## DESCENT METHODS

### General principle

1/ Choose an initial point $X_0 \in \mathbb{R}^n$
$t = 0$

While not happy (while $f$ not minimize enough)

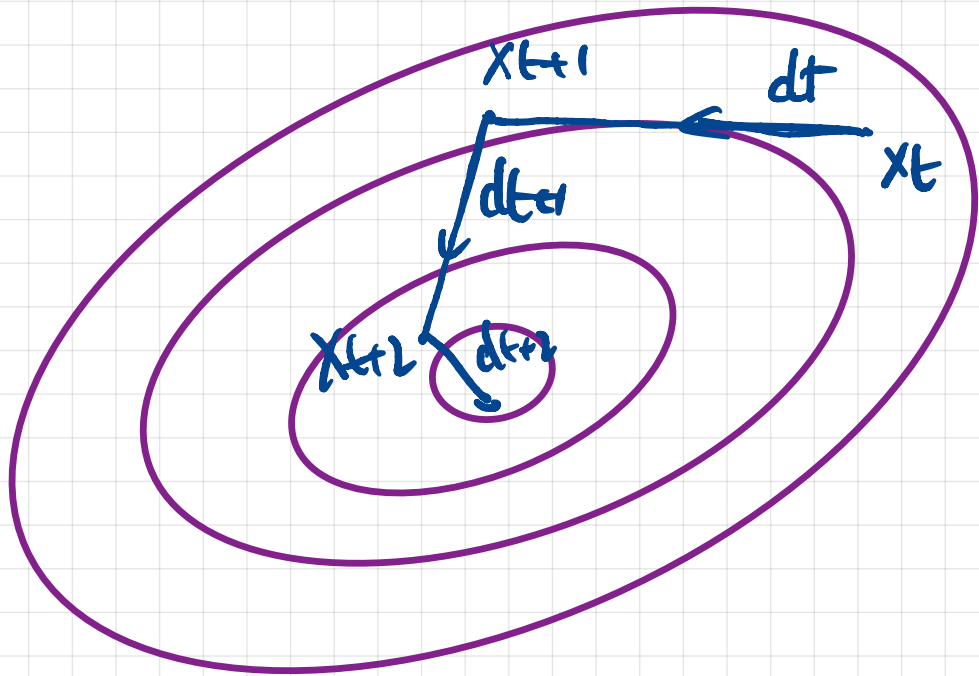- choose a descent direction $d_t \neq 0$ $d_t \in \mathbb{R}^n$
- line search
  - → choose step-size $\sigma_t > 0$
  - → set $X_{t+1} = X_t + \sigma_t d_t$

set $t = t+1$



descent direction

$X_t$ $d_t$

$X_{t+1} = X_t + \sigma_t d_t$

Remains to decide :- what is dt ?

— how to choose $\sigma t$.

How to choose a descent direction ?

Proposition

$$dt = -\nabla f(x_t)$$

For a small enough step-size $\sigma$, if I follow $-Df(x_t)$ then $f\left(x_t + \sigma\left(-Df(x_t)\right)\right) < f(x_t)$

If $\sigma$ is small enough

$$f\left(x_t - \sigma\, Df(x_t)\right) < f(x_t) \text{ and } Df(x_t) \neq 0$$

## Taylor Formula    $f$. Differentiable

$$f(x+h) = f(x) + h^T Df(x) + O\left(\|h\|^2\right)$$

$$\left(\underline{\text{Also true}} : \quad f(x+h) = f(x) + h^T Df(x) + o\left(\|h\|\right)\right)$$

$$\frac{O\left(\|h\|^2\right)}{\|h\|^2} < C \quad \left(\begin{array}{c}\text{Bounded}\\ \text{for } h \text{ small}\end{array}\right) \qquad \frac{o\left(\|h\|\right)}{\|h\|} \to 0 \quad \|h\| \to 0$$

