

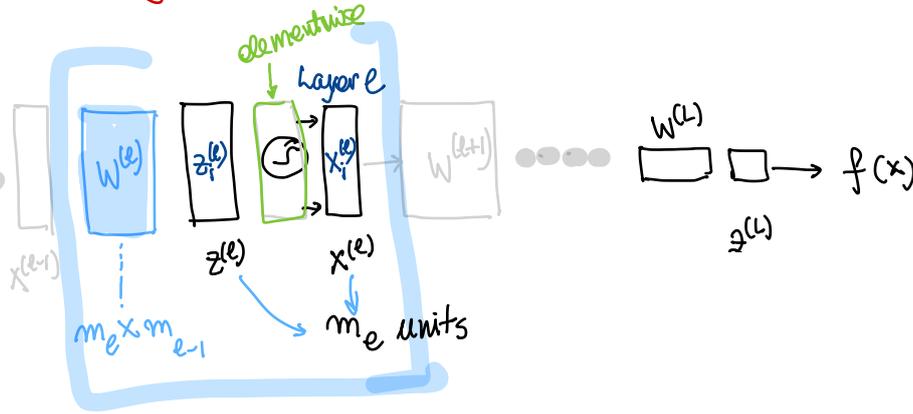
# Lecture 12

HW5 W7 due W8  
[HW6 W8]  
Q2 3/5 or W9  
HW7 W10 due W11  
Q3 W11  
[HW8 W11]

Where we left off before reading week ...

Backpropagation

Training a multi-layer network



Parameters  $W = \{ W^{(1)}, \dots, W^{(L)} \}$   $\therefore W \in \mathbb{R}^P$

Data  $\mathcal{D} = \{ (x^i, y^i), i=1:n \}$

TRAIN: by G. Descent on loss  $\mathcal{L}$

1. Init  $W$  with small random values

$$W_{ij}^{(e)} \sim N(0, \sigma^2 \frac{1}{m_e})$$

$\sigma^2$  small

2. for  $t=1, 2, \dots$

$$\left[ W \leftarrow W - \eta \frac{\partial \mathcal{L}(W^t)}{\partial W} \right] \Leftrightarrow \left[ W_j^{t+1} \leftarrow W_j^t - \eta \frac{\partial \mathcal{L}}{\partial W_j} \right]$$

