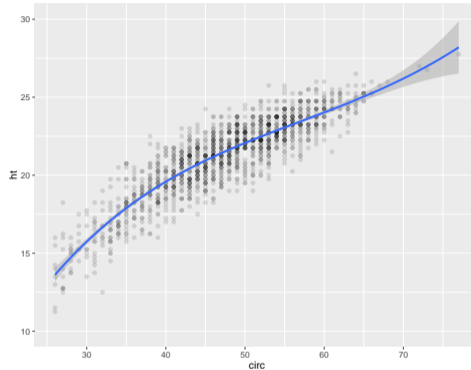
# Which Degree?



## Models

- Increasing degree = increasing complexity and better fit on the data

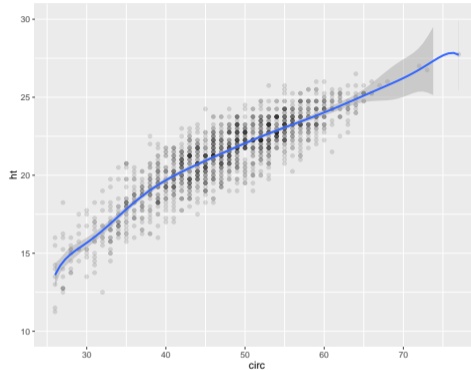
# Which Degree?



## Models

- Increasing degree = increasing complexity and better fit on the data

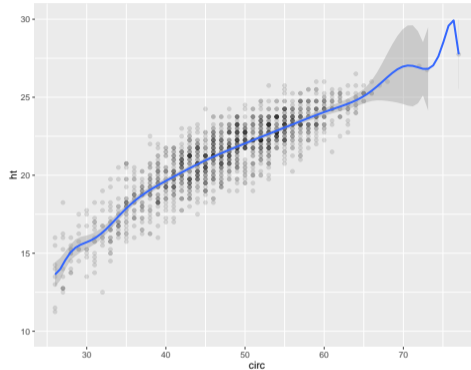
# Which Degree?



## Models

- Increasing degree = increasing complexity and better fit on the data

# Which Degree?



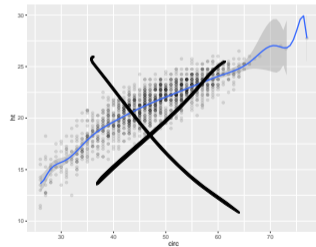
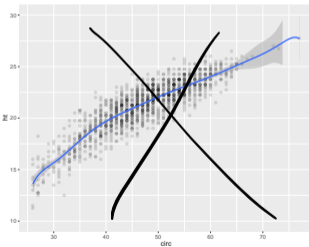
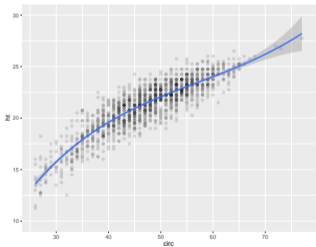
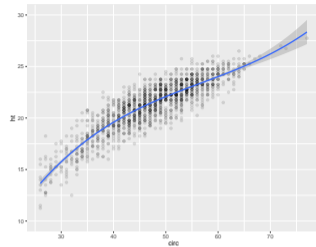
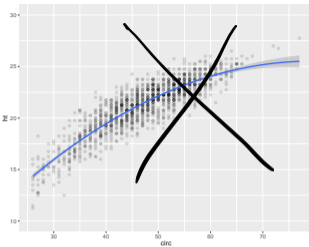
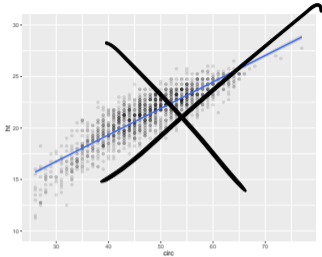
## Models

- Increasing degree = increasing complexity and better fit on the data

# Which Degree?

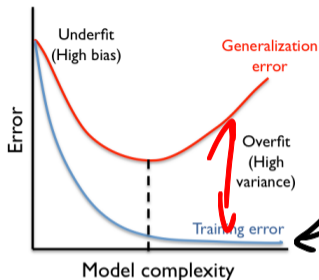
$$RMSE = \sqrt{\frac{1}{n} \sum (y_i - \hat{f}(x_i))^2}$$

A Better Point of View



## Best Degree?

- How to choose among those solution?



RMSE sur les données  
utilisées

## Error behavior

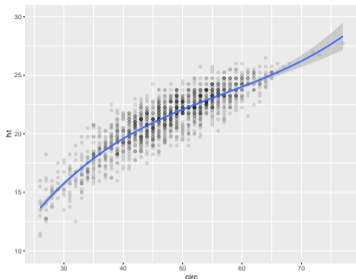
- Empirical risk (error made on the training set) decays when the complexity of the model increases.
- Quite different behavior when the error is computed on new observations (true risk / generalization error).
- Overfit for complex models: parameters learned are too specific to the learning set!
- General situation! (Think of polynomial fit...)
- Need to use an other criterion than the training error!

## Two directions

- **How to estimate** the generalization error in a different way?
- Find a way to **correct** the empirical error?

## Two Approaches

- **Cross validation:** Estimate the error on a different dataset:
  - Very efficient (and almost always used in practice!)
  - Need more data for the error computation.
- **Penalization approach:** Correct the optimism of the empirical error:
  - Require to find the correction (penalty).



## Questions

- How to build a model?
- How to fit a model to the data?
- How to assess its quality?
- How to select a model among a collection?
- How to guaranty the quality of the selected model?



- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 **A Better Point of View**
  - The Example of Univariate Linear Regression
  - **Supervised Learning**
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

## Supervised Learning Framework

- Input measurement  $\underline{X} \in \mathcal{X}$
- Output measurement  $Y \in \mathcal{Y}$ .
- $(\underline{X}, Y) \sim \mathbb{P}$  with  $\mathbb{P}$  unknown.
- **Training data** :  $\mathcal{D}_n = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$  (i.i.d.  $\sim \mathbb{P}$ )
- Often
  - $\underline{X} \in \mathbb{R}^d$  and  $Y \in \{-1, 1\}$  (classification)  $\rightarrow$  *discrimination*
  - or  $\underline{X} \in \mathbb{R}^d$  and  $Y \in \mathbb{R}$  (regression).
- A **predictor** is a function in  $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \text{ meas.}\}$

## Goal

- Construct a **good** predictor  $\hat{f}$  from the training data.
- Need to specify the meaning of good.
- Classification and regression are almost the **same** problem!

## Loss function for a generic predictor

- **Loss function:**  $\ell(Y, f(\underline{X}))$  measures the goodness of the prediction of  $Y$  by  $f(\underline{X})$
- Examples:
  - Prediction loss:  $\ell(Y, f(\underline{X})) = \mathbf{1}_{Y \neq f(\underline{X})}$
  - Quadratic loss:  $\ell(Y, f(\underline{X})) = |Y - f(\underline{X})|^2$

## Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim \mathbb{P}} [\ell(Y, f(\underline{X}))]$$

- Examples:
  - Prediction loss:  $\mathbb{E} [\ell(Y, f(\underline{X}))] = \mathbb{P} (Y \neq f(\underline{X}))$
  - Quadratic loss:  $\mathbb{E} [\ell(Y, f(\underline{X}))] = \mathbb{E} [|Y - f(\underline{X})|^2]$

- **Beware:** As  $\hat{f}$  depends on  $\mathcal{D}_n$ ,  $\mathcal{R}(\hat{f})$  is a random variable!

- The best solution  $f^*$  (which is independent of  $\mathcal{D}_n$ ) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E} [\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[ \mathbb{E}_{Y|\underline{X}} [\ell(Y, f(\underline{X}))] \right]$$

## Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

**Issue:** Solution requires to **know**  $\mathbb{E}[Y|\underline{X}]$  for all values of  $\underline{X}$ !

$$E_{y|x}[\ell(y, \beta(x))] = \ell(1, \beta(x)) \times P(Y=1|x) \\ + \ell(-1, \beta(x)) \times P(Y=-1|x)$$

$$= \begin{array}{l|l} P(Y=-1|x) & \text{si } \beta(x)=1 \\ P(Y=1|x) & \text{si } \beta(x)=-1 \end{array}$$

## Machine Learning

- Learn a rule to construct a **predictor**  $\hat{f} \in \mathcal{F}$  from the training data  $\mathcal{D}_n$  s.t. **the risk**  $\mathcal{R}(\hat{f})$  is **small on average** or with high probability with respect to  $\mathcal{D}_n$ .
- In practice, the rule should be an algorithm!

## Canonical example: Empirical Risk Minimizer

- One restricts  $f$  to a subset of functions  $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the empirical loss

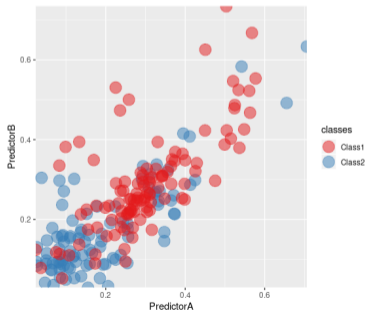
$$\hat{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Examples:
  - Linear regression
  - Linear discrimination with

$$\mathcal{S} = \{\underline{x} \mapsto \operatorname{sign}\{\underline{x}^\top \beta + \beta^{(0)}\} / \beta \in \mathbb{R}^d, \beta^{(0)} \in \mathbb{R}\}$$

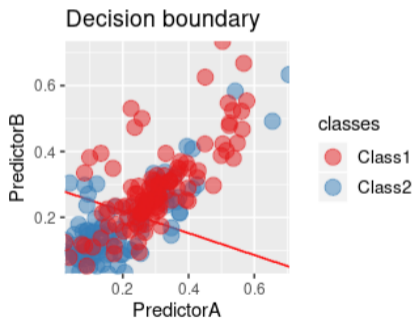
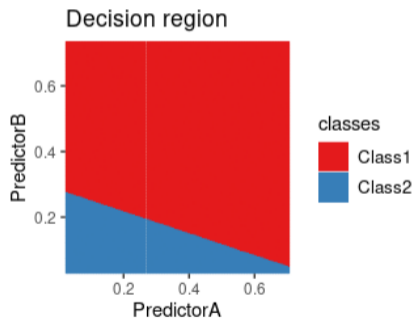
## Synthetic Dataset

- Two features/covariates.
- Two classes.
- Dataset from *Applied Predictive Modeling*, M. Kuhn and K. Johnson, Springer
- Numerical experiments with R and the caret package.



# Example: Linear Discrimination

Logistic

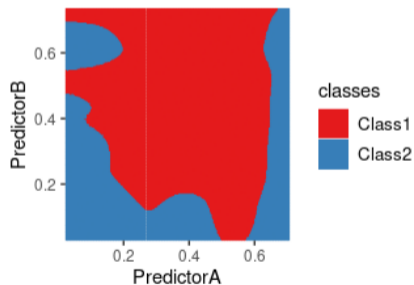




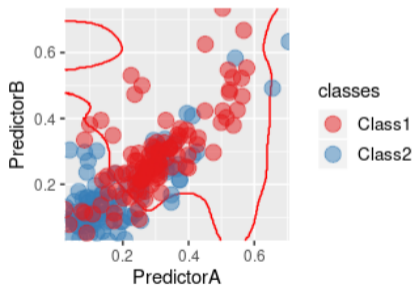
# Example: More Complex Model

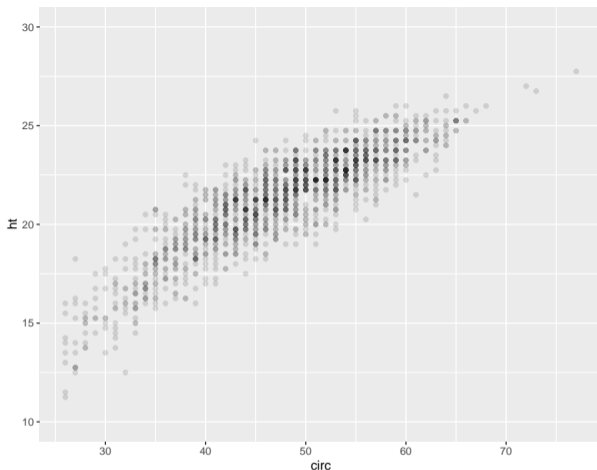
Naive Bayes with kernel density estimates

Decision region



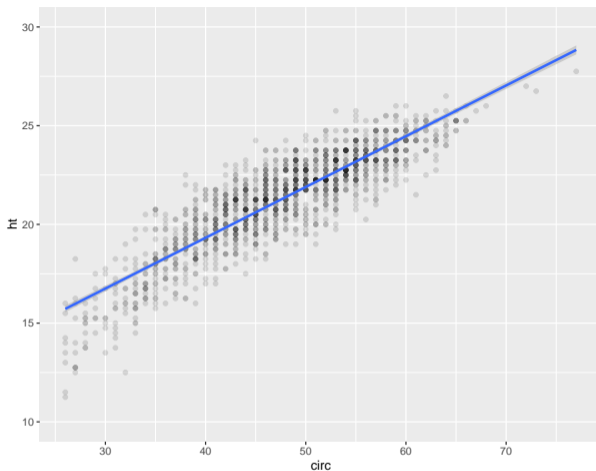
Decision boundary





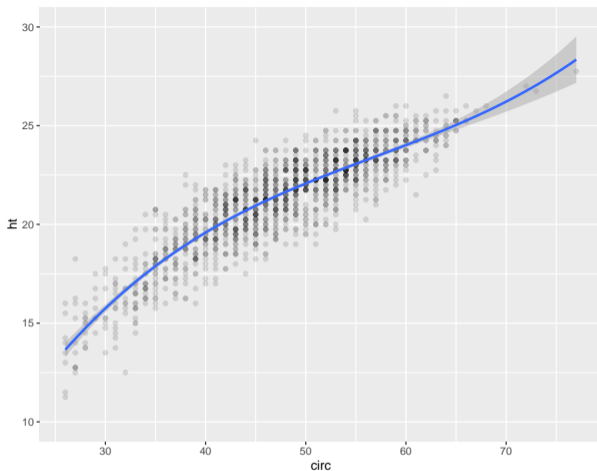
## Dataset - P.A. Cornillon

- Real dataset of 1429 eucalyptus obtained by P.A. Cornillon:
  - $\underline{X}$ : circumference /  $\underline{Y}$ : height



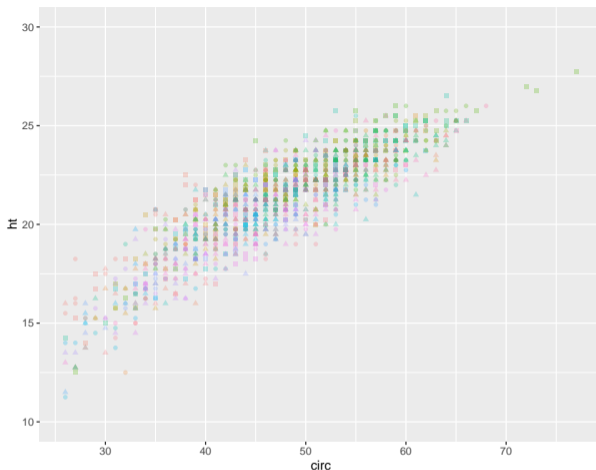
## Dataset - P.A. Cornillon

- Real dataset of 1429 eucalyptus obtained by P.A. Cornillon:
  - $\underline{X}$ : circumference /  $\underline{Y}$ : height



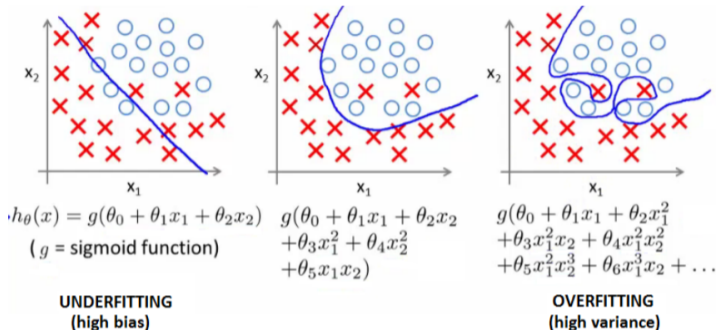
## Dataset - P.A. Cornillon

- Real dataset of 1429 eucalyptus obtained by P.A. Cornillon:
  - $\underline{X}$ : circumference /  $\underline{Y}$ : height



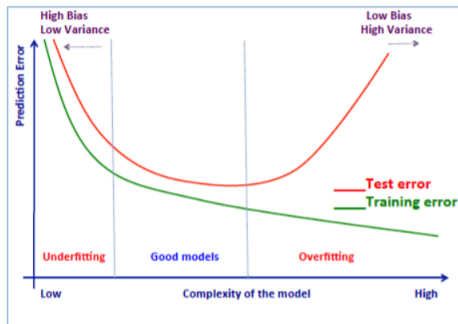
## Dataset - P.A. Cornillon

- Real dataset of 1429 eucalyptus obtained by P.A. Cornillon:
  - $\underline{X}$ : circumference, block, clone / Y: height



## Model Complexity Dilemma

- What is best a simple or a complex model?
- Too simple to be good? Too complex to be learned?

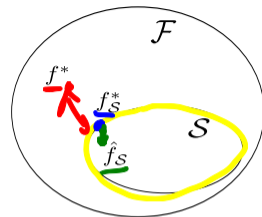


## Under-fitting / Over-fitting

- **Under-fitting:** simple model are too simple.
- **Over-fitting:** complex model are too specific to the training set.

# Bias-Variance Dilemma

- General setting:
  - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
  - Best solution:  $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
  - Class  $\mathcal{S} \subset \mathcal{F}$  of functions
  - Ideal target in  $\mathcal{S}$ :  $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
  - Estimate in  $\mathcal{S}$ :  $\hat{f}_S$  obtained with some procedure



## Approximation error and estimation error (Bias/Variance)

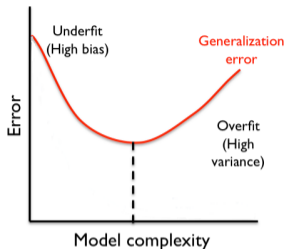
$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Approx. error can be large if the model  $\mathcal{S}$  is not suitable.
- Estimation error can be large if the model is complex.

## Agnostic approach

- No assumption (so far) on the law of  $(\underline{X}, Y)$ .





- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

**Bias-variance trade-off**  $\iff$  avoid **overfitting** and **underfitting**

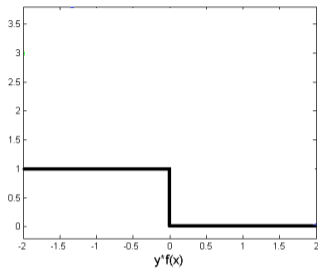
- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.

## Statistical Learning Analysis

- Error decomposition:

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Bound on the approximation term: approximation theory.
  - Probabilistic bound on the estimation term: probability theory!
  - **Goal: Agnostic bounds**, i.e. bounds that do not require assumptions on  $\mathbb{P}$ !  
(Statistical Learning?)
- 
- Often need mild assumptions on  $\mathbb{P}$ ... (Nonparametric Statistics?)



## Empirical Risk Minimizer

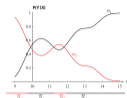
$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

*Handwritten notes:*  $\ell \in \{-1, 1\}$  (with a checkmark) and  $\epsilon \in \{0, 1\}$  (with an arrow pointing to the loss function).

- Classification loss:  $\ell^{0/1}(y, f(\underline{x})) = \mathbf{1}_{y \neq f(\underline{x})}$
- Not convex and not smooth!

# Probabilistic Point of View

## Ideal Solution and Estimation



- The best solution  $f^*$  (which is independent of  $\mathcal{D}_n$ ) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E} [\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[ \mathbb{E}_{Y|\underline{X}} [\ell(Y, f(\underline{x}))] \right]$$

### Bayes Predictor (explicit solution)

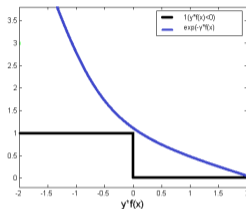
In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Issue:** Solution requires to **know**  $\mathbb{E}[Y|\underline{X}]$  for all values of  $\underline{X}$ !
- **Solution:** Replace it by an estimate.

# Optimization Point of View

## Loss Convexification



$$\beta: \mathcal{X} \rightarrow \mathbb{R}$$
$$\hat{\ell}: \text{using } (\beta(x))$$

### Minimizer of the risk

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

- **Issue:** Classification loss is not convex or smooth.
- **Solution:** Replace it by a convex majorant.

# Probabilistic and Optimization Framework

How to find a good function  $f$  with a *small* risk

$$\mathcal{R}(f) = \mathbb{E} [\ell(Y, f(\underline{X}))] \quad ?$$

**Canonical approach:**  $\hat{f}_S = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i))$

## Problems

- How to choose  $\mathcal{S}$ ?
- How to compute the minimization?

## A Probabilistic Point of View

**Solution:** For  $\underline{X}$ , estimate  $Y|\underline{X}$  plug this estimate in the Bayes classifier:  
**(Generalized) Linear Models, Kernel methods,  $k$ -nn, Naive Bayes, Tree, Bagging...**

## An Optimization Point of View

**Solution:** If necessary replace the loss  $\ell$  by an upper bound  $\ell'$  and minimize the empirical loss: **SVR, SVM, Neural Network, Tree, Boosting...**

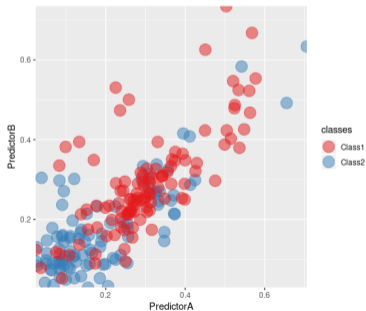
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice**
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 **Error Estimation and Method Choice**
  - **Cross Validation**
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



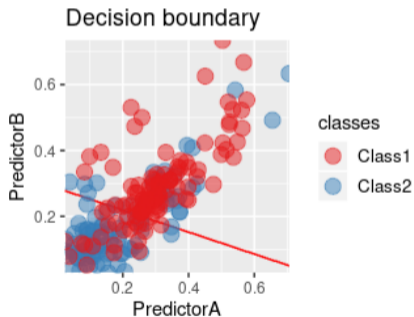
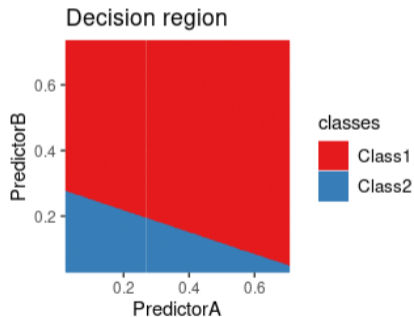
## Synthetic Dataset

- Two features/covariates.
- Two classes.
- Dataset from *Applied Predictive Modeling*, M. Kuhn and K. Johnson, Springer
- Numerical experiments with R and the caret package.



# Example: Linear Discrimination

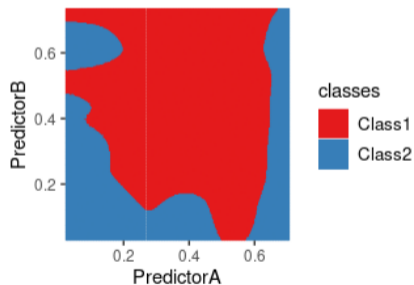
Logistic



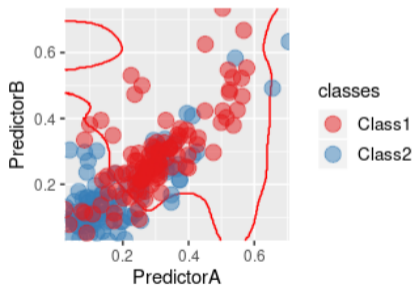
# Example: More Complex Model

Naive Bayes with kernel density estimates

Decision region

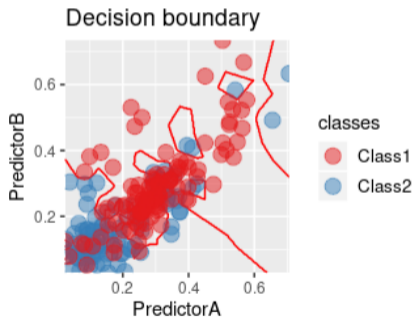
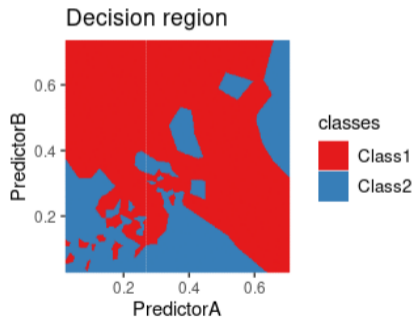


Decision boundary



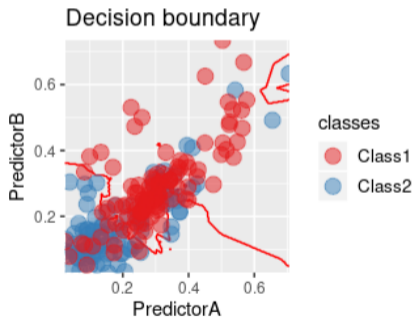
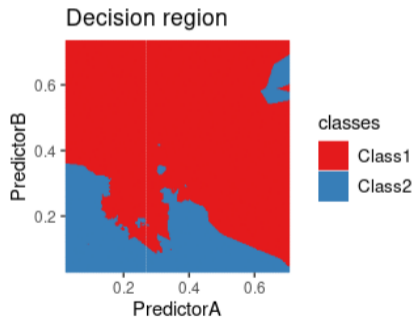
# Example: KNN

k-NN with  $k=1$



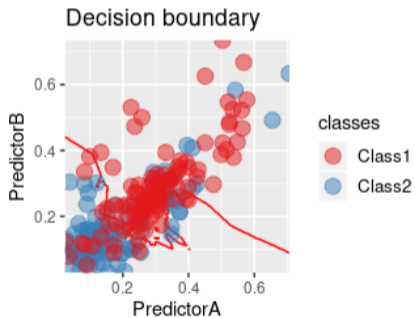
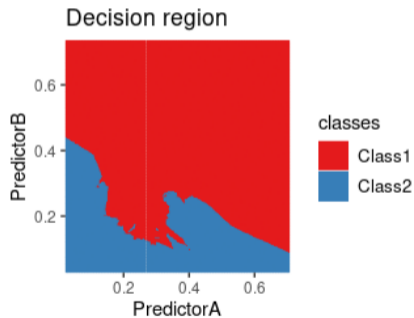
# Example: KNN

k-NN with  $k=5$



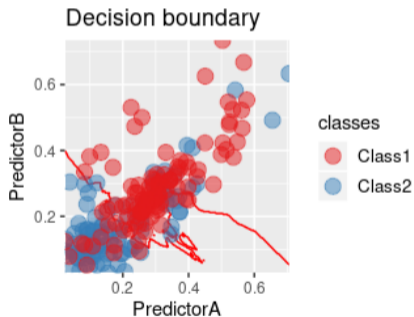
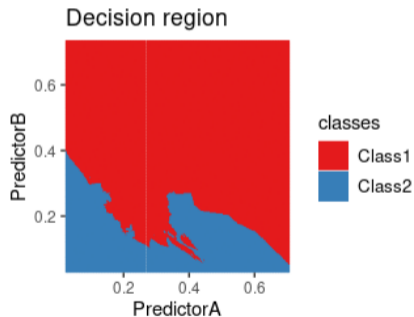
# Example: KNN

k-NN with  $k=9$



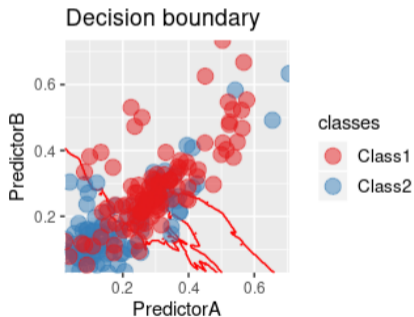
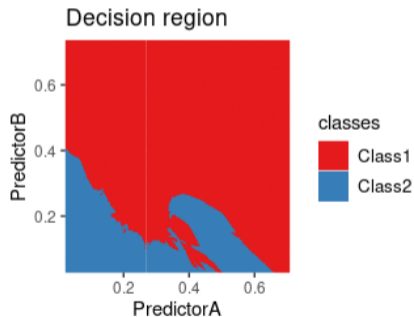
# Example: KNN

k-NN with  $k=13$



# Example: KNN

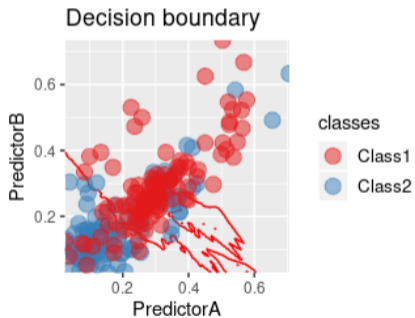
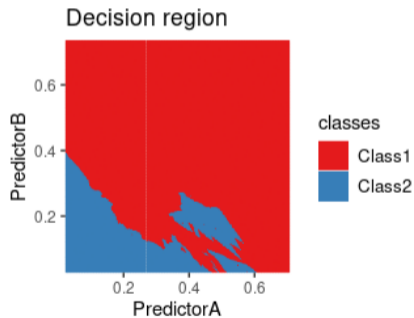
k-NN with  $k=17$





