

# Lecture October 7

In principle we compute sample expectation values

$$\begin{aligned}\text{Example MSE} &= C(\beta) \\ &= E[(y - x\beta)^2] \\ &= E[(y - \tilde{y}(x, \beta))^2]\end{aligned}$$

$$\nabla_{\beta} C(\beta) = E\left[\left(\frac{\partial C}{\partial \beta}\right)\right]$$

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2$$

$$\frac{\partial C}{\partial \beta_j} = -\frac{2}{n} \sum_{i=0}^{n-1} x_{ij} (y_i - x_{ij}\beta_j)$$

SGD (mini-batches)

Algo

initialize  $\mu_k$  (learning)

11.9.02

— 1 —  $\beta$

while stopping criterion not true

DO one epoch

- sample a mini-batch of examples from the data set
- $D = \{ (x_0, y_0), \dots, (x_{m-1}, y_{m-1}) \}$   
in total  $m$ -minibatches
- compute gradients  $g$
- update

$$\beta \leftarrow \beta - \eta_k g$$

end while

it is common to "decay" the learning rate until a saturation  $\tau$

$$\eta_k = (1 - \alpha) \eta_0 + \alpha \eta_{\tau}$$

$$\alpha = \frac{k}{\tau}$$

$\tau$  is chosen to

correspond to  
the # of iterations  
needed to make  
a few hundred  
passes through  
the training set  
 $\eta_T \sim 1\% \text{ of } \eta_0$

### Momentum SGD

$$C(\hat{\beta}) \approx C(\beta^{(n)}) - \underbrace{(\beta - \beta^{(n)})^T}_{b^T} g + \frac{1}{2} b^T H b \quad \underbrace{\quad}_{\mu g^T}$$

physics inspired:

Newton's law

position  $x(t)$  (1-dim)

Force  $f(t)$ ,  $m = 1$

$$f(t) = \frac{d^2 x}{dt^2} \Rightarrow$$

$$v(t) = \frac{dx}{dt} \quad \wedge \quad f(t) = \frac{dv}{dt}$$

at

$$v(t+\Delta t) \approx v(t) + \Delta t f(t)$$

$$x(t+\Delta t) \approx x(t) + (\Delta t)v(t)$$

Our force is proportional to  
the negative gradient  
of the cost function

$$x \rightarrow \beta \xrightarrow{\text{diag force}} \eta$$
$$v \leftarrow dv - (\Delta t) \nabla_{\beta} C(\beta)$$

$$\beta \leftarrow \beta + v$$

Momentum algo

Define  $\eta$  and momentum  
parameter  $\alpha$

initialize  $\beta$  and  $v$   
— — epochs and  
batches

while stopping criterion not  
met

DO

- sample  $m$  - mini-batch
- compute  $g \leftarrow \frac{1}{m} \nabla_{\beta} C(\beta)$
- compute  $v$   
 $v \leftarrow \alpha v - \eta g$
- update  
 $\beta \leftarrow \beta + v$

end DO

- Adagrad, RMS-prop  
Adam

Newton-Raphson

$$\beta \leftarrow \beta - \underbrace{H^{-1}}_{\text{Hessian}} \underbrace{g}_{\text{gradient}}$$

MSE example

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2$$

$$= \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i(x_i, \beta))^2$$

$$\frac{\partial C(\beta)}{\partial \beta_j} = - \frac{2}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i) \times \frac{\partial \hat{y}_i}{\partial \beta_j}$$

$$\begin{aligned} \frac{\partial^2 C}{\partial \beta_j \partial \beta_k} &= \frac{2}{n} \sum_{i=0}^{n-1} \frac{\partial \hat{y}_i}{\partial \beta_j} \frac{\partial \hat{y}_i}{\partial \beta_k} \\ &\quad - \frac{2}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i) \frac{\partial^2 \hat{y}_i}{\partial \beta_j \partial \beta_k} \end{aligned}$$

$$\approx \frac{2}{n} \sum_{i=0}^{n-1} \frac{\partial \hat{y}_i}{\partial \beta_j} \frac{\partial \hat{y}_i}{\partial \beta_k}$$

$$\left( \frac{\partial \hat{y}_i}{\partial \beta_j} \frac{\partial \hat{y}_i}{\partial \beta_j} \right) \sim \eta_j^{-1}$$

..

