

## First and second order methods

① First order methods : involving only the gradient information.

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$$

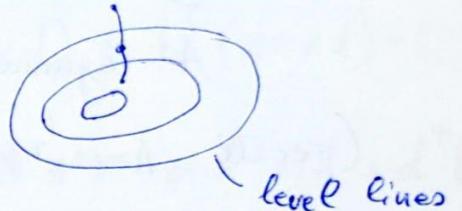
↑  
descent direction

$$f(x^{(k+1)}) < f(x^{(k)})$$

① Gradient descent = steepest descent

Denote  $g^{(k)} = \nabla f(x^{(k)})$ . Then in the method

we have  $d^{(k)} = -\frac{g^{(k)}}{\|g^{(k)}\|}$



\* What is the actual path?

Ex..  $f(x) = x_1 x_2$ .  $\nabla f(x) = (x_2^2, 2x_1 x_2)$

$$x^{(k)} = (1, 2) \Rightarrow \nabla f(x^{(k)}) = (2^2, 2 \cdot 1 \cdot 2) = (4, 4)$$

$$\Rightarrow d^{(k)} = -\frac{1}{\sqrt{4^2+4^2}} (4, 4) = \left(-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}\right)$$

\* What is the actual path? Suppose we do the exact line search, that is at each step we solve

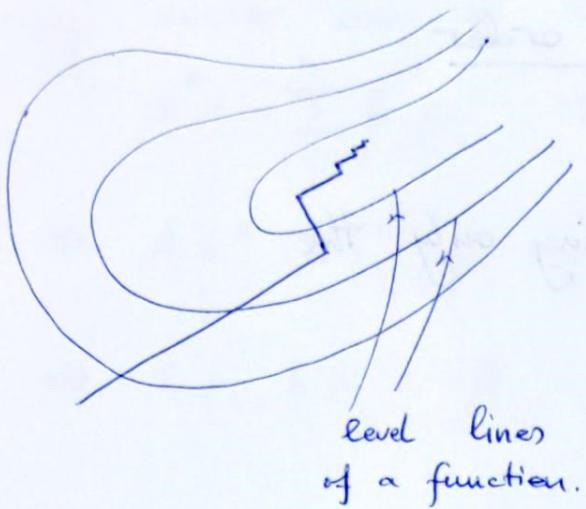
$$\alpha^{(k)} = \arg \min_{\alpha} f(x^{(k)} + \alpha d^{(k)})$$

$$\Rightarrow \frac{d}{d\alpha} f(x^{(k)} + \alpha d^{(k)}) \Big|_{\alpha=\alpha^{(k)}} = 0 \Rightarrow \nabla f(x^{(k)} + \alpha^{(k)} d^{(k)})^T \cdot d^{(k)} = 0$$

$$\text{But } d^{(k+1)} = -\frac{\nabla f(x^{(k)} + \alpha^{(k)} d^{(k)})}{\|\nabla f(x^{(k)} + \alpha^{(k)} d^{(k)})\|} \Rightarrow (d^{(k+1)})^T \cdot d^{(k)} = \langle d^{(k+1)}, d^{(k)} \rangle = 0 \Rightarrow$$

\* Note: Approximate line search doesn't satisfy this

$$d^{(k+1)} \perp d^{(k)}$$



(\*) Hence, it might take many steps to progress in a 'long valley'

## ② Conjugate gradient descent

This method overcomes (\*): suppose we want to optimize a quadratic function

$$\min_x f(x) = \frac{1}{2} x^T A x + b^T x + c,$$

$A$  symmetric, positive definite

(recall,  $A > 0 \Leftrightarrow x^T A x > 0 \quad \forall x \in \mathbb{R}^n \setminus \{0\}$ )

(also recall, that  $A$  pos. def.  $\Leftrightarrow$

For non-quadratic  
we apply using  
quadratic approxi-  
mation

$\Rightarrow \det A > 0 \Rightarrow A$  invertible

Sylvester's criterion.

$$\nabla f(x) = Ax + b = 0 \Rightarrow x^* = -A^{-1}b.$$

$\Rightarrow x^*$  is global minimum.

$$\nabla^2 f(x) = Hf(x) = A > 0$$

- We can define a scalar product:

$$\langle x, y \rangle_A := x^T A y = \langle x, Ay \rangle \leftarrow \text{since it is symmetric, order of vectors doesn't matter.}$$

- Two vectors  $x, y$  are conjugate iff  $\langle x, y \rangle_A = 0 \Leftrightarrow x^T A y = 0$ .
- Take a basis in this space:  $\{p_1, \dots, p_n\}$

Any vector can be expressed in this basis: in particular,

$$x^* = \sum_{i=1}^n \beta_i p_i \Rightarrow$$

$$\Rightarrow Ax^* = \sum_{i=1}^n \beta_i A p_i \Rightarrow b = \sum_{i=1}^n \beta_i A p_i \Rightarrow$$

$$\Rightarrow \langle p_i, b \rangle = \beta_i \langle A p_i, p_i \rangle \Rightarrow \boxed{\beta_i = \frac{\langle p_i, b \rangle}{\langle p_i, p_i \rangle_A}}.$$