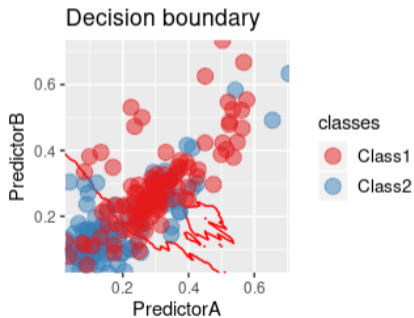
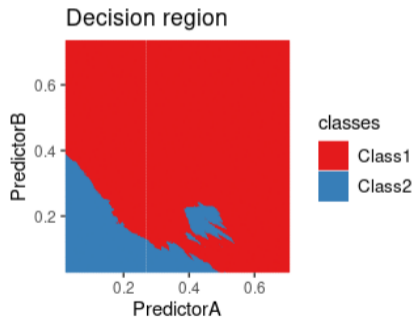
# Example: KNN

k-NN with  $k=21$



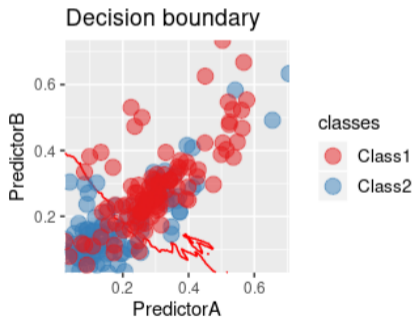
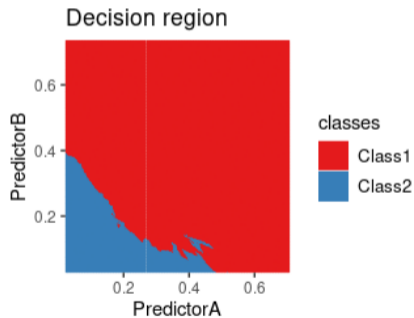
# Example: KNN

k-NN with  $k=25$



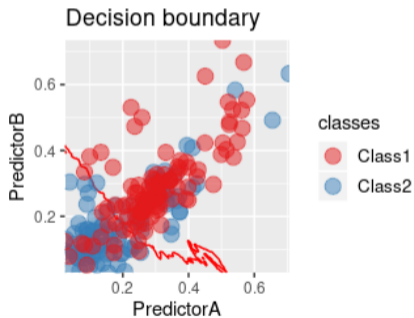
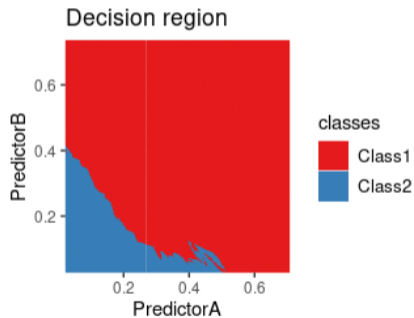
# Example: KNN

k-NN with  $k=29$



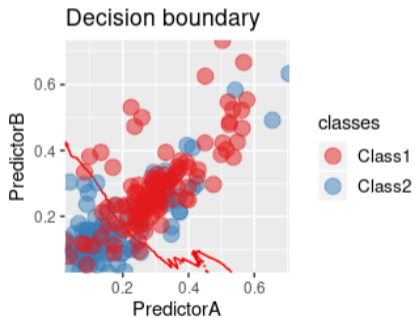
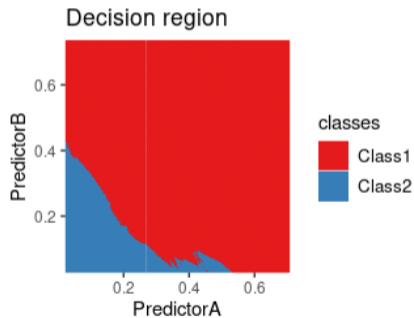
# Example: KNN

k-NN with  $k=33$



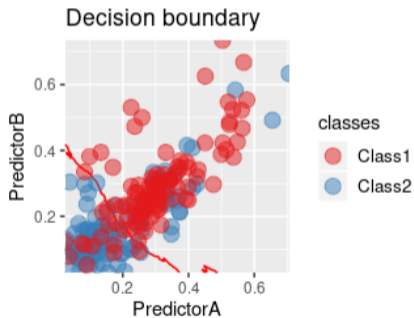
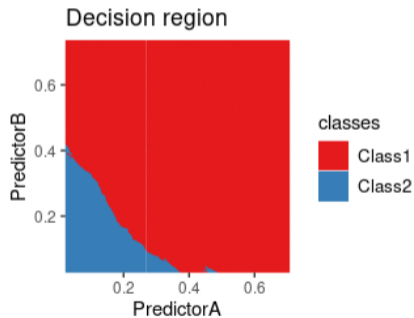
# Example: KNN

k-NN with  $k=37$



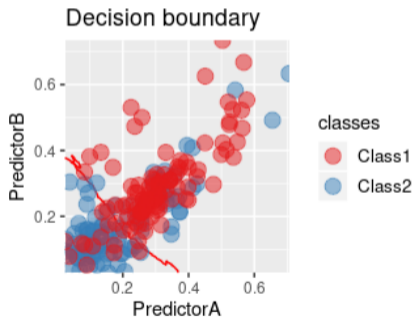
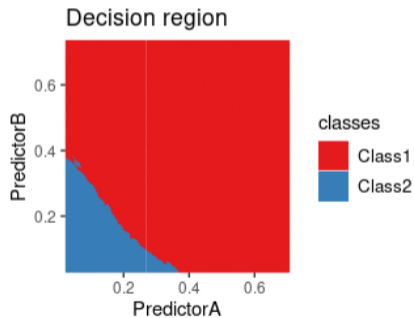
# Example: KNN

k-NN with k=45



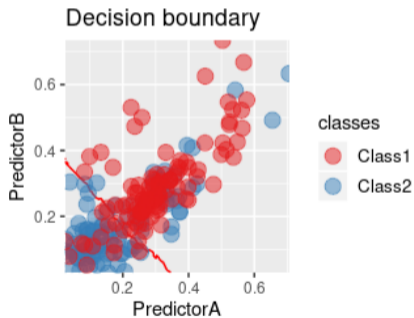
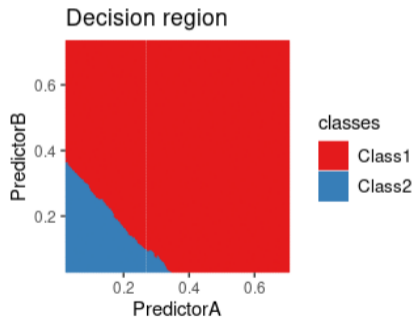
# Example: KNN

k-NN with  $k=53$



# Example: KNN

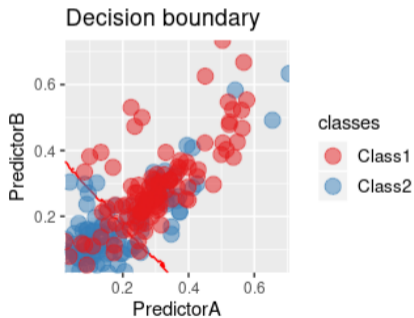
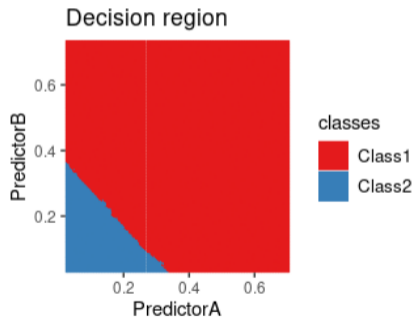
k-NN with  $k=61$





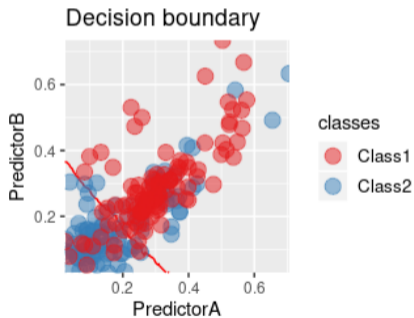
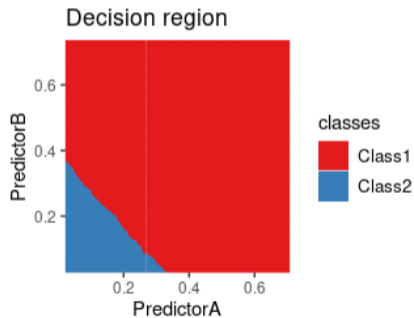
# Example: KNN

k-NN with k=69



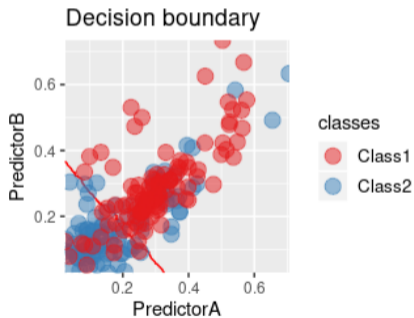
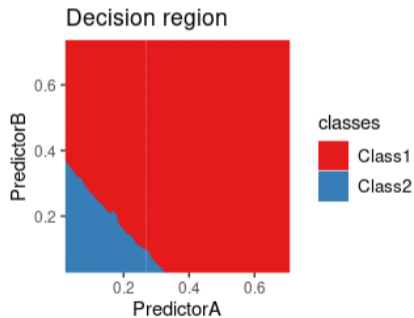
# Example: KNN

k-NN with  $k=77$



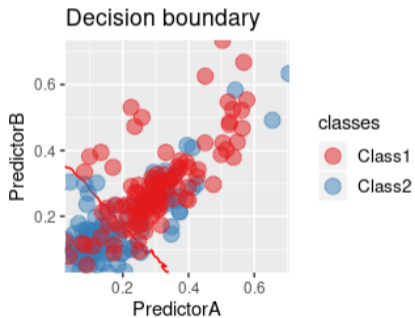
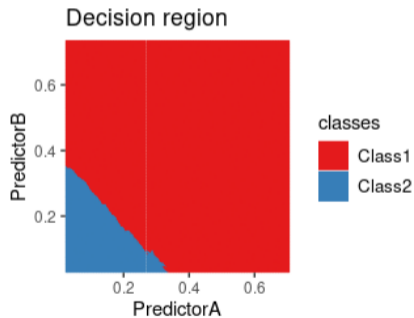
# Example: KNN

k-NN with k=85



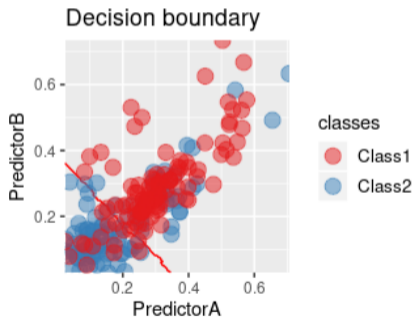
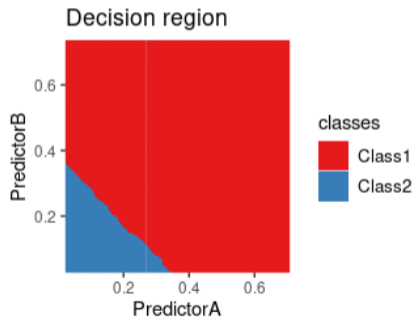
# Example: KNN

k-NN with  $k=101$



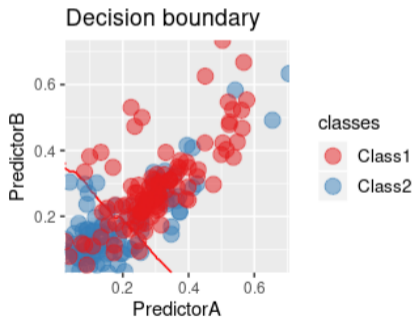
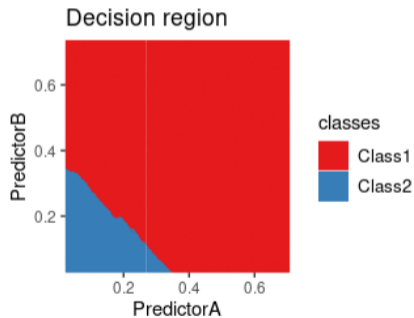
# Example: KNN

k-NN with  $k=109$



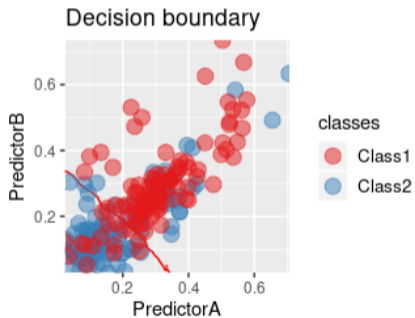
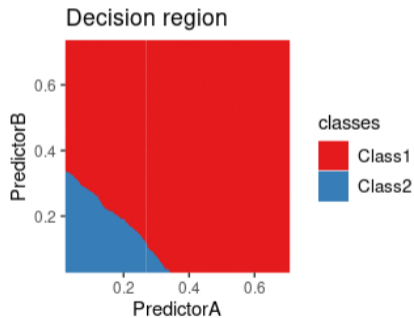
# Example: KNN

k-NN with  $k=117$



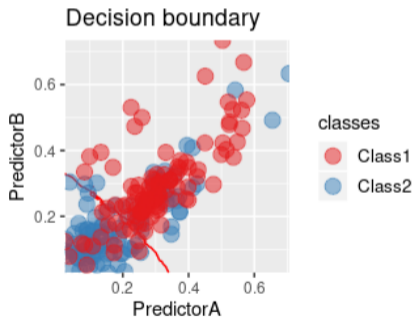
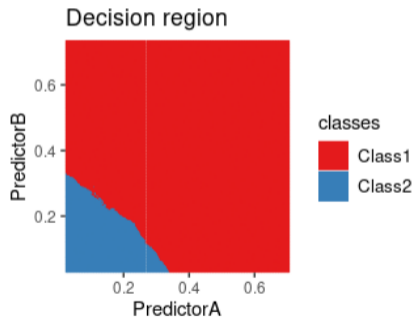
# Example: KNN

k-NN with  $k=125$



# Example: KNN

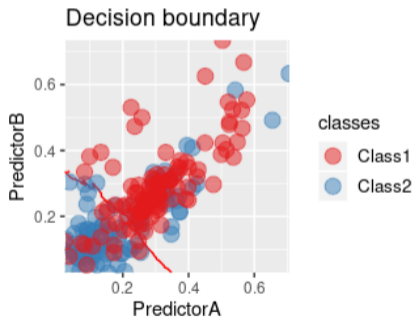
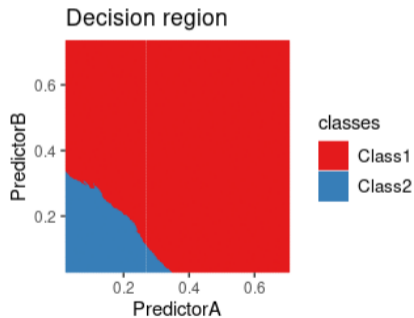
k-NN with  $k=133$





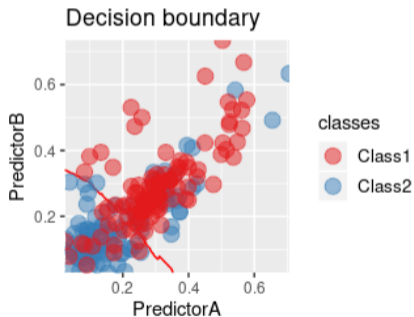
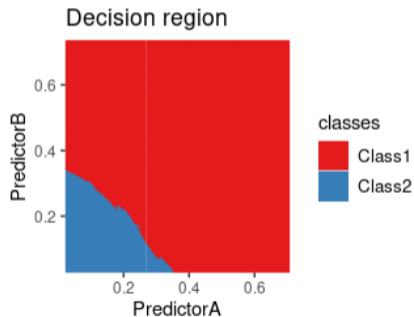
# Example: KNN

k-NN with  $k=141$



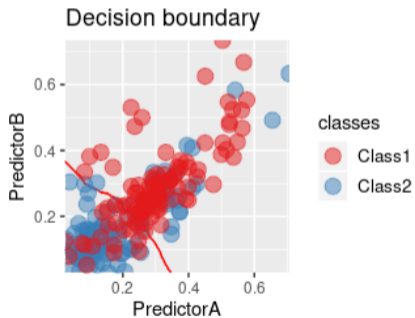
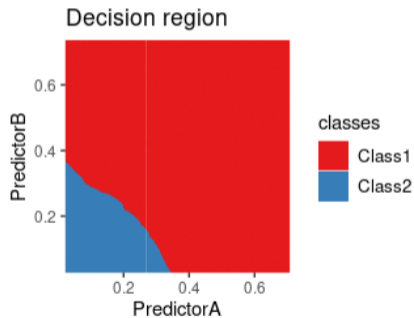
# Example: KNN

k-NN with  $k=149$



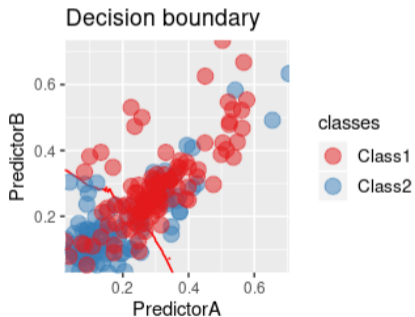
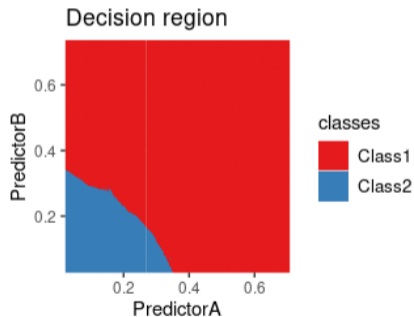
# Example: KNN

k-NN with  $k=157$



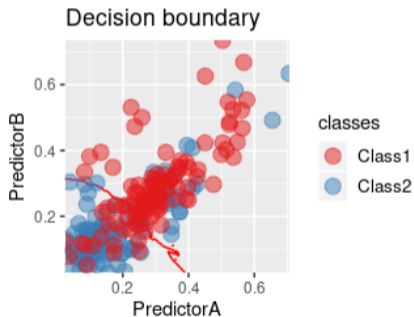
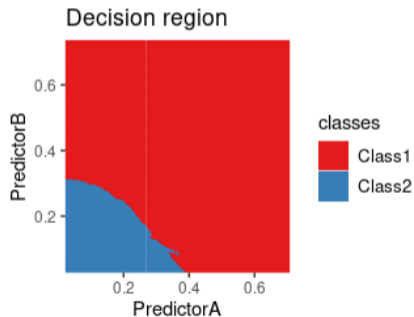
# Example: KNN

k-NN with k=165



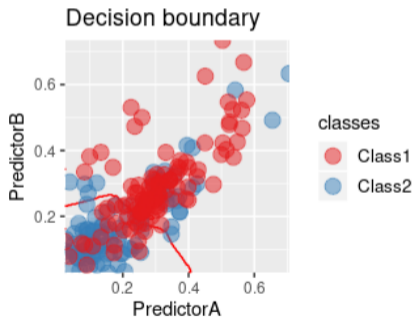
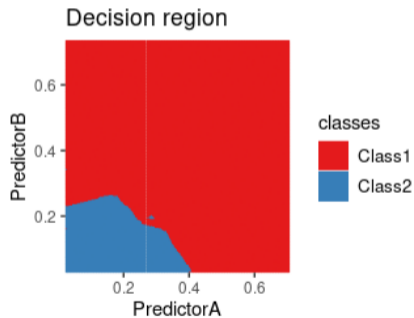
# Example: KNN

k-NN with  $k=173$



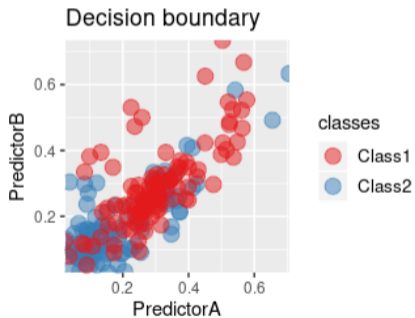
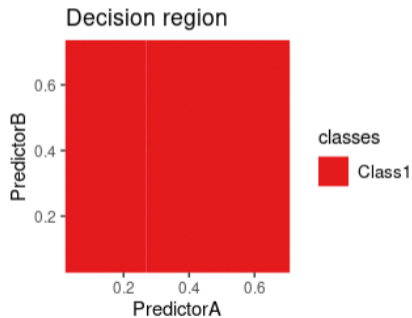
# Example: KNN

k-NN with  $k=181$



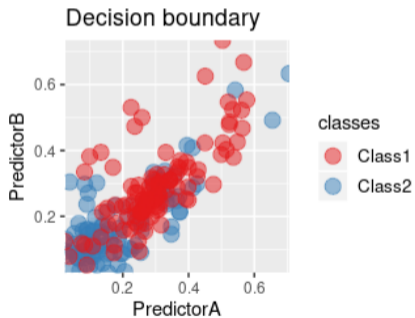
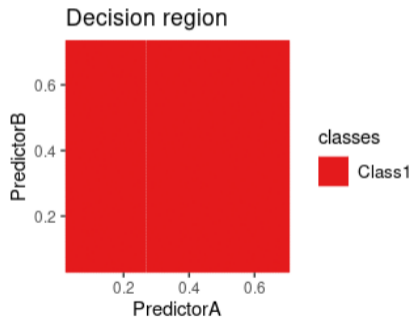
# Example: KNN

k-NN with  $k=189$



# Example: KNN

k-NN with  $k=197$



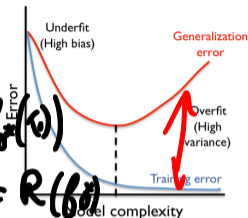


# Training Error Issue

$$\theta^* = \min \mathbb{E}[l(y, \beta_\theta(x))]$$

$$\frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\hat{\theta}}(x_i)) \leq \frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\theta^*}(x_i))$$

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\hat{\theta}}(x_i)) \right] \leq \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\theta^*}(x_i)) \right] = R(\beta_{\theta^*})$$



$$(x_i, y_i) \rightarrow \hat{\theta}$$

$$\frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\hat{\theta}}(x_i)) \rightarrow \mathbb{E}[l(y, \beta_{\hat{\theta}}(x))]$$

## Error behaviour

- Learning/training error (error made on the learning/training set) decays when the complexity of the **method** increases.
- Quite different behavior when the error is computed on new observations (generalization error).
- Overfit for complex methods: parameters learned are too specific to the learning set!
- General situation! (Think of polynomial fit...)
- Need to use a different criterion than the training error!

$$\frac{1}{n} \sum_{i=1}^n l(y_i, \beta_{\hat{\theta}}(x_i)) \rightarrow R(\beta_{\hat{\theta}})$$

## Predictor Error Estimation

- **Goal:** Given a predictor  $f$  assess its quality.
- **Method:** Hold-out error computation (/ Error correction).
- **Usage:** Compute an estimate of the error of a selected  $f$  using a **test set** to be used to monitor it in the future.
- Basic block very well understood.

## Method Selection

- **Goal:** Given a ML method assess its quality.
- **Method:** Cross Validation (/ Error correction)
- **Usage:** Compute error estimates for several ML methods using **training/validation sets** to choose the most promising one.
- Estimates can be pointwise or better intervals.
- Multiple test issues in method selection.

## Two Approaches

- **Cross validation:** Use empirical loss criterion but on independent data, very efficient (and almost always used in practice!) but slightly biased as it target uses only a fraction of the data.
- **Correction approach:** use empirical loss criterion but *correct* it with a term increasing with the complexity of  $\mathcal{S}$

$$R_n(\hat{f}_S) \rightarrow R_n(\hat{f}_S) + \text{cor}(\mathcal{S})$$

and choose the method with the smallest corrected risk.

## Which loss to use?

- The loss used in the risk: most natural!
- The loss used to estimate  $\hat{\theta}$ : penalized estimation!



- **Very simple idea:** use a second learning/verification set to compute a verification error.
- Sufficient to remove the dependency issue!
- Implicit random design setting...

## Cross Validation

- Use  $(1 - \epsilon) \times n$  observations to train and  $\epsilon \times n$  to verify!
- Possible issues:
  - Validation for a learning set of size  $(1 - \epsilon) \times n$  instead of  $n$  ?
  - Unstable error estimate if  $\epsilon n$  is too small ?
- Most classical variations:
  - Hold Out,
  - Leave One Out,
  - $V$ -fold cross validation.

## Principle

- Split the dataset  $\mathcal{D}$  in 2 sets  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$  of size  $n \times (1 - \epsilon)$  and  $n \times \epsilon$ .
- Learn  $\hat{f}^{HO}$  from the subset  $\mathcal{D}_{\text{train}}$ .
- Compute the empirical error on the subset  $\mathcal{D}_{\text{test}}$ :

$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

## Predictor Error Estimation

- Use  $\hat{f}^{HO}$  as predictor.
- Use  $\mathcal{R}_n^{HO}(\hat{f}^{HO})$  as an estimate of the error of this estimator.

## Method Selection by Cross Validation

- Compute  $\mathcal{R}_n^{HO}(\hat{f}_S^{HO})$  for all the considered methods,
- Select the method with the smallest CV error,
- Reestimate the  $\hat{f}_S$  with all the data.

## Principle

- Split the dataset  $\mathcal{D}$  in 2 sets  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$  of size  $n \times (1 - \epsilon)$  and  $n \times \epsilon$ .
- Learn  $\hat{f}^{HO}$  from the subset  $\mathcal{D}_{\text{train}}$ .
- Compute the empirical error on the subset  $\mathcal{D}_{\text{test}}$ :

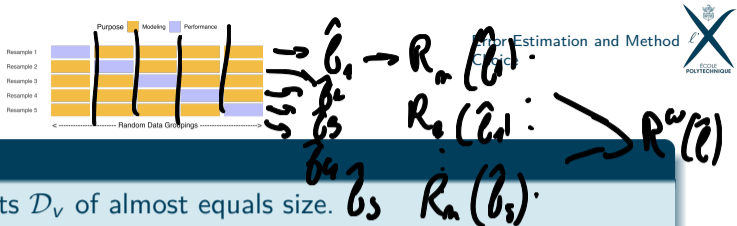
$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

- Only possible setting for error estimation.

## Hold Out Limitation for Method Selection

- Biased toward simpler method as the estimation does not use all the data initially.
- Learning variability of  $\mathcal{R}_n^{HO}(\hat{f}^{HO})$  not taken into account.

# V-fold Cross Validation



## Principle

- Split the dataset  $\mathcal{D}$  in  $V$  sets  $\mathcal{D}_v$  of almost equals size.
- For  $v \in \{1, \dots, V\}$ :
  - Learn  $\hat{f}^{-v}$  from the dataset  $\mathcal{D}$  minus the set  $\mathcal{D}_v$ .
  - Compute the empirical error:

$$\mathcal{R}_n^{-v}(\hat{f}^{-v}) = \frac{1}{n_v} \sum_{(\underline{X}_j, Y_j) \in \mathcal{D}_v} \ell(Y_j, \hat{f}^{-v}(\underline{X}_j))$$

- Compute the average empirical error:

$$\mathcal{R}_n^{CV}(\hat{f}) = \frac{1}{V} \sum_{v=1}^V \mathcal{R}_n^{-v}(\hat{f}^{-v})$$

- Estimation of the quality of a method not of a given predictor.
- Leave One Out :  $V = n$ .

## Analysis (when $n$ is a multiple of $V$ )

- The  $\mathcal{R}_n^{-v}(\hat{f}^{-v})$  are identically distributed variable but are not independent!
- Consequence:

$$\begin{aligned}\mathbb{E} \left[ \mathcal{R}_n^{CV}(\hat{f}) \right] &= \mathbb{E} \left[ \mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ \text{Var} \left[ \mathcal{R}_n^{CV}(\hat{f}) \right] &= \frac{1}{V} \text{Var} \left[ \mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ &\quad + \left(1 - \frac{1}{V}\right) \text{Cov} \left[ \mathcal{R}_n^{-v}(\hat{f}^{-v}), \mathcal{R}_n^{-v'}(\hat{f}^{-v'}) \right]\end{aligned}$$

- Average risk for a sample of size  $(1 - \frac{1}{V})n$ .
- Variance term much more complex to analyze!
- Fine analysis shows that the larger  $V$  the better...

$$C + \frac{1}{V}$$

- Accuracy/Speed tradeoff:  $V = 5$  or  $V = 10$ !



- Leave One Out =  $V$  fold for  $V = n$ : very expensive in general.

## A fast LOO formula for the linear regression

- **Prop:** for the least squares linear regression,

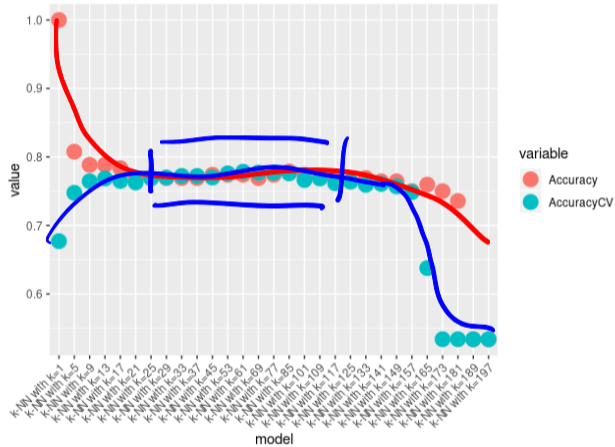
$$\hat{f}^{-i}(\underline{X}_i) = \frac{\hat{f}(\underline{X}_i) - h_{ii} Y_i}{1 - h_{ii}}$$

with  $h_{ii}$  the  $i$ th diagonal coefficient of the **hat** (projection) matrix.

- Proof based on linear algebra!
- Leads to a fast formula for LOO:

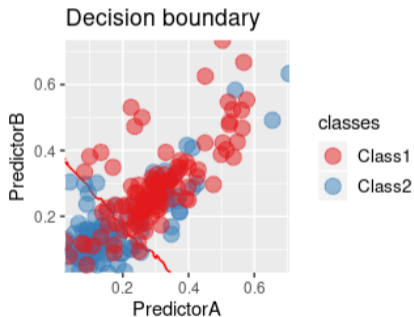
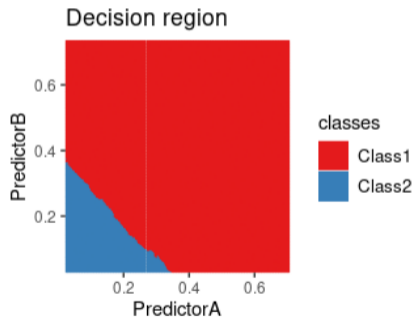
$$\mathcal{R}_n^{LOO}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i - \hat{f}(\underline{X}_i)|^2}{(1 - h_{ii})^2}$$

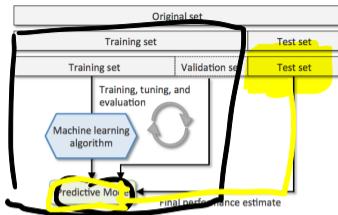
# Cross Validation



# Example: KNN ( $\hat{k} = 61$ using cross-validation)

k-NN with k=61





- **Selection Bias Issue:**

- After method selection, the cross validation is biased.
- Furthermore, it qualifies the method and not the final predictor.
- Need to (re)estimate the error of the final predictor.

## (Train/Validation)/Test strategy

- **Split** the dataset in two a (Train/Validation) and Test.
- Use **CV** with the (Train/Validation) to **select a method**.
- Train this method on (Train/Validation) to **obtain a single predictor**.
- Estimate the **performance of this predictor** on Test.
- Every choice made from the data is part of the method!

- Empirical loss of an estimator computed on the dataset used to choose is biased!
- Empirical loss is an optimistic estimate of the true loss.

## Risk Correction Heuristic

- Estimate an upper bound of this optimism for a given family.
- Correct the empirical loss by adding this upper bound.
- **Rk:** Finding such an upper bound can be complicated!
- Correction often called a **penalty**.

## Penalized Loss

- Minimization of

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\theta}(\underline{X}_i)) + \operatorname{pen}(\theta)$$

where  $\operatorname{pen}(\theta)$  is an error correction (penalty).

## Penalties

- Upper bound of the optimism of the empirical loss
- Depends on the loss and the framework!

