

OPTIONAL LECTURE 2, 26/1/2021

minimise $R_S(h)$
over H

ERROR
Empirical
Risk
Minimization

In most cases h is parameterized
by a vector of parameters $w \in \mathbb{R}^d$

$$\text{ex. } h(x) = \sum_{i=1}^d w_i x_i$$

$$h(w, x) = \sum_{i=1}^d w_i x_i$$

$$R_S(h) = \sum_{i=1}^{|S|} \ell(h, (x_i, y_i))$$

$$\text{ex. } = \sum_{i=1}^{|S|} (h(w, x_i) - y_i)^2 = \sum_{i=1}^{|S|} f(w, i)$$

$$F(w) = R_S(h)$$

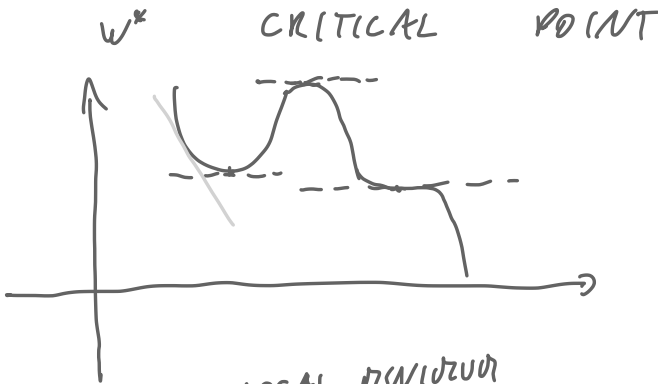
minimise $F(w)$
over $w \in \mathbb{R}^d$

minimise $F(w)$
over $w \in W \subset \mathbb{R}^d$

minimize $F(w)$
 $w \in \mathbb{R}^d$

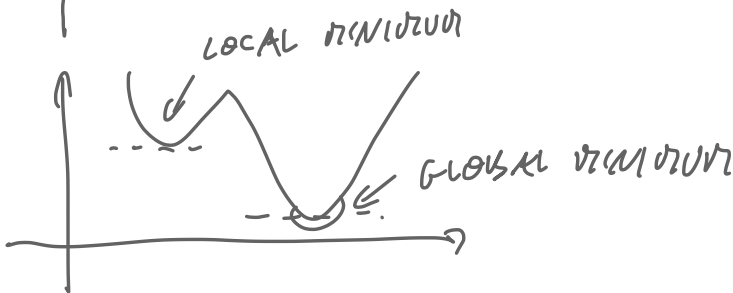
$d = 1$

$$F'(w^*) = 0$$



$$F''(w^*) > 0$$

$\Rightarrow w^*$ is a
 minimum



$d > 1$

$$\nabla F(w) = \begin{bmatrix} \frac{\partial F}{\partial w_1} \\ \frac{\partial F}{\partial w_2} \\ \vdots \\ \frac{\partial F}{\partial w_d} \end{bmatrix}$$

$$\frac{\partial F}{\partial w_1}$$

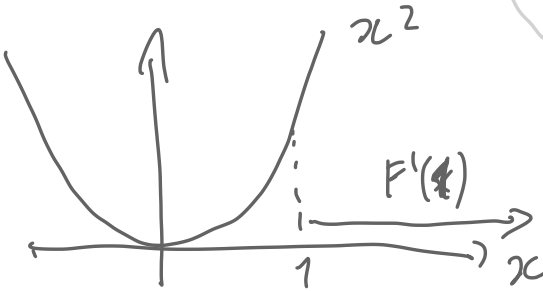
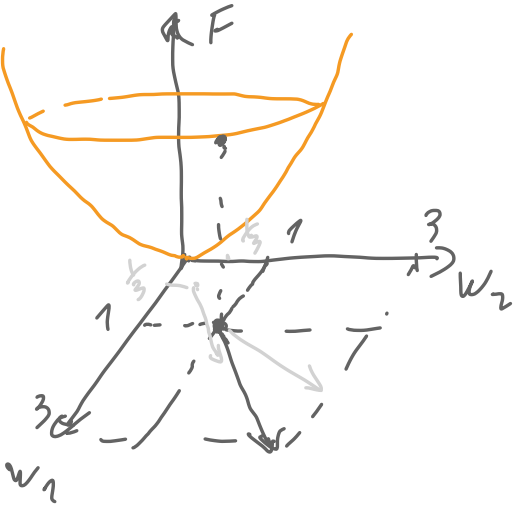
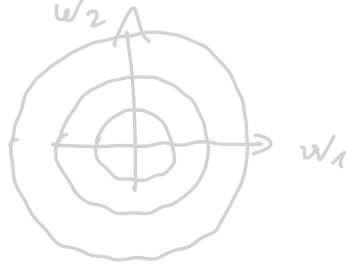
$$F(w) = w_1^2 + w_2^2$$