- We want to find it iteratively:

(1) Start with  $d^{(1)} = -g^{(1)}$

(2) Compute  $\alpha$  exactly:

Suppose  $\min_{\alpha} f(x + \alpha d)$ :

$$\frac{\partial f(x + \alpha d)}{\partial \alpha} = \frac{\partial}{\partial \alpha} \left[ \frac{1}{2} (x + \alpha d)^T A (x + \alpha d) + b^T (x + \alpha d) + c \right]$$

$$= d^T A (x + \alpha d) + d^T b = d^T A x + \alpha d^T A d + d^T b$$

$$= d^T (A x + b) + \alpha d^T A d = 0 \Rightarrow$$

$$\Rightarrow \boxed{\alpha = -\frac{d^T (A x + b)}{d^T A d}}$$

$$\alpha^{(k)} = -\frac{d^{(k)T} g^{(k)}}{\langle d^{(k)}, d^{(k)} \rangle_A} = -\frac{\langle d^{(k)}, g^{(k)} \rangle}{\langle d^{(k)}, d^{(k)} \rangle_A}$$

$$\Rightarrow x^{(2)} = x^{(1)} + \alpha^{(1)} d^{(1)}$$

Same for  $x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$

(3)  $\{d^{(k+1)} = -g^{(k+1)} + \beta^{(k)} d^{(k)}\}$ , where  $\beta^{(k)}$  is

chosen so that  $\langle d^{(k+1)}, d^{(k)} \rangle_A = 0$  (conjugate direction)

$$0 = d^{(k+1)T} A d^{(k)} = (-g^{(k+1)T} + \beta^{(k)} d^{(k)T}) A d^{(k)} =$$

$$= -g^{(k+1)T} A d^{(k)} + \beta^{(k)} d^{(k)T} A d^{(k)} \Rightarrow \beta^{(k)} = \frac{g^{(k+1)T} A d^{(k)}}{\langle d^{(k)T} A d^{(k)} \rangle}$$

$$\Leftrightarrow \boxed{\beta^{(k)} = \frac{\langle g^{(k+1)}, d^{(k)} \rangle_A}{\langle d^{(k)}, d^{(k)} \rangle_A}} \quad \text{Recall: } \langle g^{(k+1)}, d^{(k)} \rangle_A = 0$$

- Fletcher - Reeves:

$$\beta^{(k)} = \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k-1)T} \mathbf{g}^{(k-1)}}$$

- Polak - Ribière:

$$\beta^{(k)} = \frac{\mathbf{g}^{(k)T} (\mathbf{g}^{(k)} - \mathbf{g}^{(k-1)})}{\mathbf{g}^{(k-1)T} \mathbf{g}^{(k-1)}}$$

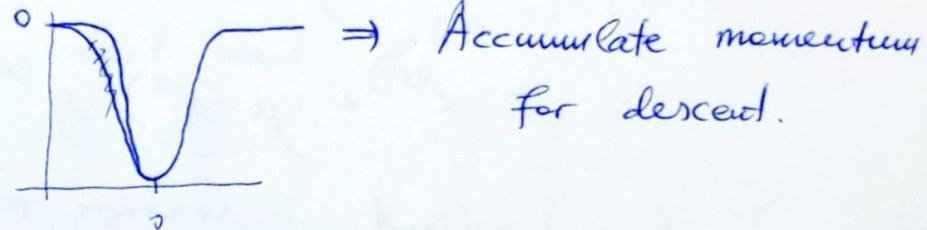
Thm: For quadratic func's, the conjugate gradient descent converges in at most  $n$  (= dimension) steps.

w/o proof

If we don't know / want to compute Hessian, we continue as follows:

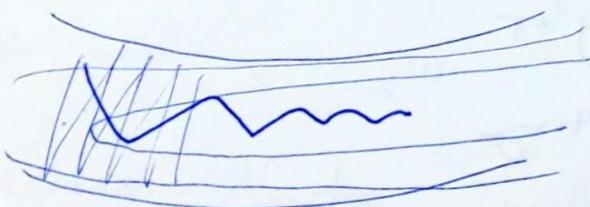
### ③ Momentum:

Plato's:  $-e^{-x^2}$



$$\begin{cases} \mathbf{v}^{(k+1)} = \beta \mathbf{v}^{(k)} - \alpha \mathbf{g}^{(k)} \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{v}^{(k+1)} \end{cases}$$

Cumulative momentum?



### ④ Nesterov momentum:

$$\begin{cases} \mathbf{v}^{(k+1)} = \beta \mathbf{v}^{(k)} - \alpha \nabla f(\underbrace{\mathbf{x}^{(k)} + \beta \mathbf{v}^{(k)}}_{\text{future position}}) \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{v}^{(k+1)} \end{cases}$$

### ⑤ Adagrad (adaptive gradient)

For sparse gradients:

$$s_i^{(k)} = \sum (g_i^{(j)})^2$$

$$x_i^{(k+1)} = x_i^{(k)} - \frac{\alpha}{\epsilon + \sqrt{s_i^{(k)}}} g_i^{(k)}$$

↑  
infrequent updates.  
-#-

$\alpha \cdot 10^{-8}$   
non-decreasing  
Learning rate drops.