## Instantiation

- Mallows Cp: Least Squares with  $\operatorname{pen}(\theta) = 2\frac{d}{n}\sigma^2$ .
- AIC Heuristics: Maximum Likelihood with  $\operatorname{pen}(\theta) = \frac{d}{n}$ .
- BIC Heuristics: Maximum Likelihood with  $\operatorname{pen}(\theta) = \log(n)\frac{d}{n}$ .
- Structural Risk Minimization: Pred. loss and clever penalty.

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice**
  - Cross Validation
  - Cross Validation and Test**
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

## Means

- **Setting:** r.v.  $e_i^{(l)}$  with  $1 \leq i \leq n_l$  and  $l \in \{1, 2\}$  and their means

$$\overline{e^{(l)}} = \frac{1}{n_l} \sum_{i=1}^{n_l} e_i^{(l)}$$

- **Question:** are the means  $\overline{e^{(l)}}$  statistically different?

## Classical i.i.d setting

- **Assumption:**  $e_i^{(l)}$  are i.i.d. for each  $l$ .
- **Test formulation:** Can we reject the null hypothesis that  $\mathbb{E} [e^{(1)}] = \mathbb{E} [e^{(2)}]$ ?
- **Methods:**
  - Gaussian (Student) test using asymptotic normality of a mean.
  - Non parametric permutation test.

- Gaussian approach is linked to confidence intervals.
- The larger  $n_l$  the smaller the confidence intervals.

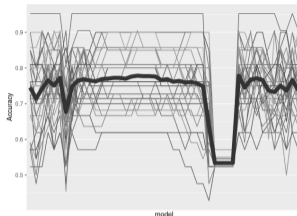


## Non i.i.d. case

- **Assumption:**  $e_i^{(l)}$  are i.d. for each  $l$  but not necessarily independent.
- **Test formulation:** Can we reject the null hypothesis that  $\mathbb{E} [e^{(1)}] = \mathbb{E} [e^{(2)}]$ ?
- **Methods:**
  - Gaussian (Student) test using asymptotic normality of a mean but variance is hard to estimate.
  - Non parametric permutation test but no confidence intervals.
- Setting for Cross Validation (other than holdout).
- Much more complicated than the i.i.d. case

## Several means

- **Assumption:**  $e_i^{(l)}$  are i.i.d. for each  $l$  but not necessarily independent.
- **Tests formulation:**
  - Can we reject the null hypothesis that the  $\mathbb{E} [e^{(l)}]$  are different?
  - Is the smaller mean statistically smaller than the second one?
- **Methods:**
  - Gaussian (Student) test using asymptotic normality of a mean with multiple tests correction.
  - Non parametric permutation test but no confidence intervals.
- Setting for Cross Validation (other than holdout).
- The more models one compares:
  - the larger the confidence intervals
  - the most probable the best model is a lucky winner
- Justify the fallback to the simplest model that could be the best one.



## CV Error, Methods and Predictors

- Cross-Validation errors: estimate of the average risk of a ML method.
- No error bound on the predictor obtained in practice.

## Probably-Approximately-Correct (PAC) Approach

- Replace the control on the average risk by a probabilistic bound

$$\mathbb{P} \left( \mathbb{E} \left[ \ell(Y, \hat{f}(X)) \right] > R \right) \leq \epsilon$$

- Requires to estimate quantiles of the risk.

# Cross Validation and Confidence Interval

- How to replace pointwise estimation by a confidence interval?
- Can we use the variability of the CV estimates?
- **Negative result:** No unbiased estimate of the variance!

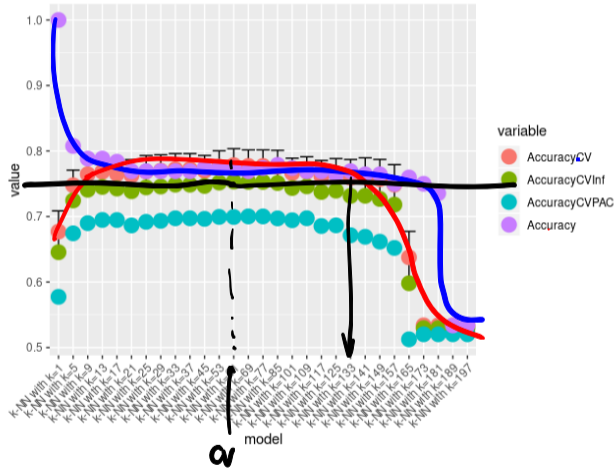
## Gaussian Interval (Comparison of the means and $\sim$ indep.)

- Compute the empirical variance and divide it by the number of folds to construct an asymptotic Gaussian confidence interval,
- Select the simplest model whose values falls into the confidence interval of the model having the smallest CV error.

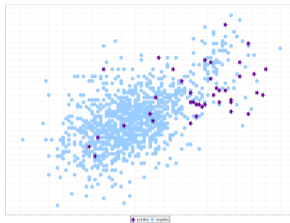
## PAC approach (Quantile, $\sim$ indep. and small risk estim. error)

- Compute the raw medians (or a larger raw quantiles)
- Select the model having the smallest quantiles to ensure a small risk with high probability.
- Always reestimate the chosen model with all the data.
- To obtain an unbiased error estimate of the final predictor: hold out error on untouched test data.

# Cross Validation



- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice**
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights**
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

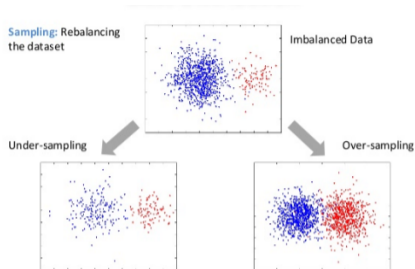


## Unbalanced Class

- **Setting:** One of the class is much more present than the other.
- **Issue:** Classifier *too attracted* by the majority class!

## Rebalanced Dataset

- **Setting:** Class proportions are different in the training and testing set (stratified sampling)
- **Issue:** Training errors are not estimate of testing errors.



## Resampling

- Modify the training dataset so that the classes are more balanced.
- Two flavors:
  - Sub-sampling which spoils data,
  - Over-sampling which needs to create *new* examples.
- **Issues:** Training data is not anymore representative of testing data
- **Hard to do it right!**



## Testing

- Testing class prob.:  $\pi_t(k)$
- Testing error target:

$$\mathbb{E}_{\pi_t} [\ell(Y, f(\underline{X}))] = \sum_k \pi_t(k) \mathbb{E} [\ell(Y, f(\underline{X})) | Y = k]$$

## Training

- Training class prob.:  $\pi_{tr}(k)$
- Training Error target:

$$\mathbb{E}_{\pi_{tr}} [\ell(Y, f(\underline{X}))] = \sum_k \pi_{tr}(k) \mathbb{E} [\ell(Y, f(\underline{X})) | Y = k]$$

## Implicit Testing Error Using the Training One

- Amounts to use a weighted loss:

$$\begin{aligned} \mathbb{E}_{\pi_{tr}} [\ell(Y, f(\underline{X}))] &= \sum_k \pi_{tr}(k) \mathbb{E} [\ell(Y, f(\underline{X})) | Y = k] \\ &= \sum_k \pi_t(k) \mathbb{E} \left[ \frac{\pi_{tr}(k)}{\pi_t(k)} \ell(Y, f(\underline{X})) \middle| Y = k \right] \\ &= \mathbb{E}_{\pi_t} \left[ \frac{\pi_{tr}(Y)}{\pi_t(Y)} \ell(Y, f(\underline{X})) \right] \end{aligned}$$

- Put more weight on less probable classes!

- In unbalanced situation, often the **cost** of misprediction is not the same for all classes (e.g. medical diagnosis, credit lending...)
- Much better to use this explicitly than to do blind resampling!

## Weighted Loss

- **Weighted loss:**

$$\ell(Y, f(\underline{X})) \rightarrow C(Y)\ell(Y, f(\underline{X}))$$

- Weighted error target:

$$\mathbb{E} [C(Y)\ell(Y, f(\underline{X}))]$$

- **Rk:** Strong link with  $\ell$  as  $C$  is independent of  $f$ .
- Often allow to reuse algorithm constructed for  $\ell$ .
- $C$  may also depends on  $\underline{X}$ ...

- The Bayes classifier is now:

$$f^* = \operatorname{argmin} \mathbb{E} [C(Y)\ell(Y, f(\underline{X}))] = \operatorname{argmin} \mathbb{E}_{\underline{X}} \left[ \mathbb{E}_{Y|\underline{X}} [C(Y)\ell(Y, f(\underline{X}))] \right]$$

## Bayes Predictor

- For  $\ell^{0/1}$  loss,

$$f^*(\underline{X}) = \operatorname{argmax}_k C(k)\mathbb{P}(Y = k|\underline{X})$$

- Same effect than a threshold modification for the binary setting!
- Allow to put more emphasis on some classes than others.

## Cost and Proportions

- Testing error target:

$$\mathbb{E}_{\pi_t} [C_t(Y)\ell(Y, f(\underline{X}))] = \sum_k \pi_t(k)C_t(k)\mathbb{E} [\ell(Y, f(\underline{X}))|Y = k]$$

- Training error target

$$\mathbb{E}_{\pi_{tr}} [C_{tr}(Y)\ell(Y, f(\underline{X}))] = \sum_k \pi_{tr}(k)C_{tr}(k)\mathbb{E} [\ell(Y, f(\underline{X}))|Y = k]$$

- **Coincide if**

$$\pi_t(k)C_t(k) = \pi_{tr}(k)C_{tr}(k)$$

- Lots of flexibility in the choice of  $C_t$ ,  $C_{tr}$  or  $\pi_{tr}$ !

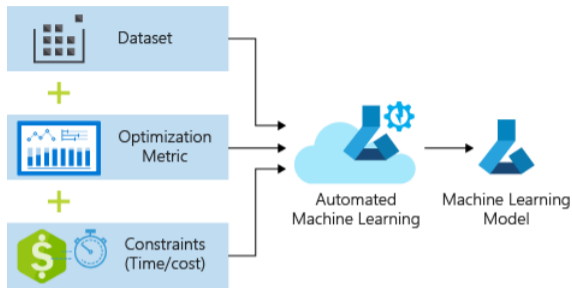
## Weighted Loss and Resampling

- **Weighted loss:** choice of a weight  $C_t \neq 1$ .
- **Resampling:** use a  $\pi_{tr} \neq \pi_t$ .
- Stratified sampling may be used to reduced the size of a dataset without loosing a low probability class!

## Combining Weights and Resampling

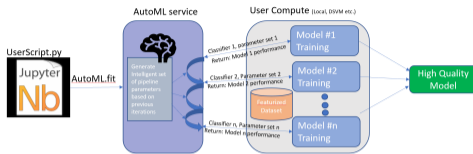
- **Weighted loss:** use  $C_{tr} = C_t$  as  $\pi_{tr} = \pi_t$ .
- **Resampling:** use an implicit  $C_t(k) = \pi_{tr}(k)/\pi_t(k)$ .
- **Combined:** use  $C_{tr}(k) = C_t(k)\pi_t(k)/\pi_{tr}(k)$
- Most ML methods allow such weights!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice**
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - **Auto ML**
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



## Auto ML

- Automatically propose a good predictor
- Rely heavily on error evaluations
- **Pros:** easy way to obtain an excellent baseline
- **Cons:** black box that can be abused. . .

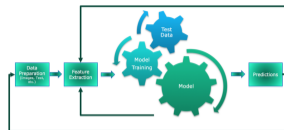


## Auto ML Task

- Input:
  - a dataset  $\mathcal{D} = (\underline{X}_i, Y_i)$
  - a loss function  $\ell(Y, f(\underline{X}))$
  - a set of possible predictors  $f_{l,h,\theta}$  corresponding to a method  $l$  in a list, with hyperparameters  $h$  and parameters  $\theta$
- Output:
  - a predictor  $f$  equal to  $f_{\hat{l},\hat{h},\hat{\theta}}$  or combining several such functions.

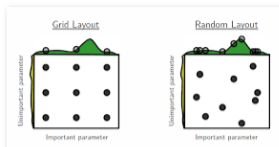


A Standard Machine Learning Pipeline



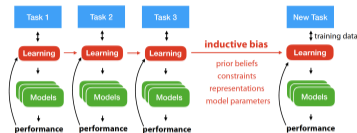
## Predictors, a.k.a fitted pipelines

- Preprocessing:
    - Feature design: normalization, coding, kernel...
    - Missing value strategy
    - Feature selection method
  - ML Method:
    - Method itself
    - Hyperparameters and architecture
    - Fitted parameters (includes optimization algorithm)
- 
- Quickly amounts to 20 to 50 design decisions!
  - **Bruteforce exploration impossible!**



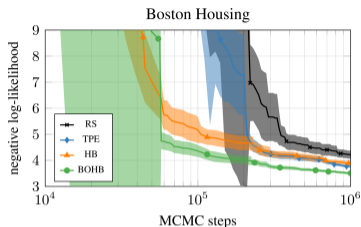
## Most Classical Approach of Auto ML

- Task rephrased as an optimization on the discrete/continuous space of methods/hyperparameters/parameters.
- Parameters obtained by classical minimization.
- Optimization of methods/hyperparameters much more challenging.
- Approaches:
  - Bruteforce: Grid search and random search
  - Clever exploration: Evolutionary algorithm
  - Surrogate based: Bayesian search and Reinforcement learning



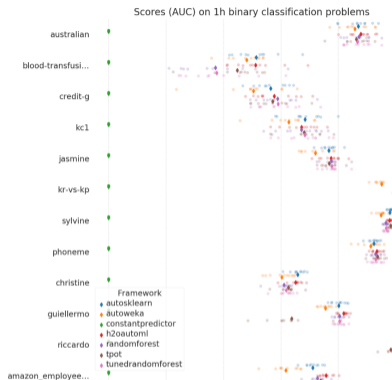
## Learn from other Learning Tasks

- Consider the choice of the method from a dataset and a metric as a learning task.
- Requires a way to describe the problems (or to compute a similarity).
- Descriptor often based on a combination of dataset properties and fast method results.
- May output a list of candidates instead of a single method.
  
- Promising but still quite experimental!



## How to obtain a good result with a time constraint?

- Brute force: Time out and methods screening with Meta-Learning (less exploration at the beginning)
- Surrogate based: Bayesian optimization (exploration/exploitation tradeoff)
- Successive elimination: Fast but not accurate performance evaluation at the beginning to eliminate the worst models (more exploration at the beginning)
- Combined strategy: Bandit strategy to obtain a more accurate estimate of errors only for the promising models (exploration/exploitation tradeoff)



## Benchmark

- Not always (much) better than a good random forest or gradient boosting predictor.
- Worth the try!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

## Logistic Regression

- Let  $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$  with  $\theta = (\beta, \beta^{(0)})$ .
- Let  $\mathbb{P}_{\theta}(Y = 1|\underline{X}) = e^{-f_{\theta}(\underline{X})} / (1 + e^{f_{\theta}(\underline{X})})$
- Estimate  $\theta$  by  $\hat{\theta}$  using a Maximum Likelihood.
- Classify using  $\mathbb{P}_{\hat{\theta}}(Y = 1|\underline{X}) > 1/2$

## $k$ Nearest Neighbors

- For any  $\underline{X}'$ , define  $\mathcal{V}_{\underline{X}'}$  as the  $k$  closest samples  $X_i$  from the dataset.
- Compute a score  $g_k = \sum_{X_i \in \mathcal{V}_{\underline{X}'}} \mathbf{1}_{Y_i=k}$
- Classify using  $\arg \max g_k$  (majority vote).

## Quadratic Discriminant Analysis

- For each class, estimate the mean  $\mu_k$  and the covariance matrix  $\Sigma_k$ .
- Estimate the proportion  $\mathbb{P}(Y = k)$  of each class.
- Compute a score  $\ln(\mathbb{P}(\underline{X}|Y = k)) + \ln(\mathbb{P}(Y = k))$

$$g_k(\underline{X}) = -\frac{1}{2}(\underline{X} - \mu_k)^\top \Sigma_k^{-1}(\underline{X} - \mu_k) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) + \ln(\mathbb{P}(Y = k))$$

- Classify using  $\arg \max g_k$
- Those three methods rely on a similar heuristic: the probabilistic point of view!



- The best solution  $f^*$  (which is independent of  $\mathcal{D}_n$ ) is

$$f^* = \arg \min_{f \in \mathcal{F}} R(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E} [\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[ \mathbb{E}_{Y|\underline{X}} [\ell(Y, f(\underline{X}))] \right]$$

## Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

**Issue:** Explicit solution requires to **know**  $Y|\underline{X}$  (or  $\mathbb{E}[Y|\underline{X}]$ ) for all values of  $\underline{X}$ !

- **Idea:** Estimate  $Y|\underline{X}$  by  $\widehat{Y|\underline{X}}$  and plug it the Bayes classifier.

## Plugin Bayes Predictor

- In binary classification with 0 – 1 loss:

$$\widehat{f}(\underline{X}) = \begin{cases} +1 & \text{if } \overline{\mathbb{P}(Y = +1|\underline{X})} \geq \overline{\mathbb{P}(Y = -1|\underline{X})} \\ & \Leftrightarrow \overline{\mathbb{P}(Y = +1|\underline{X})} \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$\widehat{f}(\underline{X}) = \mathbb{E} \left[ \widehat{Y|\underline{X}} \right]$$

- **Rk:** Direct estimation of  $\mathbb{E} [Y|\underline{X}]$  by  $\overline{\mathbb{E} [Y|\underline{X}]}$  also possible. . .

- How to estimate  $Y|\underline{X}$ ?

## Three main heuristics

- **Parametric Conditional modeling:** Estimate the law of  $Y|\underline{X}$  by a **parametric** law  $\mathcal{L}_\theta(\underline{X})$ : *(generalized) linear regression...*
- **Non Parametric Conditional modeling:** Estimate the law of  $Y|\underline{X}$  by a **non parametric** estimate: *kernel methods, loess, nearest neighbors...*
- **Fully Generative modeling:** Estimate the law of  $(\underline{X}, Y)$  and use the **Bayes formula** to deduce an estimate of  $Y|\underline{X}$ : *LDA/QDA, Naive Bayes...*
- **Rk:** Direct estimation of  $\mathbb{E}[Y|\underline{X}]$  by  $\widehat{\mathbb{E}[Y|\underline{X}]}$  also possible...

- **Input:** a data set  $\mathcal{D}_n$   
Learn  $Y|\underline{X}$  or equivalently  $\mathbb{P}(Y = k|\underline{X})$  (using the data set) and plug this estimate in the Bayes classifier
- **Output:** a classifier  $\hat{f} : \mathbb{R}^d \rightarrow \{-1, 1\}$

$$\hat{f}(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(\widehat{Y} = 1|\underline{X}) \geq \mathbb{P}(\widehat{Y} = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- Can we guaranty that the classifier is good if  $Y|\underline{X}$  is well estimated?

## Theorem

- If  $\hat{f} = \text{sign}(2\hat{p}_{+1} - 1)$  then

$$\begin{aligned}\mathbb{E} \left[ \ell^{0,1}(Y, \hat{f}(\underline{X})) \right] - \mathbb{E} \left[ \ell^{0,1}(Y, f^*(\underline{X})) \right] \\ \leq \mathbb{E} \left[ \|\widehat{Y|\underline{X}} - Y|\underline{X}\|_1 \right] \\ \leq \left( \mathbb{E} \left[ 2\text{KL}(Y|\underline{X}, \widehat{Y|\underline{X}}) \right] \right)^{1/2}\end{aligned}$$

- If one estimates  $\mathbb{P}(Y = 1|\underline{X})$  well then one estimates  $f^*$  well!
- Link between a *conditional density estimation* task and a *classification* one!
- **Rk:** In general, the conditional density estimation task is more complicated as one should be good for all values of  $\mathbb{P}(Y = 1|\underline{X})$  while the classification task focus on values around 1/2 for the 0/1 loss!
- In **regression**, (often) direct control of the quadratic loss...

- Let us denote  $p_1(\underline{X}) = \mathbb{P}(Y = 1|\underline{X})$
- Step 1: Let  $\tilde{f}(\underline{X}) = \text{sign}(2\tilde{p}_1(\underline{X}) - 1)$

$$\begin{aligned} \mathbb{E} \left[ \ell^{0/1}(Y, \tilde{f}(\underline{X})) \right] &= \mathbb{E}_{\underline{X}} \left[ p_1(\underline{X}) \mathbf{1}_{\tilde{f}(\underline{X})=-1} + (1 - p_1(\underline{X})) \mathbf{1}_{\tilde{f}(\underline{X})=1} \right] \\ &= \mathbb{E}_{\underline{X}} \left[ (1 - p_1(\underline{X})) + (2p_1(\underline{X}) - 1) \mathbf{1}_{\tilde{f}(\underline{X})=-1} \right] \end{aligned}$$

- Step 2:

$$\begin{aligned} &\mathbb{E} \left[ \ell^{0/1}(Y, \tilde{f}(\underline{X})) \right] - \mathbb{E} \left[ \ell^{0/1}(Y, \tilde{f}^*(\underline{X})) \right] \\ &= \mathbb{E}_{\underline{X}} \left[ (2p_1(\underline{X}) - 1) (\mathbf{1}_{\tilde{f}(\underline{X})=-1} - \mathbf{1}_{\tilde{f}^*(\underline{X})=-1}) \right] \end{aligned}$$

using the definition of  $\tilde{f}^* = \text{sign}(2p(\underline{X}) - 1)$

$$= \mathbb{E}_{\underline{X}} \left[ |2p_1(\underline{X}) - 1| \mathbf{1}_{\tilde{f}^*(\underline{X}) \neq \tilde{f}(\underline{X})} \right]$$

and using the fact that  $\tilde{f}^*(\underline{X}) \neq \tilde{f}(\underline{X})$  implies that  $\hat{p}(\underline{X})$  and  $p(\underline{X})$  are not on the same side with respect to 1/2

$$\leq 2\mathbb{E}_{\underline{X}} [ |p_1(\underline{X}) - \hat{p}_1(\underline{X})| ] = \mathbb{E}_{\underline{X}} [ \|p(\underline{X}) - \hat{p}(\underline{X})\|_1 ]$$

using  $\|P - Q\|_1 \leq \sqrt{2\text{KL}(P, Q)}$  and Jensen

$$\leq \mathbb{E}_{\underline{X}} \left[ \sqrt{2\text{KL}(p(\underline{X}), \hat{p}(\underline{X}))} \right] \leq (\mathbb{E}_{\underline{X}} [2\text{KL}(p(\underline{X}), \hat{p}(\underline{X}))])^{1/2}$$

## 1 Introduction

- Machine Learning
- Motivation

## 2 A Practical View

- Method or Models
- Interpretability
- Metric Choice

## 3 A Better Point of View

- The Example of Univariate Linear Regression
- Supervised Learning

## 4 Error Estimation and Method Choice

- Cross Validation
- Cross Validation and Test
- Cross Validation and Weights
- Auto ML

## 5 A Probabilistic Point of View

### ● Parametric Conditional Density Modeling

- Non Parametric Conditional Density Modeling
- Generative Modeling

## 6 Optimization Point of View

- SVM
- Penalization
- (Deep) Neural Networks
- Tree Based Methods
- Ensemble Methods

## 7 Empirical Risk Minimization

- Empirical Risk Minimization
- ERM and PAC Bayesian Analysis
- Hoeffding and Finite Class
- McDiarmid and Rademacher Complexity
- VC Dimension
- Structural Risk Minimization

## 8 References

- **Idea:** Estimate directly  $Y|\underline{X}$  by a parametric conditional density  $\mathbb{P}_\theta (Y|\underline{X})$ .

## Maximum Likelihood Approach

- Classical choice for  $\theta$ :

$$\hat{\theta} = \operatorname{argmin}_{\theta} - \sum_{i=1}^n \log \mathbb{P}_\theta (Y_i|\underline{X}_i)$$

- **Goal:** Minimize the Kullback-Leibler divergence between the conditional law of  $Y|\underline{X}$  and  $\mathbb{P}_\theta (Y|\underline{X})$

$$\mathbb{E} [\text{KL} (Y|\underline{X}, \mathbb{P}_\theta (Y|\underline{X}))]$$

- **Rk:** This is often not (exactly) the learning task!
- Large choice for the family  $\{\mathbb{P}_\theta (Y|\underline{X})\}$  but depends on  $\mathcal{Y}$  (and  $\mathcal{X}$ ).
- **Regression:** One can also model directly  $\mathbb{E} [Y|\underline{X}]$  by  $f_\theta(\underline{X})$  and estimate it with a least square criterion. . .



## Linear Models

- **Classical choice:**  $\theta = (\theta', \varphi)$

$$\mathbb{P}_{\theta}(Y|\underline{X}) = \mathbb{P}_{\underline{X}^{\top}\beta, \varphi}(Y)$$

- **Very strong assumption!**
- Classical examples:
  - Binary variable: logistic, probit...
  - Discrete variable: multinomial logistic regression...
  - Integer variable: Poisson regression...
  - Continuous variable: Gaussian regression...

## Plugin Linear Discrimination

- Model  $\mathbb{P}(Y = +1|\underline{X})$  by  $h(\underline{X}^\top \beta + \beta^{(0)})$  with  $h$  non decreasing.
- $h(\underline{X}^\top \beta + \beta^{(0)}) > 1/2 \Leftrightarrow \underline{X}^\top \beta + \beta^{(0)} - h^{-1}(1/2) > 0$
- Linear Classifier:  $\text{sign}(\underline{X}^\top \beta + \beta^{(0)} - h^{-1}(1/2))$

## Plugin Linear Classifier Estimation

- Classical choice for  $h$ :

$$h(t) = \frac{e^t}{1 + e^t} \quad \text{logit or logistic}$$

$$h(t) = F_{\mathcal{N}}(t) \quad \text{probit}$$

$$h(t) = 1 - e^{-e^t} \quad \text{log-log}$$

- Choice of the *best*  $\beta$  from the data.

## Probabilistic Model

- By construction,  $Y|\underline{X}$  follows  $\mathcal{B}(\mathbb{P}(Y = +1|\underline{X}))$
- Approximation of  $Y|\underline{X}$  by  $\mathcal{B}(h(\underline{x}^\top \beta + \beta^{(0)}))$
- *Natural* probabilistic choice for  $\beta$ : maximum likelihood estimate.
- *Natural* probabilistic choice for  $\beta$ :  $\beta$  approximately minimizing a distance between  $\mathcal{B}(h(\underline{x}^\top \beta))$  and  $\mathcal{B}(\mathbb{P}(Y = 1|\underline{X}))$ .

## Maximum Likelihood Approach

- Minimization of the opposite of the log-likelihood:

$$-\frac{1}{n} \sum_{i=1}^n \left( \mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right)$$

- Minimization possible if  $h$  is regular...

## KL Distance and log-likelihood

- *Natural* distance: Kullback-Leibler divergence

$$\begin{aligned} & \text{KL}(\mathcal{B}(\mathbb{P}(Y = 1|\underline{X})), \mathcal{B}(h(\underline{X}^\top \beta))) \\ &= \mathbb{E}_{\underline{X}} \left[ \mathbb{P}(Y = 1|\underline{X}) \log \frac{\mathbb{P}(Y = 1|\underline{X})}{h(\underline{X}^\top \beta)} \right. \\ & \quad \left. + (1 - \mathbb{P}(Y = 1|\underline{X})) \log \frac{1 - \mathbb{P}(Y = 1|\underline{X})}{1 - h(\underline{X}^\top \beta)} \right] \\ &= \mathbb{E}_{\underline{X}} \left[ -\mathbb{P}(Y = 1|\underline{X}) \log(h(\underline{X}^\top \beta)) \right. \\ & \quad \left. - (1 - \mathbb{P}(Y = 1|\underline{X})) \log(1 - h(\underline{X}^\top \beta)) \right] + C_{\underline{X}, Y} \end{aligned}$$

- Empirical counterpart = opposite of the log-likelihood:

$$-\frac{1}{n} \sum_{i=1}^n \left( \mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right)$$

## Logistic Regression and Odd

- Logistic model:  $h(t) = \frac{e^t}{1+e^t}$  (most *natural* choice...)

- The Bernoulli law  $\mathcal{B}(h(t))$  satisfies then

$$\frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = -1)} = e^t \Leftrightarrow \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = -1)} = t$$

- Interpretation in term of odd.
- Logistic model: linear model on the logarithm of the odd

$$\log \frac{\mathbb{P}(Y = 1|\underline{X})}{\mathbb{P}(Y = -1|\underline{X})} = \underline{X}^\top \beta$$

## Associated Classifier

- Plugin strategy:

$$f_\beta(\underline{X}) = \begin{cases} 1 & \text{if } \frac{e^{\underline{X}^\top \beta}}{1+e^{\underline{X}^\top \beta}} > 1/2 \Leftrightarrow \underline{X}^\top \beta > 0 \\ -1 & \text{otherwise} \end{cases}$$

## Likelihood Rewriting

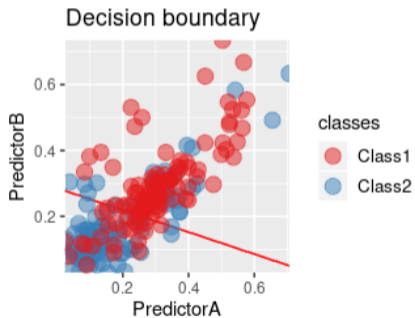
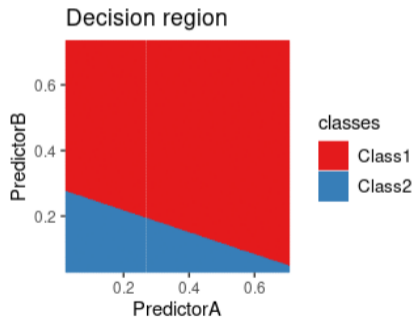
- Opposite of the log-likelihood:

$$\begin{aligned} & -\frac{1}{n} \sum_{i=1}^n \left( \mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right) \\ &= -\frac{1}{n} \sum_{i=1}^n \left( \mathbf{1}_{Y_i=1} \log \frac{e^{\underline{X}_i^\top \beta}}{1 + e^{\underline{X}_i^\top \beta}} + \mathbf{1}_{Y_i=-1} \log \frac{1}{1 + e^{\underline{X}_i^\top \beta}} \right) \\ &= \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-Y_i(\underline{X}_i^\top \beta)} \right) \end{aligned}$$

- Convex and smooth function of  $\beta$
- Easy optimization.