$$\frac{\partial F}{\partial w_1} = 2w_1$$

$$\frac{\partial F}{\partial w_2} = 2w_2$$

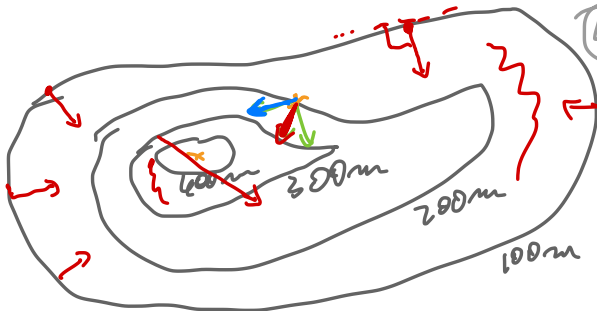
$$F(w) = w_1^2 + w_2^2$$

$$\nabla F(w) = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$



∇F → points in the direction where the function is growing the most

$\|\nabla F\|$ tells us how fast is growing in that direction



ISO LEVEL CURVES

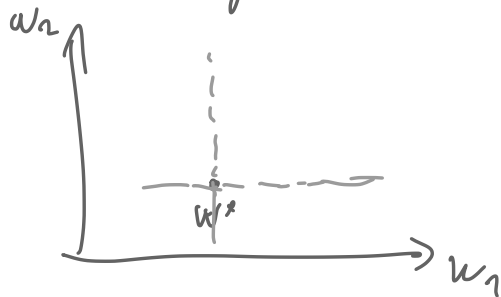
THIS FUNCT. SHOULD BE INCREASING
POSSIBLE?

∇F gives the direction of maximum increase of F

↓

it is orthogonal to the level curves of the function

$$w^* \in \arg\min F(w) \Rightarrow \nabla F(w^*) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$



$$\frac{\partial F}{\partial w_1} = 0$$

$$\frac{\partial F}{\partial w_2} = 0$$

HESSIAN OF F

$$F(w_1, w_2)$$

$$\frac{\partial^2}{\partial w_2 \partial w_1} F(w_1, w_2) =$$

$$= \frac{\partial}{\partial w_2} \left(\frac{\partial}{\partial w_1} F(w_1, w_2) \right)$$

$$\begin{aligned} \text{EX} &= \frac{\partial}{\partial w_2} \left(\frac{\partial}{\partial w_1} w_1^2 + w_2^2 \right) = \\ &= \frac{\partial}{\partial w_2} (2w_1) = 0 \end{aligned}$$

$H_F(w)$ is a matrix

$$[H_F(w)]_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} F(w)$$

ex $H_F(w) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

$$\frac{\partial}{\partial w_1} \frac{\partial}{\partial w_1} F(w) = 2$$

$$F''(w) \iff H_F(w)$$

$$F''(w) \geq 0 \Rightarrow \begin{array}{l} \text{?} \\ \text{?} \end{array} \begin{array}{l} \cancel{\det(H_F(w)) > 0} \\ H_F(w) \text{ is SEMIDEFINITE} \\ \text{POSITIVE} \end{array}$$

$$\forall x \in \mathbb{R}^d, \quad x^T H_F(w) x \geq 0$$

$$F''(w) > 0 \Rightarrow H_F(w) \text{ is DEFINITE POSITIVE}$$

$$\forall x \in \mathbb{R}^d - \{0\}, \quad x^T H_F(w) x > 0$$

$H_F(w)$ IS POSITIVE SEMIDEFINITE

$$\forall x \in \mathbb{R}^d \quad x^T H_F(w) x \geq 0$$

$$d=1 \quad H_F(w) = \left[\frac{\partial^2 H_F}{\partial w_1^2} \right] = H_F''(w)$$

$$x H_F''(w) x = x^2 \times H_F''(w) \geq 0$$

$$x^2 \geq 0 \Rightarrow H_F''(w) \geq 0$$

$$H_F(w)$$

COMPUTE SINGULAR VALUES
←————→

A IS DIAGONALIZABLE

$$A = U \Lambda U^T$$

$$U^T U = I$$

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & & \\ 0 & 0 & \ddots & \\ & & & \lambda_d \end{pmatrix}$$

