

Analytical Functions

Example: 1-D

$$f_1(x) = a(x - x_0)^2 + b$$

where $x, x_0, b \in \mathbb{R}, a \in \mathbb{R}$

Generalization:

convex quadratic function

$$f_2(x) = \frac{1}{2}(x - x_0)^T A (x - x_0) + b$$

where $x, x_0 \in \mathbb{R}^n, b \in \mathbb{R}, A \in \mathbb{R}^{\{n \times n\}}$
and A symmetric positive definite (SPD)

Exercise:

What is the minimum of $f_2(x)$?

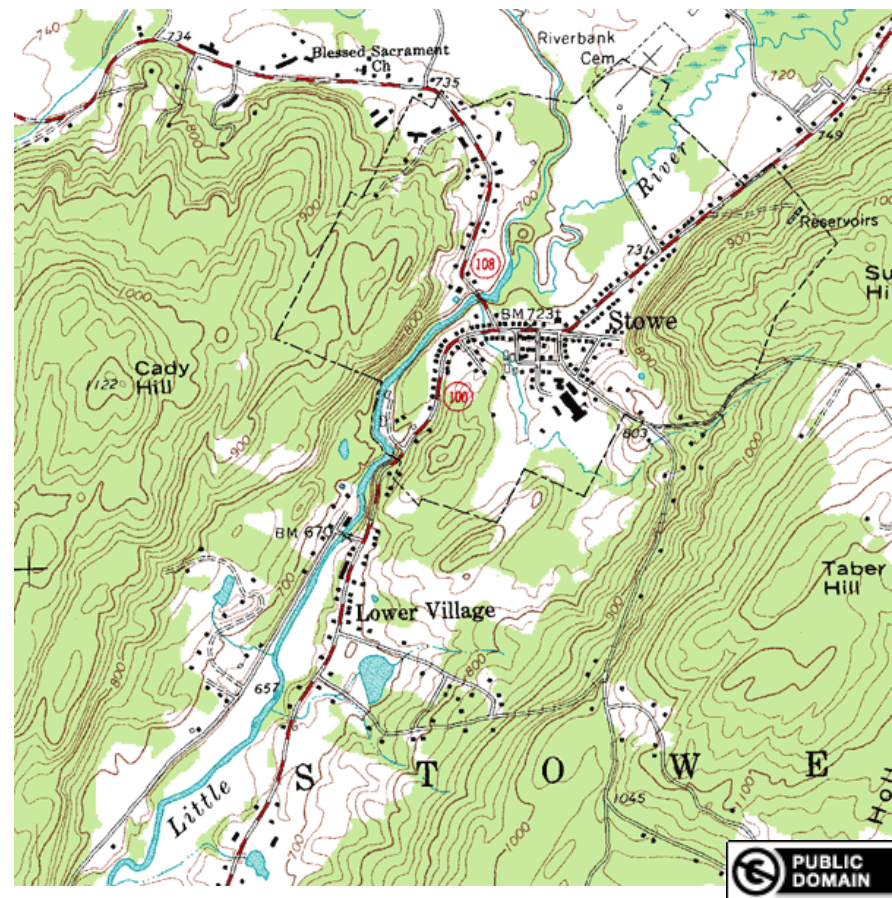
Levels Sets of Convex Quadratic Functions

Continuation of exercise:
What are the level sets of f_2 ?

Reminder: level sets of a function

$$L_c = \{x \in \mathbb{R}^n \mid f(x) = c\}$$

(similar to topography lines /
level sets on a map)

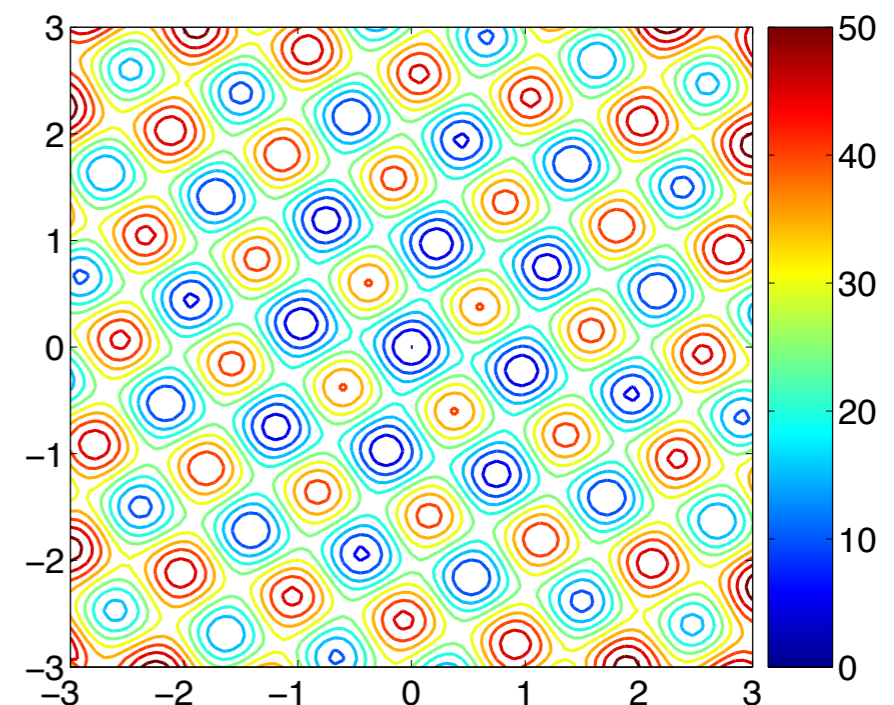
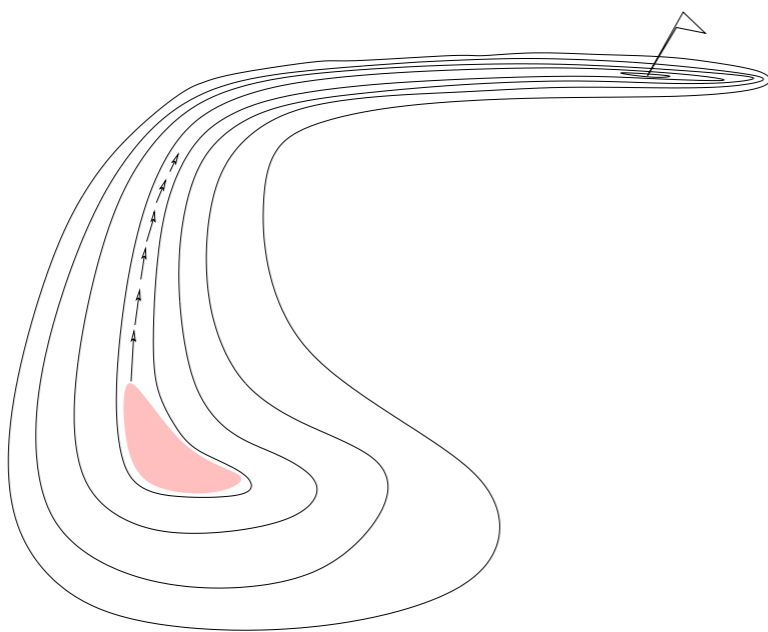


Level Sets: Visualization of a Function

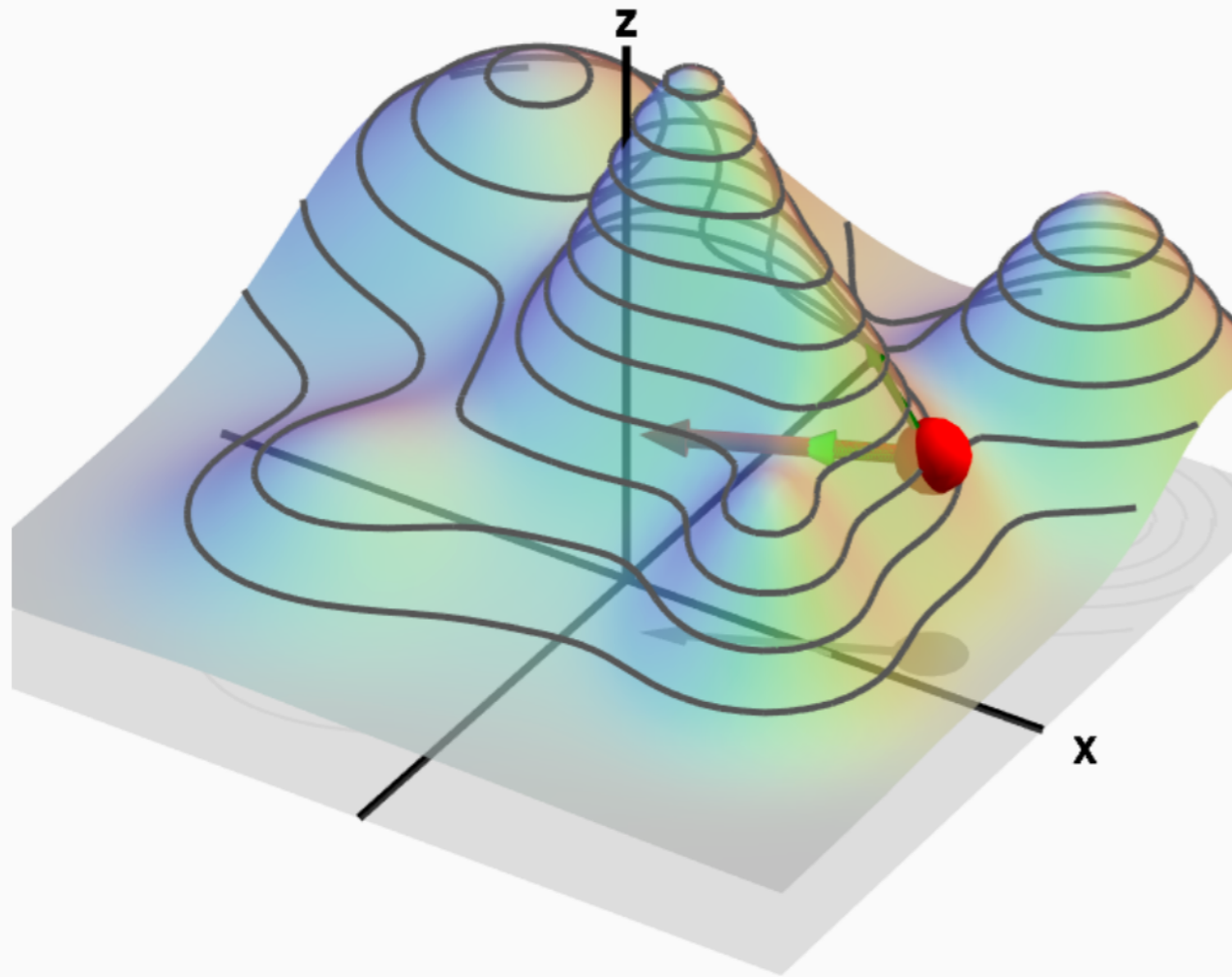
One-dimensional (1-D) representations are often misleading (as 1-D optimization is “trivial”, see slides related to curse of dimensionality), we therefore often represent **level-sets** of functions