# Example: Logistic

Logistic



## Transformed Representation

- From  $\underline{X}$  to  $\Phi(\underline{X})!$
- New description of  $\underline{X}$  leads to a different **linear** model:

$$f_{\beta}(\underline{X}) = \Phi(\underline{X})^{\top} \beta$$

## Feature Design

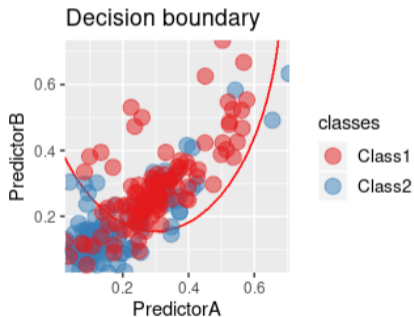
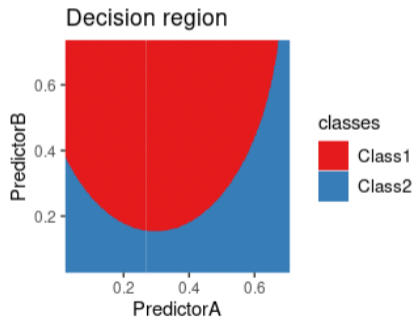
- Art of choosing  $\Phi$ .
- Examples:
  - Renormalization, (domain specific) transform
  - Basis decomposition
  - Interaction between different variables. . .

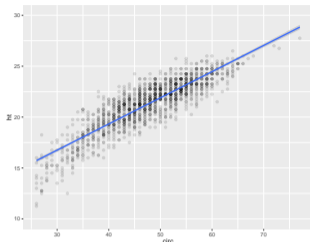


# Example: Quadratic Logistic

$$\underbrace{(Y_A, X_B)}_{\in \mathbb{R}^2} \rightarrow \underbrace{(X_A, X_B, X_A^2, X_B^2, X_A X_B)}_{\in \mathbb{R}^5}$$

Quadratic Logistic





## Gaussian Linear Model

- **Model:**  $Y|\underline{X} \sim \mathcal{N}(\underline{X}^\top \beta, \sigma^2)$  plus independence
- Probably the most classical model of all time!
- Maximum Likelihood with explicit formulas for the two parameters.
- In regression, estimation of  $\mathbb{E}[Y|\underline{X}]$  is sufficient: other/no model for the noise possible.

## Generalized Linear Model

- Model entirely characterized by its mean (up to a scalar nuisance parameter) ( $v(\mathbb{E}_\theta [Y]) = \theta$  with  $v$  invertible).
- Exponential family: Probability law family  $P_\theta$  such that the density can be written

$$f(y, \theta, \varphi) = e^{\frac{y\theta - v(\theta)}{\varphi} + w(y, \varphi)}$$

where  $\varphi$  is a nuisance parameter and  $w$  a function independent of  $\theta$ .

- Examples:
  - Gaussian:  $f(y, \theta, \varphi) = e^{-\frac{y\theta - \theta^2/2}{\varphi} - \frac{y^2/2}{\varphi}}$
  - Bernoulli:  $f(y, \theta) = e^{y\theta - \ln(1+e^\theta)}$  ( $\theta = \ln p/(1-p)$ )
  - Poisson:  $f(y, \theta) = e^{(y\theta - e^\theta) + \ln(y!)}$  ( $\theta = \ln \lambda$ )
- Linear Conditional model:  $Y|\underline{X} \sim P_{\underline{x}^\top \beta} \dots$

- ML fit of the parameters

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - **Non Parametric Conditional Density Modeling**
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- **Idea:** Estimate  $Y|\underline{X}$  or  $\mathbb{E}[Y|\underline{X}]$  directly without resorting to an explicit parametric model.

## Non Parametric Conditional Estimation

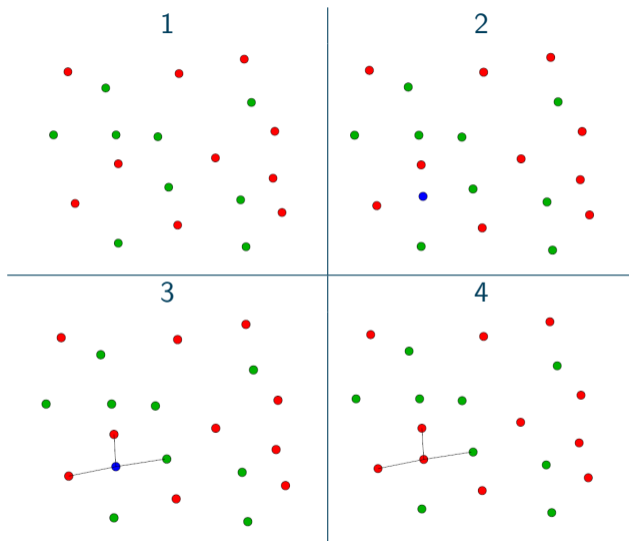
- Two heuristics:
  - $Y|\underline{X}$  (or  $\mathbb{E}[Y|\underline{X}]$ ) is almost constant (or simple) in a neighborhood of  $\underline{X}$ . (Kernel methods)
  - $Y|\underline{X}$  (or  $\mathbb{E}[Y|\underline{X}]$ ) can be approximated by a model whose dimension depends on the complexity and the number of observation. (Quite similar to parametric model plus model selection...)
- Focus on **kernel methods!**

- **Idea:** The behavior of  $Y|\underline{X}$  is locally *constant* or simple!

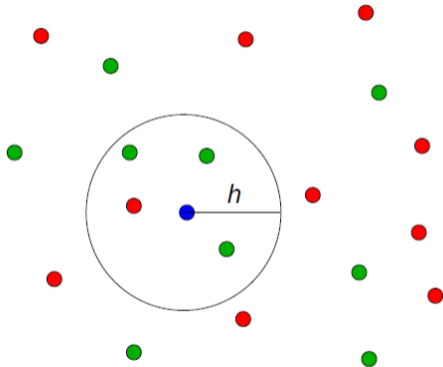
## Kernel

- Choose a kernel  $K$  (think of a weighted neighborhood).
  - For each  $\tilde{X}$ , compute a simple localized estimate of  $Y|\underline{X}$
  - Use this local estimate to take the decision
- 
- In regression, estimation of  $\mathbb{E}[Y|\underline{X}]$  is sufficient.

# Example: $k$ Nearest-Neighbors (with $k = 3$ )



# Example: $k$ Nearest-Neighbors (with $k = 4$ )





- Neighborhood  $\mathcal{V}_{\underline{x}}$  of  $\underline{x}$ :  $k$  learning samples closest from  $\underline{x}$ .

## $k$ -NN as local conditional density estimate

$$\mathbb{P}(\widehat{Y} = 1 | \underline{X}) = \frac{\sum_{\underline{X}_i \in \mathcal{V}_{\underline{X}}} \mathbf{1}_{\{Y_i = +1\}}}{|\mathcal{V}_{\underline{X}}|}$$

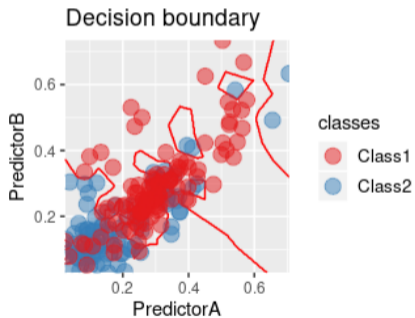
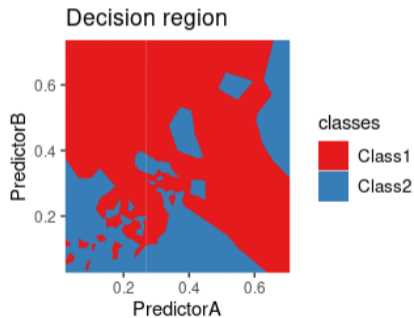
- KNN Classifier:

$$\widehat{f}_{KNN}(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(\widehat{Y} = 1 | \underline{X}) \geq \mathbb{P}(\widehat{Y} = -1 | \underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Lazy learning:** all the computations have to be done at prediction time.
- **Remark:** You can also use your favorite kernel estimator...

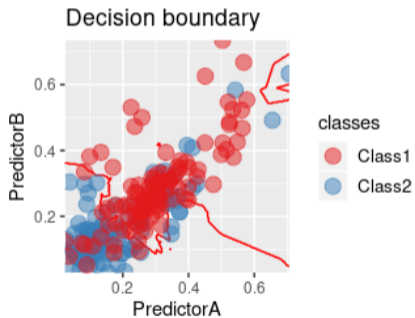
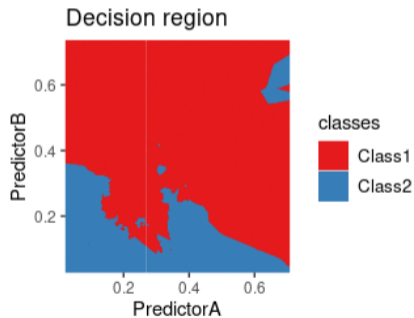
# Example: KNN

k-NN with k=1



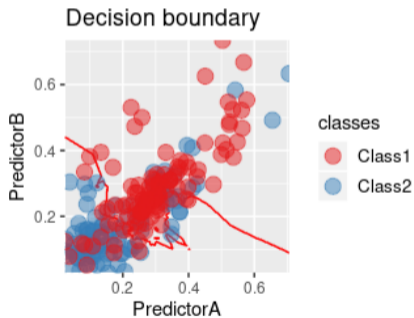
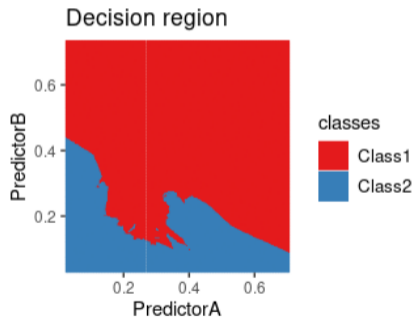
# Example: KNN

k-NN with  $k=5$



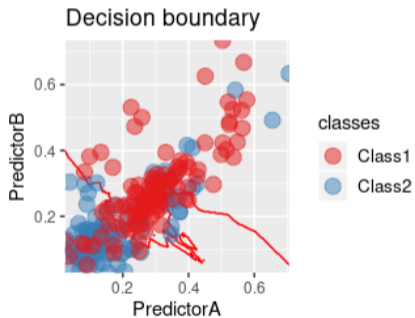
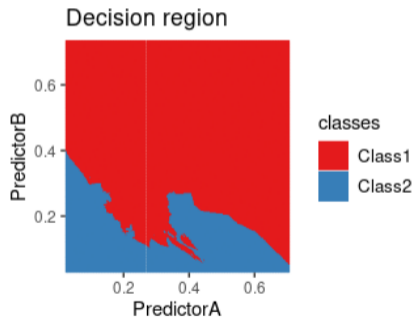
# Example: KNN

k-NN with  $k=9$



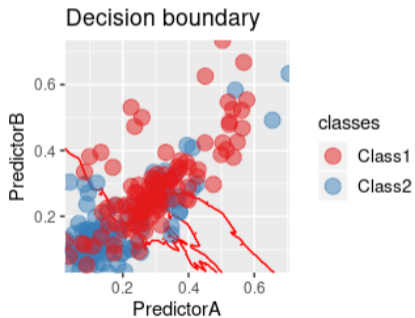
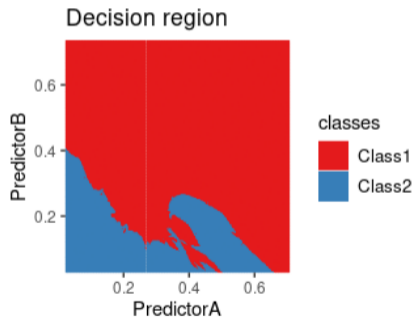
# Example: KNN

k-NN with  $k=13$



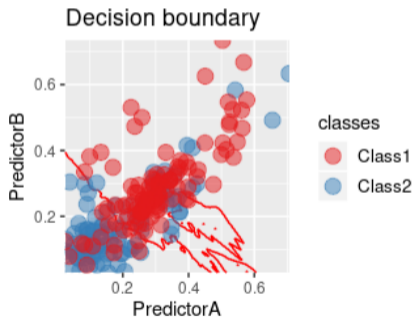
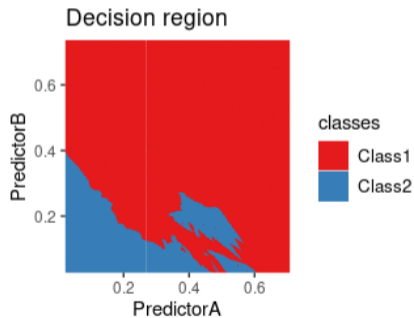
# Example: KNN

k-NN with  $k=17$



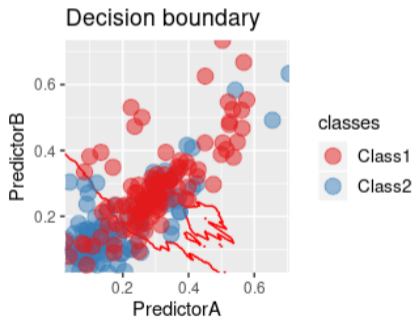
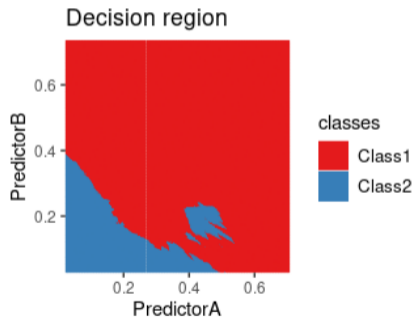
# Example: KNN

k-NN with k=21



# Example: KNN

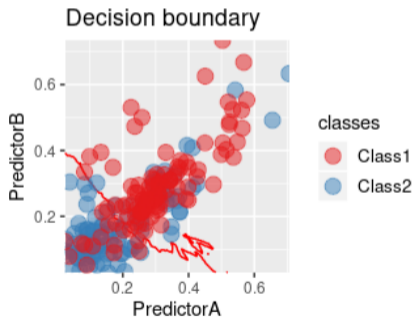
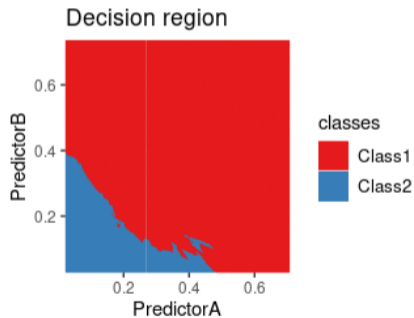
k-NN with  $k=25$





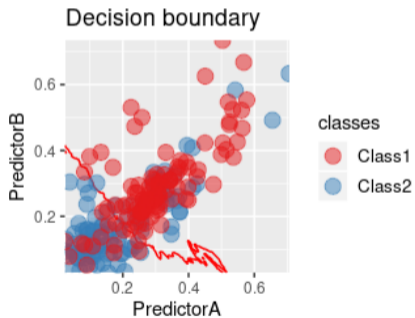
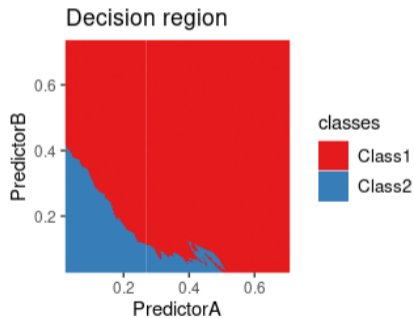
# Example: KNN

k-NN with  $k=29$



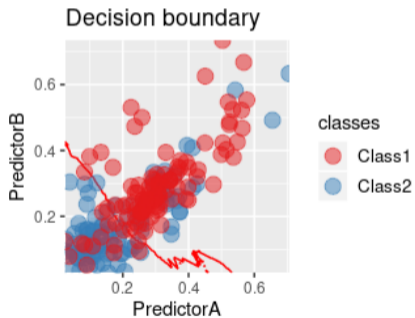
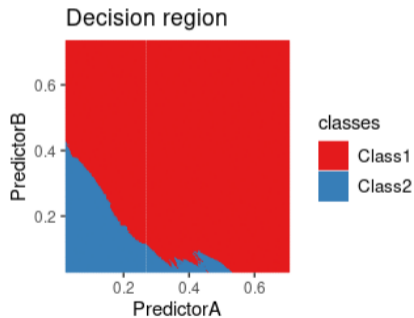
# Example: KNN

k-NN with  $k=33$



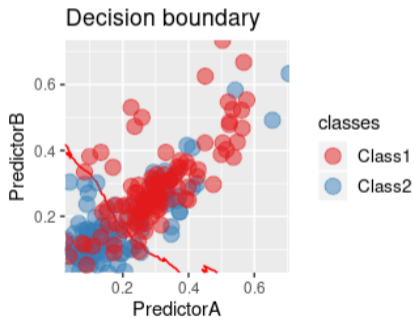
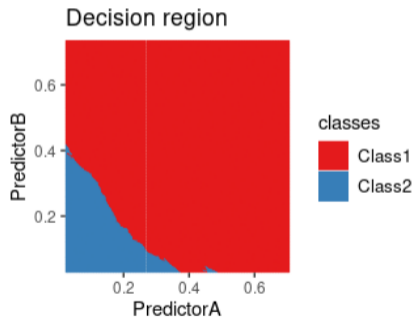
# Example: KNN

k-NN with  $k=37$



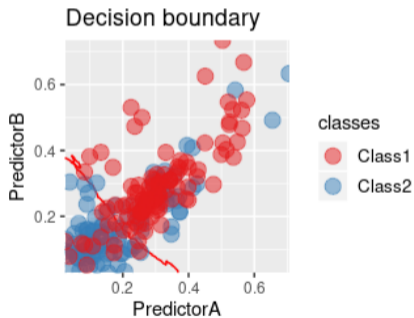
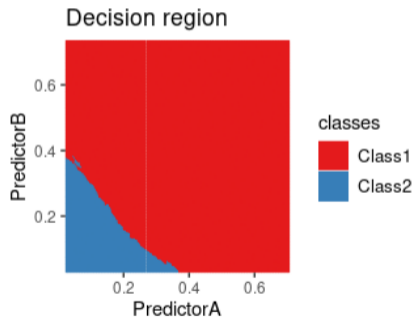
# Example: KNN

k-NN with k=45



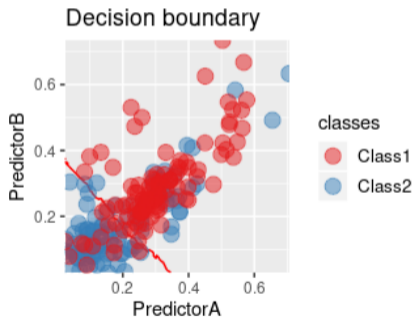
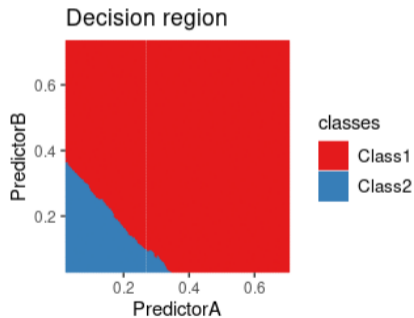
# Example: KNN

k-NN with k=53



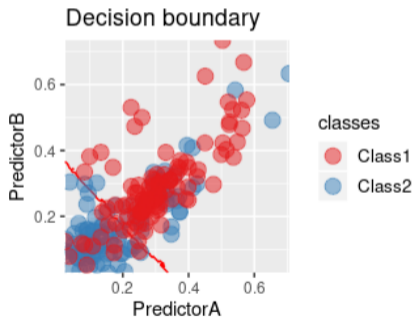
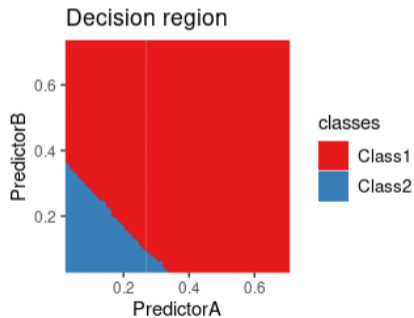
# Example: KNN

k-NN with  $k=61$



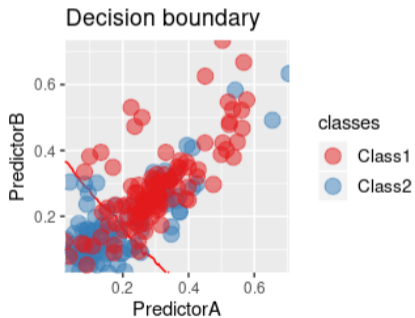
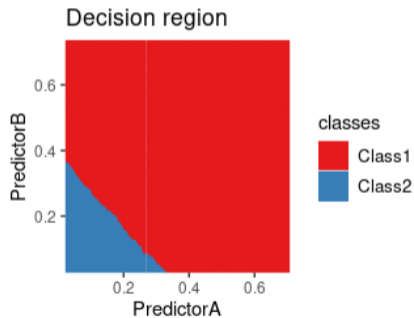
# Example: KNN

k-NN with k=69



# Example: KNN

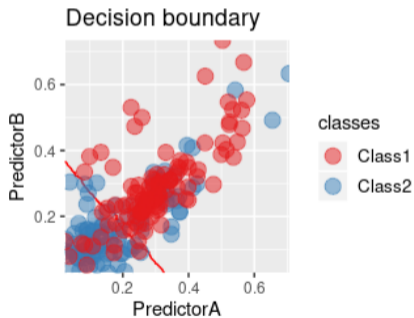
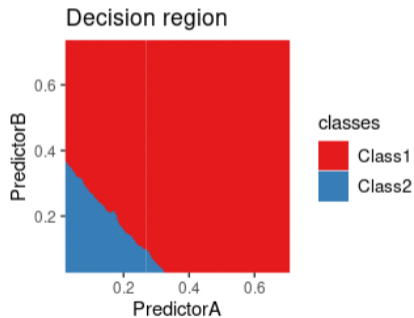
k-NN with  $k=77$





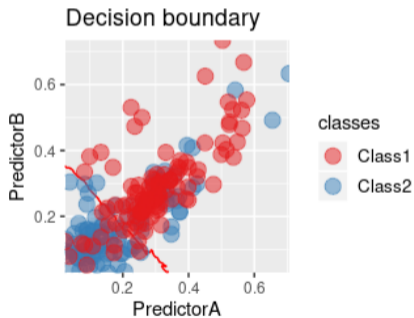
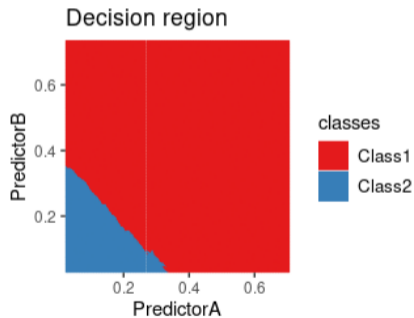
# Example: KNN

k-NN with k=85



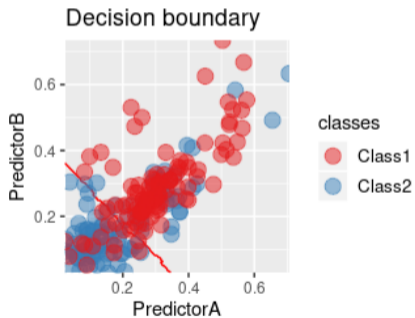
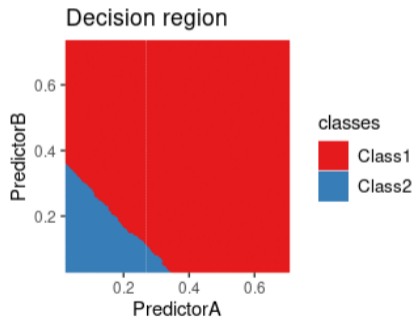
# Example: KNN

k-NN with  $k=101$



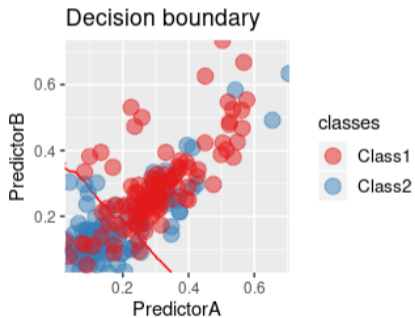
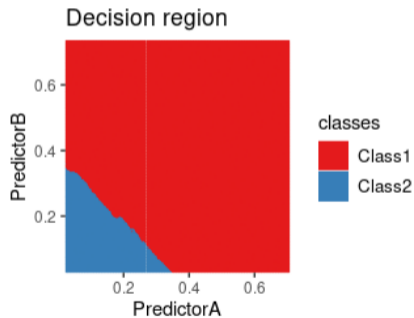
# Example: KNN

k-NN with  $k=109$



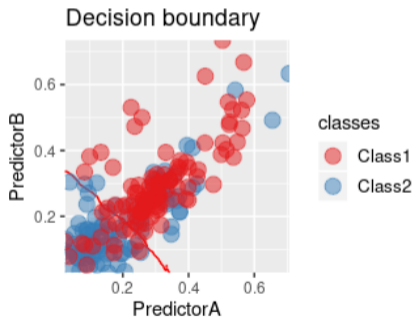
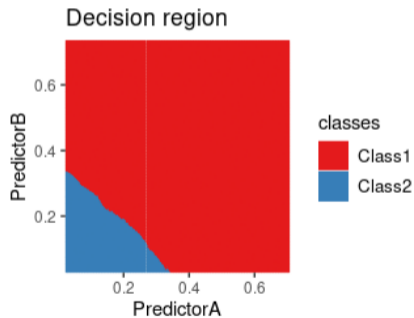
# Example: KNN

k-NN with  $k=117$



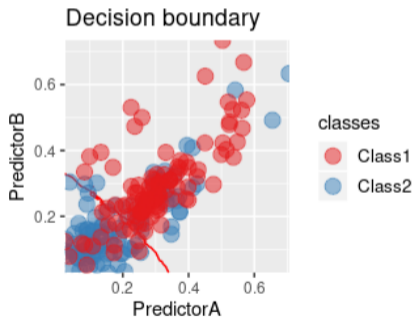
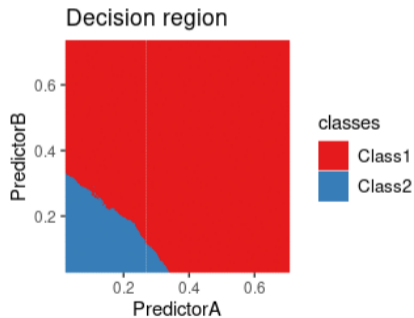
# Example: KNN

k-NN with  $k=125$



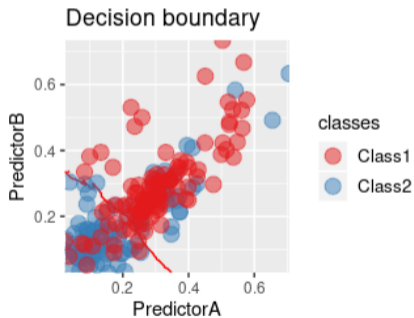
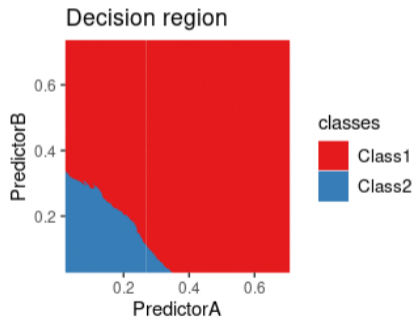
# Example: KNN

k-NN with  $k=133$



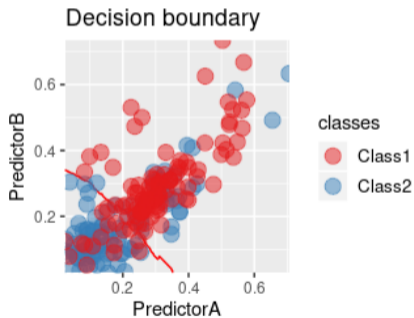
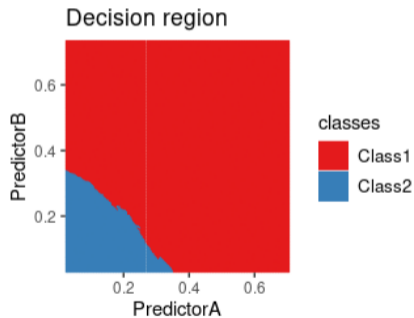
# Example: KNN

k-NN with k=141



# Example: KNN

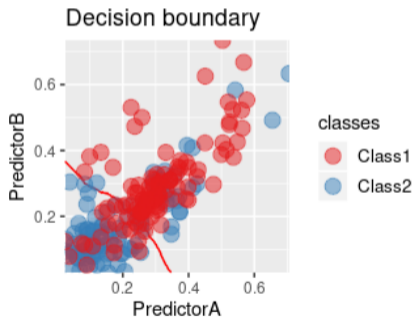
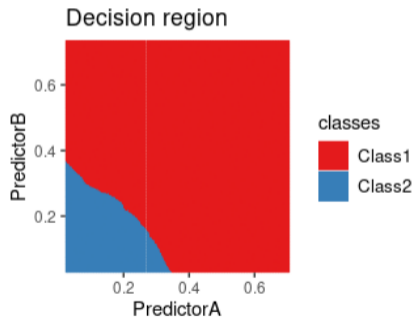
k-NN with  $k=149$





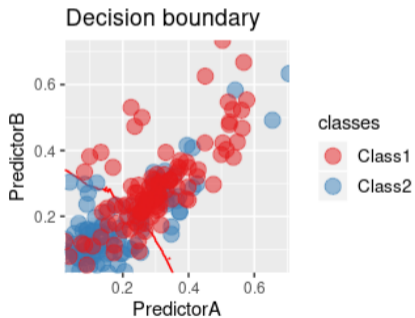
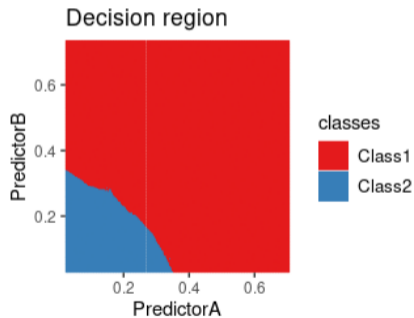
# Example: KNN

k-NN with  $k=157$



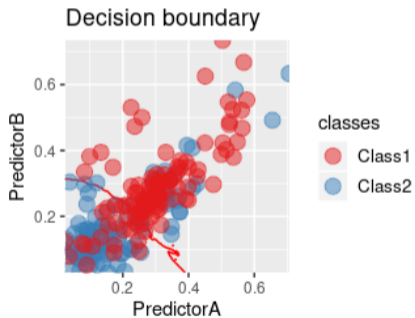
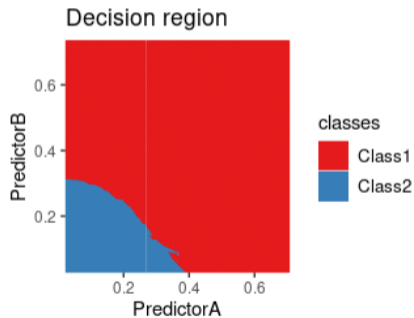
# Example: KNN

k-NN with  $k=165$



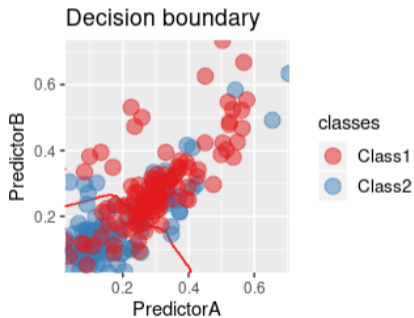
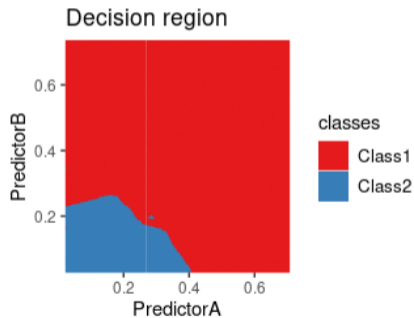
# Example: KNN

k-NN with  $k=173$



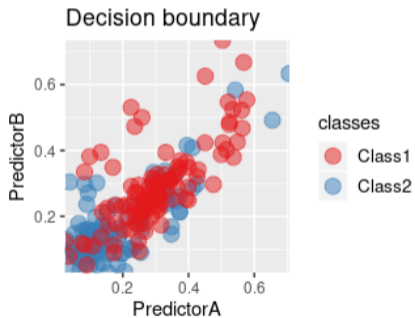
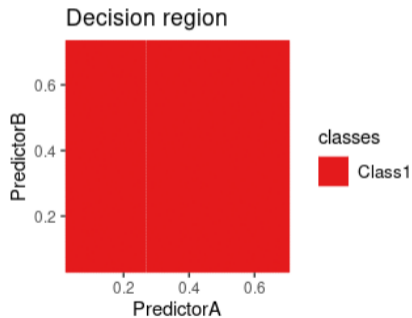
# Example: KNN

k-NN with  $k=181$



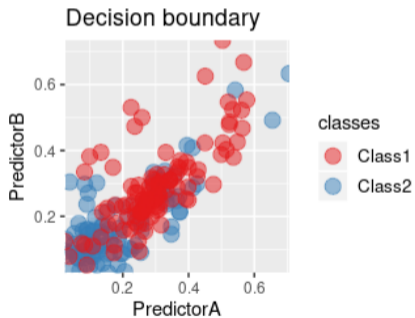
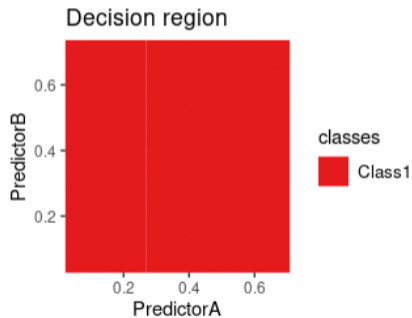
# Example: KNN

k-NN with  $k=189$



# Example: KNN

k-NN with  $k=197$



## A naive idea

- $\mathbb{E}[Y|\underline{X}]$  can be approximated by a local average:

$$\hat{f}(\underline{X}) = \frac{1}{|\{\underline{X}_i \in \mathcal{N}(\underline{X})\}|} \sum_{\underline{X}_i \in \mathcal{N}(\underline{X})} Y_i$$



where  $\mathcal{B}(\underline{X})$  is a neighborhood of  $\underline{X}$ .