$\Rightarrow$  accelerated methods. (RMS Prop, Adadelta, Adam)

adaptive momentum.

⑥ Hypergradient descent: adapt learning rate.

$\Rightarrow$  derivative of learning rate.

$$\frac{\partial f(x^{(k)})}{\partial \alpha} = -(g^{(k)})^T g^{(k-1)}$$

$$\Rightarrow \alpha^{(k+1)} = \alpha^{(k)} - \mu \frac{\partial f(x^{(k)})}{\partial \alpha} =$$

$$= \alpha^{(k)} + \mu (g^{(k)})^T g^{(k-1)} \cdot \text{Two step to remember}$$

$\mu$  = hypergradient learning rate.

Can be combined w/ previous methods.

Ex: Do 2 steps of conjugate grad. descent

$$f(x,y) = x^2 + xy + y^2 + 5$$

$(x,y) = (1,1)$  starting

Do 2 steps of hypergrad. descent

② Second order methods.: Use quadratic approx.

Newton's method:

$$x^{(k+1)} = x^{(k)} - (H^{(k)})^{-1} g^{(k)}$$

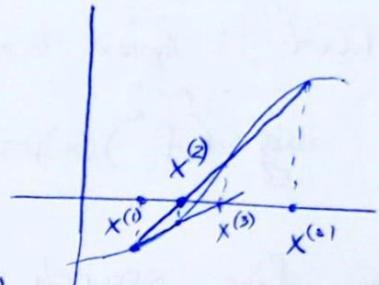
(in 1D:  $x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})}$ )

roots of  $f'$ .

- Works well near minimum  
(when  $H^{(k)} \succ 0$ )

① Secant method: don't want  $f''$   
(in 1D)

$$f''(x^{(k)}) \approx \frac{f'(x^{(k)}) - f'(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$



and  $x^{(k+1)} - x^{(k)} = \frac{x^{(k)} - x^{(k-1)}}{f'(x^{(k)}) - f'(x^{(k-1)})} f'(x^{(k)})$

We need  $x^{(0)}, x^{(1)}$

② Quasi-Newton: "multidim. secant"

[see below]

~~$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} Q^{(k)} g^{(k)}$$~~

~ inverse Hessian

• Davidon - Fletcher - Powell (DFP)  $Q^{(1)} = \frac{\text{Id}}{\gamma} \otimes (Q\gamma)^T$

~~$$\gamma^{(k+1)} = g^{(k+1)} - g^{(k)}$$~~

~~$$\delta^{(k+1)} = x^{(k+1)} - x^{(k)}$$~~

~~$$Q \leftarrow Q - \frac{Q \delta (\delta^T Q)}{\delta^T Q \delta} + \frac{\delta \delta^T}{\delta^T \delta}$$~~

1)  $Q \nabla f(x)$  remains symmetric, positive definite

2)  $f(x) = \frac{1}{2} x^T A x + b^T x + c \Rightarrow Q = A^{-1}$

- - -

(l) Quasi-Newton method: 'multidim secant'

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} H^{(k)} g^{(k)}$$

$\uparrow$   
~ inverse Hessian

There are two methods

- Davidon - Fletcher - Powell (DFP) historically first
- Broyden - Fletcher - Goldfarb - Shanno (BFGS)  
(also, L-BFGS, with limited memory)

These methods are ~~of~~ dual to each other.

$$Q_k(x) = y^{(k)} + (g^{(k)})^T (x - x^{(k)}) + \frac{1}{2} (x - x^{(k)})^T H^{(k)} (x - x^{(k)})$$

Do the line search (for the starting function):

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} (H^{(k)})^{-1} g^{(k)}$$

This gives  $x^{(k+1)}$ . Now we need to find the next  $H^{(k+1)} \sim H(x^{(k+1)}) = \nabla^2 f(x^{(k+1)})$ .

Conditions:

$$\left. \begin{array}{l} \nabla Q_{k+1}(x^{(k+1)}) = g^{(k+1)} \\ \nabla Q_{k+1}(x^{(k)}) = g^{(k)} \\ Q_{k+1}(x^{(k+1)}) = y^{(k+1)} \end{array} \right\} \Rightarrow H^{(k+1)} \underbrace{(x^{(k+1)} - x^{(k)})}_{s^{(k)}} = \underbrace{g^{(k+1)} - g^{(k)}}_{y^{(k)}}$$

$$\Rightarrow \boxed{H^{(k+1)} s^{(k)} = y^{(k)}} \quad \boxed{\text{secant condition}}$$

It is positive definite if  $(s^{(k)})^T y^{(k)} > 0$  curvature condition  
(related to convexity of  $f$ ; can be enforced)