U IS ORTHOGONAL MATRIX

$\{\lambda_i\}$ EIGENVALUES

$$A = U \Lambda U^T$$

$$AU = U \Lambda U^T U =$$

$$AU = U \Lambda$$

$$A u_i = \lambda_i \cdot u_i$$

U IS THE MATRIX
ARE EIGEN
VECTORS

$A = A^T \Rightarrow A$ IS DIAGONALIZABLE

$H_F = H_F^T$ H_F IS

$$H_F(\omega) = U \Lambda U^T$$

H_F IS POSITIVE SEMIDEFINITE

IFF $\lambda_i \geq 0 \quad \forall i$

H_F IS POSITIVE DEFINITE

IFF $\lambda_i > 0 \quad \forall i$

Imagine $\lambda_j < 0$ I pick $x = u_j$

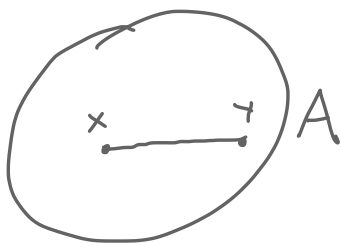
$$u_j^T H_F u_j = u_j^T \lambda_j u_j = \underbrace{\|u_j\|^2}_{>0} \lambda_j < 0$$

$\Rightarrow H_F$ IS NOT POSITIVE
SEMIDEFINITE

$$\boxed{H_F(\omega)} \longrightarrow U^2 F(\omega)$$

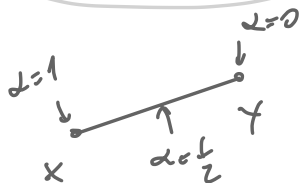
CONVEXITY

Convex set

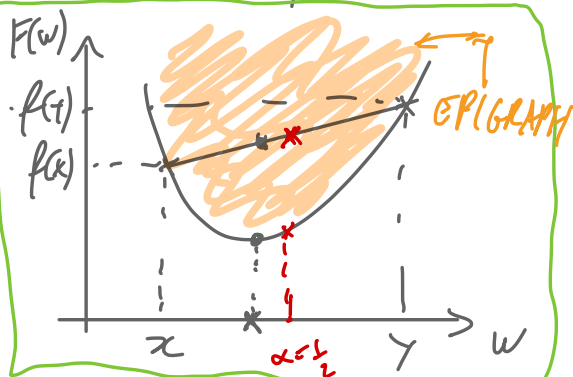


A is convex if
 $\forall x, y \in A, \forall \alpha \in [0, 1]$

$$\alpha x + (1-\alpha)y \in A$$



Convex function



F is convex if

$$\forall x, y \quad \forall \alpha \in [0, 1]$$

$$f(\alpha x + (1-\alpha)y) \leq$$

$$\alpha f(x) + (1-\alpha)f(y)$$

WHY DO WE CARE ABOUT CONVEXITY

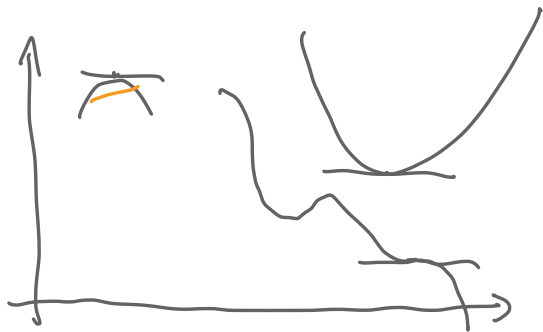
CONVEX OPT. PROBLEM

$$\begin{array}{l} \text{minimize } F(w) \\ \text{w.r.t } W \end{array}$$

W is a convex set

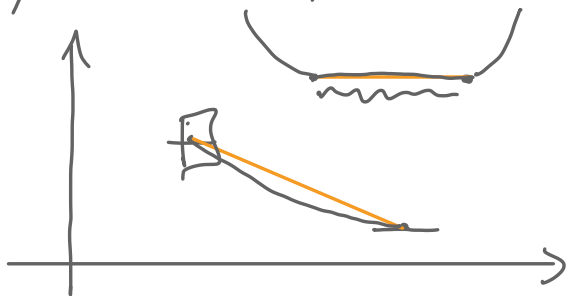
F is a convex function

WHY CONVEX OPT. PROBLEMS ARE NICE?