step size  $\eta$   $g = \text{gradient}$

$x_j^i$  ← training ex.  $i$

$x_j^{(e)}$  ← layer  $e$

$x_j^x$  units

$x_2^{(0)i}$  = input value layer 0, 2nd attribute  $i$ th example

I comp/iter ✓

NP

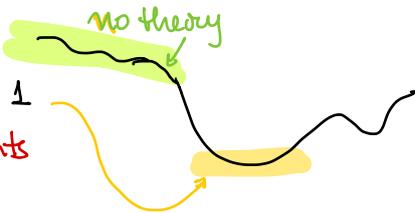
T #iterations ←

GD is order 1

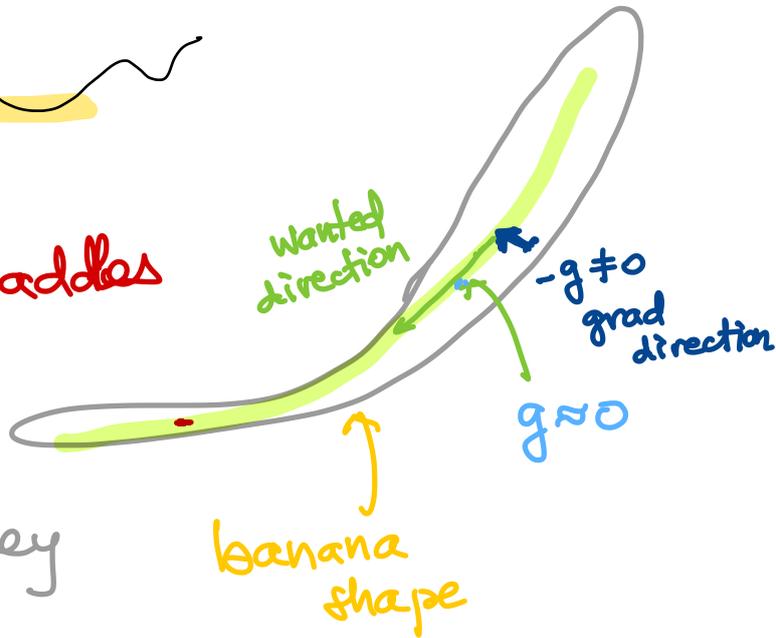
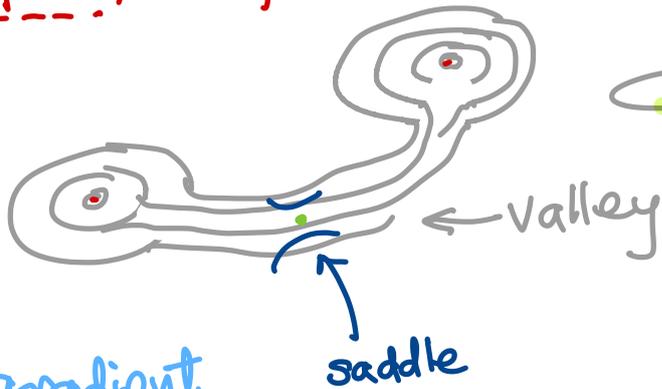
S saturation / vanishing gradients

L local minima

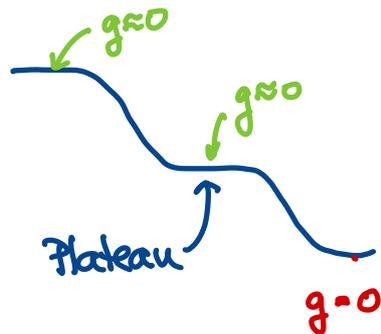
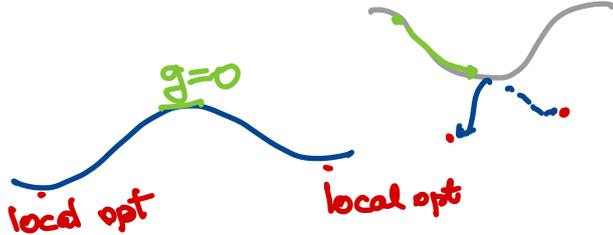
O overfitting



# Local Optima, Valleys, Plateaus, Saddles



$g^t = \text{gradient at } t$



# Stochastic Grad Descent

$$w_j^{t+1} \leftarrow w_j^t - \eta \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{\partial \ell(y_i, f(x_i; w^t))}{\partial w_j}}_{g_i \text{ contribution of } (x_i, y_i)}$$

for any  $t$ :  $g = \frac{1}{n} \sum_i g_i$

Idea sample  $i \sim \text{unif}[1:n]$   $\times n$  times  $\Rightarrow \mathcal{B} = \{(x^i, y^i), i=1:n\} \subset \mathcal{D}$   
"batch"

$$g_{\mathcal{B}} = \frac{1}{n'} \sum_{i \in \mathcal{B}} g_i \approx g$$

$$\sigma^2 = \text{Var}\{g_i\}$$

$$g = \text{mean}\{g_i\}$$

random variable:  $g_{\mathcal{B}} = g + \underbrace{\varepsilon}_{\text{"noise"}}$

$$\begin{aligned} \text{mean}(\varepsilon) &= 0 \\ \text{Var}(\varepsilon) &= \frac{\sigma^2}{n'} \\ \text{std}(\varepsilon) &= \frac{\sigma}{\sqrt{n'}} \end{aligned}$$