To maintain symmetry & positive definiteness,

we choose

$$H^{(k+1)} = H^{(k)} + \alpha uu^T + \beta vv^T, \quad \alpha, \beta > 0$$

↑                      ↓  
matrix

\* Note that:  $w^T(uu^T)w = \langle u, w \rangle^2 \geq 0$  iff  $w \neq 0$  and  $u \perp w$

Choose:  $u = y^{(k)}$  and  $v = H^{(k)} s^{(k)}$

$$H^{(k+1)} s^{(k)} = y^{(k)} = H^{(k)} s^{(k)} + \alpha y^{(k)} (y^{(k)})^T s^{(k)} + \\ + \beta H^{(k)} s^{(k)} \underbrace{(s^{(k)})^T (H^{(k)})^T s^{(k)}}_{\text{some vectors.}} \Rightarrow$$

$$\Rightarrow y^{(k)} \left( 1 - \alpha \cdot \underbrace{(y^{(k)})^T \cdot s^{(k)}}_{\geq 0 \text{ (curvature cond.)}} \right) = H^{(k)} s^{(k)} \left( 1 + \beta \cdot \underbrace{(s^{(k)})^T (H^{(k)})^T s^{(k)}}_{\geq 0 \text{ by induction}} \right)$$

$H^{(k)}$  is pos. def.

$$\Rightarrow \alpha = \frac{1}{(y^{(k)})^T \cdot s^{(k)}} \quad \beta = - \frac{1}{(s^{(k)})^T (H^{(k)})^T s^{(k)}}$$

Conclusion:

$$H^{(k+1)} = H^{(k)} + \frac{y^{(k)} (y^{(k)})^T}{(y^{(k)})^T \cdot s^{(k)}} - \frac{H^{(k)} s^{(k)} (s^{(k)})^T (H^{(k)})^T}{(s^{(k)})^T (H^{(k)})^T s^{(k)}}$$

$$H^{(1)} = \text{id}$$

Exercise: Do 1 step in Newton method for

$$f(x) = f(x_1, x_2) = (x_1 + 1)^2 + (x_2 + 3)^2 + 4$$

starting at  $(0, 0)$ .

# Stochastic Gradient Descent

In steepest descent,

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \underbrace{\nabla f(x^{(k)})}_{\text{we want to change only this!}}, \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

What optimization problems we are interested in?

normalization  $\min_x \frac{1}{n} \sum_{i=1}^n f_i(x)$  (finite sum problems)

① Why so?  $(x_1, y_1), \dots, (x_m, y_m) \in \mathbb{R}^n \times Y$  training data.

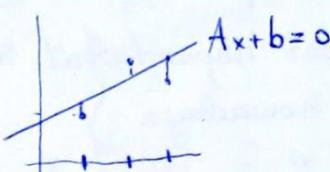
Both  $n$  and  $m$  are large!

think of images

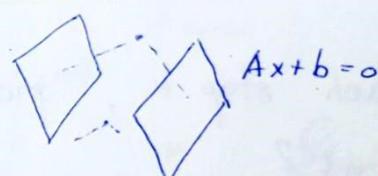
think of users in social networks.

- Yes/No Classification
- Numbers Regression

\*  $\frac{1}{n} \|Ax - b\|_2^2 = \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i)^2 = \frac{1}{n} \sum f_i(x)$  Least-squares



or minimize distance to hyperplanes



\*  $\frac{1}{n} \|Ax - b\|_2^2 + \lambda \|x\|_1 = \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i)^2 + \lambda \sum_{j=1}^d |x_j|$

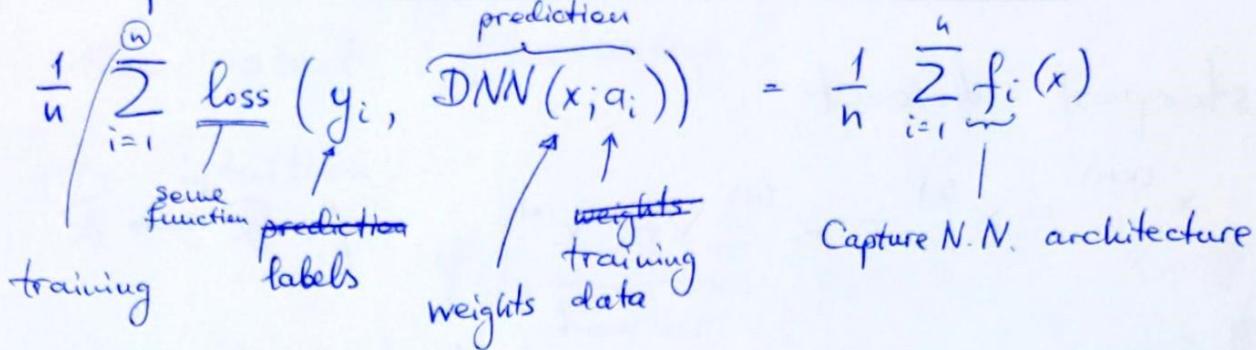
coming from constrained optimization

Lasso  $\ell_1$ -least squares  
(we will see in constraint optimization)  
least absolute shrinkage and selection operator