IN A CONVEX OPT PROBLEM

1) IF $\nabla F(w^*) = 0 \Rightarrow w^*$ IS A LOCAL MINIMUM



2) IF w^* IS A LOCAL MINIMUM IT IS ALSO A GLOBAL MINIMUM



THERE ARE EFFICIENT ALGORITHMS TO SOLVE CONVEX OPTIMIZATION PROBLEMS

CLASSIFICATION PROBLEM

$$\ell(h, (x_i, y_i)) = \mathbb{1}_{h(x_i) \neq y_i} \quad \text{NATURAL CHOICE}$$

CORRON CHOICE

CROSS-ENTROPY LOSS

$$\ell(h, (x, y)) = -y \log h(x) - (1-y) \log(1-h(x))$$

$$h(x) = w^T x$$

THIS LOSS IS DIFFERENTIABLE

THIS LOSS IS CONVEX IN w

SURROGATE LOSS

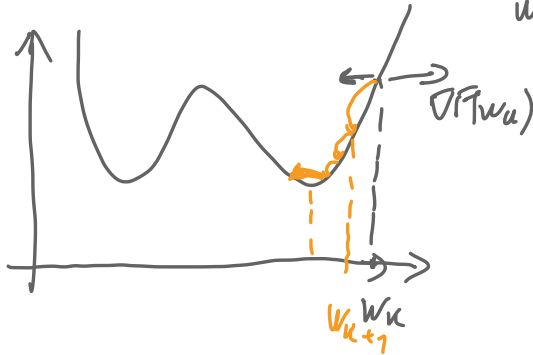
NATURAL CHOICE FOR $\ell_{\text{NAT}}(h, (x, y))$

YOU PICK $\ell_{\text{SUR}}(h, (x, y))$

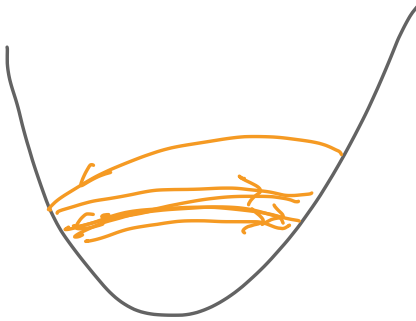
$$1) \ell_{\text{NAT}}(h, (x, y)) \leq \ell_{\text{SUR}}(h, (x, y))$$

$$2) \ell_{\text{SUR}}(h, (x, y)) \text{ IS CONVEX}$$

minimize $F(w)$
 $w \in \mathbb{R}^d$



$$w_{k+1} = w_k - \underbrace{\alpha_k}_{\text{LEARNING RATE}} \nabla F(w_k)$$



α const

$$\nabla F(w_k) \approx 0$$

$$F(w) = \sum_{i=1}^{|S|} f(w, i)$$

$$w_{k+1} = w_k - \alpha \nabla F(w_k)$$

FULL-BATCH
 GRADIENT METHOD

$$\nabla F(w_k) = \nabla \left(\sum_i f(w, i) \right) = \sum_i \nabla f(w, i)$$

BATCH = DATASET

AT EVERY ITERATION YOU TAKE ONLY
ONE POINT IN THE DATASET
UNIFORMLY AT RANDOM ξ

$$w_{k+1} = w_k - \eta \nabla f(w, \xi)$$

STOCHASTIC GRADIENT METHOD
DESCENT

YOU CAN PICK A RANDOM SUBSET ξ_k
OF SAMPLES FROM THE DATASET

$$\frac{1}{n_k} \sum_{i=1}^{n_k} \nabla f(w, \xi_{k,i})$$

$$\xi_k = \{ \xi_{k,1}, \xi_{k,2}, \dots, \xi_{k,n_k} \}$$

MINI-BATCH GRADIENT METHOD

$$n_k = 1 \quad \equiv \quad \text{SGD}$$

$$n_k = |S| \quad \equiv \quad \text{FULL-BATCH}$$

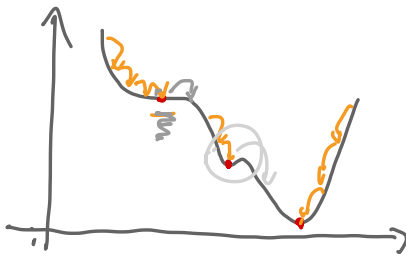
minimize $F(w) = \sum_{i=1}^{|S|} f(w, i)$

WHY NOT ALWAYS FULL-BATCH?