- **Heuristic:**

- If  $\underline{X} \rightarrow \mathbb{E}[Y|\underline{X}]$  is regular then

$$\mathbb{E}[Y|\underline{X}] \simeq \mathbb{E}[\mathbb{E}[Y|\underline{X}'] | \underline{X}' \in \mathcal{N}(\underline{X})] = \mathbb{E}[Y|\underline{X}' \in \mathcal{N}(\underline{X})]$$

- Replace an expectation by an empirical average:

$$\mathbb{E}[Y|\underline{X}' \in \mathcal{N}(\underline{X})] \simeq \frac{1}{|\{\underline{X}_i \in \mathcal{N}(\underline{X})\}|} \sum_{\underline{X}_i \in \mathcal{N}(\underline{X})} Y_i$$

## Neighborhood and Size

- Most classical choice:  $\mathcal{N}(\underline{X}) = \{ \underline{X}', \|\underline{X} - \underline{X}'\| \leq h \}$  where  $\|\cdot\|$  is a (pseudo) norm and  $h$  a size (bandwidth) parameter.
- In principle, the norm and  $h$  could vary with  $\underline{X}$ , and the norm can be replaced by a (pseudo) distance.
- Focus here on a fixed distance with a fixed bandwidth  $h$  cased.

## Bandwidth Heuristic

- A **large bandwidth** ensures that the average is taken on many samples and thus the **variance is small**...
- A **small bandwidth** is thus that the approximation  $\mathbb{E}[Y|\underline{X}] \simeq \mathbb{E}[Y|\underline{X}' \in \mathcal{N}(\underline{X})]$  is more accurate (**small bias**).



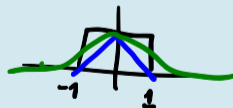
## Weighted Local Average

- Replace the neighborhood  $\mathcal{N}(\underline{X})$  by a decaying **window function**  $w(\underline{X}, \underline{X}')$ .
- $\mathbb{E}[Y|\underline{X}]$  can be approximated by a **weighted local average**:

$$\hat{f}(\underline{X}) = \frac{\sum_i w(\underline{X}, \underline{X}'_i) Y_i}{\sum_i w(\underline{X}, \underline{X}'_i)}.$$

## Kernel

- Most classical choice:  $w(\underline{X}, \underline{X}') = K\left(\frac{\underline{X} - \underline{X}'}{h}\right)$  where  $h$  the bandwidth is a scale parameter.
- Examples:
  - **Box kernel**:  $K(t) = \mathbf{1}_{\|t\| \leq 1}$  (Neighborhood)
  - **Triangular kernel**:  $K(t) = \max(1 - \|t\|, 0)$ .
  - **Gaussian kernel**:  $K(t) = e^{-t^2/2}$
- **Rk**:  $K$  and  $\lambda K$  yields the same estimate.



## Nadaraya-Watson Heuristic

- Provided all the **densities** exist

$$\mathbb{E}[Y|\underline{X}] = \frac{\int Y p(\underline{X}, Y) dY}{\int p(Y, \underline{X}) dY} = \frac{\int Y p(\underline{X}, Y) dY}{p(\underline{X})}$$

- Replace the unknown densities by their **estimates**:

$$\hat{p}(\underline{X}) = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i)$$

$$\hat{p}(\underline{X}, Y) = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i) K'(Y - Y_i)$$

- Now if  $K'$  is a kernel such that  $\int Y K'(Y) dY = 0$  then

$$\int Y \hat{p}(\underline{X}, Y) dY = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i) Y_i$$

## Nadaraya-Watson

- Resulting estimator of  $\mathbb{E}[Y|\underline{X}]$

$$\hat{f}(\underline{X}) = \frac{\sum_{i=1}^n Y_i K_h(\underline{X} - \underline{X}_i)}{\sum_{i=1}^n K_h(\underline{X} - \underline{X}_i)}$$

- Same **local weighted average** estimator!

## Bandwidth Choice

- Bandwidth  $h$  of  $K$  allows to **balance between bias and variance**.
- Theoretical analysis of the error is possible.
- The smoother the densities the easier the estimation but the optimal bandwidth depends on the unknown regularity!

## Another Point of View on Kernel

- Nadaraya-Watson estimator:

$$\hat{f}(\underline{X}) = \frac{\sum_{i=1}^n Y_i K_h(\underline{X} - \underline{X}_i)}{\sum_{i=1}^n K_h(\underline{X} - \underline{X}_i)}$$

- Can be view as a **minimizer** of

$$\sum_{i=1}^n |Y_i - \beta|^2 K_h(\underline{X} - \underline{X}_i)$$

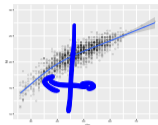
- **Local regression** of order 0!



## Local Linear Model

- Estimate  $\mathbb{E}[Y|\underline{X}]$  by  $\hat{f}(\underline{X}) = \phi(\underline{X})^\top \hat{\beta}(\underline{X})$  where  $\phi$  is any function of  $\underline{X}$  and  $\hat{\beta}(\underline{X})$  is the minimizer of

$$\sum_{i=1}^n |Y_i - \phi(\underline{X}_i)^\top \beta|^2 K_h(\underline{X} - \underline{X}_i).$$



## 1D Nonparametric Regression

- Assume that  $\underline{X} \in \mathbb{R}$  and let  $\phi(\underline{X}) = (1, \underline{X}, \dots, \underline{X}^d)$ .
- **LOESS estimate:**  $\hat{f}(\underline{X}) = \sum_{j=0}^d \hat{\beta}(\underline{X}^{(j)}) \underline{X}^j$  with  $\hat{\beta}(\underline{X})$  minimizing

$$\sum_{i=1}^n |Y_i - \sum_{j=0}^d \beta^{(j)} \underline{X}_i^j|^2 K_h(\underline{X} - \underline{X}_i).$$

- Most classical kernel used: Tricubic kernel

$$K(t) = \max(1 - |t|^3, 0)^3$$

- Most classical degree: 2...
- Local bandwidth choice such that a proportion of points belongs to the window.

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - **Generative Modeling**
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- **Idea:** If one knows the law of  $(\underline{X}, Y)$  everything is easy!

## Bayes formula

- With a slight abuse of notation,

$$\begin{aligned}\mathbb{P}(Y|\underline{X}) &= \frac{\mathbb{P}((\underline{X}, Y))}{\mathbb{P}(\underline{X})} \\ &= \frac{\mathbb{P}(\underline{X}|Y)\mathbb{P}(Y)}{\mathbb{P}(\underline{X})}\end{aligned}$$

- **Generative Modeling:**

- Propose a model for  $(\underline{X}, Y)$  (or equivalently  $\underline{X}|Y$  and  $Y$ ),
- Estimate it as a density estimation problem,
- Plug the estimate in the Bayes formula
- Plug the conditional estimate in the Bayes *classifier*.
- **Rk:** Require to estimate  $(\underline{X}, Y)$  rather than only  $Y|\underline{X}$ !
- Great flexibility in the model design but may lead to complex computation.

- Simpler setting in classification!

## Bayes formula

$$\mathbb{P}(Y = k|\underline{X}) = \frac{\mathbb{P}(\underline{X}|Y = k)\mathbb{P}(Y = k)}{\mathbb{P}(\underline{X})}$$

- Binary Bayes classifier (the best solution)

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = 1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Heuristic:** Estimate those quantities and plug the estimations.
- By using different models/estimators for  $\mathbb{P}(\underline{X}|Y)$ , we get different classifiers.
- **Rk:** No need to renormalize by  $\mathbb{P}(\underline{X})$  to take the decision!



## Discriminant Analysis (Gaussian model)

- The densities are modeled as multivariate normal, i.e.,

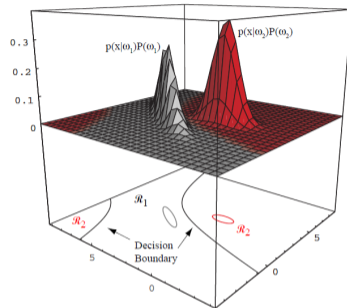
$$\mathbb{P}(\underline{X}|Y = k) \sim \mathcal{N}_{\mu_k, \Sigma_k}$$

- Discriminant functions:  $g_k(\underline{X}) = \ln(\mathbb{P}(\underline{X}|Y = k)) + \ln(\mathbb{P}(Y = k))$

$$g_k(\underline{X}) = -\frac{1}{2}(\underline{X} - \mu_k)^\top \Sigma_k^{-1}(\underline{X} - \mu_k) \\ - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) + \ln(\mathbb{P}(Y = k))$$

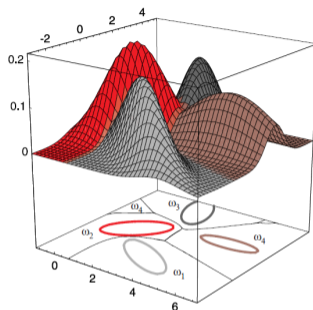
- QDA (different  $\Sigma_k$  in each class) and LDA ( $\Sigma_k = \Sigma$  for all  $k$ )

- **Beware: this model can be false but the methodology remains valid!**



## Quadratic Discriminant Analysis

- The probability densities are Gaussian
- The effect of any decision rule is to divide the feature space into some decision regions  $\mathcal{R}_1, \mathcal{R}_2$
- The regions are separated by decision boundaries



## Quadratic Discriminant Analysis

- The probability densities are Gaussian
- The effect of any decision rule is to divide the feature space into some decision regions  $\mathcal{R}_1, \mathcal{R}_2, \dots, \mathcal{R}_c$
- The regions are separated by decision boundaries

## Estimation

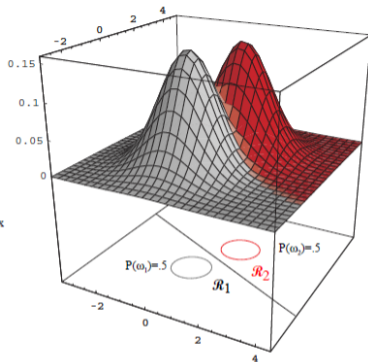
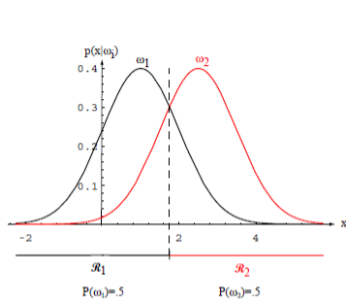
In practice, we will need to estimate  $\mu_k$ ,  $\Sigma_k$  and  $\mathbb{P}_k := \mathbb{P}(Y = k)$

- The estimate proportion  $\mathbb{P}(\widehat{Y} = k) = \frac{n_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i = k\}}$
- Maximum likelihood estimate of  $\widehat{\mu}_k$  and  $\widehat{\Sigma}_k$  (explicit formulas)

- DA classifier

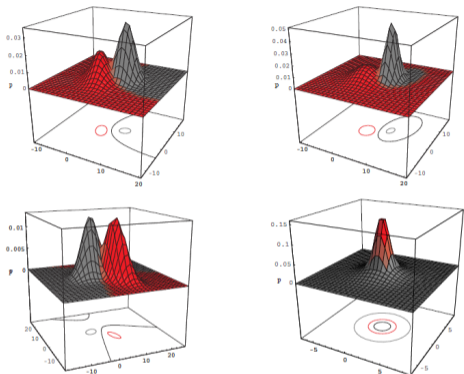
$$\widehat{f}_G(\underline{X}) = \begin{cases} +1 & \text{if } \widehat{g}_{+1}(\underline{X}) \geq \widehat{g}_{-1}(\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- Decision boundaries: quadratic = degree 2 polynomials.
- If one imposes  $\Sigma_{-1} = \Sigma_1 = \Sigma$  then the decision boundaries is a linear hyperplane.



## Linear Discriminant Analysis

- $\Sigma_{\omega_1} = \Sigma_{\omega_2} = \Sigma$
- The decision boundaries are linear hyperplanes

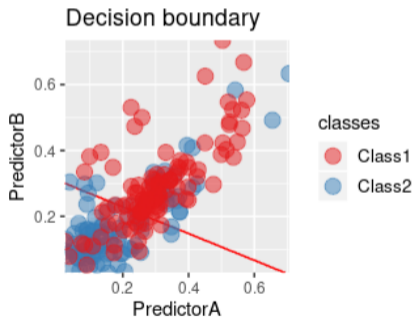
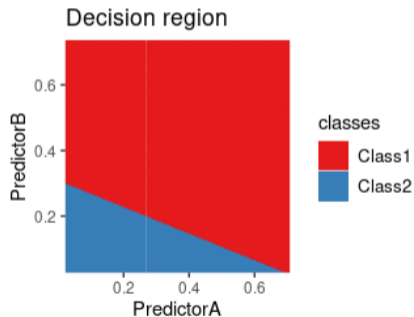


## Quadratic Discriminant Analysis

- $\Sigma_{\omega_1} \neq \Sigma_{\omega_2}$
- Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general quadratics.

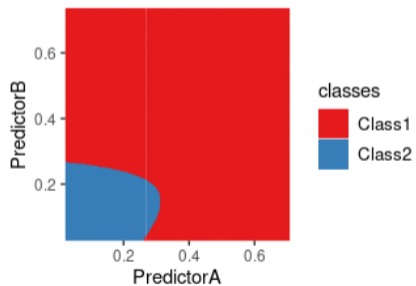
# Example: LDA

## Linear Discriminant Analysis

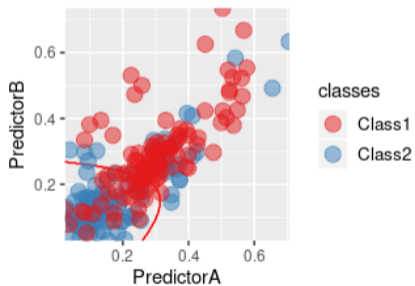


## Quadratic Discriminant Analysis

Decision region



Decision boundary





## Naive Bayes

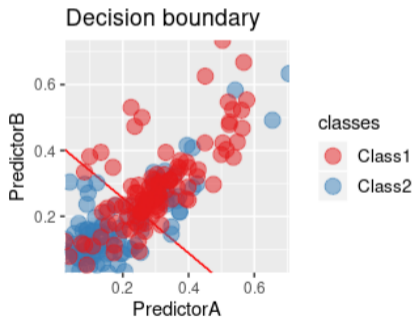
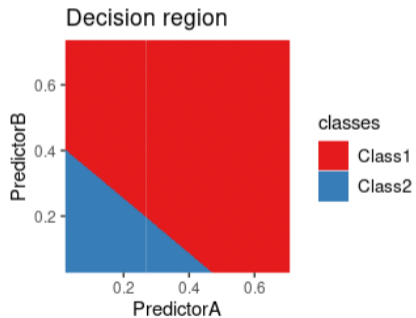
- Classical algorithm using a crude modeling for  $\mathbb{P}(\underline{X}|Y)$ :
  - Feature **independence** assumption:

$$\mathbb{P}(\underline{X}|Y) = \prod_{l=1}^d \mathbb{P}(X^{(l)}|Y)$$

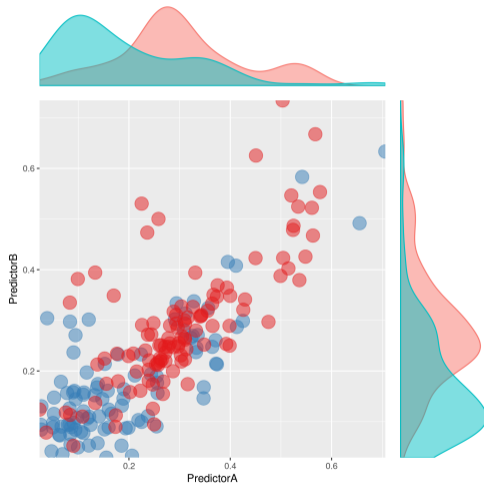
- Simple featurewise model: binomial if binary, multinomial if finite and Gaussian if continuous
- If all features are continuous, similar to the previous Gaussian but with a **diagonal covariance matrix!**
- Very simple learning even in **very high dimension!**

# Example: Naive Bayes

Naive Bayes with Gaussian model



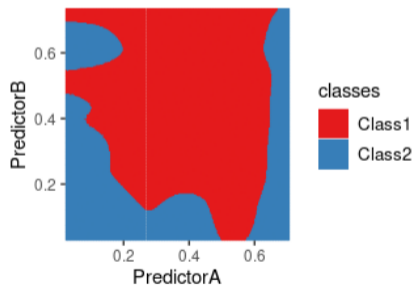
# Naive Bayes with Density Estimation



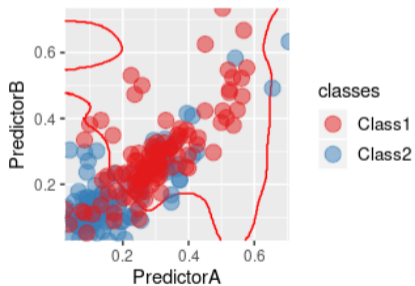
# Example: Naive Bayes

Naive Bayes with kernel density estimates

Decision region



Decision boundary



- Other models of the world!

## Bayesian Approach

- Generative Model plus prior on the parameters
- Inference thanks to the Bayes formula

## Graphical Models

- Markov type models on Graphs

## Gaussian Processes

- Multivariate Gaussian models

- ...

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

# Probabilistic and Optimization Framework

How to find a good function  $f$  with a *small* risk

$$\mathcal{R}(f) = \mathbb{E} [\ell(Y, f(\underline{X}))] \quad ?$$

**Canonical approach:**  $\hat{f}_S = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i))$

## Problems

- How to choose  $\mathcal{S}$ ?
- How to compute the minimization?

## A Probabilistic Point of View

**Solution:** For  $\underline{X}$ , estimate  $Y|\underline{X}$  plug this estimate in the Bayes classifier:  
**(Generalized) Linear Models, Kernel methods,  $k$ -nn, Naive Bayes, Tree, Bagging...**

## An Optimization Point of View

**Solution:** If necessary replace the loss  $\ell$  by an upper bound  $\ell'$  and minimize the empirical loss: **SVR, SVM, Neural Network, Tree, Boosting...**

## Penalized Logistic Regression

- Let  $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$  with  $\theta = (\beta, \beta^{(0)})$ .
- Find  $\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-Y_i f_{\theta}(\underline{X}_i)}) + \lambda \|\beta\|_1$
- Classify using  $\text{sign}(f_{\hat{\theta}})$

## ~~Penalized Logistic Regression~~

Support Vector Machine

- Let  $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$  with  $\theta = (\beta, \beta^{(0)})$ .
- Find  $\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i f_{\theta}(\underline{X}_i), 0) + \lambda \|\beta\|_2^2$
- Classify using  $\text{sign}(f_{\hat{\theta}})$

→ Examine RKHs  
Choix du noyau



## Deep Learning

- Let  $f_{\theta}(\underline{X})$  with  $f$  a feed forward neural network outputting two values with a softmax layer as a last layer.

- Optimize by gradient descent the cross-entropy  $\frac{1}{n} \sum_{i=1}^n \log (f_{\theta}(\underline{X}_i)^{(Y_i)})$

- Classify using  $\text{sign}(f_{\hat{\theta}})$

- Those three methods rely on a similar heuristic: the optimization point of view!

- The best solution  $f^*$  is the one minimizing

$$f^* = \arg \min R(f) = \arg \min \mathbb{E} [\ell(Y, f(\underline{X}))]$$

## Empirical Risk Minimization

- One restricts  $f$  to a subset of functions  $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the average empirical loss

$$\hat{f} = \hat{f}_\theta = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Intractable for the  $\ell^{0/1}$  loss!

## Risk Convexification

- Replace the loss  $\ell(Y, f_\theta(\underline{X}))$  by a convex upperbound  $\ell'(Y, f_\theta(\underline{X}))$  (surrogate loss).
- Minimize the average of the surrogate empirical loss

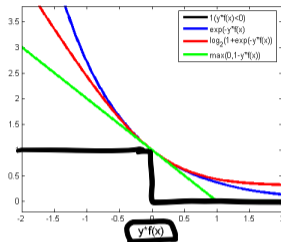
$$\tilde{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell'(Y_i, f_\theta(\underline{X}_i))$$

$$b: \mathcal{X} \rightarrow \{-1, 1\} \quad | \quad \tilde{b}: \mathcal{X} \rightarrow \mathbb{R}$$

- Use  $\hat{f} = \operatorname{sign}(\tilde{f})$
- Much easier optimization.

## Instantiation

- Logistic (Revisited)
- Support Vector Machine
- (Deep) Neural Network
- Boosting



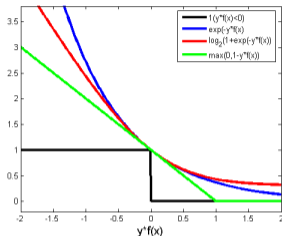
## Convexification

- Replace the loss  $\ell^{0/1}(Y, f(\underline{X}))$  by

$$\ell'(Y, f(\underline{X})) = l(Yf(\underline{X}))$$

with  $l$  a convex function.

- **Further mild assumption:**  $l$  is decreasing, differentiable at 0 and  $l'(0) < 0$ .



## Classical convexification

- Logistic loss:  $\ell'(Y, f(\underline{X})) = \log(1 + e^{-Yf(\underline{X})})$  (Logistic / NN)
- Hinge loss:  $\ell'(Y, f(\underline{X})) = (1 - Yf(\underline{X}))_+$  (SVM)
- Exponential loss:  $\ell'(Y, f(\underline{X})) = e^{-Yf(\underline{X})}$  (Boosting. . .)

## The Target is the Bayes Classifier

- The minimizer of

$$\mathbb{E} [\ell'(Y, f(\underline{X}))] = \mathbb{E} [I(Yf(\underline{X}))]$$

is the Bayes classifier  $f^* = \text{sign}(2\eta(\underline{X}) - 1)$

## Control of the Excess Risk

- It exists a convex function  $\Psi$  such that

$$\begin{aligned} \Psi \left( \mathbb{E} \left[ \ell^{0/1}(Y, \text{sign}(f(\underline{X}))) \right] - \mathbb{E} \left[ \ell^{0/1}(Y, f^*(\underline{X})) \right] \right) \\ \leq \mathbb{E} [\ell'(Y, f(\underline{X}))] - \mathbb{E} [\ell'(Y, f^*(\underline{X}))] \end{aligned}$$

- Theoretical guarantee!

- By definition,

$$\mathbb{E} [l(Yf)|\underline{X}] = \eta(\underline{X})l(f) + (1 - \eta(\underline{X}))l(-f)$$

- Let  $H(f, \eta) = \eta l(f) + (1 - \eta)l(-f)$ , the optimal value for  $\tilde{f}$  satisfies

$$\delta H(\tilde{f}, \eta) = -\eta \delta l(\tilde{f}) + (1 - \eta) \delta l(-\tilde{f}) \ni 0.$$

- With a slight abuse of notation, we denote by  $\delta l(\tilde{f})$  and  $\delta l(-\tilde{f})$  the two subgradients such that

$$\eta \delta l(\tilde{f}) - (1 - \eta) \delta l(-\tilde{f}) = 0$$

- Now we discuss the sign of  $\tilde{f}$ :

- If  $\tilde{f} > 0$ ,  $\delta l(-\tilde{f}) < \delta l(\tilde{f})$  and thus  $\eta > (1 - \eta)$ , i.e.  $2\eta - 1 > 0$ .
- Conversely, if  $\tilde{f} < 0$  then  $2\eta - 1 < 0$

- Thus  $\text{sign}(\tilde{f}) = \text{sign}(2\eta - 1)$  i.e. the minimizer of  $\mathbb{E} [l(yf)|\underline{X}]$  is  $f^*(\underline{X}) = \text{sign}(2\eta(\underline{X}) - 1)$

- We define  $H(\eta) = \inf_f H(f, \eta) = \inf_f (\eta l(f) + (1 - \eta)l(-f))$ . By construction,  $H$  is a concave function satisfying  $H(1/2 + x) = H(1/2 - x)$ .
- Furthermore, one verify that if we consider the minimum over the *wrong sign classifiers*,  $\inf_{f, f(2\eta-1) < 0} H(f, \eta) = l(0)$ .
- Indeed,

$$\begin{aligned} & \inf_{f, f(2\eta-1) < 0} H(f, \eta) \\ &= \inf_{f, f(2\eta-1) < 0} (\eta l(f) + (1 - \eta)l(-f)) \\ &\geq \inf_{f, f(2\eta-1) < 0} (\eta(l(0) + l'(0)f) + (1 - \eta)(l(0) - l'(0)f)) \\ &\geq l(0) + \inf_{f, f(2\eta-1) < 0} l'(0)f(2\eta - 1) = l(0) \end{aligned}$$



- Furthermore,

$$\begin{aligned}\mathbb{E} [\ell'(Y, f(\underline{X}))] &= \mathbb{E}_{\underline{X}} [H(f, \eta(\underline{X}))] \\ \mathbb{E} [\ell'(Y, f^*(\underline{X}))] &= \mathbb{E}_{\underline{X}} [H(\eta(\underline{X}))]\end{aligned}$$

- We define then

$$\Psi(\theta) = l(0) - H\left(\frac{1+\theta}{2}\right)$$

which is thus a convex function satisfying  $\Psi(0) = 0$  and  $\Psi(\theta) > 0$  for  $\theta > 0$ .

- Recall that

$$\begin{aligned}\mathbb{E} \left[ \ell^{0/1}(Y, \text{sign}(f(\underline{X}))) \right] &- \mathbb{E} \left[ \ell^{0/1}(Y, f^*(\underline{X})) \right] \\ &= \mathbb{E}_{\underline{X}} \left[ |2\eta(\underline{X}) - 1| \mathbf{1}_{f^*(\underline{X}) \neq \text{sign}(f(\underline{X}))} \right]\end{aligned}$$

- Using Jensen inequality, we derive

$$\begin{aligned}\Psi \left( \mathbb{E} \left[ \ell^{0/1}(Y, \text{sign}(f(\underline{X}))) \right] - \mathbb{E} \left[ \ell^{0/1}(Y, f^*(\underline{X})) \right] \right) \\ \leq \mathbb{E}_{\underline{X}} \left[ \Psi \left( |2\eta(\underline{X}) - 1| \mathbf{1}_{f^*(\underline{X}) \neq \text{sign}(f(\underline{X}))} \right) \right]\end{aligned}$$

- Using  $\Psi(0) = 0$  and the symmetry of  $H$ ,

$$\begin{aligned}
 & \Psi \left( \mathbb{E} \left[ \ell^{0/1}(Y, \text{sign}(f(\underline{X}))) \right] - \mathbb{E} \left[ \ell^{0/1}(Y, f^*(\underline{X})) \right] \right) \\
 & \leq \mathbb{E}_{\underline{X}} \left[ \left( l(0) - H \left( \left( \frac{1 + |2\eta(\underline{X}) - 1|}{2} \right) \right) \right) \mathbf{1}_{f^*(\underline{X}) \neq \text{sign}(f(\underline{X}))} \right] \\
 & \leq \mathbb{E}_{\underline{X}} \left[ (l(0) - H(\eta(\underline{X}))) \mathbf{1}_{f^*(\underline{X}) \neq \text{sign}(f(\underline{X}))} \right] \\
 & \leq \mathbb{E}_{\underline{X}} \left[ (l(0) - H(\eta(\underline{X}))) \mathbf{1}_{f(\underline{X})(2\eta(\underline{X})-1) < 0} \right]
 \end{aligned}$$

- Using the property of the wrong sign classifiers

$$\begin{aligned}
 & \Psi \left( \mathbb{E} \left[ \ell^{0/1}(Y, \text{sign}(f(\underline{X}))) \right] - \mathbb{E} \left[ \ell^{0/1}(Y, f^*(\underline{X})) \right] \right) \\
 & \leq \mathbb{E}_{\underline{X}} \left[ (H(f, \eta(\underline{X})) - H(f^*, \eta(\underline{X}))) \mathbf{1}_{f(\underline{X})(2\eta(\underline{X})-1) < 0} \right] \\
 & \leq \mathbb{E}_{\underline{X}} \left[ (H(f, \eta(\underline{X})) - H(f^*, \eta(\underline{X}))) \right] \\
 & \leq \mathbb{E} \left[ \ell'(Y, f(\underline{X})) \right] - \mathbb{E} \left[ \ell'(Y, f^*(\underline{X})) \right]
 \end{aligned}$$

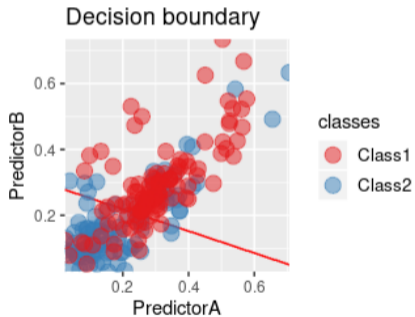
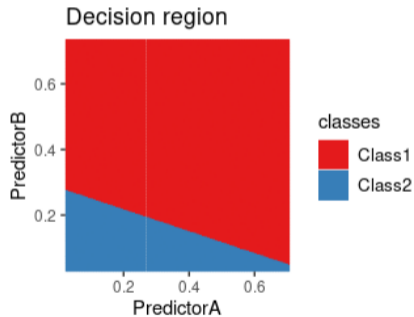
- Ideal solution:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

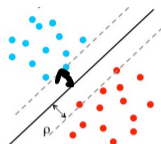
## Logistic regression

- Use  $f(\underline{X}) = \underline{X}^\top \beta + \beta^{(0)}$ .
- Use the logistic loss  $\ell(y, f) = \log_2(1 + e^{-yf})$ , i.e. the -log-likelihood.
- Different vision than the statistician but same algorithm!