in theory  $n'=1$

in practice  $\rightarrow n' > 1$

## Benefits

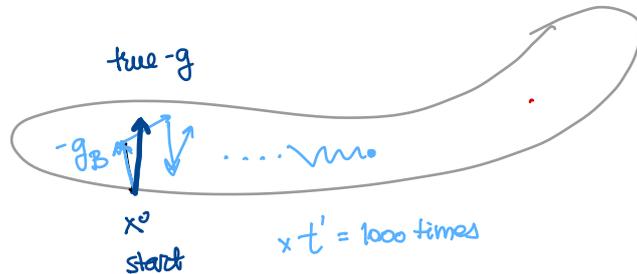
- 1) Faster computation / iteration **(I)**  $n'p$  vs  $np$
- 2)  $\text{---||---}$  overall  $\underline{T_{n'p}}$  vs  $\underline{T_{np}}$
- 3) Randomness: avoid escape saddles, shallow local minima,

# Why SGD "faster" than GD?

$$n = 10^6$$

$$n' = 1000$$

$$t' = \frac{n}{n'} = 1000 = \#\mathcal{B}'_s \text{ in } \mathcal{D}$$



## Theoretical convergence

near  $W^*$  optimum

- increase  $n' \rightarrow n$
- average  $g_{\mathcal{B}^t}^t$  over last  $\tilde{T}$  iterations
- $\eta \rightarrow 0$  (like  $\sim \frac{1}{t}$ )

Momentum — can we mimic 2nd order methods cheaply?

GD: 1st order approx

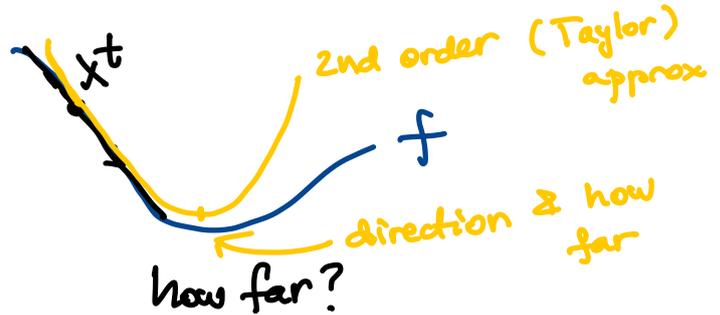
$$g \in \mathbb{R}^p$$

Newton: 2nd order

$$\left[ \frac{\partial^2 L}{\partial w_j \partial w_{j'}} \right]_{j, j' = 1:p} \in \mathbb{R}^{p \times p}$$

Hessian

TOO EXPENSIVE



Heavy Ball/Momentum

$$\gamma \in (0, 1)$$

$$w_j^{t+1} \leftarrow w_j^t - \gamma \eta \underbrace{\frac{\partial L^t}{\partial w_j}}_{g_j^t} + (1-\gamma) \underbrace{(w_j^t - w_j^{t-1})}_{\text{previous step}}$$

$$w_j^{t+1} \leftarrow w_j^t - \eta g_j^t + (1-\eta)(w_j^t - w_j^{t-1})$$

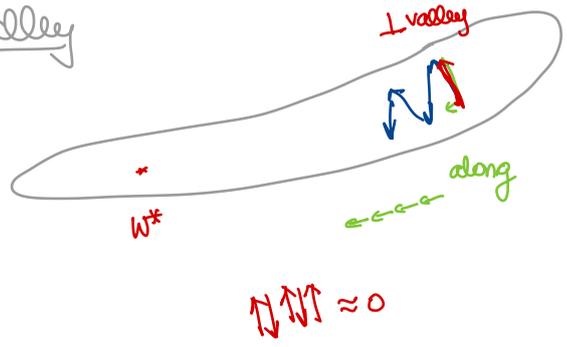
$$w_j^t - w_j^{t-1} = \eta g_j^{t-1} + (1-\eta)(w_j^{t-1} - w_j^{t-2})$$

$$w_j^{t-1} - w_j^{t-2} = \eta g_j^{t-2} + (1-\eta)(\dots)$$

Ex

$$w_j^{t+1} \leftarrow w_j^t - \eta \left\{ g_j^t + (1-\eta)g_j^{t-1} + (1-\eta)^2 g_j^{t-2} + \dots \right\}$$

Valley



Plateau