$$\mathcal{L}_c = \{x \in \mathbb{R}^n \mid f(x) = c, \}, c \in \mathbb{R}$$

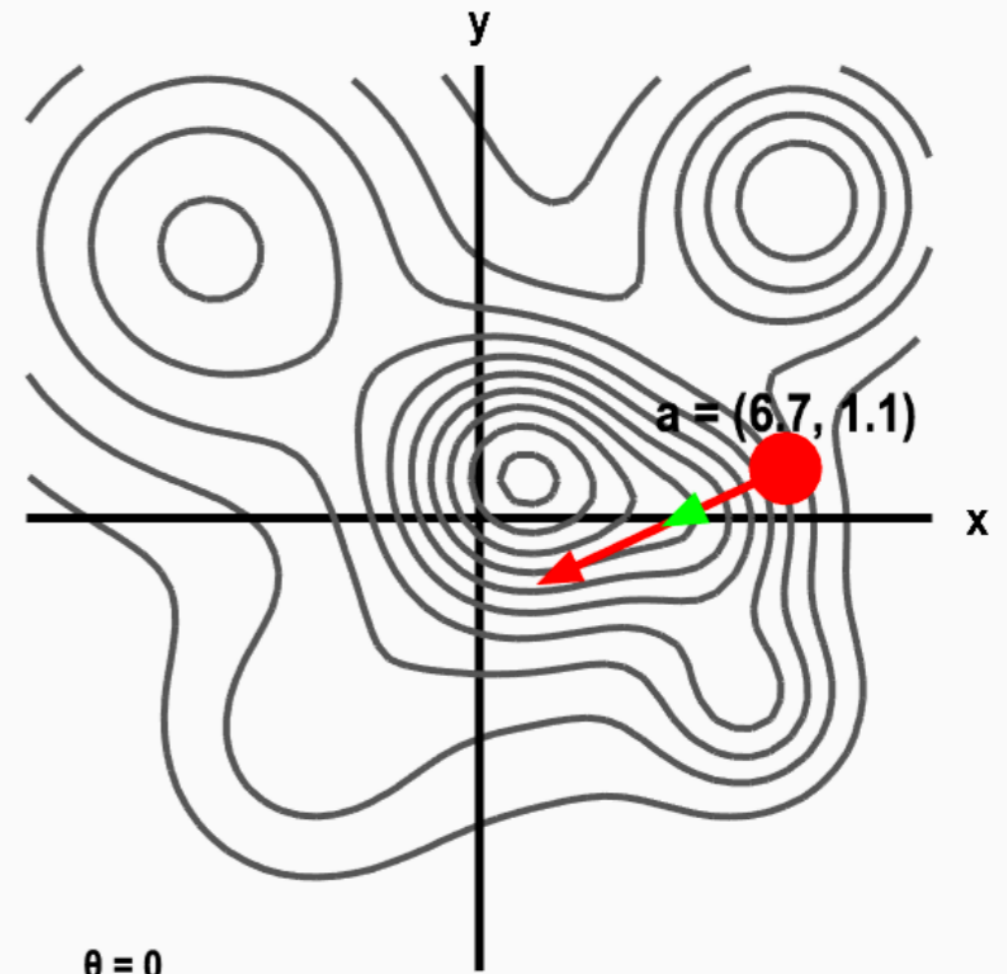
Examples of level sets in 2D



Level Sets: Visualization of a Function



$\theta = 0$
 $u = (-0.91, -0.42)$
 $a = (6.7, 1.1)$
 $\nabla f(a) = (-1.81, -0.85)$
 $D_u f(a) = 2.00$
 $|\nabla f(a)| = 2.00$



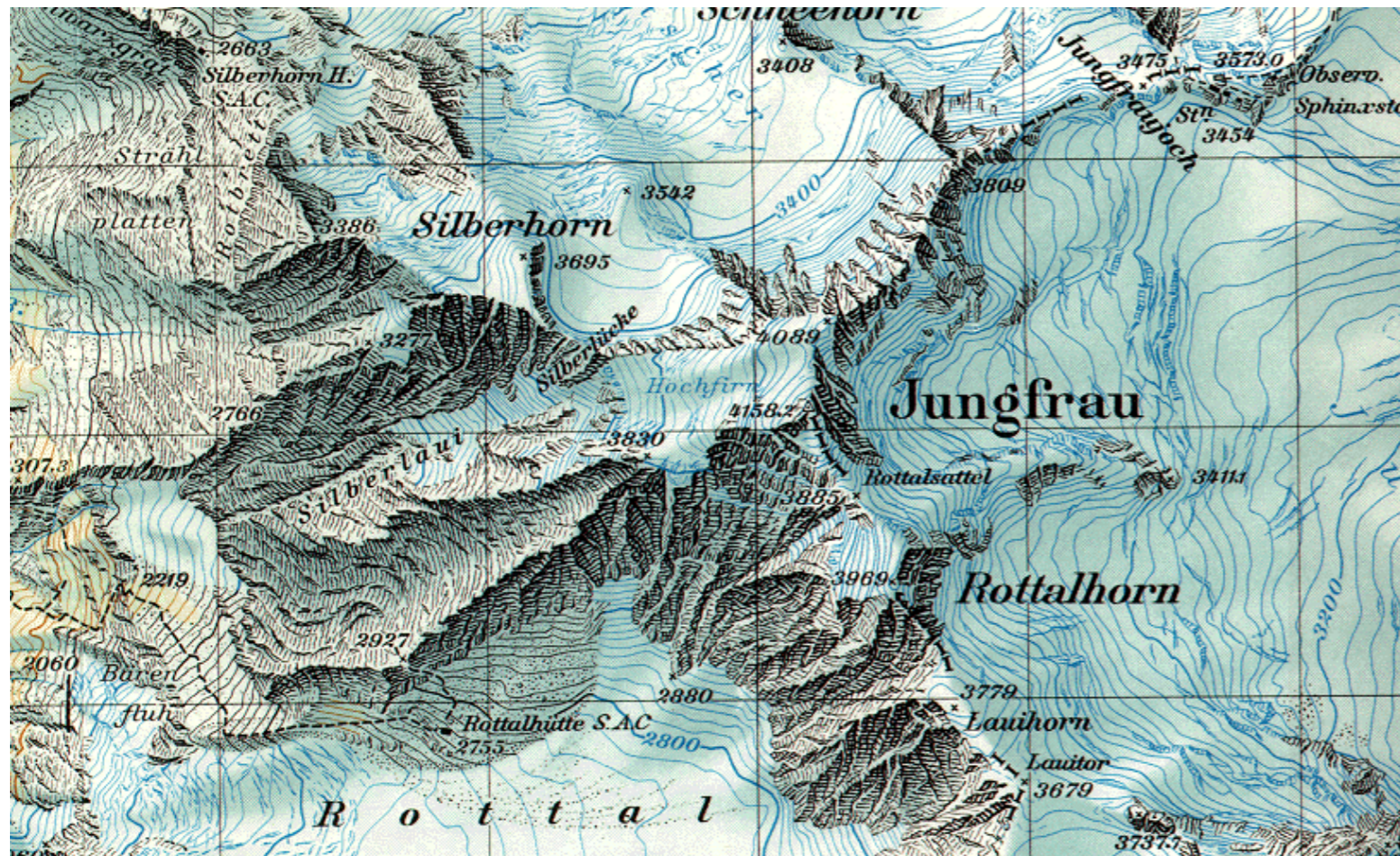
$\theta = 0$
 $u = (-0.91, -0.42)$
 $\nabla f(a) = (-1.81, -0.85)$
 $a = (6.7, 1.1)$
 $D_u f(a) = 2.00$
 $|\nabla f(a)| = 2.00$
 $f(a) = 4.87$