\*  $\frac{1}{2} \|x\|_2^2 + \frac{C}{n} \sum_{i=1}^n \max\{0, 1 - y_i (x^T a_i + b)\}$  SVM

$$\begin{array}{c:c} \cdot & \cdot \\ \cdot & \cdot \\ \vdots & \vdots \\ i & -1 \end{array} \quad \begin{array}{c} \text{P.S.} \\ \{-1, +1\} \end{array}$$

## \* Deep Neural Networks



Capture N.N. architecture

## \* Maximum likelihood estimations, etc...

$$\rightsquigarrow \min_x \underbrace{\frac{1}{n} \sum_{i=1}^n f_i(x)}_f$$

If we use <sup>fast</sup> gradient descent, then

$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}) = \\ &= x^{(k)} - \alpha^{(k)} \underbrace{\frac{1}{n} \sum_{i=1}^n \nabla f_i(x^{(k)})}_{\text{Huge sum! Computing one step is long}} \end{aligned}$$

Huge sum! Computing one step is long

rightsquigarrow take only 1 component  $\rightsquigarrow$

rightsquigarrow stochastic gradient descent. (the most important method nowadays)

② What if at each step  $k$ , pick randomly an integer

$$i(k) \in \{1, 2, \dots, n\}?$$

And perform the update

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \underbrace{\nabla f_{i(k)}(x^{(k)})}_{\text{STOCHASTIC GRADIENT}}$$

*n times faster than using  $\nabla f$*

Convergence, if  
 $\sum_{k=1}^{\infty} \alpha^{(k)} = \infty$  and  $\sum_k \alpha^{(k)} < \infty$   
 $\alpha^{(k)} \sim \frac{C}{k}$

[Robbins, Monro, 1951] ~ 12250 citations (1 citation every 2 days)

DEMO :  $(1-x_1)^2 + 100(x_2 - x^2)^2$  It is not descent as such!

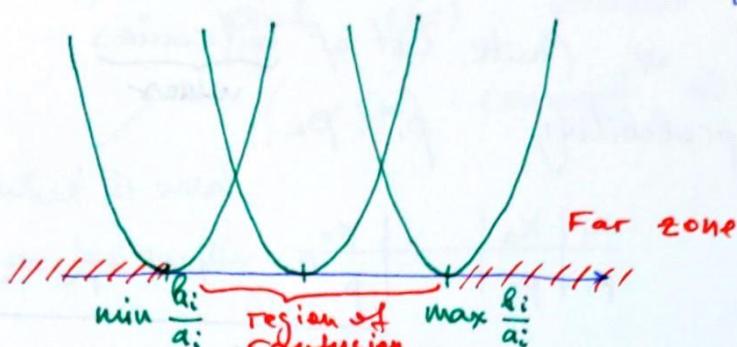
• Best step size??

③ Works really well at the beginning, and then start fluctuating?

1D intuition:

$$\min f(x) = \frac{1}{2} \sum_{i=1}^n (a_i x - b_i)^2, \quad x \in \mathbb{R}$$

$$a_i \in \mathbb{R}$$



$$\nabla f = 0 \Rightarrow \sum a_i (a_i x - b_i) = 0 \Rightarrow$$

$$\Rightarrow x^* = \frac{\sum a_i b_i}{\sum a_i^2}$$

For components, the minimum of  $f_i(x) = \frac{1}{2} (a_i x - b_i)^2$

$$\text{is } x_i^* = \frac{b_i}{a_i}$$

Now let  $R = [\underbrace{\min_i x_i^*}_{x_{\min}^*}, \underbrace{\max_i x_i^*}_{x_{\max}^*}]$ . Therefore,

$$\forall i, x_{\min} \cdot a_i \cdot b_i \leq x_{\max}^* \cdot a_i \Rightarrow$$

$$\Rightarrow x^* \leq x_{\max}^* \frac{\sum a_i^2}{\sum a_i^2} = x_{\max}^*. \text{ Similarly for } x_{\min}.$$

Therefore,  $x^* \in R$ .

If  $x$  lies outside  $R$ , then

$$\nabla f_i(x) = a_i (a_i x - b_i)$$

have the same

$$\nabla f(x) = \sum_{i=1}^n a_i (a_i x - b_i)$$

signs!

$\Rightarrow$  using S.G.  $\nabla f_i(x)$  instead of  $\nabla f(x)$  also gives progress!

Outside we have a component in the direction of true gradient!  
Quickly do reasonable progress is sometimes only what is needed! !