## Logistic



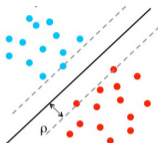
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



- Linear classifier:  $\text{sign}(\underline{X}^\top \beta + \beta^{(0)})$
- Separable case:  $\exists(\beta, \beta^{(0)}), \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) > 0!$

How to choose  $(\beta, \beta^{(0)})$  so that the separation is maximal?

- Strict separation:  $\exists(\beta, \beta^{(0)}), \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1$
- Distance between  $\underline{X}^\top \beta + \beta^{(0)} = 1$  and  $\underline{X}^\top \beta + \beta^{(0)} = -1$ :  
$$\frac{2}{\|\beta\|}$$
- Maximizing this distance is equivalent to minimizing  $\frac{1}{2}\|\beta\|^2$ .



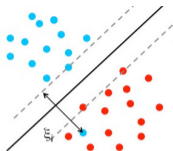
## Separable SVM

- Constrained optimization formulation:

$$\min \frac{1}{2} \|\beta\|^2 \quad \text{with} \quad \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1$$

- Quadratic Programming setting.
- Efficient solver available. . .





- What about the non separable case?

## SVM relaxation

- Relax the assumptions

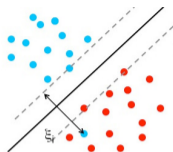
$$\forall i, Y_i(\underline{X}_i^T \beta + \beta^{(0)}) \geq 1 \quad \text{to} \quad \forall i, Y_i(\underline{X}_i^T \beta + \beta^{(0)}) \geq 1 - s_i$$

with the **slack variables**  $s_i \geq 0$

- Keep those slack variables as small as possible by minimizing

$$\frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i$$

where  $C > 0$  is the **goodness-of-fit strength**



## SVM

- Constrained optimization formulation:

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- **Hinge Loss** reformulation:

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \underbrace{\max(0, 1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}))}_{\text{Hinge Loss}}$$

- Constrained convex optimization algorithms vs gradient descent algorithms.

- Convex relaxation:

$$\operatorname{argmin} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0)$$

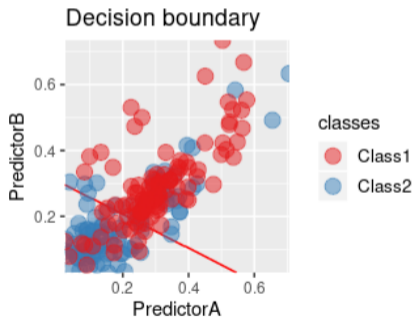
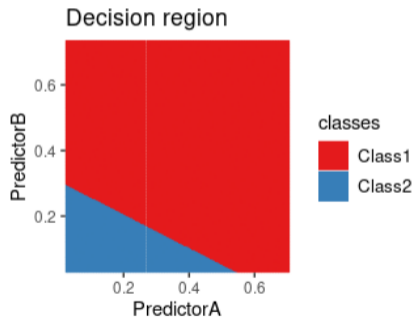
$$= \operatorname{argmin} \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0) + \frac{1}{Cn} \frac{1}{2} \|\beta\|^2$$

- **Prop:**  $\ell^{0/1}(Y_i, \operatorname{sign}(\underline{X}_i^\top \beta + \beta^{(0)})) \leq \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0)$

## Penalized convex relaxation (Tikhonov!)

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, \operatorname{sign}(\underline{X}_i^\top \beta + \beta^{(0)})) \\ & \leq \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0) + \frac{1}{Cn} \frac{1}{2} \|\beta\|^2 \end{aligned}$$

## Support Vector Machine



## Constrained Minimization

- Goal:

$$\min_x f(x)$$

$$\text{with } \begin{cases} h_j(x) = 0, & j = 1, \dots, p \\ g_i(x) \leq 0, & i = 1, \dots, q \end{cases}$$

- or rather with argmin!

## Different Setting

- $f, h_j, g_i$  **differentiable**
- $f$  **convex**,  $h_j$  **affine** and  $g_i$  **convex**.

## Feasibility

- $x$  is **feasible** if  $h_j(x) = 0$  and  $g_i(x) \leq 0$ .
- **Rk:** The set of feasible points may be empty

## Constrained Minimization

- Goal:

$$p^* = \min_x f(x) \quad \text{with} \quad \begin{cases} h_j(x) = 0, & j = 1, \dots, p \\ g_i(x) \leq 0, & i = 1, \dots, q \end{cases}$$

## Lagrangian

- **Def:**

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{j=1}^p \lambda_j h_j(x) + \sum_{i=1}^q \mu_i g_i(x)$$

with  $\lambda \in \mathbb{R}^p$  and  $\mu \in (\mathbb{R}^+)^q$ .

- The  $\lambda_j$  and  $\mu_i$  are called the dual (or Lagrange) variables.

- **Prop:**

$$\max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \mathcal{L}(x, \lambda, \mu) = \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ +\infty & \text{otherwise} \end{cases}$$

$$\min_x \max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \mathcal{L}(x, \lambda, \mu) = p^*$$

## Lagrangian

- **Def:**

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{j=1}^p \lambda_j h_j(x) + \sum_{i=1}^q \mu_i g_i(x)$$

with  $\lambda \in \mathbb{R}^p$  and  $\mu \in (\mathbb{R}^+)^q$ .

## Lagrangian Dual

- Lagrangian dual function:

$$Q(\lambda, \mu) = \min_x \mathcal{L}(x, \lambda, \mu)$$

- **Prop:**

$$Q(\lambda, \mu) \leq f(x), \text{ for all feasible } x$$
$$\max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} Q(\lambda, \mu) \leq \min_{x \text{ feasible}} f(x)$$

## Primal

- Primal:

$$p^* = \min_{x \in \mathcal{X}} f(x) \text{ with } \begin{cases} h_j(x) = 0, & j = 1, \dots, p \\ g_i(x) \leq 0, & i = 1, \dots, q \end{cases}$$

## Dual

- Dual:

$$q^* = \max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} Q(\lambda, \mu) = \max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \min_x \mathcal{L}(x, \lambda, \mu)$$

## Duality

- Always **weak duality**:

$$q^* \leq p^*$$
$$\max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \min_x \mathcal{L}(x, \lambda, \mu) \leq \min_x \max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \mathcal{L}(x, \lambda, \mu)$$

- Not always strong duality  $q^* = p^*$ .



## Strong Duality

- **Strong duality:**

$$q^* = p^*$$
$$\max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \min_x \mathcal{L}(x, \lambda, \mu) = \min_x \max_{\lambda \in \mathbb{R}^p, \mu \in (\mathbb{R}^+)^q} \mathcal{L}(x, \lambda, \mu)$$

- Allow to compute the solution of one problem from the other.
- Requires some assumptions!

## Strong Duality under Convexity and Slater's Condition

- $f$  **convex**,  $h_j$  **affine** and  $g_i$  **convex**.
- **Slater's condition:** it exists a feasible point such that  $h_j(x) = 0$  for all  $j$  and  $g_i(x) < 0$  for all  $i$ .
- Sufficient to prove **strong duality**.
- **Rk:** If the  $g_i$  are affine, it suffices to have  $h_j(x) = 0$  for all  $j$  and  $g_i(x) \leq 0$  for all  $i$ .

## Karush-Kuhn-Tucker Condition

- Stationarity:

$$\nabla_x \mathcal{L}(x^*, \lambda, \mu) = \nabla f(x^*) + \sum_j \lambda_j \nabla h(x^*) + \sum_i \mu_i \nabla g(x^*) = 0$$

- Primal admissibility:

$$h_j(x^*) = 0 \quad \text{and} \quad g_i(x^*) \leq 0$$

- Dual admissibility:

$$\mu_i \geq 0$$

- Complementary slackness:

$$\mu_i g_i(x^*) = 0$$

## KKT Theorem

- If  $f$  **convex**,  $h_j$  **affine** and  $g_i$  **convex**, all are differentiable and **strong duality** holds then  $x^*$  is a **solution** of the primal problem **if and only if** the **KKT condition holds**

## SVM

- Constrained optimization formulation:

$$\min \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

## SVM Lagrangian

- Lagrangian:

$$\begin{aligned} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \\ & + \sum_i \alpha_i (1 - s_i - Y_i(\underline{X}_i^\top \beta + \beta^{(0)})) - \sum_i \mu_i s_i \end{aligned}$$

## KKT Optimality Conditions

- Stationarity:

$$\nabla_{\beta} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = \beta - \sum_i \alpha_i Y_i \underline{X}_i = 0$$

$$\nabla_{\beta^{(0)}} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = - \sum_i \alpha_i = 0$$

$$\nabla_{s_i} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = C - \alpha_i - \mu_i = 0$$

- Primal and dual admissibility:

$$(1 - s_i - Y_i(\underline{X}_i^{\top} \beta + \beta^{(0)})) \leq 0, \quad s_i \geq 0, \quad \alpha_i \geq 0, \quad \text{and} \quad \mu_i \geq 0$$

- Complementary slackness:

$$\alpha_i(1 - s_i - Y_i(\underline{X}_i^{\top} \beta + \beta^{(0)})) = 0 \quad \text{and} \quad \mu_i s_i = 0$$

## Consequence

- $\beta^* = \sum_i \alpha_i Y_i \underline{X}_i$  and  $0 \leq \alpha_i \leq C$ .
- If  $\alpha_i \neq 0$ ,  $\underline{X}_i$  is called a **support vector** and either
  - $s_i = 0$  and  $Y_i(\underline{X}_i^{\top} \beta + \beta^{(0)}) = 1$  (margin hyperplane),
  - or  $\alpha_i = C$  (outliers).
- $\beta^{(0)*} = Y_i - \underline{X}_i^{\top} \beta^*$  for any support vector with  $0 < \alpha_i < C$ .

## SVM Lagrangian Dual

- Lagrangian Dual:

$$Q(\alpha, \mu) = \min_{\beta, \beta^{(0)}, s} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu)$$

- Prop:

- if  $\sum_i \alpha_i Y_i \neq 0$  or  $\exists i, \alpha_i + \mu_i \neq C$ ,

$$Q(\alpha, \mu) = -\infty$$

- if  $\sum_i \alpha_i Y_i = 0$  and  $\forall i, \alpha_i + \mu_i = C$ ,

$$Q(\alpha, \mu) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \underline{X}_i^\top \underline{X}_j$$

## SVM Dual problem

- Dual problem is a Quadratic Programming problem:

$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \underline{X}_i^\top \underline{X}_j$$

- Involves the  $\underline{X}_i$  only through their scalar products.

## Mercer Representation Theorem

- For any loss  $\ell$  and any increasing function  $\Phi$ , the minimizer in  $\beta$  of

$$\sum_{i=1}^n \ell(Y_i, \underline{X}_i^\top \beta + \beta^{(0)}) + \Phi(\|\beta\|_2)$$

is a linear combination of the input points  $\beta^* = \sum_{i=1}^n \alpha'_i \underline{X}_i$ .

- Minimization problem in  $\alpha'$ :

$$\sum_{i=1}^n \ell(Y_i, \sum_j \alpha'_j \underline{X}_i^\top \underline{X}_j + \beta^{(0)}) + \Phi(\|\beta\|_2)$$

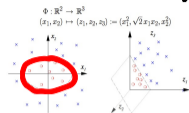
involving only the scalar product of the data.

- Optimal predictor requires only to compute scalar products.

$$\hat{f}^*(\underline{X}) = \underline{X}^\top \beta^* + \beta^{(0),*} = \sum_i \alpha'_i \underline{X}_i^\top \underline{X}$$

- Transform a problem in dimension  $\dim(\mathcal{X})$  in a problem in dimension  $n$ .
- Direct minimization in  $\beta$  can be more efficient...

$$(x_1, x_2) \rightarrow (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$

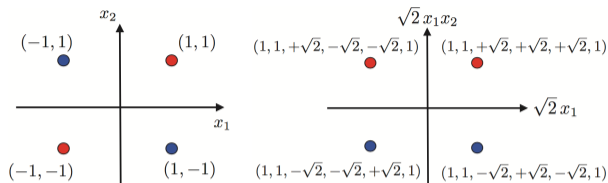


## Feature Engineering

- Art of creating **new features** from the existing one  $\underline{X}$ .
- Example: add monomials  $(\underline{X}^{(j)})^2, \underline{X}^{(j)} \underline{X}^{(j')}, \dots$
- Adding feature increases the dimension.

## Feature Map

- Application  $\phi: \mathcal{X} \rightarrow \mathbb{H}$  with  $\mathbb{H}$  an Hilbert space.
- Linear decision boundary in  $\mathbb{H}$ :  $\phi(\underline{X})^\top \beta + \beta^{(0)} = 0$  is **not an hyperplane anymore** in  $\mathcal{X}$ .
- **Heuristic:** Increasing dimension allows to make data almost linearly separable.



## Polynomial Mapping of order 2

- $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$

$$\phi(\underline{X}) = \left( (\underline{X}^{(1)})^2, (\underline{X}^{(2)})^2, \sqrt{2}\underline{X}^{(1)}\underline{X}^{(2)}, \sqrt{2}\underline{X}^{(1)}, \sqrt{2}\underline{X}^{(2)}, 1 \right)$$

- Allow to solve the XOR classification problem with the *hyperplane*  $\underline{X}^{(1)}\underline{X}^{(2)} = 0$ .

## Polynomial Mapping and Scalar Product

- **Prop:**

$$\phi(\underline{X})^\top \phi(\underline{X}') = (1 + \underline{X}^\top \underline{X}')^2$$



## Primal, Lagrangian and Dual

- Primal:

$$\min \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\phi(\underline{X}_i)^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- Lagrangian:

$$\begin{aligned} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) &= \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \\ &\quad + \sum_i \alpha_i (1 - s_i - Y_i(\phi(\underline{X}_i)^\top \beta + \beta^{(0)})) - \sum_i \mu_i s_i \end{aligned}$$

- Dual:

$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \phi(\underline{X}_i)^\top \phi(\underline{X}_j)$$

- Optimal  $\phi(\underline{X})^\top \beta^* + \beta^{(0),*} = \sum_i \alpha_i Y_i \phi(\underline{X})^\top \phi(\underline{X}_i)$

- Only need to know to compute  $\phi(\underline{X})^\top \phi(\underline{X}')$  to obtain the solution.

- Many algorithms (e.g. SVM) require only to be able to compute the scalar product  $\phi(\underline{X})^\top \phi(\underline{X}')$ .

## Kernel

- Any application

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

is called a **kernel** over  $\mathcal{X}$ .

## Kernel Trick

- Computing directly the **kernel**  $k(x, x') = \phi(\underline{X})^\top \phi(\underline{X}')$  may be easier than computing  $\phi(\underline{X})$ ,  $\phi(\underline{X}')$  and then the scalar product.
- Here  $k$  is defined from  $\phi$ .
- Under some assumption on  $k$ ,  $\phi$  can be implicitly *defined* from  $k$ !

## Positive Definite Symmetric Kernels

- A kernel  $k$  is PDS if and only if
  - $k$  is symmetric, i.e.

$$k(\underline{X}, \underline{X}') = k(\underline{X}', \underline{X})$$

- for any  $N \in \mathbb{N}$  and any  $(\underline{X}_1, \dots, \underline{X}_N) \in \mathcal{X}^N$ ,

$$\mathbf{K} = [k(\underline{X}_i, \underline{X}_j)]_{1 \leq i, j \leq N}$$

is positive semi-definite, i.e.  $\forall u \in \mathbb{R}^N$

$$u^T \mathbf{K} u = \sum_{1 \leq i, j \leq N} u^{(i)} u^{(j)} k(\underline{X}_i, \underline{X}_j) \geq 0$$

or equivalently all the eigenvalues of  $\mathbf{K}$  are non-negative.

- The matrix  $\mathbf{K}$  is called the **Gram matrix** associated to  $(\underline{X}_1, \dots, \underline{X}_N)$ .

## Mercer Theorem

- For any PDS kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , it exists a Hilbert space  $\mathbb{H} \subset \mathbb{R}^{\mathcal{X}}$  with a scalar product  $\langle \cdot, \cdot \rangle_{\mathbb{H}}$  such that

- it exists a mapping  $\phi : \mathcal{X} \rightarrow \mathbb{H}$  satisfying

$$k(\underline{X}, \underline{X}') = \langle \phi(\underline{X}), \phi(\underline{X}') \rangle_{\mathbb{H}}$$

- the **reproducing property** holds, i.e. for any  $h \in \mathbb{H}$  and any  $\underline{X} \in \mathcal{X}$

$$h(\underline{X}) = \langle h, k(\underline{X}, \cdot) \rangle_{\mathbb{H}}.$$

$$\begin{aligned} \mathbb{H} : \alpha_i \sum_i \alpha_i k(\underline{X}_i, \underline{x}) & \Bigg| \Bigg| \sum \alpha_i k(\underline{X}_i, \underline{z}) \Bigg|_{\mathbb{H}} \\ & = \sum \alpha_i \alpha_j k(\underline{X}_i, \underline{X}_j) \end{aligned}$$

- By def.,  $\mathbb{H}$  is a **reproducing kernel Hilbert space** (RKHS).
- $\mathbb{H}$  is called the **feature space** associated to  $k$  and  $\phi$  the **feature mapping**.
- No unicity in general.
- **Rk:** if  $k(\underline{X}, \underline{X}') = \phi'(\underline{X})^T \phi'(\underline{X}')$  with  $\phi' : \mathcal{X} \rightarrow \mathbb{R}^p$  then
  - $\mathbb{H}$  can be chosen as  $\{\underline{X} \mapsto \phi'(\underline{X})^T \beta, \beta \in \mathbb{R}^p\}$  and  $\|\underline{X} \mapsto \phi'(\underline{X})^T \beta\|_{\mathbb{H}}^2 = \|\beta\|_2^2$ .
  - $\phi(\underline{X})(\underline{X}') = \underline{X}^T \underline{X}'$ .

## Separable Kernel

- For any function  $\Psi : \mathcal{X} \rightarrow \mathbb{R}$ ,  $k(\underline{X}, \underline{X}') = \Psi(\underline{X})\Psi(\underline{X}')$  is PDS.

## Kernel Stability

- For any PDS kernels  $k_1$  and  $k_2$ ,  $k_1 + k_2$  and  $k_1 k_2$  are PDS kernels.
- For any sequence of PDS kernels  $k_n$  converging pointwise to a kernel  $k$ ,  $k$  is a PDS kernel.
- For any PDS kernel  $k$  such that  $|k| \leq r$  and any power series  $\sum_n a_n z^n$  with  $a_n \geq 0$  and a convergence radius larger than  $r$ ,  $\sum_n a_n k^n$  is a PDS kernel.
- For any PDS kernel  $k$ , the renormalized kernel  $k'(\underline{X}, \underline{X}') = \frac{k(\underline{X}, \underline{X}')}{\sqrt{k(\underline{X}, \underline{X})k(\underline{X}', \underline{X}')}} is a PDS kernel.$
- Cauchy-Schwartz for  $k$  PDS:  $k(\underline{X}, \underline{X}')^2 \leq k(\underline{X}, \underline{X})k(\underline{X}', \underline{X}')$

## PDS Kernels

- Vanilla kernel:

$$k(\underline{X}, \underline{X}') = \underline{X}^\top \underline{X}'$$

- Polynomial kernel:

$$k(\underline{X}, \underline{X}') = (1 + \underline{X}^\top \underline{X}')^k$$

- Gaussian RBF kernel:

$$k(\underline{X}, \underline{X}') = \exp\left(-\gamma \|\underline{X} - \underline{X}'\|^2\right)$$

- Tanh kernel:

$$k(\underline{X}, \underline{X}') = \tanh(a \underline{X}^\top \underline{X}' + b)$$

- Most classical is the Gaussian RBF kernel...
- Lots of freedom to construct kernel for non classical data.

## Representer Theorem

- Let  $k$  be a PDS kernel and  $\mathbb{H}$  its corresponding RKHS, for any increasing function  $\Phi$  and any function  $L : \mathbb{R}^n \rightarrow \mathbb{R}$ , the optimization problem

$$\operatorname{argmin}_{h \in \mathbb{H}} L(h(\underline{X}_1), \dots, h(\underline{X}_n)) + \Phi(\|h\|)$$

admits only solutions of the form

$$\sum_{i=1}^n \alpha'_i k(\underline{X}_i, \cdot).$$

- Examples:
  - (kernelized) SVM
  - (kernelized) Penalized Logistic Regression (Ridge)
  - (kernelized) Penalized Regression (Ridge)

## Primal

- Constrained Optimization:

$$\min_{f \in \mathbb{H}, \beta^{(0)}, s} \|f\|_{\mathbb{H}}^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(f(\underline{X}_i) + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- Hinge loss:

$$\min_{f \in \mathbb{H}, \beta^{(0)}} \|f\|_{\mathbb{H}}^2 + C \sum_{i=1}^n \max(0, 1 - Y_i(f(\underline{X}_i) + \beta^{(0)}))$$

- Representer:

$$\min_{\alpha', \beta^{(0)}} \sum_{i,j} \alpha'_i \alpha'_j k(\underline{X}_i, \underline{X}_j) + C \sum_{i=1}^n \max(0, 1 - Y_i(\sum_j \alpha'_j k(\underline{X}_j, \underline{X}_i) + \beta^{(0)}))$$

## Dual

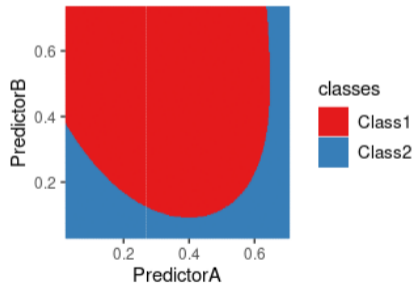
- Dual:

$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j k(\underline{X}_i, \underline{X}_j)$$

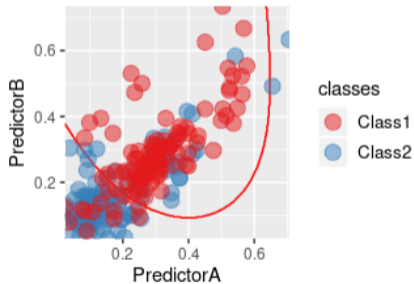


## Support Vector Machine with polynomial kernel

Decision region

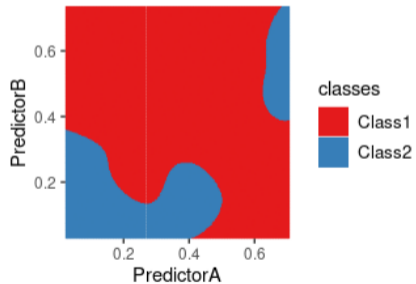


Decision boundary

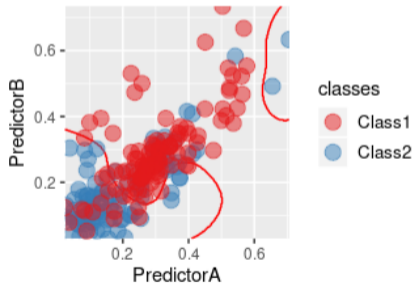


## Support Vector Machine with Gaussian kernel

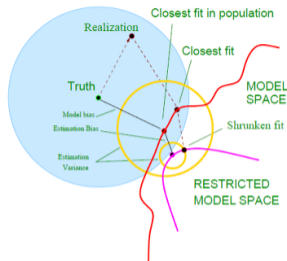
Decision region



Decision boundary



- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - **Penalization**
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



## Bias-Variance Issue

- Most complex models may not be the best ones due to the variability of the estimate.
- Naive idea: can we *simplify* our model without losing too much?
  - by using only a subset of the variables?
  - by forcing the coefficients to be small?
- Can we do better than exploring all possibilities?

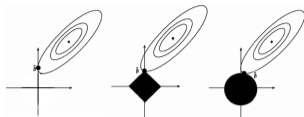
- **Setting:** Gen. linear model = prediction of  $Y$  by  $h(\underline{x}^T \beta)$ .

## Model coefficients

- Model entirely specified by  $\beta$ .
- Coefficientwise:
  - $\beta^{(i)} = 0$  means that the  $i$ th covariate is not used.
  - $\beta^{(i)} \sim 0$  means that the  $i$ th covariate as a *low* influence. . .
- If some covariates are useless, better use a simpler model. . .

## Submodels

- *Simplify* the model through a constraint on  $\beta$ !
- Examples:
  - Support: Impose that  $\beta^{(i)} = 0$  for  $i \notin I$ .
  - Support size: Impose that  $\|\beta\|_0 = \sum_{i=1}^d \mathbf{1}_{\beta^{(i)} \neq 0} < C$
  - Norm: Impose that  $\|\beta\|_p < C$  with  $1 \leq p$  (Often  $p = 2$  or  $p = 1$ )



## Sparsity

- $\beta$  is sparse if its number of non-zero coefficients ( $l_0$ ) is small. . .
- Easy interpretation in term of dimension/complexity.

## Norm Constraint and Sparsity

- Sparsest solution obtained by definition with the  $l_0$  norm.
- No induced sparsity with the  $l_2$  norm. . .
- Sparsity with the  $l_1$  norm (can even be proved to be the same than with the  $l_0$  norm under some assumptions).
- Geometric explanation.

## Constrained Optimization

- Choose a constant  $C$ .
- Compute  $\beta$  as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d, \|\beta\|_p \leq C} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(\underline{x}_i^\top \beta))$$

## Lagrangian Reformulation

- Choose  $\lambda$  and compute  $\beta$  as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(\underline{x}_i^\top \beta)) + \lambda \|\beta\|_p^{p'}$$

with  $p' = p$  except if  $p = 0$  where  $p' = 1$ .

- Easier calibration... but no explicit model  $\mathcal{S}$ .
- **Rk:**  $\|\beta\|_p$  is not scaling invariant if  $p \neq 0$ ...
- Initial rescaling issue.

## Penalized Linear Model

- Minimization of

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, h(\underline{x}_i^\top \beta)) + \operatorname{pen}(\beta)$$

where  $\operatorname{pen}(\beta)$  is a (sparsity promoting) penalty

- Variable selection if  $\beta$  is sparse.

## Classical Penalties

- AIC:  $\operatorname{pen}(\beta) = \lambda \|\beta\|_0$  (non convex / sparsity)
- Ridge:  $\operatorname{pen}(\beta) = \lambda \|\beta\|_2^2$  (convex / no sparsity)
- Lasso:  $\operatorname{pen}(\beta) = \lambda \|\beta\|_1$  (convex / sparsity)
- Elastic net:  $\operatorname{pen}(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$  (convex / sparsity)

- Easy optimization if  $\operatorname{pen}$  (and the loss) is convex...
- **Need to specify  $\lambda$  to define a ML method!**



## Classical Examples

- Penalized Least Squares
  - Penalized Logistic Regression
  - Penalized Maximum Likelihood
  - SVM
  - Tree pruning
- 
- Sometimes used even if the parametrization is not linear. . .

## Penalized $\ell^{0/1}$ loss (Structural Risk Minimization)

- Minimization of

$$\operatorname{argmin}_{f_m, m \in \mathcal{M}, f_m \in \mathcal{S}_m} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f_m(\underline{X}_i)) + \operatorname{pen}(m)$$

where  $\operatorname{pen}(m)$  is a complexity driven penalty...

- No easy optimization here!

## Classical Penalties

- Finite class:  $\operatorname{pen}(m) = \lambda \sqrt{\frac{\log |\mathcal{M}|}{n}}$
- Finite VC Dimension:  $\operatorname{pen}(m) = \lambda \sqrt{\frac{d_{VC}(\mathcal{S}_m) \log\left(\frac{en}{d_{VC}(\mathcal{S}_m)}\right)}{n}}$
- **Need to specify  $\lambda$ !**

## Penalized convexified $\ell$ loss

- Minimization of

$$\operatorname{argmin}_{f_m, m \in \mathcal{M}, f_m \in \mathcal{S}_m} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_m(\underline{X}_i)) + \operatorname{pen}(m)$$

where  $\operatorname{pen}(m)$  is a complexity driven penalty. . .

- Easy optimization here!
- **Reuse the previous  $\operatorname{pen}(m)$ !**
- **Need to specify  $\lambda$ !**
- SVM case:
  - $d_{VC} \sim \|\beta\|^2$  which advocates for a penalty in  $\lambda\|\beta\|$ . . .
  - A penalty in  $\lambda'\|\beta\|^2$  is more convenient numerically and there is a correspondence between the two problems. . .