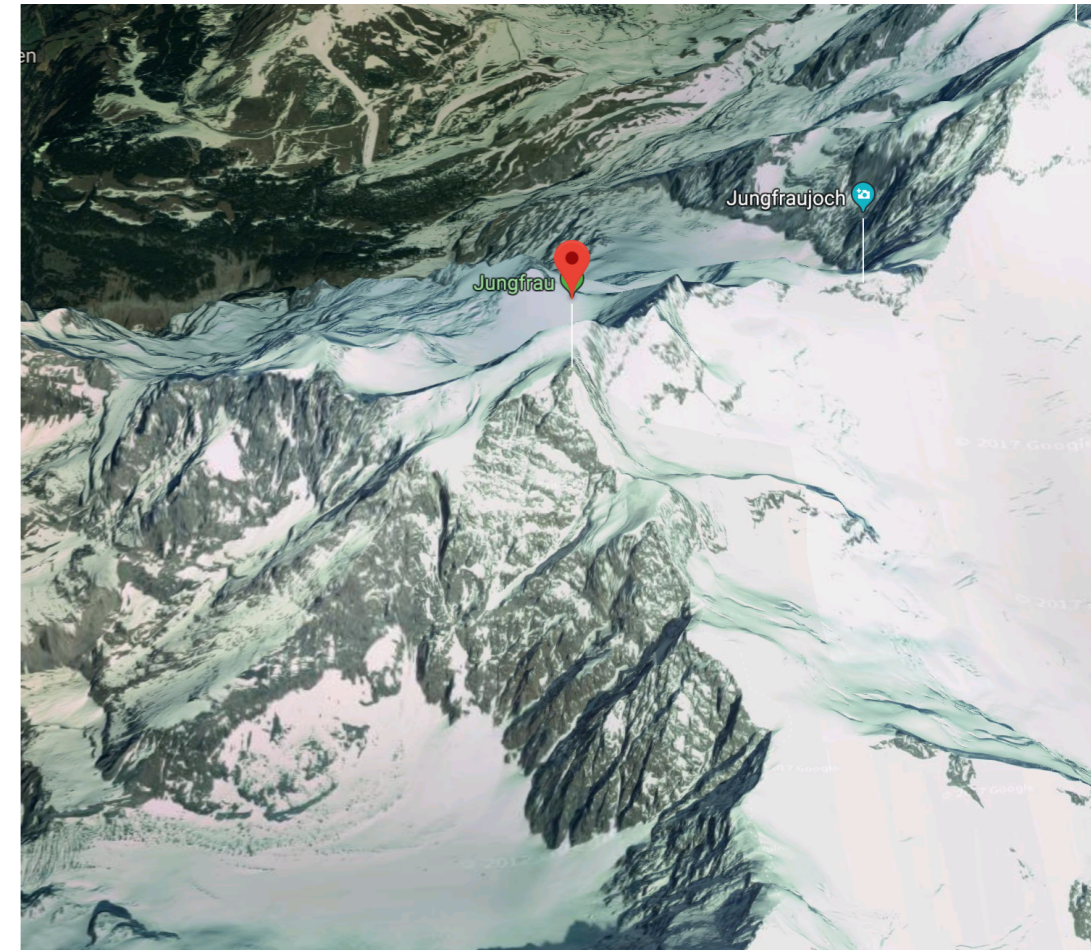
Source: Nykamp DQ, "Directional derivative on a mountain." From *Math Insight*. http://mathinsight.org/applet/directional_derivative_mountain

Level Sets: Topographic Map

The function is the altitude



Topographic map



3-D picture

Levels Sets of Convex Quadratic Functions

Continuation of exercise:

What are the level sets of f_2 ?

$$f_2(x) = \frac{1}{2} (x-x_0)^T A (x-x_0) + b$$

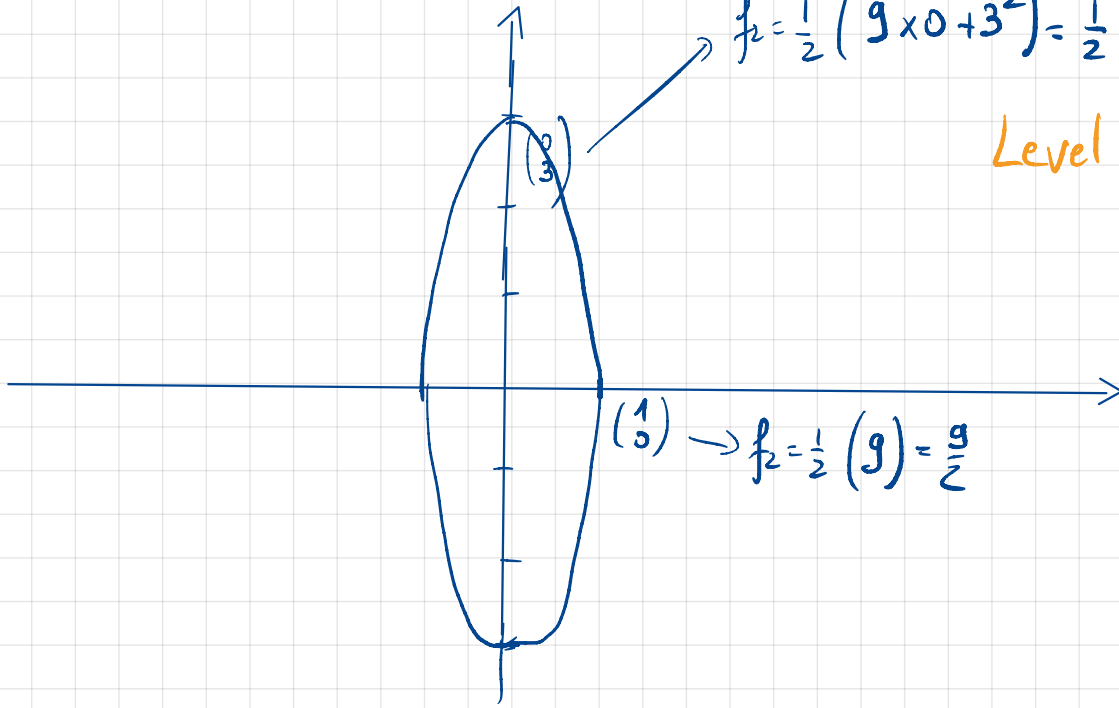
A SPD

- Probably too complicated in general, thus an example here
 - Consider $A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$, $b = 0$, $n = 2$
 - a) Compute $f_2(x)$.
 - b) Plot the level sets of $f_2(x)$.
 - c) More generally, for $n = 2$, if A is SPD with eigenvalues $\lambda_1 = 9$ and $\lambda_2 = 1$, what are the level sets of $f_2(x)$?
- Not necessarily diagonal*

$$A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$$

$$f_2(x) = \frac{1}{2} (9x_1^2 + x_2^2)$$

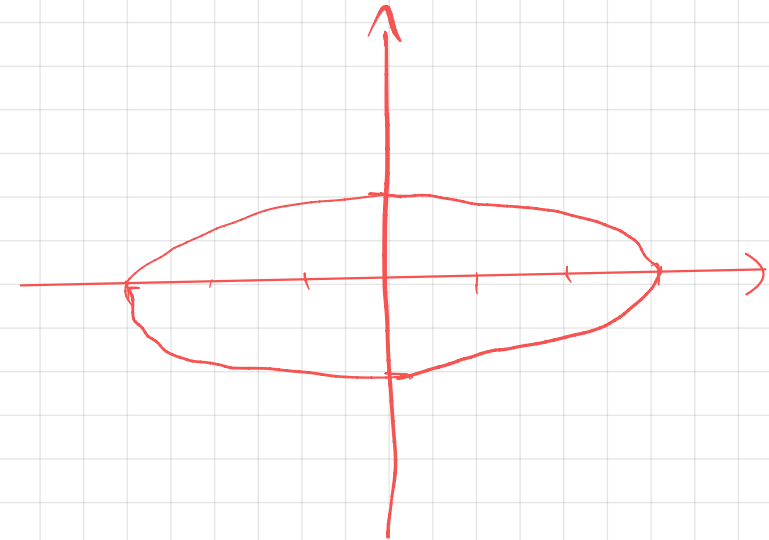
$$x = (x_1, x_2)$$



$$f_2 = \frac{1}{2} (9 \times 0 + 3^2) = \frac{1}{2} (9) = \frac{9}{2}$$

Level sets are ellipsoid, long-axis: y-axis
small axis: x-axis

$$\text{If } A = \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix}, \quad f_2(x) = \frac{1}{2} (x_1^2 + 9x_2^2) \rightarrow$$



A is symmetric, positive, definite:

$$A = P D P^T$$

from the spectral theorem.

P is orthogonal

P contains the eigenvectors of A

$$f_2(x) = \frac{1}{2} x^T A x = \frac{1}{2} x^T P D P^T x$$

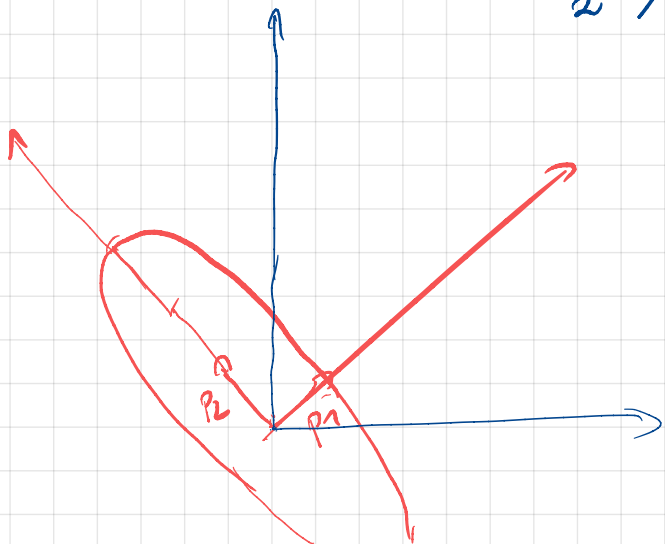
$$= \frac{1}{2} \underbrace{(P^T x)^T}_{x^T} D \underbrace{P^T x}_y$$

$$= \frac{1}{2} y^T D y$$

$$= \frac{1}{2} y^T \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix} y = \frac{1}{2} (9y_1^2 + y_2^2)$$