- ④ Stochastic gradient descent uses stochastic gradients'  $\hat{g}(x)$  such that  $x^{(k+1)} = x^{(k)} - \alpha g^{(k)}$  where  $g = \hat{g}$
- $E[\hat{g}(x)] = \nabla f(x)$  • Unbiased estimate of true gradient
- ↑  
expectation (mean, average)

\* Digression:

For a random variable  $X$  w/ finite list of outcomes  $x_1, \dots, x_k$  that occur w/ probability  $p_1, \dots, p_k$ ,

$$E[X] = x_1 p_1 + \dots + x_k p_k$$

$$\begin{array}{c|c|c|c} x_1 & x_2 & \dots & x_k \\ \hline p_1 & p_2 & & p_k \end{array}, \quad p_1 + \dots + p_k = 1$$



$X$  = we see  $\Rightarrow$  # of dots.

$$E[X] = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + \dots + 6 \cdot \frac{1}{6} = 3.5$$

Properties:

$$E[\text{const}] = \text{const}$$

$$E[aX + bY] = aE[X] + bE[Y]$$

$$X \leq Y \Rightarrow E[X] \leq E[Y]$$

If we pick  $\hat{x}$



$\hat{X}$  = # of heads in row (conseq. heads)



$$\begin{array}{c|c|c|c} 0 & 1 & 2 & 3 \\ \hline 1 & 8 & \frac{1}{8} & \end{array}$$

If  $\hat{X} = 0$  :: you win 10 €

If  $\hat{X} = 1$  :: you loose 5 €

$X$  = Expected profit? The rest: loose 1 €

~~$\begin{array}{c|c|c|c} 1 & 2 & 3 & 4 \\ \hline 1 & 8 & \frac{1}{8} & \end{array}$~~

$$\begin{array}{c|c|c} 10 & -5 & -1 \\ \hline \frac{1}{6} & \frac{2}{6} & \frac{3}{6} \end{array}$$

$$E[X] = 10 \cdot \frac{1}{6} - 5 \cdot \frac{2}{6} - 1 \cdot \frac{3}{6} = -\frac{1}{2}$$

For us, we choose  $i(k)$  uniformly at random

$$\Leftrightarrow P(i(k) = s) = \frac{1}{n} \Rightarrow \text{horst average}$$

$$E[\nabla f_{i(k)}(x)] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x) = \nabla f(x).$$

Therefore,  $\nabla f_{i(k)}(x)$  is a stochastic gradient.

We also need to care about variance

$$\text{Var}(X) = E[(X - E(X))^2]$$

(examples above)

Should be small

$$\text{exercist} \quad \leftarrow \text{Var}[\nabla f_{i(k)}(x)] = \frac{1}{n} \sum (\nabla f_i - \nabla f)^2(x)$$

⑤ Practice and theory: which of the two versions to use?

- pick  $i(k)$  uniformly at random

$\nearrow$  theory (  $i(k)$  is independent of  $i(k-1), i(k-2), \dots$  )

- pick  $i(k)$  without replacement

$\nearrow$  (  $i(k)$  cannot be equal to  $i(k-1), i(k-2), \dots$  )

what is used

(  $k$  is smaller  $n$  )

Mini-Batch idea:

$$x^{(k+1)} = x^{(k)} - \frac{\alpha^{(k)}}{|I_k|} \sum_{j \in I_k} \nabla f_j(x^{(k)})$$

$\underbrace{\text{mini batch}}_{\text{s.g.}}$  → decreases variance.  
→ how to select them?

Ex.: Suppose  $|I_k|=2$ , show that variance is smaller

$\nwarrow$  compute variance.

Assume that we pick a mini batch uniformly at random.

\* Pro

⑥ SGD w/ momentum:

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} g(x^{(k)}) + \beta (x^{(k)} - x^{(k-1)})$$

$\uparrow$  s.g.       $\uparrow$  for stabilization.

⑦ Why does SGD work?

Let us start w/ regular gradient descent.

(a) Assume  $|u^\top Hf(x) \cdot u| \leq L \|u\|^2 \forall x$

By Taylor,  $\exists \xi_k$

$$\begin{aligned}
 f(x^{(k+1)}) &= f(x^{(k)} - \alpha \nabla f(x^{(k)})) = \\
 &= f(x^{(k)}) - \alpha \nabla f(x^{(k)})^T \cdot \nabla f(x^{(k)}) + \frac{1}{2} (\alpha \nabla f(x^{(k)}))^T \nabla^2 f(\xi_k) \\
 &\leq f(x^{(k)}) - \alpha \|\nabla f(x^{(k)})\|^2 + \frac{\alpha^2 L}{2} \|\nabla f(x^{(k)})\|^2 \\
 &= f(x^{(k)}) - \alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(x^{(k)})\|^2
 \end{aligned}$$

$\left. \begin{array}{l} \text{Descent!} \\ \text{Lemma!} \end{array} \right\}$

If  $\alpha$  small enough (e.g.,  $\alpha \leq \frac{1}{L}$ ), then

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \frac{1}{2} \alpha \|\nabla f(x^{(k)})\|^2 \quad ! \quad (\#)$$

decreases.