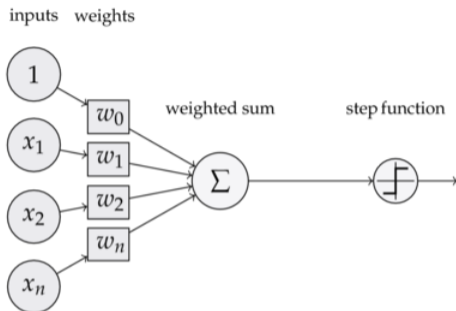
## Practical Selection Methodology

- Choose a penalty family  $\text{pen}_\lambda$ .
  - Compute a CV error for the penalty  $\text{pen}_\lambda$  for all  $\lambda \in \Lambda$ .
  - Determine  $\hat{\lambda}$  the  $\lambda$  minimizing the CV error.
  - Compute the final model with the penalty  $\text{pen}_{\hat{\lambda}}$ .
- CV allows to select a ML method, penalized estimation with a penalty  $\text{pen}_{\hat{\lambda}}$ , not a single predictor hence the need of a final reestimation.

## Why not using CV on a grid? *on $\beta$ ?*

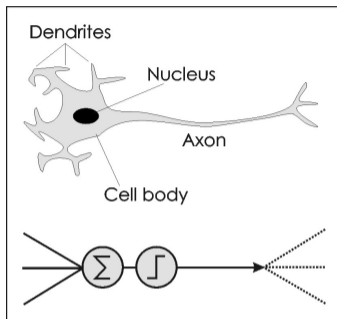
- Grid size scales exponentially with the dimension!
- **If the penalized minimization is easy**, much cheaper to compute the CV error for all  $\lambda \in \Lambda$ ...
- CV performs best when the set of candidates is not too big (or is structured...)

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - **(Deep) Neural Networks**
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



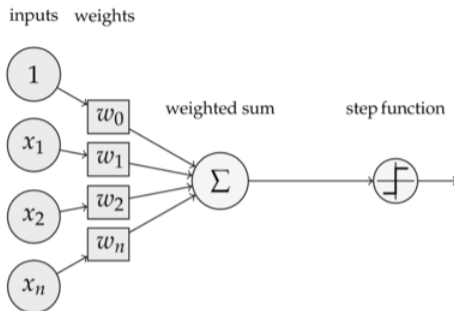
## Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



## Perceptron (Rosenblatt 1957)

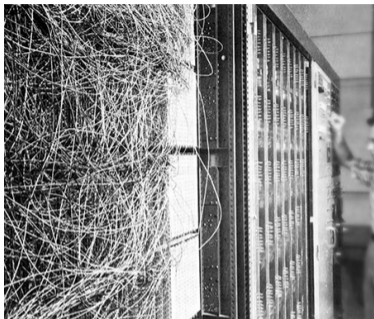
- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



## Perceptron (Rosenblatt 1957)

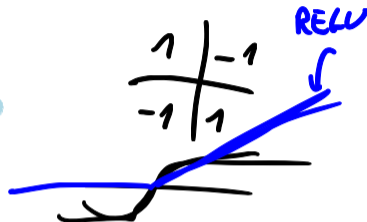
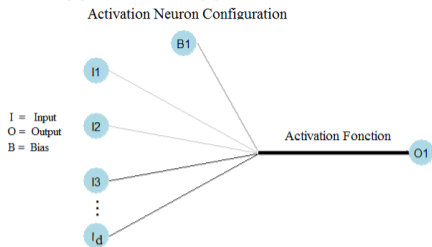
- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.





## Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



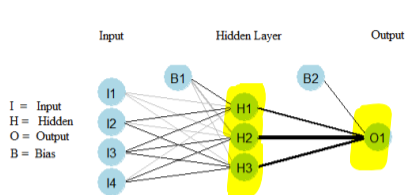
## Artificial neuron

- Structure:
  - Mix inputs with a **weighted sum**,
  - Apply a (non linear) **activation function** to this sum,
  - Possibly threshold the result to make a decision.
- Weights learned by minimizing a loss function.

## Logistic unit

- Structure:
  - Mix inputs with a **weighted sum**,
  - Apply the **logistic function**  
 $\sigma(t) = e^t / (1 + e^t)$ ,
  - Threshold at 1/2 to make a decision!
- Logistic weights learned by minimizing the -log-likelihood.

- Equivalent to linear regression when using a linear activation function!

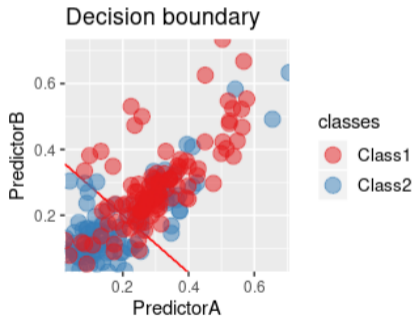
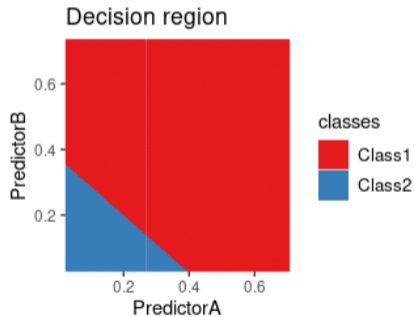


$$\frac{\partial}{\partial w_i'} \beta(z) = O_u(H_u(z)) = \sum_{k=1}^n \frac{\partial O_u}{\partial z_k} (H_{w'}(z)) \times \frac{\partial H_{w'}(z)}{\partial w_i'}$$

## MLP (Rumelhart, McClelland, Hinton - 1986)

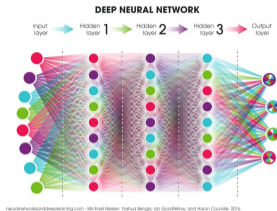
- Multilayer Perceptron: cascade of layers of artificial neuron units.
- Optimization through a gradient descent algorithm with a clever implementation (**Backprop**).
- Construction of a function by composing simple units.
- MLP corresponds to a specific direct acyclic graph structure.
- Non convex optimization problem!

## Neural Network



## Universal Approximation Theorem (Hornik, 1991)

- A **single hidden layer neural network** with a linear output unit can **approximate** any continuous function **arbitrarily well** given enough hidden units.
- Valid for most activation functions.
- No bounds on the number of required units. . . (Asymptotic flavor)
- A single hidden layer is sufficient but more may require less units.

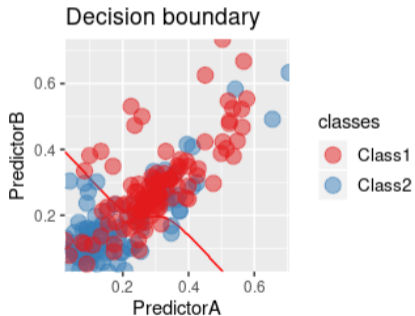
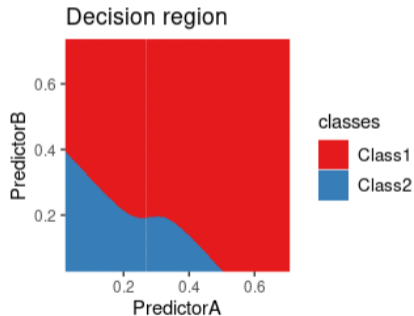


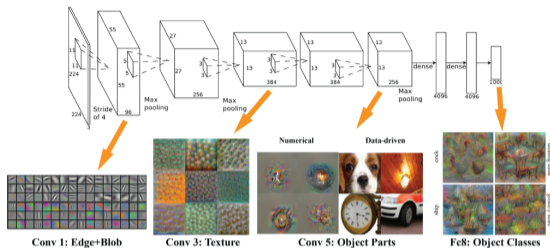
$$\mathcal{I} \rightarrow H_1(\mathcal{I}) \rightarrow H_2(H_1(\mathcal{I}))$$
$$\vdots$$
$$H_n(H_{n-1}(\dots H_1(\mathcal{I})))$$

## Deep Neural Network structure

- Deep cascade of layers!
- No conceptual novelty. . .
- But a **lot of tricks** allowing to obtain a good solution: clever initialization, better activation function, weight regularization, accelerated stochastic gradient descent, early stopping. . .
- Use of GPU and a lot of data. . .
- Very impressive results!

H2O NN





## Family of Machine Learning algorithm combining:

- a (deep) multilayered structure,
  - a clever optimization including initialization and regularization.
- 
- Examples: Deep NN, AutoEncoder, Recursive NN, GAN, Transformer...
  - Interpretation as a **Representation Learning**.
  - **Transfer learning**: use as initialization a pretrained net.
  - Very efficient and still evolving!



PROC. OF THE IEEE, NOVEMBER 1998

7

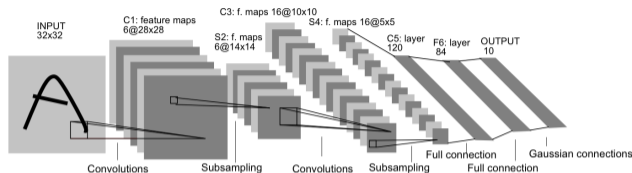
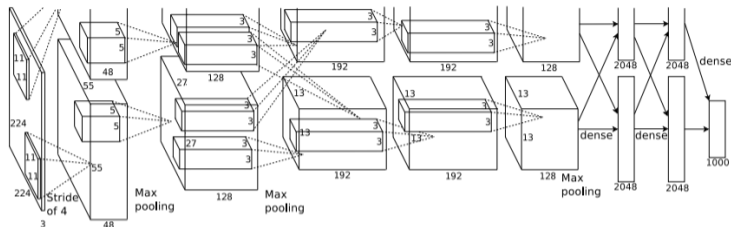


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

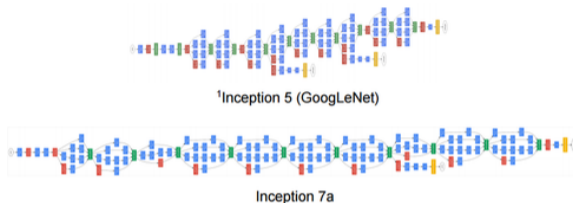
## Le Net - Y. LeCun (1989)

- 6 hidden layer architecture.
- Drastic reduction of the number of parameters through a translation invariance principle (convolution).
- Required 3 days of training for 60 000 examples!
- Tremendous improvement.
- Representation learned through the task.



## Alexnet - A. Krizhevsky, I. Sutskever, G. Hinton (2012)

- Bigger and deeper layers and thus much more parameters.
- Clever initialization scheme, RELU, renormalization and use of GPU.
- 6 days of training for 1.2 millions images.
- Tremendous improvement. . .



<sup>1</sup>Inception 5 (GoogLeNet)

Inception 7a

<sup>1</sup>Going Deeper with Convolutions, [C. Szegedy et al, CVPR 2015]



## Trends

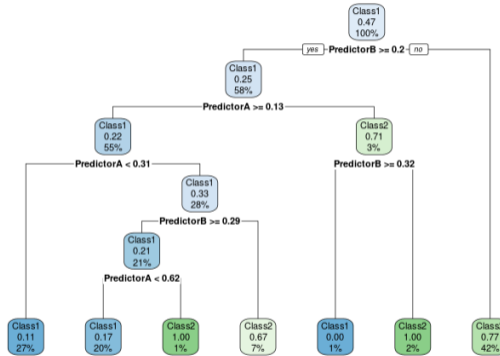
- Bigger and bigger networks! (GoogLeNet / Residual Neural Network / Transformers. . . )
- More computational power to learn better representation.
- Work in Progress!

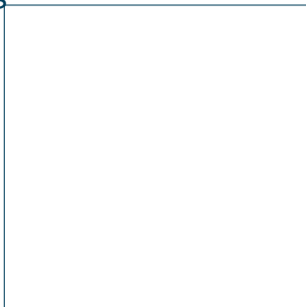
- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - **Tree Based Methods**
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



## Tree principle (CART by Breiman (85) / ID3 by Quinlan (86))

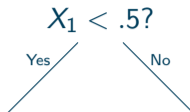
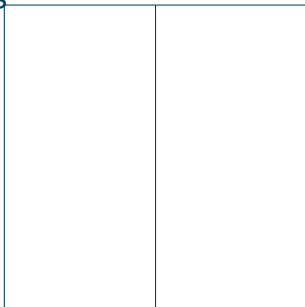
- Construction of a recursive partition through a tree structured set of questions (splits around a given value of a variable)
- For a given partition, probabilistic approach **and** optimization approach yields the same classifier!
- A simple majority vote/averaging in each leaf
- Quality of the prediction depends on the tree (the partition).
- **Intuitively:**
  - small leaves lead to low bias but large variance
  - large leaves lead to large bias but low variance. . .
- **Issue:** Minim. of the (penalized) empirical error is NP hard!
- Practical tree construction are all based on two steps:
  - a top-down step in which branches are created (branching)
  - a bottom-up in which branches are removed (pruning)





## Greedy top-bottom approach

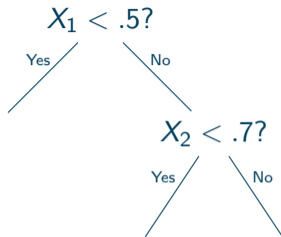
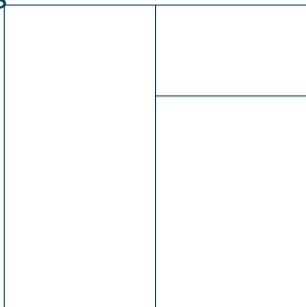
- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .



## Greedy top-bottom approach

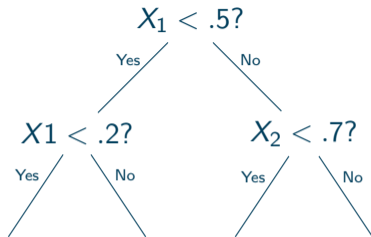
- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .





## Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible...



## Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .

## Various definition of *inhomogeneous*

- **CART:** empirical loss based criterion (least squares/prediction error)

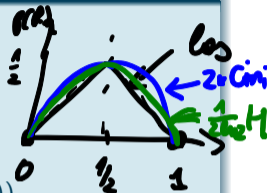
$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} \ell(y_i, y(R)) + \sum_{\underline{x}_i \in \bar{R}} \ell(y_i, y(\bar{R}))$$

- **CART:** Gini index (Classification)

$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} p(R)(1 - p(R)) + \sum_{\underline{x}_i \in \bar{R}} p(\bar{R})(1 - p(\bar{R}))$$

- **C4.5:** entropy based criterion (Information Theory)

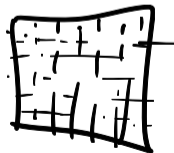
$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} H(R) + \sum_{\underline{x}_i \in \bar{R}} H(\bar{R})$$

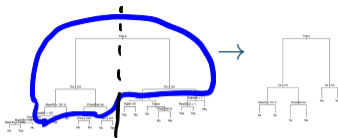


- CART with Gini is probably the most used technique. . .
- Other criterion based on  $\chi^2$  homogeneity or based on different local predictors (generalized linear models. . .)

## Choice of the split in a given region

- Compute the criterion for **all features and all possible splitting points** (necessarily among the data values in the region)
- Choose the one **minimizing** the criterion
- **Variations:** split at all categories of a categorical variables using a clever category ordering (ID3), split at a fixed position (median/mean)
- **Stopping rules:**
  - when a leaf/region contains less than a prescribed number of observations
  - when the region is sufficiently homogeneous. . .
- May lead to a quite complex tree: over-fitting possible!
- Additional pruning often use.





- **Model selection** within the (rooted) subtrees of previous tree!
- Number of subtrees can be quite large but the tree structure allows to find the best model efficiently.

## Key idea

- The predictor in a leaf depends only on the values in this leaf.
- **Efficient bottom-up (dynamic programming) algorithm** if the criterion used satisfies an additive property

$$C(\mathcal{T}) = \sum_{\mathcal{L} \in \mathcal{T}} c(\mathcal{L})$$

- Example: AIC / CV.

## Examples of criterion satisfying this assumptions

- AIC type criterion:

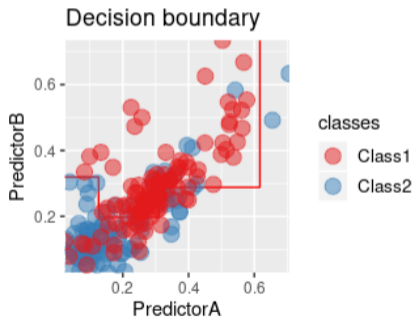
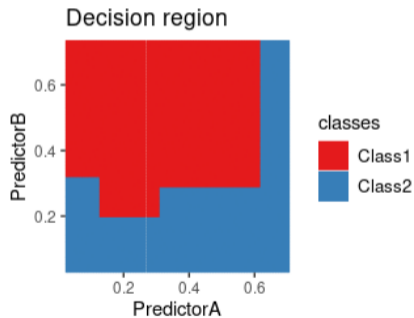
$$\sum_{i=1}^n \ell'(y_i, f_{\mathcal{L}(\underline{x}_i)}(\underline{x}_i)) + \lambda |\mathcal{T}| = \sum_{\mathcal{L} \in \mathcal{T}} \left( \sum_{\underline{x}_i \in \mathcal{L}} \ell'(y_i, f_{\mathcal{L}}(\underline{x}_i)) + \lambda \right)$$

- Simple cross-Validation (with  $(\underline{x}'_i, y'_i)$  a different dataset):

$$\sum_{i=1}^{n'} \ell'(y'_i, f_{\mathcal{L}(\underline{x}'_i)}(\underline{x}'_i)) = \sum_{\mathcal{L} \in \mathcal{T}} \left( \sum_{\underline{x}'_i \in \mathcal{L}} \ell'(y'_i, f_{\mathcal{L}}(\underline{x}'_i)) \right)$$

- Limit over-fitting for a single tree.
- **Rk:** almost never used when combining several trees. . .

## CART



## Pros

- Leads to a easily interpretable model
- Fast computation of the prediction
- Easily deals with categorical features (and missing values)

## Cons

- Greedy optimization
- Hard decision boundaries
- Lack of stability



- Lack of robustness for single trees.
- How to combine trees?

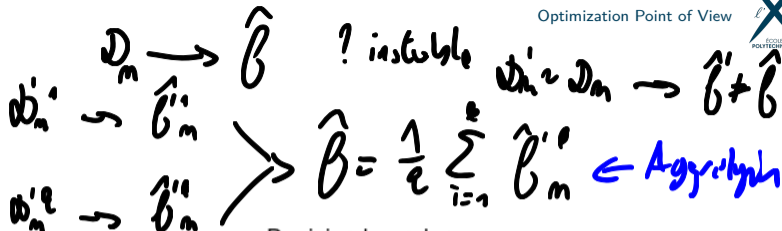
## Parallel construction

- Construct several trees from bootstrapped samples and average the responses (**Bagging**)
- Add more randomness in the tree construction (**Random Forests**)

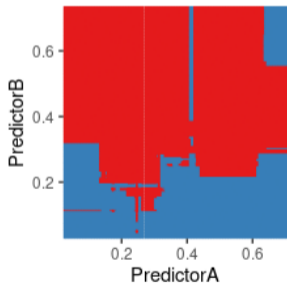
## Sequential construction

- Construct a sequence of trees by reweighting sequentially the samples according to their difficulties (**AdaBoost**)
- Reinterpretation as a stagewise additive model (**Boosting**)

# Ensemble methods



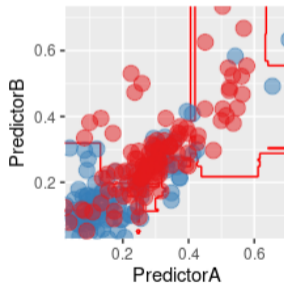
Decision region



classes

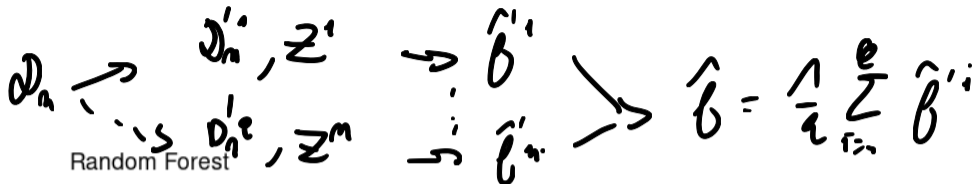
- Class1
- Class2

Decision boundary

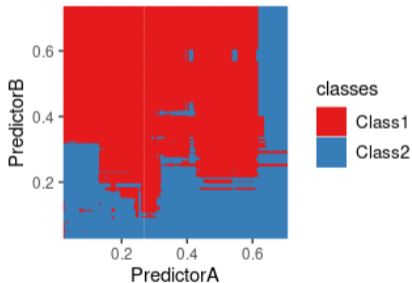


classes

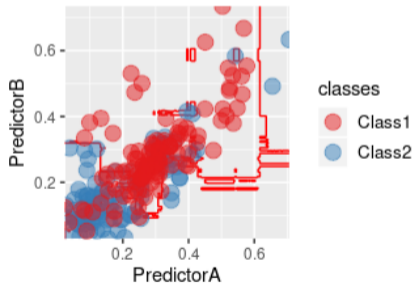
- Class1
- Class2



Decision region



Decision boundary



# Ensemble methods

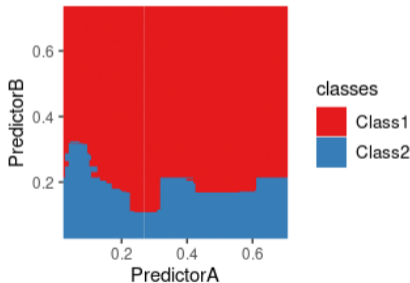
$$\mathcal{D}_n = (x, y) : x \rightarrow f_1(x) \in \mathcal{Y} \rightarrow \hat{\mathcal{D}}_n = (x, y - b_0(x)) : b_2(x) \approx y - b_1(x)$$

$$b_1(x) = \gamma \ln \frac{1}{b_1} \hat{\mathcal{D}}_n(x, y - b_0(x) - b_2(x))$$

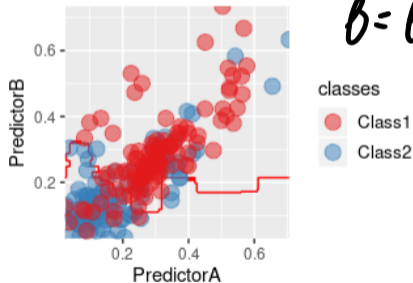
$$\hat{b} = b_1(x), b_2(x) \rightarrow b_3(x)$$

AdaBoost

Decision region



Decision boundary



## 1 Introduction

- Machine Learning
- Motivation

## 2 A Practical View

- Method or Models
- Interpretability
- Metric Choice

## 3 A Better Point of View

- The Example of Univariate Linear Regression
- Supervised Learning

## 4 Error Estimation and Method Choice

- Cross Validation
- Cross Validation and Test
- Cross Validation and Weights
- Auto ML

## 5 A Probabilistic Point of View

- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling

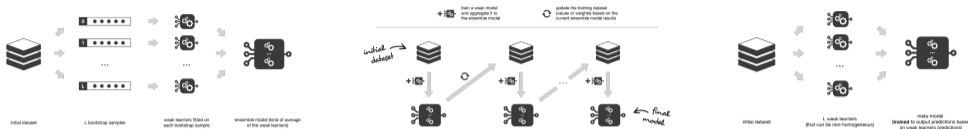
## 6 Optimization Point of View

- SVM
- Penalization
- (Deep) Neural Networks
- Tree Based Methods
- **Ensemble Methods**

## 7 Empirical Risk Minimization

- Empirical Risk Minimization
- ERM and PAC Bayesian Analysis
- Hoeffding and Finite Class
- McDiarmid and Rademacher Complexity
- VC Dimension
- Structural Risk Minimization

## 8 References



## Ensemble Methods

- **Averaging:** combine several models by averaging (bagging, random forests, ...)
  - **Boosting:** construct a sequence of (weak) classifiers (XGBoost, Lightgbm)
  - **Stacking:** use the outputs of several models as features (tpot...)
- 
- Loss of interpretability but gain in performance
  - Beware of overfitting with stacking: the second learning step should be done with fresh data.
  - No end to end optimization as in deep learning!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization**
  - **Empirical Risk Minimization**
  - **ERM and PAC Bayesian Analysis**
  - **Hoeffding and Finite Class**
  - **McDiarmid and Rademacher Complexity**
  - **VC Dimension**
  - **Structural Risk Minimization**
- 8 References

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 **Empirical Risk Minimization**
  - **Empirical Risk Minimization**
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References



## Empirical Risk Minimizer (ERM)

- For any loss  $\ell$  and function class  $\mathcal{S}$ ,

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i)) = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}_n(f)$$

- Key property:

$$\mathcal{R}_n(\hat{f}) \leq \mathcal{R}_n(f), \forall f \in \mathcal{S}$$

- **Minimization not always tractable in practice!**
- Focus on the  $\ell^{0/1}$  case:
  - only algorithm is to try all the functions,
  - not feasible if there are many functions
  - but interesting hindsight!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization**
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis**
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- Theoretical control of the random (error estimation) term:

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*)$$

## Probably Almost Correct Analysis

- **Theoretical guarantee** that

$$\mathbb{P} \left( \mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \epsilon_S(\delta) \right) \geq 1 - \delta$$

for a suitable  $\epsilon_S(\delta) \geq 0$ .

- Implies:

- $\mathbb{P} \left( \mathcal{R}(\hat{f}) - \mathcal{R}(f^*) \leq \mathcal{R}(f_S^*) - \mathcal{R}(f^*) + \epsilon_S(\delta) \right) \geq 1 - \delta$

- $\mathbb{E} \left[ \mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \right] \leq \int_0^{+\infty} \delta_S(\epsilon) d\epsilon$

- The result should hold without any assumption on the law ***P***!

- By construction:

$$\begin{aligned} \mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) &= \mathcal{R}(\hat{f}) - \mathcal{R}_n(\hat{f}) + \mathcal{R}_n(\hat{f}) - \mathcal{R}_n(f_S^*) + \mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*) \\ \text{???} \quad &\Rightarrow \underbrace{\mathcal{R}(\hat{f}) - \mathcal{R}_n(\hat{f})}_{\text{LGN}} + \underbrace{\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*)}_{\text{LGN}} \\ &\leq (\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*)) - (\mathcal{R}_n(\hat{f}) - \mathcal{R}_n(f_S^*)) \end{aligned}$$

## Four possible upperbounds

- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} ((\mathcal{R}(f) - \mathcal{R}(f_S^*)) - (\mathcal{R}_n(f) - \mathcal{R}_n(f_S^*)))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) + (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) + \sup_{f \in \mathcal{S}} (\mathcal{R}_n(f) - \mathcal{R}(f))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq 2 \sup_{f \in \mathcal{S}} |\mathcal{R}(f) - \mathcal{R}_n(f)|$

- Supremum of centered random variables!
- **Key:** Concentration of each variable...

- By construction, for any  $f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') = \mathcal{R}_n(f') + (\mathcal{R}(f') - \mathcal{R}_n(f'))$$

## A uniform upper bound for the error

- Simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f))$$

- Supremum of centered random variables!
- **Key:** Concentration of each variable. . .
- Can be interpreted as a justification of the ERM!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 **Empirical Risk Minimization**
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - **Hoeffding and Finite Class**
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- Empirical loss:

$$\mathcal{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

## Properties

- $\ell^{0/1}(Y_i, f(\underline{X}_i))$  are i.i.d. random variables in  $[0, 1]$ .

## Concentration

$$\mathbb{P}(\mathcal{R}(f) - \mathcal{R}_n(f) \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

$$\mathbb{P}(\mathcal{R}_n(f) - \mathcal{R}(f) \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

$$\mathbb{P}(|\mathcal{R}_n(f) - \mathcal{R}(f)| \leq \epsilon) \geq 1 - 2e^{-2n\epsilon^2}$$

- Concentration of sum of bounded independent variables!
- Hoeffding theorem.

## Theorem

- Let  $Z_i$  be a sequence of ind. centered r.v. supported in  $[a_i, b_i]$  then

$$\mathbb{P} \left( \sum_{i=1}^n Z_i \geq \epsilon \right) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- Proof ingredients:

- Chernov bounds:

$$\mathbb{P} \left( \sum_{i=1}^n Z_i \geq \epsilon \right) \leq \frac{\mathbb{E} \left[ e^{\lambda \sum_{i=1}^n Z_i} \right]}{e^{\lambda \epsilon}} \leq \frac{\prod_{i=1}^n \mathbb{E} \left[ e^{\lambda Z_i} \right]}{e^{\lambda \epsilon}}$$

- Exponential moment bounds:  $\mathbb{E} \left[ e^{\lambda Z_i} \right] \leq e^{\frac{\lambda^2 (b_i - a_i)^2}{8}}$
- Optimization in  $\lambda$

- Prop:**

$$\mathbb{E} \left[ e^{\lambda \sum_{i=1}^n Z_i} \right] \leq e^{\frac{\lambda^2 \sum_{i=1}^n (b_i - a_i)^2}{8}}.$$



- Exponential moment function:

$$\Psi_Z(\lambda) = \log \mathbb{E} \left[ e^{\lambda Z} \right]$$

$$\Psi_Z(0) = 0$$

$$\Psi'_Z(\lambda) = \frac{\mathbb{E} \left[ Z e^{\lambda Z} \right]}{\mathbb{E} \left[ e^{\lambda Z} \right]}$$

$$\Psi'_Z(0) = \frac{\mathbb{E} [Z]}{\Psi_Z(0)} = 0$$

$$\begin{aligned} \Psi''_Z(\lambda) &= \frac{\mathbb{E} \left[ Z^2 e^{\lambda Z} \right]}{\mathbb{E} \left[ e^{\lambda Z} \right]} - \frac{\left( \mathbb{E} \left[ Z e^{\lambda Z} \right] \right)^2}{\left( \mathbb{E} \left[ e^{\lambda Z} \right] \right)^2} \\ &= \mathbb{E} \left[ Z^2 e^{\lambda Z - \psi(\lambda)} \right] - \mathbb{E} \left[ Z e^{\lambda Z - \psi(\lambda)} \right]^2 \\ &= \text{Var} [Z'] \end{aligned}$$

with  $Z'$  a random variable with density  $e^{\lambda Z - \psi(\lambda)}$  with respect to  $dZ$ .