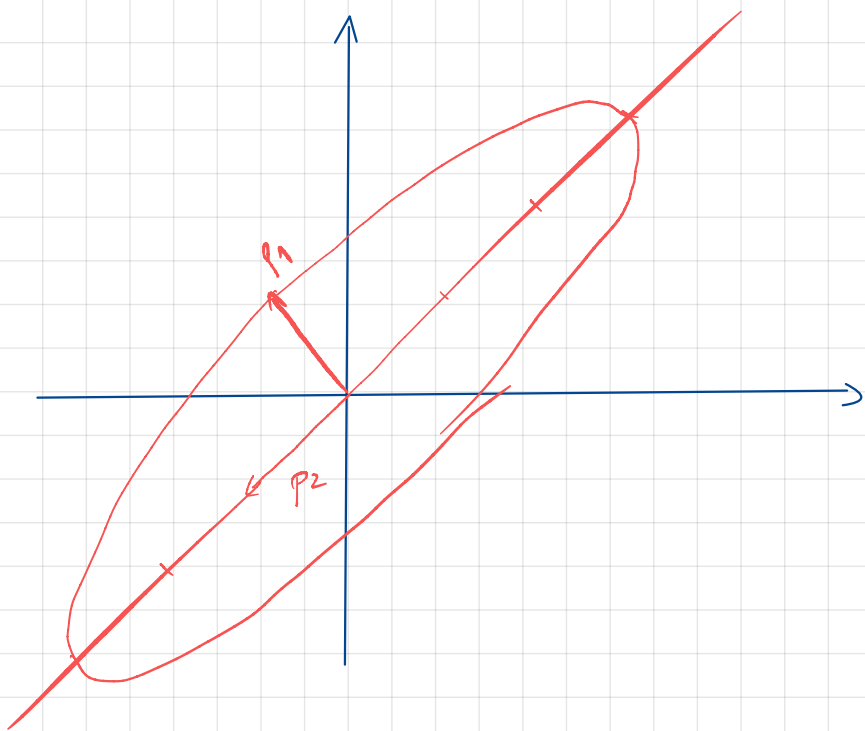
$$y = P^T x$$

$$D = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$$



p_1, p_2 eigenvectors of A associated to $\lambda_1 = 9, \lambda_2 = 1$

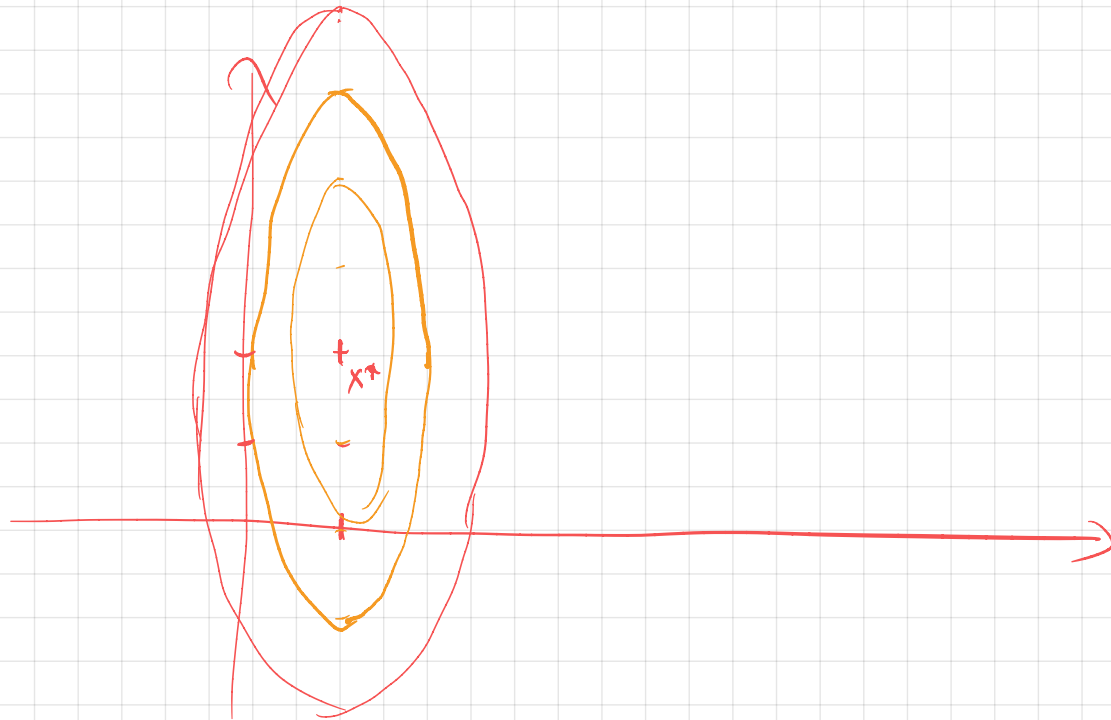
"Same" ellipsoid than before but rotated
The main axis of ellipsoid are the eigenvectors of A .



We have assumed before $x^* = 0$, if $x^* = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and we consider

$$f(x) = \frac{1}{2} (x - x^*)^T \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix} (x - x^*)$$

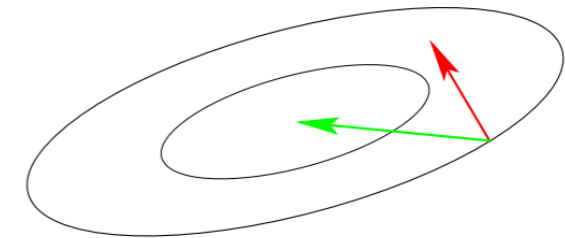
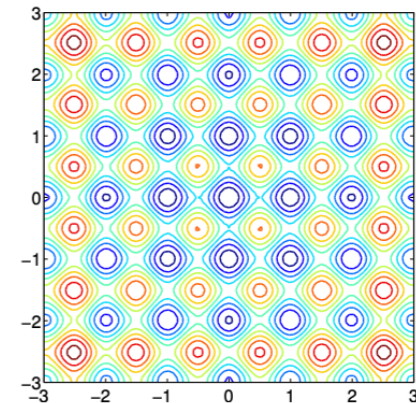
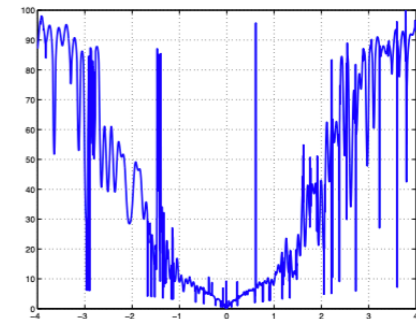
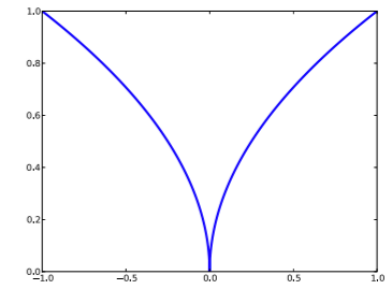
then the optimum of f is in x^* and the ellipsoid are centered around x^* , i.e.



What Makes a Function Difficult to Solve?

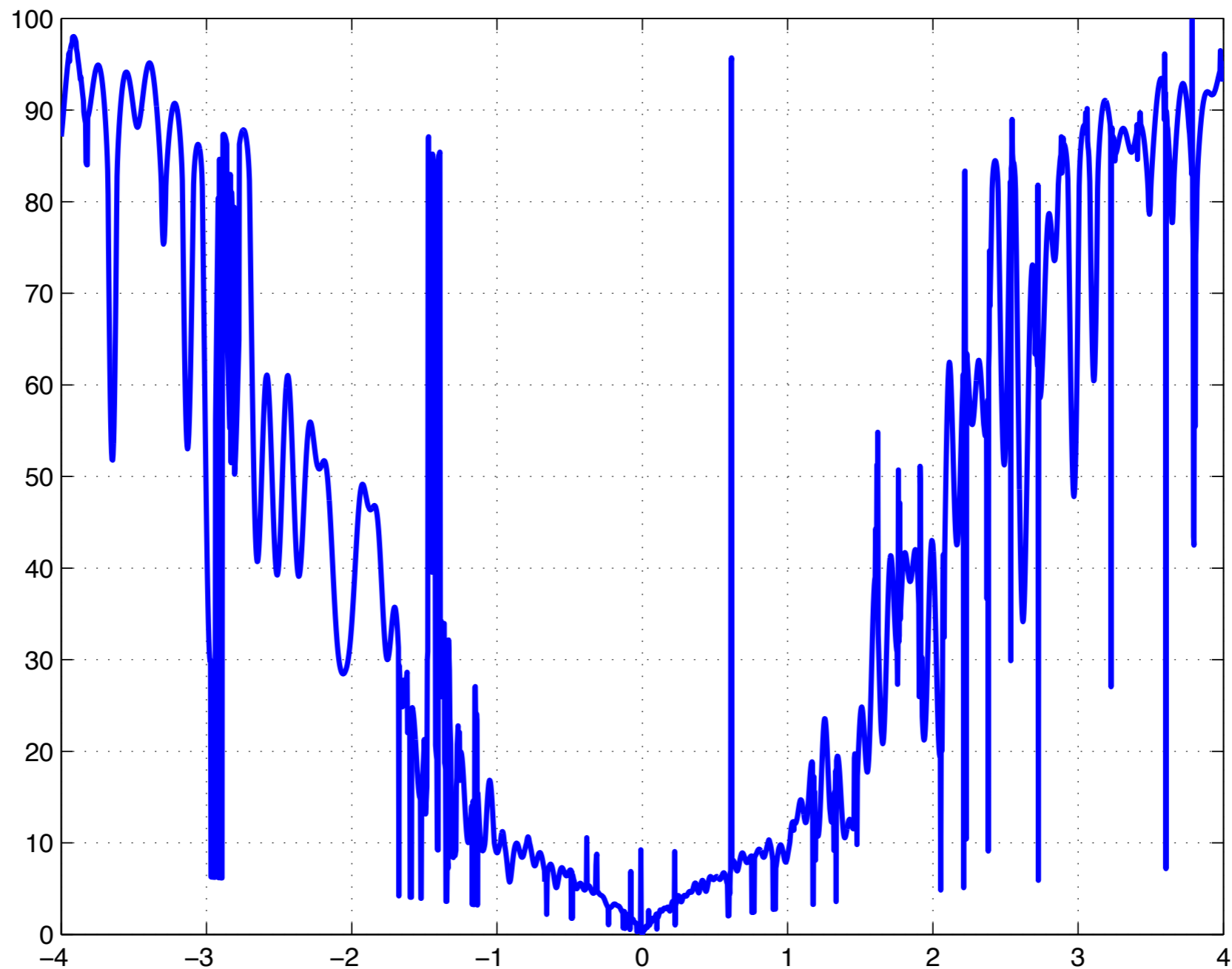
- ▶ non-linear, non-quadratic, non-convex
 - on linear and quadratic functions much better search policies are available
- ▶ ruggedness
 - non-smooth, discontinuous, multimodal, and/or noisy function
- ▶ dimensionality (size of search space)
 - (considerably) larger than three
- ▶ non-separability
 - dependencies between the objective variables
- ▶ ill-conditioning

~~Why stochastic search?~~



gradient direction Newton direction

Ruggedness



A cut of a 4-D function that can easily be solved with the
CMA-ES algorithm

Why is Optimization a non-trivial Problem?