Apply Taylor Formula to $f(x_t - \sigma Df(x_t))$

$$f\left(x_t \underbrace{- \sigma Df(x_t)}_{h}\right) = f(x_t) + \left(-\sigma Df(x_t)\right)^T Df(x_t) + O\left(\|\sigma Df(x_t)\|^2\right)$$

$$= f(x_t) - \sigma \| Df(x_t)\|^2 + \sigma^2 O\left(\| Df(x_t)\|^2\right)$$

For $\sigma$ small enough

$$f(x_t - \sigma Df(x_t)) \approx f(x_t) \underbrace{- \sigma \| Df(x_t)\|^2}_{< 0}$$

$$< f(x_t)$$

$\hookrightarrow$ $Df(x_t)$ is a descent direction.

$$O\left(\sigma^2 \| Df(x)\|^2\right) = \sigma^2 O\left(\| Df(x)\|^2\right)$$

$$O\left(\sigma \|h\|\right) = \sigma \, O\left(\|h\|\right)$$

# CHOICE OF THE STEP-SIZE

$$f\left(x_t - \sigma_t \; x_t\right)$$

$$\sigma_t = \frac{\sigma}{cond\left(\nabla^2 f(x_t)\right)}$$

$\rightarrow$ constant ?

$\hookrightarrow$ small enough.

What is the optimal step-size?



$$f\left(x_t - \sigma \, \nabla f(x_t)\right)$$

$$\sigma \overset{g}{\longmapsto} f\left(x_t - \sigma \, \nabla f(x_t)\right)$$

$\mathbb{R}$

I can minimize

$$\sigma \longmapsto f\left(x_t - \sigma \nabla f(x_t)\right)$$

$\sigma \in \mathbb{R}$

1D optimization

minimize $f$ along-the gradient direction starting from $x_t$.

## optimal step. size

$$\sigma_t = \underset{\sigma \geq 0}{\arg\min} \; f(x_t - \sigma \nabla f(x_t))$$

Typically too expensive to do those 1D optimization perfectly

There exists different rules to approximate optimal step-size. One widely used is called Armijo rule.
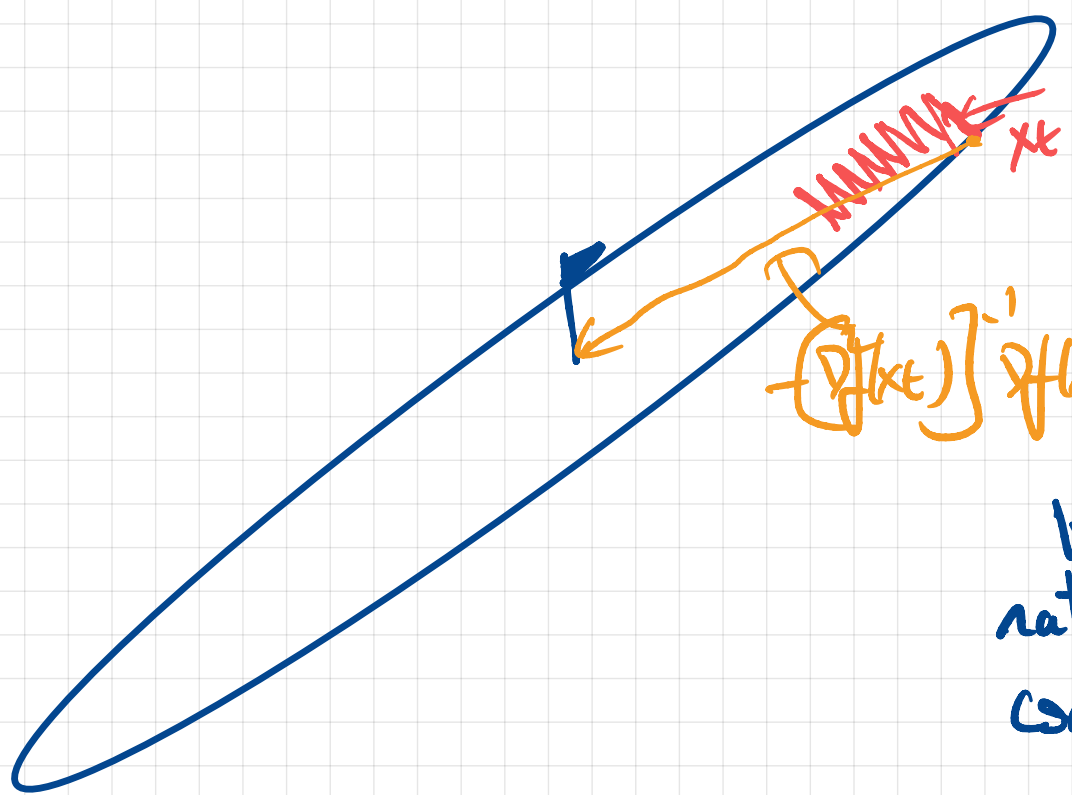
## Stopping criterion:

- If $|f(x_{t+1}) - f(x_t)| < \varepsilon$
- Stop after a certain number of iterations
  $\hookrightarrow$ no guarantee on whether the algorithm has converged.