- Now as  $Z' \in [a, b]$ ,

$$\begin{aligned}\text{Var} [Z'] &= \mathbb{E} [(Z' - \mathbb{E} [Z'])^2] \\ &\leq \mathbb{E} [(Z' - (a + b)/2)^2] \leq (b - a)^2/4\end{aligned}$$

- As  $\Psi_Z(0) = 0$ ,  $\Psi'_Z(0) = 0$  and  $\Psi''_Z(\lambda) \leq (b - a)^2/4$ , Taylor formula ensures that  $\exists \theta \in [0, \lambda]$  such that

$$\Psi_Z(\lambda) = \frac{1}{2} \Psi''_Z(\theta) \lambda^2 \leq \frac{(b - a)^2}{8} \lambda^2$$

- This yields

$$\mathbb{E} [e^{\lambda Z}] \leq e^{\frac{(b-a)^2}{8} \lambda^2}$$

- Optimization in  $\lambda$ :

$$\frac{\prod_{i=1}^n \mathbb{E} \left[ e^{\lambda Z_i} \right]}{e^{\lambda \epsilon}} \leq e^{\frac{\lambda^2}{8} \sum_{i=1}^n (b_i - a_i)^2 - \lambda \epsilon}$$

is minimal for  $\lambda = 4\epsilon / (\sum_{i=1}^n (b_i - a_i)^2)$  for which

$$e^{\frac{\lambda^2}{8} \sum_{i=1}^n (b_i - a_i)^2 - \lambda \epsilon} = e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}.$$

## Theorem

- Let  $Z_i$  be a sequence of independent centered random variables supported in  $[a_i, b_i]$  then

$$\mathbb{P} \left( \sum_{i=1}^n Z_i \geq \epsilon \right) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- $Z_i = \frac{1}{n} \left( \mathbb{E} \left[ \ell^{0/1}(Y, f(\underline{X})) \right] - \ell^{0/1}(Y_i, f(\underline{X}_i)) \right)$
- $\mathbb{E} [Z_i] = 0$  and  $Z_i \in \left[ \frac{1}{n} \left( \mathbb{E} \left[ \ell^{0/1}(Y, f(\underline{X})) \right] - 1 \right), \frac{1}{n} \mathbb{E} \left[ \ell^{0/1}(Y, f(\underline{X})) \right] \right]$
- Concentration:

$$\mathbb{P} (\mathcal{R}(f) - \mathcal{R}_n(f) \geq \epsilon) \leq e^{-2n\epsilon^2}$$

- By symmetry,

$$\mathbb{P} (\mathcal{R}_n(f) - \mathcal{R}(f) \geq \epsilon) \leq e^{-2n\epsilon^2}$$

- Combining the two yields

$$\mathbb{P} (|\mathcal{R}_n(f) - \mathcal{R}(f)| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$$

## Concentration

- If  $\mathcal{S}$  is finite of cardinality  $|\mathcal{S}|$ ,

$$\mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) \leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \geq 1 - \delta$$

$$\mathbb{P} \left( \sup_f |\mathcal{R}_n(f) - \mathcal{R}(f)| \leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \geq 1 - 2\delta$$

- Control of the supremum by a quantity depending on the cardinality and the probability parameter  $\delta$ .
- Simple combination of Hoeffding and a union bound.

$$|\mathcal{I}| e^{-2n\epsilon^2} \leq \delta$$

$$\mathbb{P} \left( \sup_{i \in \mathcal{I}} z_i \geq \epsilon \right) \leq \sum_i \mathbb{P}(z_i \geq \epsilon) = |\mathcal{I}| \mathbb{P}(z \geq \epsilon)$$

$$e^{-2n\epsilon^2}$$

## PAC Bounds

- If  $\mathcal{S}$  is finite of cardinality  $|\mathcal{S}|$ , with proba greater than  $1 - 2\delta$

$$\begin{aligned}\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) &\leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}} \\ &\leq 2\sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}}\end{aligned}$$

- If  $\mathcal{S}$  is finite of cardinality  $|\mathcal{S}|$ , with proba greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\begin{aligned}\mathcal{R}(f') &\leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \\ &\leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}\end{aligned}$$

- Using Hoeffding,

- $\forall f \in \mathcal{S}$ ,

$$\mathbb{P} \left( (\mathcal{R}(f) - \mathcal{R}_n(f)) \geq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \leq e^{-2n \left( \frac{\log |\mathcal{S}| + \log(1/\delta)}{2n} \right)} = \frac{\delta}{|\mathcal{S}|}$$

- We also have

$$\mathbb{P} \left( (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*)) \geq \sqrt{\frac{\log(1/\delta)}{2n}} \right) \leq e^{-2n \left( \frac{\log(1/\delta)}{2n} \right)} = \delta$$

- Using a union bound strategy

$$\begin{aligned} & \mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) \geq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \\ & \leq \sum_{f \in \mathcal{S}} \mathbb{P} \left( (\mathcal{R}_n(f) - \mathcal{R}(f)) \geq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \\ & \leq \sum_{f \in \mathcal{S}} \frac{\delta}{|\mathcal{S}|} = \delta \end{aligned}$$



- Now

$$\begin{aligned} \mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) + (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*)) \geq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}} \right) \\ \leq \mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) \geq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \\ + \mathbb{P} \left( (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*)) \geq \sqrt{\frac{\log(1/\delta)}{2n}} \right) \leq \delta + \delta = 2\delta \end{aligned}$$

- The second bound can be obtained directly by bounding  $2 \sup_f |\mathcal{R}(f) - \mathcal{R}_n(f)|$ .

## PAC Bounds

- If  $\mathcal{S}$  is finite of cardinality  $|\mathcal{S}|$ , with proba greater than  $1 - 2\delta$

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_{\mathcal{S}}^*) \leq \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- If  $\mathcal{S}$  is finite of cardinality  $|\mathcal{S}|$ , with proba greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Risk increases with the cardinality of  $\mathcal{S}$ .
- Similar issue in cross-validation!
- No direct extension for an infinite  $\mathcal{S}$ ...

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 **Empirical Risk Minimization**
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - **McDiarmid and Rademacher Complexity**
  - VC Dimension
  - Structural Risk Minimization
- 8 References

- Supremum of Empirical losses:

$$\begin{aligned}\Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}_n) &= \sup_{f \in \mathcal{S}} \mathcal{R}(f) - \mathcal{R}_n(f) \\ &= \sup_{f \in \mathcal{S}} \left( \mathbb{E} \left[ \ell^{0/1}(Y, f(\underline{X})) \right] - \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i)) \right)\end{aligned}$$

## Properties

- Bounded difference:

$$|\Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}_j, \dots, \underline{X}_n) - \Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}'_j, \dots, \underline{X}_n)| \leq 1/n$$

## Concentration

$$\mathbb{P}(\Delta_n(\mathcal{S}) - \mathbb{E}[\Delta_n(\mathcal{S})] \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

- Concentration of bounded difference function.
- Generalization of Hoeffding theorem: McDiarmid Theorem.

$$\sup_{\beta} \left[ \mathbb{E}[\ell(\gamma, \beta(x)) - \frac{1}{n} \sum \ell(\gamma_i, \beta(x_i))] \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[ \mathbb{E}[\ell(\gamma_i, \beta(x_i))] - \ell(\gamma_i, \beta(x_i)) \right]$$

$$= \frac{1}{n} \sum_{i \neq 0} \mathbb{E}[\ell(\gamma_i, \beta(x_i)) - \ell(\gamma_i, \beta(x_i))] +$$

$$+ \frac{1}{n} \left[ \mathbb{E}[\ell(\gamma_0, \beta(x_0)) - \ell(\gamma_0, \beta(x_0))] \right]$$

$$- \frac{1}{n} \left[ \mathbb{E}[\ell(\gamma_0, \beta(x_0)) - \ell(\gamma_0, \beta(x_0))] \right]$$

$$+ \frac{1}{n} \left[ \mathbb{E}[\ell(\gamma_0, \beta(x_0)) - \ell(\gamma_0, \beta(x_0))] \right]$$

$$\leq \frac{1}{n} \left[ \ell(\gamma_0, \beta(x_0)) - \ell(\gamma_0, \beta(x_0)) \right]$$

## Bounded difference function

- $g : \mathcal{X}^n \rightarrow \mathbb{R}$  is a bounded difference function if it exist  $c_i$  such that

$$\forall (\underline{X}_i)_{i=1}^n, (\underline{X}'_i)_{i=1}^n \in \mathbb{R},$$

$$\left| g(\underline{X}_1, \dots, \underline{X}_j, \dots, \underline{X}_n) - g(\underline{X}_1, \dots, \underline{X}'_j, \dots, \underline{X}_n) \right| \leq c_j$$

## Theorem

- If  $g$  is a bounded difference function and  $\underline{X}_i$  are independent random variables then

$$\mathbb{P} (g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n)] \geq \epsilon) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}}$$

$$\mathbb{P} (\mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n)] - g(\underline{X}_1, \dots, \underline{X}_n) \geq \epsilon) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}}$$

- Proof ingredients:
  - Chernov bounds
  - Martingale decomposition...

- Let  $g_i = \mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n) | \underline{X}_1, \dots, \underline{X}_i]$  so that  $g_0 = \mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n)]$  and  $g_n = g(\underline{X}_1, \dots, \underline{X}_n)$ .
- By construction,

$$g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n)] = \sum_{i=0}^{n-1} g_{n-i} - g_{n-i-1}.$$

- Now

$$\begin{aligned} & \mathbb{E} \left[ e^{\lambda(g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E} [g(\underline{X}_1, \dots, \underline{X}_n)])} \right] \\ &= \mathbb{E} \left[ e^{\lambda(\sum_{i=0}^{n-1} g_{n-i} - g_{n-i-1})} \right] \\ &= \mathbb{E} \left[ \prod_{i=0}^{n-1} e^{\lambda(g_{n-i} - g_{n-i-1})} \right] \end{aligned}$$

by conditioning we have

$$= \mathbb{E} \left[ \mathbb{E} \left[ \prod_{i=0}^{n-1} e^{\lambda(g_{n-i} - g_{n-i-1})} \mid \underline{X}_2, \dots, \underline{X}_n \right] \right]$$

and using the fact that, for  $i \geq 1$ ,  $g_i$  is constant conditionally to  $\underline{X}_2, \dots, \underline{X}_n$

$$= \mathbb{E} \left[ \mathbb{E} \left[ e^{\lambda(g_1 - g_0)} \mid \underline{X}_2, \dots, \underline{X}_n \right] \prod_{i=1}^{n-1} e^{\lambda(g_{i+1} - g_i)} \right]$$

Now  $g_1 - g_0$  is by construction a centered random variable bounded in absolute value by  $c_n$  and thus

$$\leq e^{\frac{\lambda^2 c_n^2}{8}} \mathbb{E} \left[ \prod_{i=1}^{n-1} e^{\lambda(g_i - g_{i+1})} \right]$$



- Reusing the same technique recursively, we obtain

$$\mathbb{E} \left[ e^{\lambda(g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E}[g(\underline{X}_1, \dots, \underline{X}_n)])} \right] \leq e^{-\frac{\lambda^2 \sum_{i=1}^n c_n^2}{8}}$$

- Optimizing the corresponding Chernov bound in  $\lambda$  as in Hoeffding proof yields the result.

## Theorem

- If  $g$  is a bounded difference function and  $\underline{X}_i$  are independent random variables then

$$\mathbb{P}(g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E}[g(\underline{X}_1, \dots, \underline{X}_n)] \geq \epsilon) \leq e^{\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}}$$

- Using  $g = \Delta_n(\mathcal{S})$  for which  $c_i = 1/n$  yields immediately

$$\mathbb{P}(\Delta_n(\mathcal{S}) - \mathbb{E}[\Delta_n(\mathcal{S})] \geq \epsilon) \leq e^{\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}} = e^{-2n\epsilon^2}$$

- We derive then

$$\mathbb{P}(\Delta_n(\mathcal{S}) \geq \mathbb{E}[\Delta_n(\mathcal{S})] + \epsilon) \leq e^{\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}} = e^{-2n\epsilon^2}$$

- It remains to upperbound

$$\mathbb{E}[\Delta_n] = \mathbb{E}\left[\sup_{f \in \mathcal{S}} \mathcal{R}(f) - \mathcal{R}_n(f)\right]$$

## Theorem

- Let  $\sigma_i$  be a sequence of i.i.d. random symmetric Bernoulli variables (Rademacher variables):

$$\mathbb{E} \left[ \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq 2 \mathbb{E} \left[ \sup_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell^{0/1}(Y_i, f(\underline{X}_i)) \right]$$

## Rademacher complexity

- Let  $B \subset \mathbf{R}^n$ , the Rademacher complexity of  $B$  is defined as

$$R_n(B) = \mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right]$$

- Theorem gives an upper bound of the expectation in term of the **average Rademacher complexity of the random set**  
 $B_n(\mathcal{S}) = \{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}$ .
- Back to finite setting:** This set is at most of cardinality  $2^n$ !

- We prove a slightly more general result.
- Let  $\mathcal{H}$  be a set of  $n$ -tuple of functions of  $Z_i = (\underline{X}_i, Y_i)$ , we will prove that

$$\mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - \mathbb{E} [h_i(Z_i)]) \right] \leq 2 \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h_i(Z_i) \right]$$

- Using  $h_i(Z_i) = \ell^{0/1}(Y_i, f(\underline{X}_i))$  yields the result...
- Now, the doubling trick consists in introducing a second set of samples  $Z'_i$  with the same distribution:

$$\begin{aligned} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - \mathbb{E} [h_i(Z_i)]) \right] &= \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - \mathbb{E} [h_i(Z'_i)]) \right] \\ &= \mathbb{E}_Z \left[ \sup_{h \in \mathcal{H}} \mathbb{E}_{Z'} \left[ \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - h_i(Z'_i)) \right] \right] \end{aligned}$$

- We may now upperbound this term by exchanging the sup and the expectation to obtain

$$\mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - \mathbb{E} [h_i(Z_i)]) \right] \leq \mathbb{E}_{Z, Z'} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - h_i(Z'_i)) \right]$$

- By construction  $h_i(Z_i) - h_i(Z'_i)$  is a symmetric random variable and has thus the same law than  $\sigma_i(h_i(Z_i) - h_i(Z'_i))$  where  $\sigma_i$  is a sequence of i.i.d. Rademacher variable.

- Thus

$$\mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (h_i(Z_i) - \mathbb{E} [h_i(Z_i)]) \right] \leq \mathbb{E}_{Z, Z', \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i (h_i(Z_i) - h_i(Z'_i)) \right]$$

- Now we can split the sum in two and obtain

$$\begin{aligned}\mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h_i(Z_i) - \mathbb{E} [h_i(Z_i)] \right] &\leq \mathbb{E}_{Z, Z', \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h_i(Z_i) \right] \\ &\quad + \mathbb{E}_{Z, Z', \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i (-h_i(Z'_i)) \right] \\ &\leq \mathbb{E}_{Z, \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h_i(Z_i) \right] \\ &\quad + \mathbb{E}_{Z, \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (-\sigma_i) h_i(Z_i) \right]\end{aligned}$$

which yields by symmetry of  $\sigma_i$

$$\mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h_i(Z_i) - \mathbb{E} [h_i(Z_i)] \right] \leq 2 \mathbb{E}_{Z, \sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h_i(Z_i) \right]$$

## Theorem

- If  $B$  is finite and such that  $\forall b \in B, \frac{1}{n} \|b\|_2^2 \leq M^2$ , then

$$R_n(B) = \mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right] \leq \sqrt{\frac{2M^2 \log |B|}{n}}$$

- If  $B = B_n(\mathcal{S}) = \{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}$ , we have  $M = 1$  and thus

$$R_n(B) \leq \sqrt{\frac{2 \log |B_n(\mathcal{S})|}{n}}$$

- We obtain immediately

$$\mathbb{E} \left[ \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right].$$

- By Jensen inequality

$$\begin{aligned} e^{\lambda \mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right]} &\leq \mathbb{E} \left[ e^{\lambda \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i} \right] \\ &\leq \sum_{b \in B} \mathbb{E} \left[ e^{\lambda \frac{1}{n} \sum_{i=1}^n \sigma_i b_i} \right] \end{aligned}$$

We can now use Hoeffding inequality with  $c_i = 2b_i/n$  to obtain

$$\begin{aligned} e^{\lambda \mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right]} &\leq \sum_{b \in B} e^{\frac{4\lambda^2 \sum_{i=1}^n b_i^2}{8n^2}} \\ &\leq \sum_{b \in B} e^{\frac{M^2 \lambda^2}{2n}} = |B| e^{\frac{M^2 \lambda^2}{2n}} \end{aligned}$$

- Thus

$$\mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right] \leq \frac{\log B}{\lambda} + \frac{M^2 \lambda}{2n}$$



- The optimal value is given by  $\lambda = \sqrt{\frac{2n \log B}{M^2}}$  yielding

$$\mathbb{E} \left[ \sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right] \leq \sqrt{\frac{2 \log |B_n(\mathcal{S})|}{n}}$$

~ log |S| ?

## Theorem

- With probability greater than  $1 - 2\delta$ ,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- This is a direct consequence of the previous bound.

- The second bound is straightforward.
- For the first one,

$$\begin{aligned} \mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) + \mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*) \geq \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{2 \log(1/\delta)}{n}} \right) \\ \leq \mathbb{P} \left( \sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) \geq \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{\log(1/\delta)}{2n}} \right) \\ + \mathbb{P} \left( (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*)) \geq \sqrt{\frac{\log(1/\delta)}{2n}} \right) \leq \delta + \delta = 2\delta \end{aligned}$$

## Corollary

- If  $\mathcal{S}$  is finite then with probability greater than  $1 - 2\delta$

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_{\mathcal{S}}^*) \leq \sqrt{\frac{8 \log |\mathcal{S}|}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- If  $\mathcal{S}$  is finite then with probability greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8 \log |\mathcal{S}|}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- It suffices to notice that

$$|B_n(\mathcal{S})| = |\{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}| \leq |\mathcal{S}|$$

- Same result with Hoeffding but with **better** constants!

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Difference due to the *crude* upperbound of

$$\mathbb{E} \left[ \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right]$$

- **Why bother?:** We do not have to assume that  $\mathcal{S}$  is finite!

$$|B_n(\mathcal{S})| \leq 2^n$$

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization**
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension**
  - Structural Risk Minimization
- 8 References

## Theorem

$$\mathbb{E} \left[ \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right]$$

- Key quantity:  $\mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right]$
- Hard to control due to its structure!

## A first data dependent upperbound

$$\mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] \leq \sqrt{\frac{8 \log \mathbb{E} [|B_n(\mathcal{S})|]}{n}} \quad (\text{Jensen})$$

- Depends on the unknown  $P$ !

$$\mathbb{E} [\beta(x)] \leq \sup_x \beta(x)$$

## Shattering Coefficient (or Growth Function)

- The shattering coefficient of the class  $\mathcal{S}$ ,  $s(\mathcal{S}, n)$ , is defined as

$$s(\mathcal{S}, n) = \sup_{((\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)) \in (\mathcal{X} \times \{-1, 1\})^n} |\{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}|$$

- By construction,  $|B_n(\mathcal{S})| \leq s(\mathcal{S}, n) \leq \min(2^n, |\mathcal{S}|)$ !

## A data independent upperbound

$$\mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] \leq \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}}$$



## Theorem

- With probability greater than  $1 - 2\delta$ ,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Depends only on the class  $\mathcal{S}$ !

$$s(\mathcal{S}, d) \leq 2^d$$

## VC Dimension

- The VC dimension  $d_{VC}$  of  $\mathcal{S}$  is defined as the largest integer  $d$  such that
$$s(\mathcal{S}, d) = 2^d$$

- The VC dimension can be infinite!

## VC Dimension and Dimension

- **Prop:** If  $\text{span}(\mathcal{S})$  corresponds to the sign of functions in a linear space of dimension  $d$  then  $d_{VC} \leq d$ .
- VC dimension similar to the usual dimension.

## Sauer's Lemma

- If the VC dimension  $d_{VC}$  of  $\mathcal{S}$  is finite

$$s(\mathcal{S}, n) \leq \begin{cases} 2^n & \text{if } n \leq d_{VC} \\ \left(\frac{en}{d_{VC}}\right)^{d_{VC}} & \text{if } n > d_{VC} \end{cases}$$

- **Cor.:**  $\log s(\mathcal{S}, n) \leq d_{VC} \log\left(\frac{en}{d_{VC}}\right)$  if  $n > d_{VC}$ .

- Let  $B(\mathcal{S}, C = (c_1, \dots, c_n)) = \{(f(c_1), \dots, f(c_n)), f \in \mathcal{S}\}$ ,  
 $|B(\mathcal{S}, C)| = |\{(\ell^{0/1}(Y_i, f(c_i)))_{i=1}^n, Y_i \in \{-1, 1\}, f \in \mathcal{S}\}|$   
 so that  $s(\mathcal{S}, n) = \sup_{C, |C|=n} |B(\mathcal{S}, C)|$ .

- We say that  $\mathcal{S}$  shatters  $B$  iff  $|B(\mathcal{S}, B)| = 2^B$  or  $B = \emptyset$ .
- Lemma:**  $|B(\mathcal{S}, C)| \leq |\{A \subset C, \mathcal{S} \text{ shatters } A\}|$
- The VC dimension  $d$  is the largest size of a shattered set so that

$$|\{A \subset C, \mathcal{S} \text{ shatters } A\}| \leq \sum_{i=0}^{\min(n,d)} \binom{n}{i}$$

- If  $n \leq d$ ,  $\sum_{i=0}^n \binom{n}{i} = 2^n$ .
- If  $n \geq d$ ,

$$\left(\frac{d}{n}\right)^d \sum_{i=0}^d \binom{n}{i} \leq \sum_{i=0}^d \binom{n}{i} \left(\frac{d}{n}\right)^i \leq \sum_{i=0}^n \binom{n}{i} \left(\frac{d}{n}\right)^i = \left(1 + \frac{d}{n}\right)^n \leq e^d$$

- The lemma will be proved by induction on the size of  $C$ .
- If  $|C| = 1$  then
  - either  $\mathcal{S}$  shatters  $\{c_1\}$  and  $|B(\mathcal{S}, C)| = 2 \leq |\{\emptyset, C\}| = |\{A \subset C, \mathcal{S} \text{ shatters } A\}|$ ,
  - or  $\mathcal{S}$  does not shatter  $\{c_1\}$  and  $|B(\mathcal{S}, C)| = 1 \leq |\{\emptyset\}| = |\{A \subset C, \mathcal{S} \text{ shatters } A\}|$
- Now assume the property is true for all  $C'$  of size smaller than  $n - 1$ .
- Let  $C = \{c_1, \dots, c_n\}$ , by definition,
$$\begin{aligned} |B(\mathcal{S}, C)| &= |\{(f(c_1), \dots, f(c_n)), f \in \mathcal{S}\}| \\ &= |\{(-1, f(c_2), \dots, f(c_n)), f \in \mathcal{S}, f(c_1) = -1\}| \\ &\quad + |\{(1, f(c_2), \dots, f(c_n)), f \in \mathcal{S}, f(c_1) = 1\}| \\ &= |\{(f(c_2), \dots, f(c_n)), f \in \mathcal{S}\}| \\ &\quad + |\{(f(c_2), \dots, f(c_n)), f \in \mathcal{S}, \exists f', f(c_1) = -f'(c_1), f(c_i) = f'(c_i), i \neq 1\}| \end{aligned}$$

- Now by construction and induction,

$$\begin{aligned} |\{(f(c_2), \dots, f(c_n)), f \in \mathcal{S}\}| &= |B(\mathcal{S}, C \setminus \{c_1\})| \\ &\leq |\{A \subset C \setminus \{c_1\}, \mathcal{S} \text{ shatters } A\}| \\ &\leq |\{A \subset C, c_1 \notin A, \mathcal{S} \text{ shatters } A\}| \end{aligned}$$

- Now let  $\mathcal{S}' = \{f \in \mathcal{S}, \exists f', f(c_1) = -f'(c_1), f(c_i) = f'(c_i), i \neq 1\}$

$$\begin{aligned} &|\{(f(c_2), \dots, f(c_n)), f \in \mathcal{S}, \exists f', f(c_1) = -f'(c_1), f(c_i) = f'(c_i), i \neq 1\}| \\ &= |B(\mathcal{S}', C \setminus \{c_1\})| \leq |\{A \subset C \setminus \{c_1\}, \mathcal{S}' \text{ shatters } A\}| \\ &\leq |\{A \subset C \setminus \{c_1\}, \mathcal{S}' \text{ shatters } A \cup \{c_1\}\}| \quad \text{by definition of } \mathcal{S}' \\ &\leq |\{A \subset C, c_1 \in A, \mathcal{S}' \text{ shatters } A\}| \leq |\{A \subset C, c_1 \in A, \mathcal{S} \text{ shatters } A\}| \end{aligned}$$

- Summing those two bounds yields the lemma property for  $C$ :

$$\begin{aligned} |B(\mathcal{S}, C)| &\leq |\{A \subset C, c_1 \notin A, \mathcal{S} \text{ shatters } A\}| + |\{A \subset C, c_1 \in A, \mathcal{S} \text{ shatters } A\}| \\ &\leq |\{A \subset C, \mathcal{S} \text{ shatters } A\}| \end{aligned}$$

## PAC Bounds

- If  $\mathcal{S}$  is of VC dimension  $d_{VC}$  then if  $n > d_{VC}$
- With probability greater than  $1 - 2\delta$ ,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{8d_{VC} \log\left(\frac{en}{d_{VC}}\right)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8d_{VC} \log\left(\frac{en}{d_{VC}}\right)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- **Rk:** If  $d_{VC} = +\infty$  no uniform PAC bounds exists!

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization**
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - **Structural Risk Minimization**
- 8 References



## PAC Bounds

- Let  $\pi_f > 0$  such that  $\sum_{f \in \mathcal{S}} \pi_f = 1$
- With proba greater than  $1 - 2\delta$ ,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{\log(1/\pi_f)}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With proba greater than  $1 - \delta$ , simultaneously  $\forall f' \in \mathcal{S}$ ,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log(1/\pi_f)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Very similar proof than the uniform one!
- Much more interesting idea when combined with several models...

$$\hat{\beta} = \underset{b \in \mathcal{S}}{\text{argmin}} R_n(\hat{\beta})$$

$$\mathcal{S}_1 \dots \mathcal{S}_m$$

$$y' = U\beta_i$$

$$\hat{\beta}_i = \underset{b \in \mathcal{S}_i}{\text{argmin}} R_n(\hat{\beta}_i)$$

$$\hookrightarrow \hat{\beta}_i = \underset{i}{\text{argmin}} R_n(\hat{\beta}_i) + \sqrt{\frac{8 \log S_n(i)}{n}} + \sqrt{\frac{\log^2 \tilde{m}}{2n}}$$

- Assume we have a countable collection of set  $(\mathcal{S}_m)_{m \in \mathcal{M}}$  and let  $\pi_m$  be such that  $\sum_{m \in \mathcal{M}} \pi_m = 1$ .

## Non Uniform Risk Bound

- With probability  $1 - \delta$ , simultaneously for all  $m \in \mathcal{M}$  and all  $f \in \mathcal{S}_m$ ,

$$\mathcal{R}(f) \leq \mathcal{R}_n(f) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

## Structural Risk Minimization

- Choose  $\hat{f}$  as the minimizer over  $m \in \mathcal{M}$  and  $f \in \mathcal{S}_m$  of

$$\mathcal{R}_n(f) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}}$$

- Mimics the minimization of the integrated risk!

## PAC Bound

- If  $\hat{f}$  is the SRM minimizer then with probability  $1 - 2\delta$ ,

$$\mathcal{R}(\hat{f}) \leq \inf_{m \in \mathcal{M}} \inf_{f \in \mathcal{S}_m} \left( \mathcal{R}(f) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} \right) + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- The SRM minimizer balances the risk  $\mathcal{R}(f)$  and the upper bound on the estimation error  $\mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}}$ .
- $\mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right]$  can be replaced by an upper bound (for instance a VC based one)...

EU  $\quad U \approx \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}}$

- Let  $f_{\mathcal{S}_{m^*}}^*$  be the minimizer over  $m \in \mathcal{M}$  and  $f \in \mathcal{S}_m$  of

$$\mathcal{R}(f) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}}$$

- The previous bound is thus equivalent to the fact that with probability greater than  $1 - 2\delta$ ,

$$\begin{aligned} \mathcal{R}(\hat{f}) &\leq \mathcal{R}(f_{\mathcal{S}_{m^*}}^*) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_{m^*})|}{n}} \right] + \sqrt{\frac{\log(1/\pi_{m^*})}{2n}} \\ &\quad + \sqrt{\frac{2 \log(1/\delta)}{n}} \end{aligned}$$

- We use then that with probability  $1 - \delta$

$$\begin{aligned}\mathcal{R}(\hat{f}) &\leq \mathcal{R}_n(f) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}} \\ &\leq \mathcal{R}_n(f_{\mathcal{S}_{m^*}^*}) + \mathbb{E} \left[ \sqrt{\frac{8 \log |B_n(\mathcal{S}_{m^*})|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}\end{aligned}$$

- Combining this with the fact that with probability  $1 - \delta$

$$\mathcal{R}(f_{\mathcal{S}_{m^*}^*}) \leq \mathcal{R}_n(f_{\mathcal{S}_{m^*}^*}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

yields the result.

- 1 Introduction
  - Machine Learning
  - Motivation
- 2 A Practical View
  - Method or Models
  - Interpretability
  - Metric Choice
- 3 A Better Point of View
  - The Example of Univariate Linear Regression
  - Supervised Learning
- 4 Error Estimation and Method Choice
  - Cross Validation
  - Cross Validation and Test
  - Cross Validation and Weights
  - Auto ML
- 5 A Probabilistic Point of View
  - Parametric Conditional Density Modeling
  - Non Parametric Conditional Density Modeling
  - Generative Modeling
- 6 Optimization Point of View
  - SVM
  - Penalization
  - (Deep) Neural Networks
  - Tree Based Methods
  - Ensemble Methods
- 7 Empirical Risk Minimization
  - Empirical Risk Minimization
  - ERM and PAC Bayesian Analysis
  - Hoeffding and Finite Class
  - McDiarmid and Rademacher Complexity
  - VC Dimension
  - Structural Risk Minimization
- 8 References

# References



T. Hastie, R. Tibshirani, and J. Friedman.  
*The Elements of Statistical Learning.*  
Springer Series in Statistics, 2009



G. James, D. Witten, T. Hastie, and  
R. Tibshirani.  
*An Introduction to Statistical Learning with  
Applications in R.*  
Springer, 2014



A. Géron.  
*Hands-On Machine Learning with  
Scikit-Learn, Keras and TensorFlow (2nd ed.)*  
O'Reilly, 2019



Ch. Giraud.  
*Introduction to High-Dimensional Statistics.*  
CRC Press, 2014



M. Mohri, A. Rostamizadeh, and  
A. Talwalkar.  
*Foundations of Machine Learning.*  
MIT Press, 2012



S. Shalev-Shwartz and S. Ben-David.  
*Understanding Machine Learning.*  
Cambridge University Press, 2014



B. Schölkopf and A. Smola.  
*Learning with kernels.*  
The MIT Press, 2002

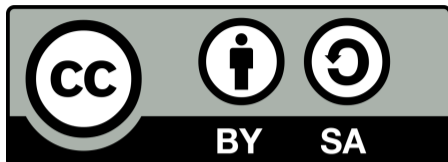


F. Chollet.  
*Deep Learning with Python.*  
Manning, 2017



F. Chollet and J.J. Allaire.  
*Deep Learning with R.*  
Manning, 2017





## Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
  - **Share:** copy and redistribute the material in any medium or format
  - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
  - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
  - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
  - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

## Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.