Curse of dimensionality

if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

Why is Optimization a non-trivial Problem?

Curse of dimensionality

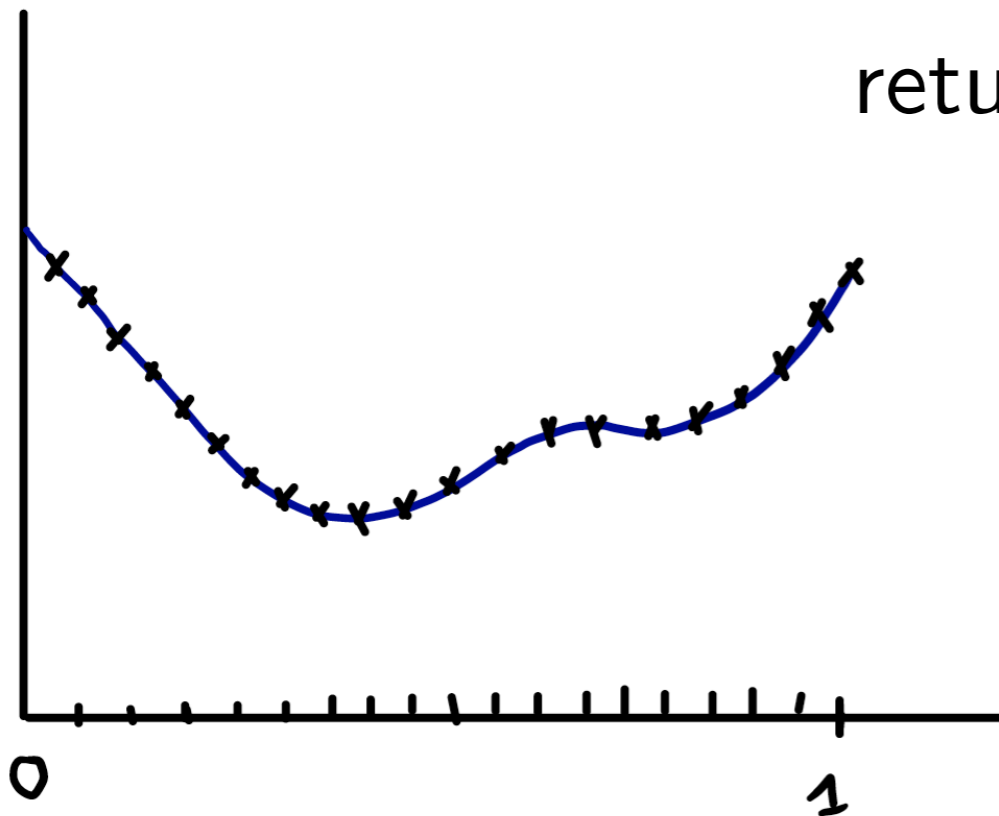
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



Why is Optimization a non-trivial Problem?

Curse of dimensionality

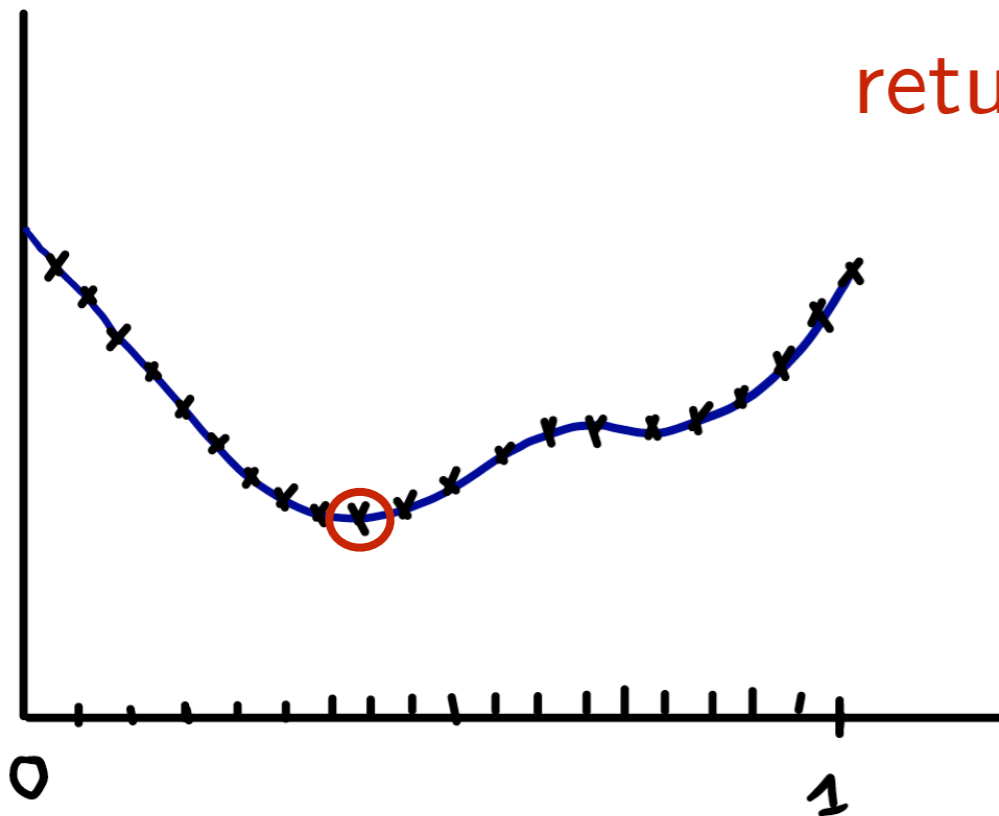
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



Why is Optimization a non-trivial Problem?

Curse of dimensionality

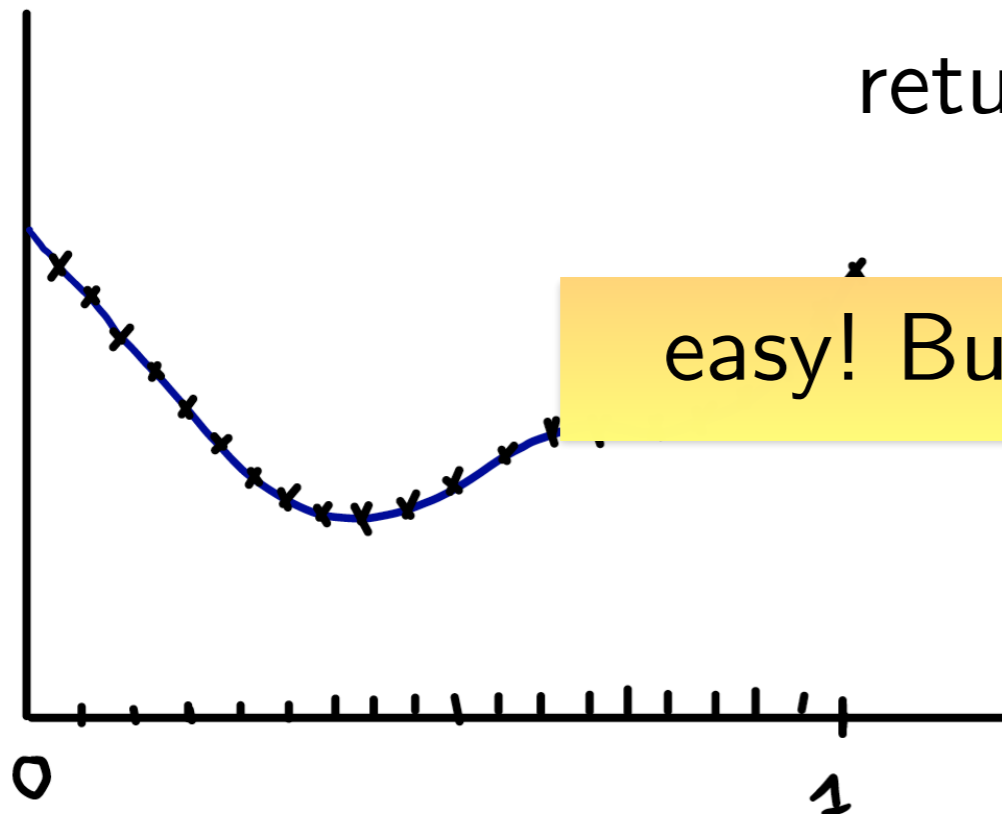
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



easy! But how does it scale when n increases?

1-D optimization is trivial

Curse of Dimensionality

The term **curse of dimensionality** (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 100 points onto a real interval, say $[0,1]$.

How many points would you need to get a similar coverage (in terms of distance between adjacent points) in dimension 10?

Curse of Dimensionality

The term **curse of dimensionality** (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 100 points onto a real interval, say $[0,1]$. To get similar coverage, in terms of distance between adjacent points, of the 10-dimensional space $[0,1]^{10}$ would require $100^{10} = 10^{20}$ points. A 100 points appear now as isolated points in a vast empty space.

Consequence: a **search policy** (e.g. exhaustive search) that is valuable in small dimensions **might be useless** in moderate or large dimensional search spaces.

Curse of Dimensionality

How long would it take to evaluate 10^{20} points?