- Stop if $\|Df(x_t)\|$ small

Remark: If you need to maximize $f$ and you only have a minimizer implemented, then just minimize $-f$.

If I implement gradient descent but need to maximize, I can implement gradient ascent.

$$x_{t+1} = x_t + \sigma t \, Df(x_t)$$

## Gradient descent is slow on ill-conditionned problems



$-Df(x_t)$

$x_t$

$-[Df(x_t)]^{-1} Df(x_t)$

On a ill-conditionned function $-Df$ typically points in the wron direction and the convergence will be very slow.

We can prove that the convergence rate is slower the larger the condition number.

Instead of $-Df(x_t)$, we could follow the Newton direction $-\left[D^2 f(x_t)\right]^{-1} Df(x_t)$

The Newton direction minimizes the locally quadratic approximation of $f$.

$$f(x + \Delta x) = \underbrace{f(x) + Df(x)^T \Delta x + \frac{1}{2} (\Delta x)^T Df(x) \Delta x}_{\text{quadratic approximation of } f}$$

Newton direction minimizes this.

If we can obtain in a "cheap" way the Newton direction, we should use it.

But often too expensive to obtain $D^2 f(x_t)$ and to invert it.

In the convex-quadratic case, the function equals its second order approximation and the Newton direction is perfect as it points towards the optimum.

For non convex-quadratic case, the Newton is typically good to follow but not point directly towards the optimum.

## QUASI-NEWTON METHOD : BFGS (70's)
L-BFGS

Broyden    Fletcher    Goldfarb  Shannon

$$x_{t+1} = x_t - \partial_t H_t \, Df(x_t)$$

approximation of $\left[ \nabla^2 f(x_t) \right]^{-1}$

Ht is updated iteratively using $Df(x_t)$ ( without computing the Hessian matrix) and it approximates $[D^2 f(x_{t-1})]^{-1}$

cf wikepedia page for update.

Large scale version of BFGS : L-BFGS

Limited-memory BFGS

# STOCHASTIC GRADIENT DESCENT

Minimize loss function of the following form:

$$P(w) = \frac{1}{N} \sum_{i=1}^{N} Q_i(w) \qquad N : \# \text{ Date}$$

$$\# \text{ Examples}$$

w can be weights of Neural Network.

How do we minimize Q?

Gradient descent: $\quad DQ(w) = \frac{1}{N} \sum_{i=1}^{N} DQ_i(w)$

$$w_{t+1} = w_t - \sigma_t \, DQ(w_t) \quad [\text{update of weights}]$$

BACK PROPAGATION algorithm to compute $DQ_i(w)$

Typically N is very large, computation of all $\nabla Q_i(w)$ $i = 1, \ldots, N$ too expensive.

Instead we approximate $\nabla Q(w)$

$$\nabla Q(w) \approx \nabla Q_i(w) \quad \left(\begin{array}{l}\text{Gradient of} \\ \text{single example}\end{array}\right)$$

Also do mini batches:

$$\nabla Q(w) \approx \frac{1}{n\text{ batches}} \sum_{r=1}^{n\text{batches}} \nabla Q_i(w)$$

$$n\text{ batches} \ll N$$

# STOCHASTIC Gradient descent :

Choose an initial vectors of parameters and a step-size $\eta$

While not happy

- Randomly shuffle examples in training set
- For $i = 1, \ldots, N$

$$w \longleftarrow w - \eta \, D\, \mathcal{Q}i\,(w)$$

(possibly mini-batches)

Not covered : - choice of step-size ( step-size adapted using "momentum techniques" ADAM )

- increase # elements of mini-batches