Note further,

$$\begin{aligned}
 \frac{1}{2} \alpha \sum_{k=0}^{T-1} \|\nabla f(x^{(k)})\|^2 &\leq \sum_{k=0}^{T-1} f(x^{(k)}) - f(x^{(k+1)}) = \\
 &\stackrel{T \text{ iterations}}{=} f(x^{(0)}) - f(x^{(T)}) \leq f(x^{(0)}) - \underset{\substack{\uparrow \\ f(x^*)}}{f^*} \\
 \Rightarrow \min \|\nabla f(x^{(k)})\|^2 &\leq \frac{1}{T} \sum_{k=0}^{T-1} \|\nabla f(x^{(k)})\|^2 \leq \underset{\substack{\uparrow \\ \text{solution}}}{\frac{2(f(x^{(0)}) - f^*)}{\alpha T}}
 \end{aligned}$$

$\Rightarrow$  we converge.

(b) Let us now look at SGD.

We assume

$$\|\nabla_i f_i(x)\| \leq G \quad (1^{\text{st}} \text{ der. bound})$$

$$\|u^\top \nabla^2 f(x) \cdot u\| \leq L \|u\|^2 \quad (2^{\text{nd}} \text{ derivative bound})$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \alpha_k^k (\nabla f_{\tilde{i}_k}^k(x^{(k)}))^T \cdot \nabla f(x^{(k)}) +$$

$$+ \frac{(\alpha^{(k)})^2 L}{2} \cdot \underbrace{\|\nabla f_{\tilde{i}_k}^k(x^{(k)})\|^2}_{\text{can be of any sign!}}$$

→ Take Expectation: Use bounded variance here!

$$\mathbb{E}[f(x^{(k+1)})] \leq \mathbb{E}[f(x^{(k)})] - \alpha^{(k)} \mathbb{E}[\nabla f_{\tilde{i}_k}^k(x^{(k)})^T \cdot \nabla f(x^{(k)})]$$

$$+ \frac{(\alpha^{(k)})^2 G^2 L}{2} \quad \text{tricky!}$$

$$\Rightarrow \mathbb{E}[\nabla f_{\tilde{i}_k}^k(x^{(k)})] = \nabla f(x^{(k)}), \text{ and hence}$$

$$\mathbb{E}[f(x^{(k+1)})] \leq \mathbb{E}[f(x^{(k)})] - \alpha^{(k)} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] + \frac{(\alpha^{(k)})^2 G^2 L}{2}$$

(Compare to #)

$$\Rightarrow \sum_{k=0}^T \alpha^{(k)} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq f(x^{(0)}) - \mathbb{E}[f(x^{(T)})] + \frac{G^2 L}{2} \sum_{k=0}^{T-1} (\alpha^{(k)})^2$$

$$\leq f(x^{(0)}) - f^* + \frac{G^2 L}{2} \sum_{k=0}^{T-1} (\alpha^{(k)})^2.$$

$$(f(x^{(T)}) \geq f^* \Rightarrow \mathbb{E}[f(x^{(T)})] \geq f^*)$$

How to make sense of this?

Run SGD for a random number of iterations  $T$

with  $P(T=t) = \frac{\alpha_t^{(t)}}{\sum_{k=0}^{T-1} \alpha^{(k)}}$  [New discrete random variable]

$$\Rightarrow \mathbb{E}[\|\nabla f(x^{(T)})\|^2] = \sum_{k=0}^{T-1} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \cdot P(T=k) =$$

$$= \frac{1}{\sum_{k=0}^{T-1} \alpha^{(k)}} \sum_{k=0}^{T-1} \alpha^{(k)} \cdot \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq \left( \sum_{k=0}^{T-1} \alpha^{(k)} \right)^{-1} \cdot (f(x^{(0)}) - f^* +$$

$$+ \frac{G^2 L}{2} \sum_{k=0}^{T-1} (\alpha^{(k)})^2)$$

→ these are better estimates for convex objectives.

○ If the learning rate is constant, then

$$\alpha^{(k)} = \alpha \Rightarrow$$

$$\Rightarrow E[\|\nabla f(x^{(T)})\|^2] \leq \frac{f(x^*) - f^*}{\alpha T} + \underbrace{\frac{\alpha L G^2}{2}}_{\downarrow 0 \text{ as } T \rightarrow \infty}$$

↑  
doesn't go to zero!

This is what we observed. { 'noise ball' term,  
'uncertainty range'

○ If  $\sum_{k=0}^{T-1} \alpha^{(k)}$  grows faster than  $\sum_{k=0}^{T-1} (\alpha^{(k)})^2$ , then

For example  $E[\|\nabla f(x^{(T)})\|^2] \xrightarrow[T \rightarrow \infty]{} 0$

Example:  $\alpha^{(k)} = \frac{1}{\sqrt{k+1}}$

$$\sum_{k=0}^{T-1} \alpha^{(k)} = \sum_{k=0}^{T-1} \frac{1}{\sqrt{k+1}} \sim \int_0^T \frac{1}{\sqrt{x}} dx = 2\sqrt{T}$$

$$\sum (\alpha^{(k)})^2 = \sum_{k=0}^{T-1} \frac{1}{k+1} \sim \int_1^T \frac{1}{x} dx = \log T$$

$$\Rightarrow E[\|\nabla f(x^{(T)})\|^2] \leq O\left(\frac{1}{\sqrt{T}}\right)$$

↑ up to approx.  
of integrals

→ There are better estimates for convex objectives.