Curse of Dimensionality

How long would it take to evaluate 10^{20} points?

```
import timeit
timeit.timeit('import numpy as np ;
np.sum(np.ones(10)*np.ones(10))', number=1000000)
> 7.0521080493927
```

7 seconds for 10^6 evaluations of $f(x) = \sum_{i=1}^{10} x_i^2$

We would need more than 10^8 days for evaluating 10^{20} points

[As a reference: origin of human species: roughly 6×10^8 days]

Separability

Given $f : x = (x_1, \dots, x_n) \in \mathbb{R}^n \mapsto f(x) \in \mathbb{R}$, let us define the 1-D functions that are cuts of f along the different coordinates:

$$f_{(x_1^i, \dots, x_n^i)}^i(y) = f(x_1^i, \dots, x_{i-1}^i, y, x_{i+1}^i, \dots, x_n^i)$$

for $(x_1^i, \dots, x_n^i) \in \mathbb{R}^{n-1}$, with $(x_1^i, \dots, x_n^i) = (x_1^i, \dots, x_{i-1}^i, x_{i+1}^i, \dots, x_n^i)$

Definition: A function f is **separable** if for all i , for all $(x_1^i, \dots, x_n^i) \in \mathbb{R}^{n-1}$, for all $(\hat{x}_1^i, \dots, \hat{x}_n^i) \in \mathbb{R}^{n-1}$

$$\operatorname{argmin}_y f_{(x_1^i, \dots, x_n^i)}^i(y) = \operatorname{argmin}_y f_{(\hat{x}_1^i, \dots, \hat{x}_n^i)}^i(y)$$

a weak definition of separability

Separability (cont)

Proposition: Let f be a **separable** then for all x_i^j

$$\operatorname{argmin} f(x_1, \dots, x_n) = \left(\operatorname{argmin}_{(x_2^1, \dots, x_n^1)} f^1(x_1), \dots, \operatorname{argmin}_{(x_1^n, \dots, x_{n-1}^n)} f^n(x_n) \right)$$

and f can be optimized using n minimization along the coordinates.

Exercise: prove the previous proposition

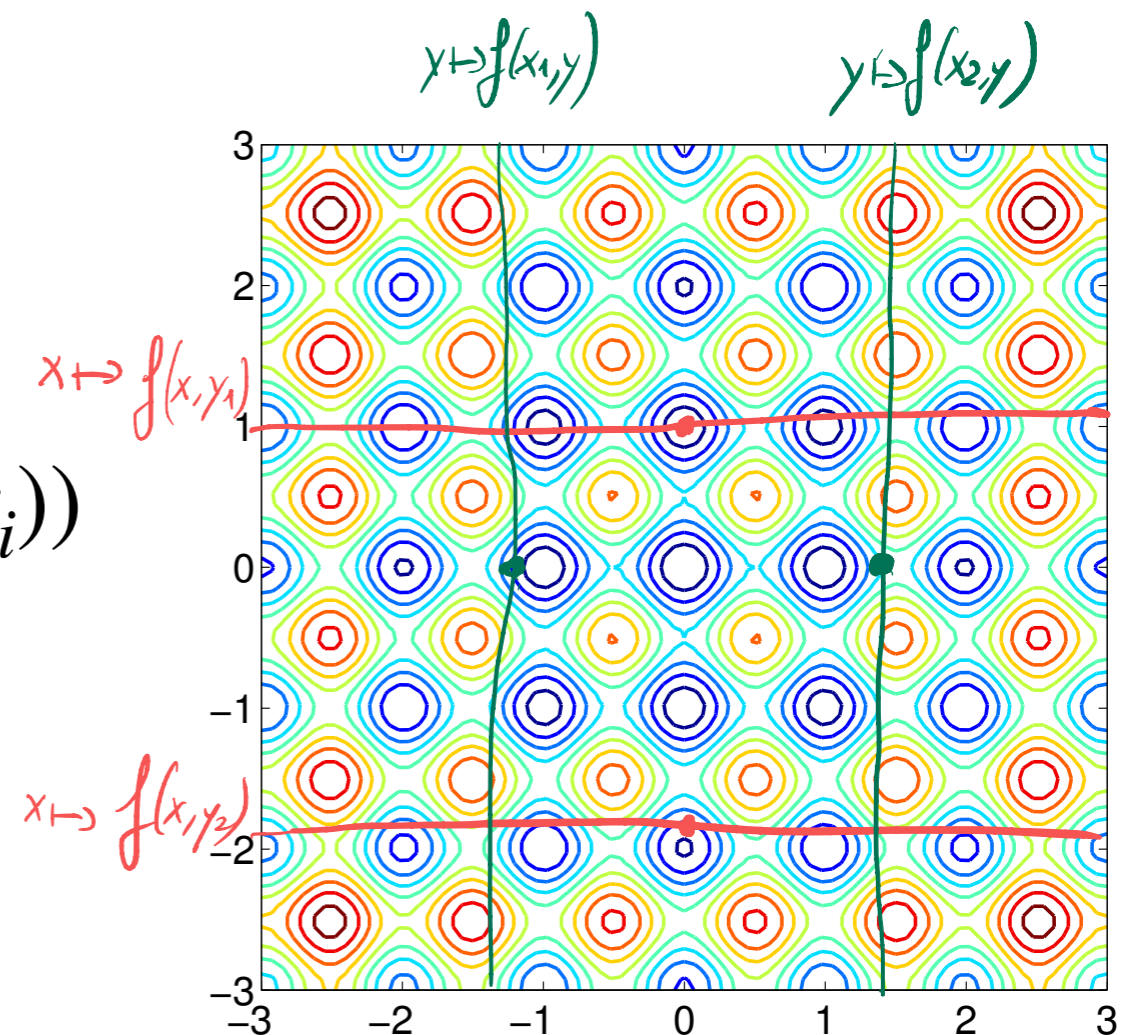
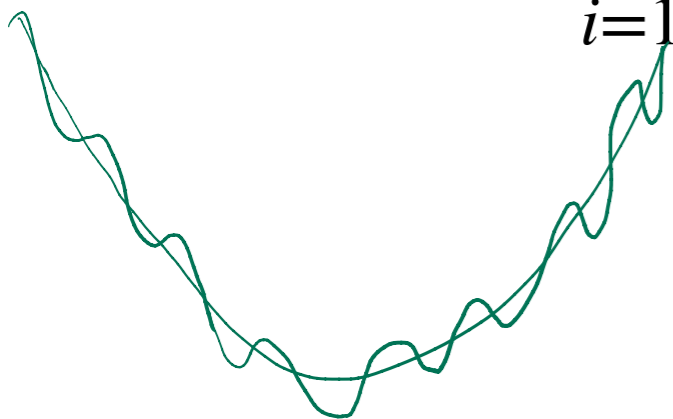
HONE EXERCICE

Example: Additively Decomposable Functions

Exercise: Let $f(x_1, \dots, x_n) = \sum_{i=1}^n h_i(x_i)$ for h_i having a unique argmin. Prove that f is separable. We say in this case that f is additively decomposable.

Example: Rastrigin function

$$f(x) = 10n + \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i))$$



Non-separable Problems

Separable problems are typically easy to optimize. Yet **difficult real-world problems are non-separable.**

One needs to be careful when evaluating optimization algorithms that not too many test functions are separable and if so that the *algorithms do not exploit separability.*

***Otherwise:** good performance on test problems will not reflect good performance of the algorithm to solve difficult problems*

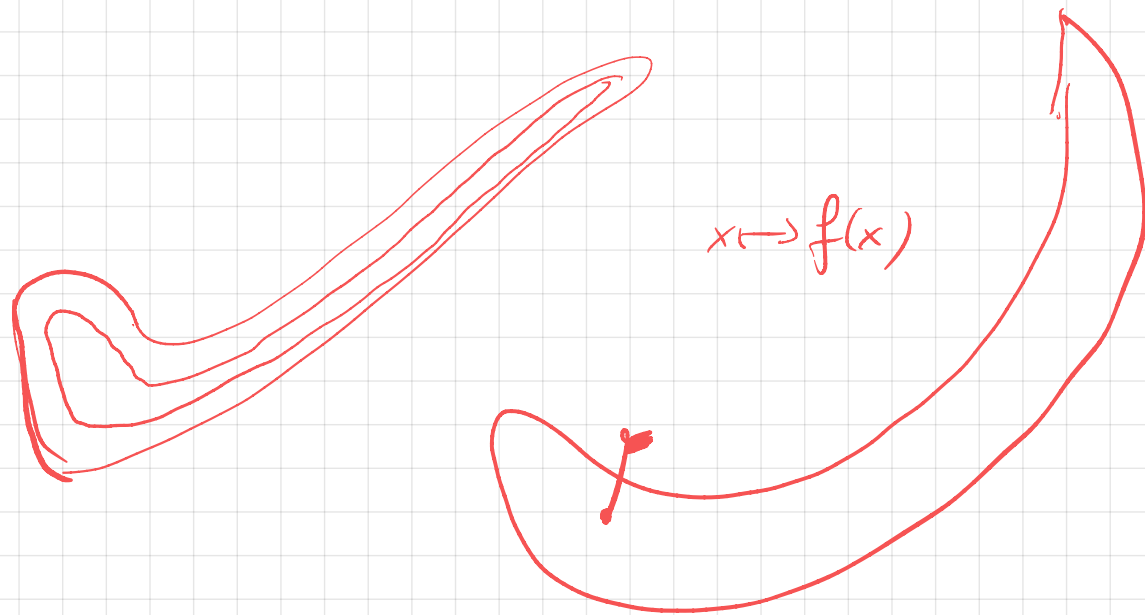
Algorithms known to exploit separability:

Many Genetic Algorithms (GA), Most Particle Swarm Optimization (PSO)

If I give you $f: \begin{cases} \mathbb{R}^n \rightarrow \mathbb{R} \\ x \mapsto f(x) \end{cases}$ which is separable

How can you build a non-separable function?

Rosenbrock function:



3?

A linear $x \mapsto f(Ax)$ is separable. No.