- I MAY NOT HAVE ENOUGH MEMORY

$\mathcal{S} \subset S$ minimize $\sum_{(x_i, y_i) \in \mathcal{S}} \ell(h, (x_i, y_i))$

- NOISE CAN BE USEFUL ~~OVERFITTING~~

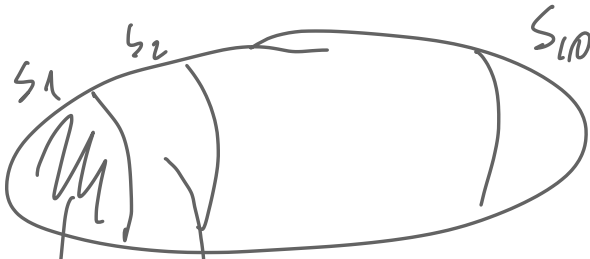


FB could converge to any of them

THE NOISE CAN HELP OF
AVOIDING TO GET STUCK IN
BAD CRITICAL POINTS.

IT DOESN'T HAPPEN FOR CONVEX
PROBLEM.

• GO BACK MEMORY



$$\boxed{\text{RAM}} \rightarrow \frac{1}{|S_1|} \sum_{i \in S_1} \nabla f(w, i) = f_1$$

$$\boxed{\text{RAM}} \rightarrow \frac{1}{|S_2|} \sum_{i \in S_2} \nabla f(w, i) = f_2$$

$$\frac{f_1 + f_2 + \dots + f_{10}}{10} = \frac{1}{|S|} \sum_{i \in S} \nabla f(w, i) \\ \equiv \text{FB}$$