Why

$$\mathbb{E} [\nabla f_{i(k)}(x^{(k)})] = \nabla f(x^{(k)}) ?$$

$$\mathbb{E}[X | Y=y] = \sum x \cdot \frac{P(X=x, Y=y)}{P(Y=y)}.$$

random var.

If  $X$  and  $Y$  are independent, then  $P(X=x, Y=y) = P(X=x) \cdot P(Y=y) \Rightarrow \mathbb{E}[X | Y=y] = \mathbb{E}[X]$

Let  $s_1, \dots, s_{k-1}$  is the sequence of indices s.t.

$$x^{(k)} = x^{(k-1)} + \alpha^{(k-1)} \nabla f_{s_{k-1}}(x^{(k-1)})$$

$$x^{(k-1)} = x^{(k-2)} + \alpha^{(k-2)} \nabla f_{s_{k-2}}(x^{(k-2)})$$

$$x^{(0)} = x^{(\infty)} + \alpha^{(\infty)} \nabla f_{s_0}(x^{(\infty)})$$

$$\Rightarrow \mathbb{E}[\nabla f_{i(k)}(x^{(k)})] = \sum_{s=1}^n \underbrace{\nabla f_s(x^{(k)})}_{\text{value}} \cdot \frac{P(i(k)=s, i(k-1) \neq i(k-2), \dots, i(0) = s_{k-1}, s_{k-2}, \dots, s_0))}{P((i(k-1), \dots, i(0)) = (s_{k-1}, \dots, s_0))}$$

$$= \sum_{s=1}^n \nabla f_s(x^{(k)}) \cdot P(i(k)=s) =$$

$$= \frac{1}{n} \sum_{s=1}^n \nabla f_s(x^{(k)}) = \nabla f(x^{(k)}).$$

Let us write  $E_i[\nabla f_{i(\alpha)}(x^{(\alpha)})]$  for expectation

Assume  $E[\|\nabla_{\nabla f_{i(\alpha)}} f(x) - \nabla f(x)\|^2] \leq \sigma^2$

$$\Leftrightarrow \frac{1}{n} \sum_{i=1}^n \|\nabla_{\nabla f_i(x)} f(x) - \nabla f(x)\|^2 \leq \sigma^2$$

↑

Variance is bounded.  $\Rightarrow E[\|\nabla f_{i(\alpha)}(x)\|^2] \leq$

Similarly, we get:  $\leq \|\nabla f(x)\|^2 + \sigma^2$

$$f^*(x^{(\alpha+1)}) \leq f^*(x^{(\alpha)}) - \alpha^{(\alpha)} \langle \nabla_{\nabla f_{i(\alpha)}} f(x^{(\alpha)}), \nabla f(x^{(\alpha)}) \rangle + \\ + \frac{(\alpha^{(\alpha)})^2 L}{2} \|\nabla f_{i(\alpha)}(x^{(\alpha)})\|^2.$$

$$\Rightarrow E[f^*(x^{(\alpha+1)}) | x^{(\alpha)}] \leq E[f^*(x^{(\alpha)}) | x^{(\alpha)}] \dots$$

$$= f^*(x^{(\alpha)}) - \alpha^{(\alpha)} \left\langle \overbrace{E[\nabla f_{i(\alpha)}(x^{(\alpha)}) | x^{(\alpha)}]}^{\nabla f(x^{(\alpha)}) \text{ as above}}, \nabla f(x^{(\alpha)}) \right\rangle + \\ + \frac{(\alpha^{(\alpha)})^2 L}{2} E[\|\nabla f_{i(\alpha)}(x^{(\alpha)})\|^2 | x^{(\alpha)}] \\ = f^*(x^{(\alpha)}) - \alpha^{(\alpha)} \|\nabla f(x^{(\alpha)})\|^2 + \frac{(\alpha^{(\alpha)})^2 L}{2} \dots$$

$$\Rightarrow E[f^*(x^{(\alpha+1)}) | x^{(\alpha)}] \leq f^*(x^{(\alpha)}) - \frac{\alpha^{(\alpha)}}{2} \left(1 - \frac{\alpha^{(\alpha)} L}{2}\right) \|\nabla f(x^{(\alpha)})\|^2 + \\ + \frac{(\alpha^{(\alpha)})^2 \sigma^2 L}{2}$$

$(\alpha^{(\alpha)} \leq \frac{1}{L})$

$$\Rightarrow E[f^*(x^{(\alpha+1)})] \leq E[f^*(x^{(\alpha)})] - \frac{\alpha^{(\alpha)}}{2} E[\|\nabla f(x^{(\alpha)})\|^2] + \frac{(\alpha^{(\alpha)})^2 \sigma^2 L}{2}.$$

(This is more parallel to GD method.)