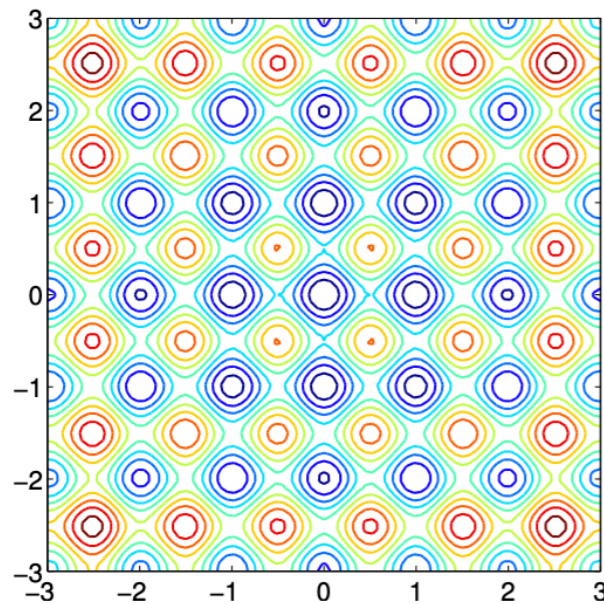
Non-separable Problems

Building a non-separable problem from a separable one

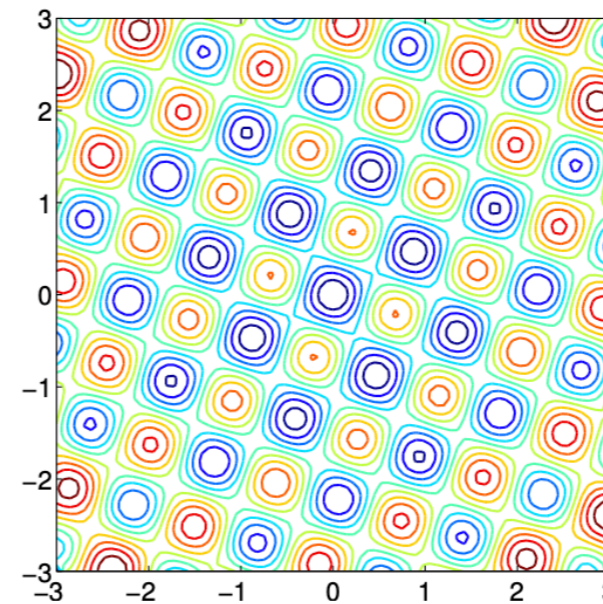
Rotating the coordinate system

- ▶ $f : \mathbf{x} \mapsto f(\mathbf{x})$ separable
- ▶ $f : \mathbf{x} \mapsto f(\mathbf{R}\mathbf{x})$ non-separable

R rotation matrix



R
→



¹ Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

² Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

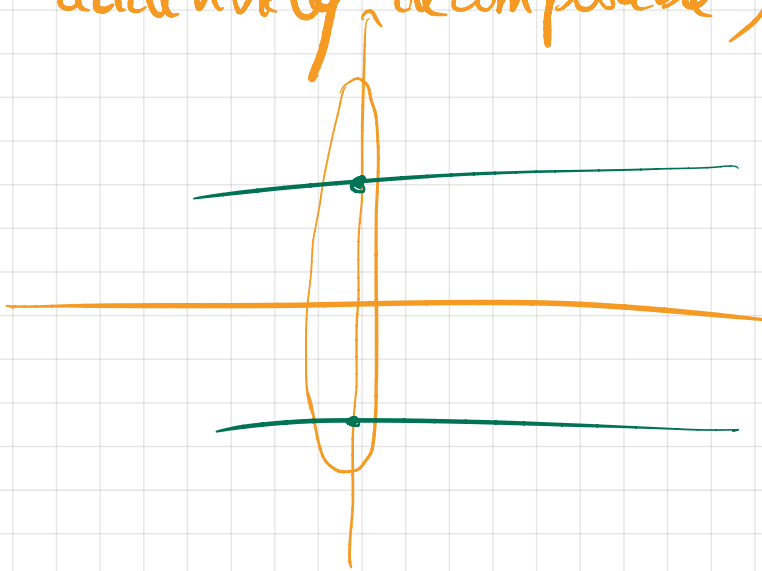
Let $f(x) = \frac{1}{2} x^T A x$ where A is symmetric positive definite.

Is f separable?

If $A = \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix}$, is f separable?

$$f(x) = \underbrace{\frac{9}{2} x_1^2}_{h_1(x_1)} + \underbrace{\frac{1}{2} x_2^2}_{h_2(x_2)} = h_1(x_1) + h_2(x_2)$$

f is then additively decomposable, so it is separable.



If A is diagonal, then f is separable.

If A is not diagonal, then f is not separable.

$$f(x) = \frac{1}{2} x^T A x \quad \text{where } A \text{ is not diagonal}$$

I can write f as the rotation of a separable function:

From the spectral theorem

$$A = P D P^T$$

↑
Diagonal

$$\begin{aligned} f(x) &= \frac{1}{2} x^T P D P^T x \\ &= \frac{1}{2} (P^T x)^T D P^T x \end{aligned}$$

$$= g(P^T x)$$

↑
separable

orthogonal
↑

Rotation

$$g(x) = \frac{1}{2} x^T D x$$

Is separable because D is diagonal

Let f be convex quadratic, i.e. $f = \frac{1}{2} (x-x_0)^T A (x-x_0) + b$
where A is SPD.

$$(f \text{ is separable}) \Leftrightarrow (A \text{ is diagonal})$$

In addition, any convex quadratic function can be written as

$$f(x) = g(Px) \quad \text{where } g \text{ is separable}$$

P is orthogonal

Ill-conditioned Problems - Case of Convex-quadratic functions

Exercise: Consider a convex-quadratic function

$f(x) = \frac{1}{2}(x - x^*)^T H (x - x^*)$ with H a symmetric, positive, definite (SPD) matrix.

~~1. why is it called a convex-quadratic function? What is the H matrix of f ?~~

The condition number of the matrix H (with respect to the Euclidean norm) is defined as

$$\text{cond}(H) = \frac{\lambda_{\max}(H)}{\lambda_{\min}(H)}$$

with $\lambda_{\max}()$ and $\lambda_{\min}()$ being respectively the largest and smallest eigenvalues.

Ill-conditioned Problems

Ill-conditioned means a high condition number of the ~~matrix~~ matrix H .

Consider now the specific case of the function $f(x) = \frac{1}{2}(x_1^2 + 9x_2^2)$

1. Compute its ~~Hessian~~ Hessian matrix, its condition number
2. Plots the level sets of f , relate the condition number to the axis ratio of the level sets of f
3. Generalize to a general convex-quadratic function

Real-world problems are often ill-conditioned.

4. Why do you think it is the case?
5. why are ill-conditioned problems difficult?

~~(see also **Exercise 2.5**)~~

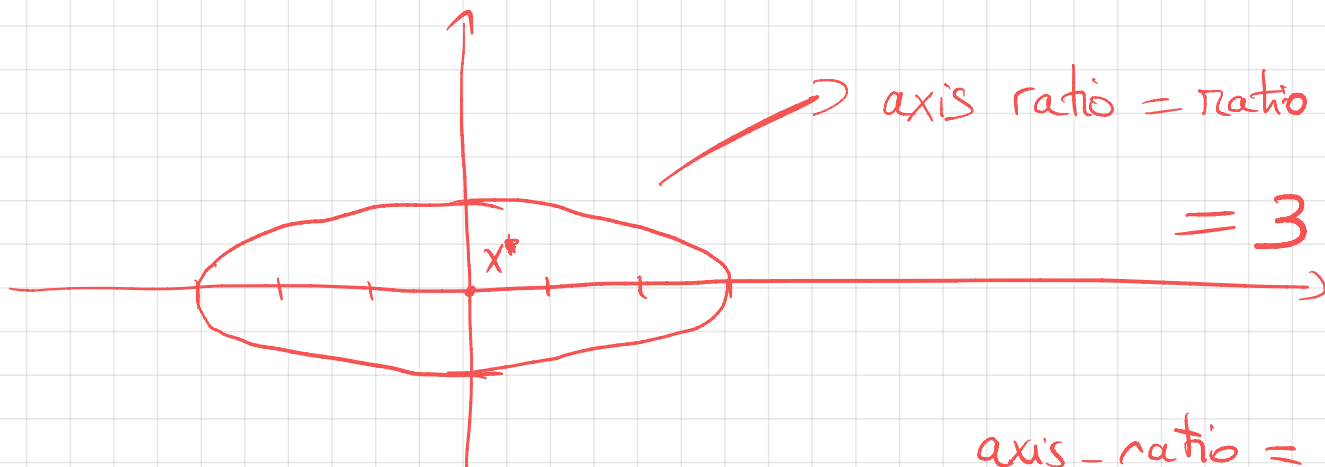
$$f(x) = \frac{1}{2} (x_1^2 + 9x_2^2) = \frac{1}{2} (x - x^*)^T H (x - x^*)$$

$$= \frac{1}{2} x^T \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix} x \quad \left(\begin{array}{l} \text{i.e. } x^* = 0 \\ H = \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix} \end{array} \right)$$

$$x = (x_1, x_2)$$

Then $H = \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix}$ (and $x^* = 0$)

$$\text{cond}(H) = \frac{\lambda_{\max}(H)}{\lambda_{\min}(H)} = \frac{9}{1} = 9$$



axis ratio = ratio between largest and smallest axis of ellipsoid

= 3

$$\text{axis-ratio} = \sqrt{\text{cond}(H)}$$

An ill-conditioned convex-quadratic problem will have a large ratio between the largest axis and smallest axis of the ellipsoid level set.

Why do we often encounter ill-conditioned problems (in the real world)?