FB

$$W_2 = W_1 - \alpha \frac{1}{|S|} \sum_{i \in S} \nabla f(w_i)$$

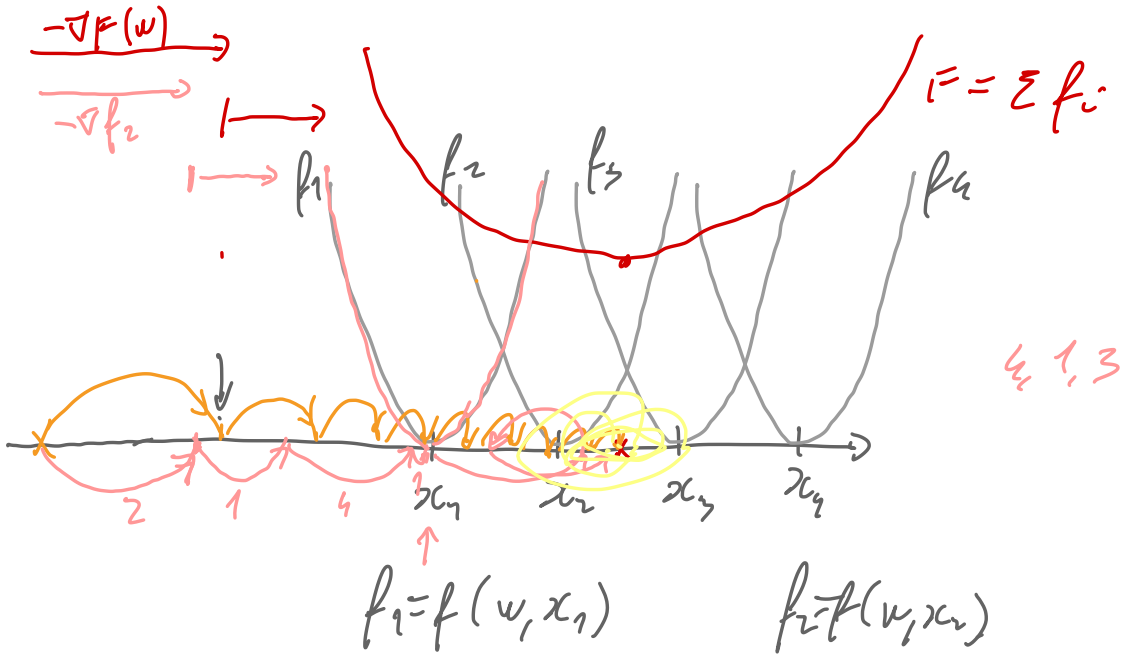
MB

$$\begin{aligned} W_2 &= W_1 - \alpha \frac{1}{|S_1|} \sum_{i \in S_1} \nabla f(w_i) \\ W_3 &= W_2 - \alpha \frac{1}{|S_2|} \sum_{i \in S_2} \nabla f(w_i) \\ &\dots \\ W_{10} &= W_0 - \alpha \frac{1}{|S_{10}|} \sum_{i \in S_{10}} \nabla f(w_i) \end{aligned}$$

10

IS IT BETTER TO DO 1 PRECISE
UPDATE (FB)

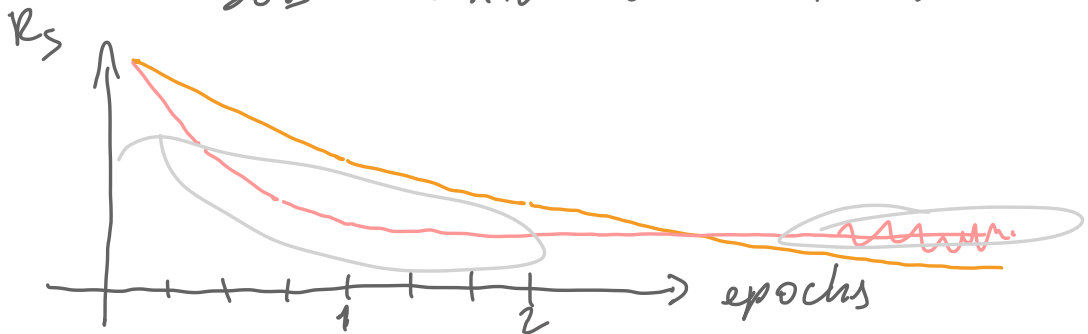
OR 10 NOT-SO PRECISE UPDATES (MB)?



← REGION OF CONFUSION →

1) WHEN WE ARE FAR FROM THE MINIMUM SGD GIVES US THE RIGHT DIRECTION

2) WHEN WE ARE CLOSE TO THE MINIMUM SGD CAN GET "CONFUSED"

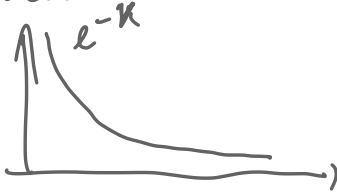


IDEA : START WITH SMALL η
AND PROGRESSIVELY INCREASE η
VARIANCE REDUCTION

Theoretical results

$$\text{FB: } R_S(w_k) - R_S(w^*) \sim \rho^k \quad \rho < 1$$

LINEAR CONVERGENCE



$$\rho^k < \varepsilon \quad \Rightarrow \quad k_\varepsilon \in O\left(\ln \frac{1}{\varepsilon}\right)$$

NUMBER OF ITERATIONS

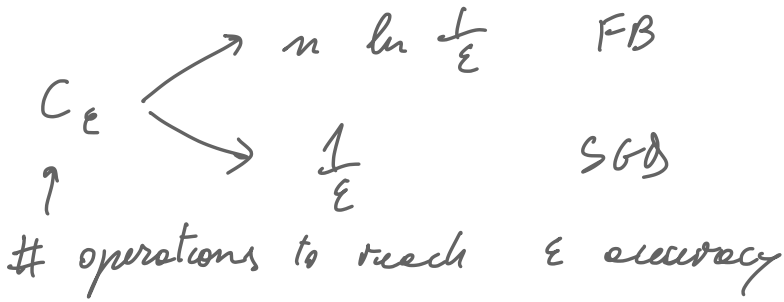
SGD

$$R_S(w_k) - R_S(w^*) \sim \frac{1}{k} \quad \text{SUBLINEAR}$$

$$k_\varepsilon \in O\left(\frac{1}{\varepsilon}\right)$$

FB does $n = |S|$ calculations for iteration

SGD does 1 calculation for iteration



If the dataset is huge
it can be $n \ln \frac{1}{\epsilon} \gg \frac{1}{\epsilon}$

T for one iteration

- Above $T \propto$ # operations
- $T = T_0 + \beta$ (# operations)
 \uparrow
 $\beta = 0$