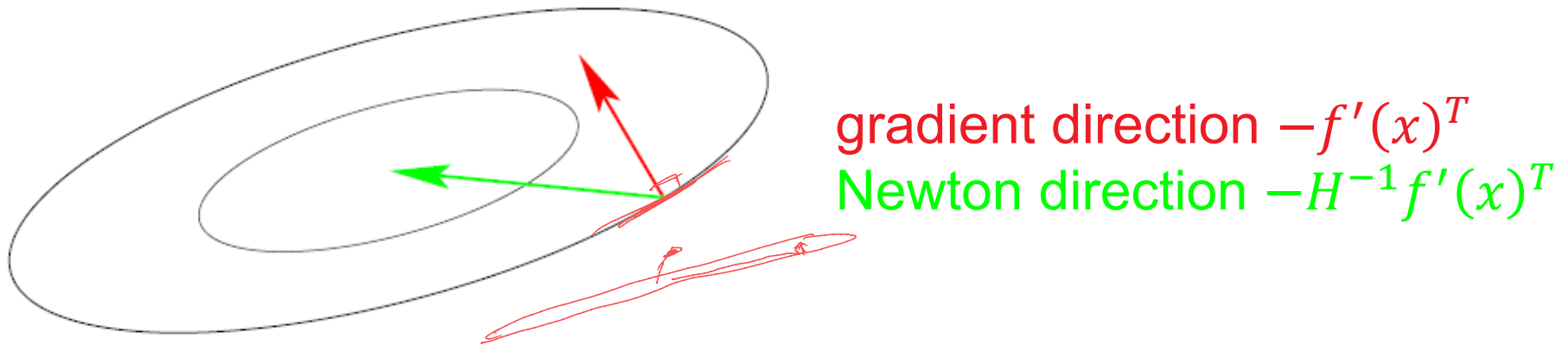
→ Because we optimize often variables that have different units/scales with different orders of magnitude.

III-Conditioned Problems: Curvature of Level Sets

Consider the convex-quadratic function

$$f(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^T H (\mathbf{x} - \mathbf{x}^*) = \frac{1}{2} \sum_i h_{i,i} x_i^2 + \frac{1}{2} \sum_{i,j} h_{i,j} x_i x_j$$

H is Hessian matrix of f and symmetric positive definite



*Ill-conditioning means **squeezed level sets** (high curvature).
Condition number equals nine here. Condition numbers up to 10^{10}
are not unusual in real-world problems.*

If $H \approx I$ (small condition number of H) first order information (e.g. the gradient) is sufficient. Otherwise **second order information** (estimation of H^{-1}) information necessary.

Reminder: Different Notions of Optimum

Unconstrained case

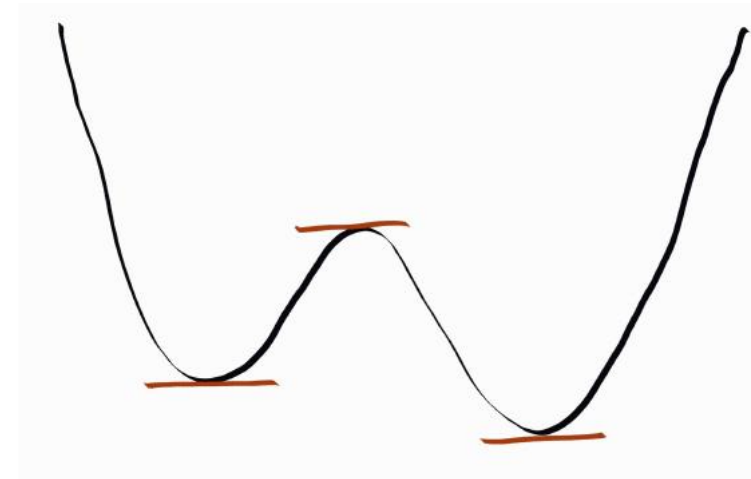
- local vs. global
 - local minimum \mathbf{x}^* : \exists a neighborhood V of \mathbf{x}^* such that $\forall \mathbf{x} \in V: f(\mathbf{x}) \geq f(\mathbf{x}^*)$
 - global minimum: $\forall \mathbf{x} \in \Omega: f(\mathbf{x}) \geq f(\mathbf{x}^*)$
- strict local minimum if the inequality is strict



Mathematical Characterization of Optima

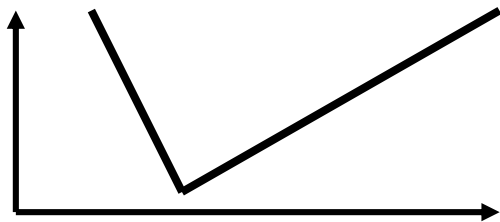
Objective: Derive general characterization of optima

Example: if $f: \mathbb{R} \rightarrow \mathbb{R}$ differentiable,
 $f'(x) = 0$ at optimal points



- generalization to $f: \mathbb{R}^n \rightarrow \mathbb{R}$?
- generalization to constrained problems?

Remark: notion of optimum independent of notion of derivability

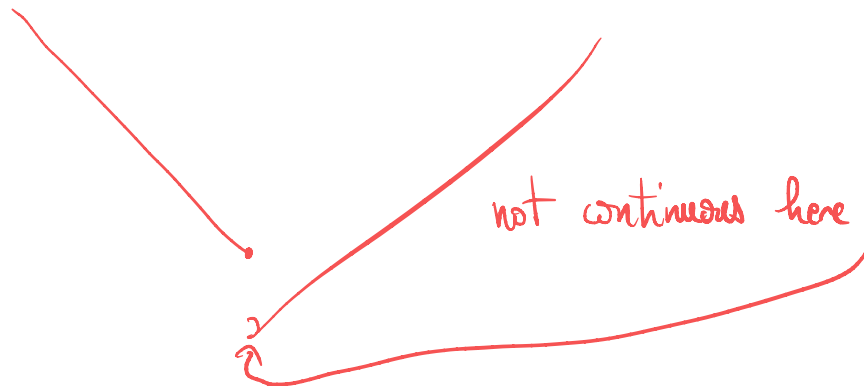
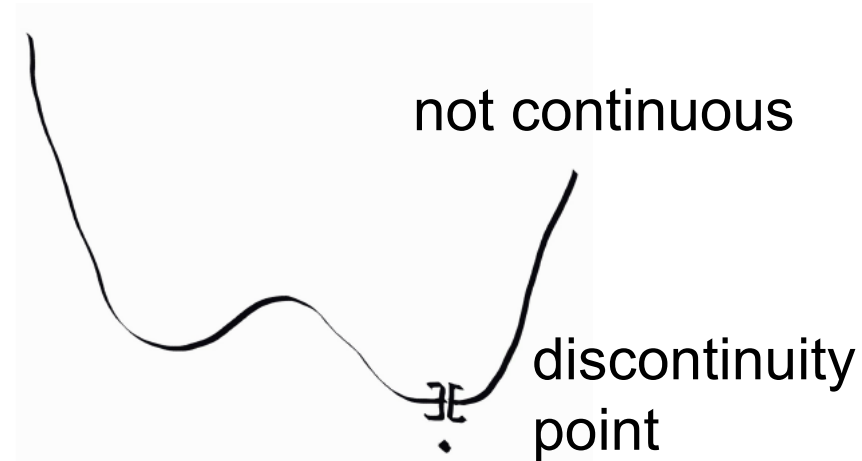
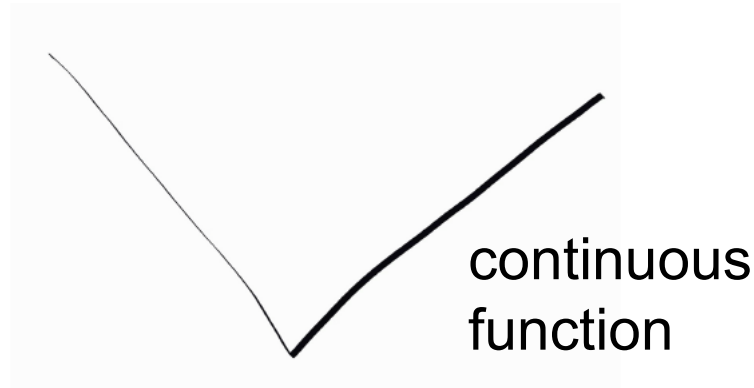


optima of such function can be easily
approached by certain type of methods

Reminder: Continuity of a Function

$f: (V, \| \cdot \|_V) \rightarrow (W, \| \cdot \|_W)$ is continuous in $x \in V$ if

$\forall \epsilon > 0, \exists \eta > 0$ such that $\forall y \in V: \|x - y\|_V \leq \eta; \|f(x) - f(y)\|_W \leq \epsilon$



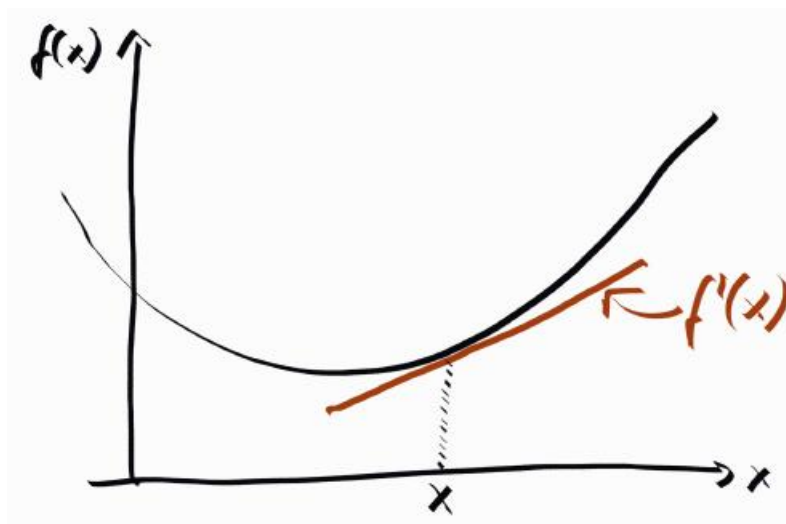
Reminder: Differentiability in 1D ($n=1$)

$f: \mathbb{R} \rightarrow \mathbb{R}$ is differentiable in $x \in \mathbb{R}$ if

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \text{ exists, } h \in \mathbb{R}$$

Notation:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$



The derivative corresponds to the slope of the tangent in x .