## Algo Adagrad

Define  $\eta$

initialize  $\beta$   
— — —  $\delta \sim 10^{-7}$

while stopping not met

DO

— sample minibatches

— compute  $g \leftarrow \frac{1}{m} \nabla_{\beta} C(\beta)$

— accumulate squared  
gradients

$$r \leftarrow r + g \odot g$$

— compute  $\Delta \beta$

$$\Delta \beta \leftarrow -\frac{\eta}{\delta + \sqrt{r}} \odot g$$

— update  $\beta \leftarrow \beta + \Delta \beta$

end DO

In RMS-prop:

$$r \leftarrow \rho r + (1 - \rho) g \odot g$$

$$\Delta \beta = - \frac{\eta}{\sqrt{S+2}} \odot g$$

go-to cues

SGD with momentum

RMS-prop ( in deep learning  
for non-convex  
 $C(\beta)$  )

Adagrad + ADAM

Automatic differentia-  
tion = Autograd

$$f' \approx \frac{f(x+\Delta x) - f(x-\Delta x)}{2\Delta x}$$

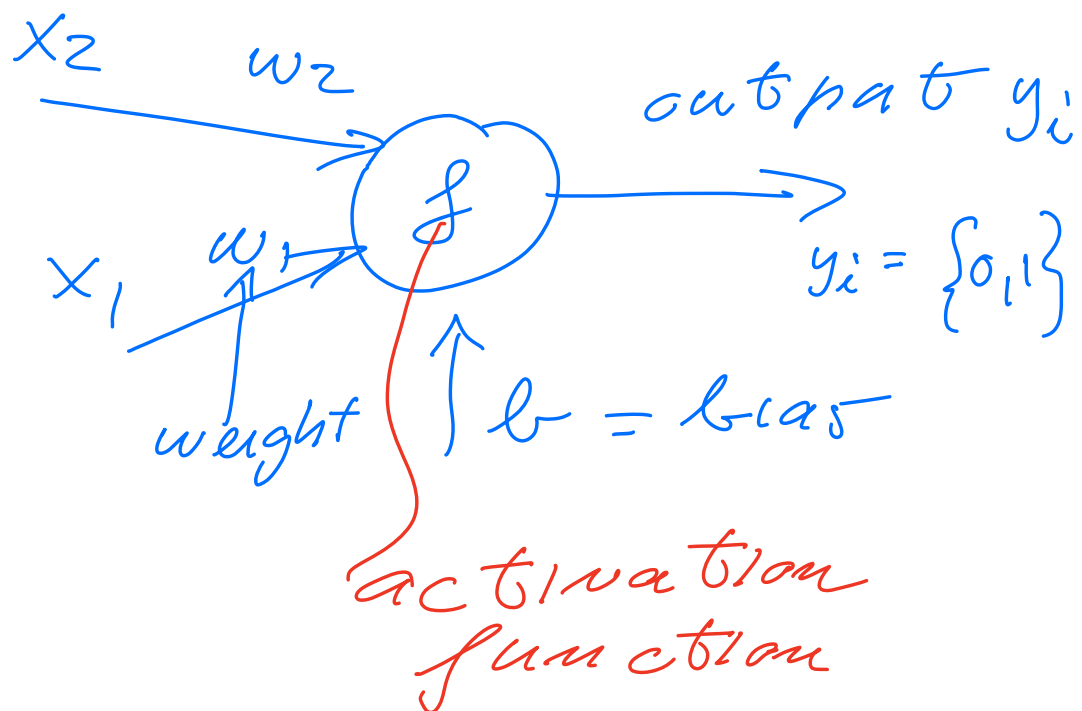
$$f'' \approx \frac{f(x+\Delta x) + f(x-\Delta x) - 2f(x)}{\Delta x^2}$$



# Deep-learning and Feed Forward Neural network (FFNN)

Basic definition:

- single perception model with one neuron / *node* / *unit*



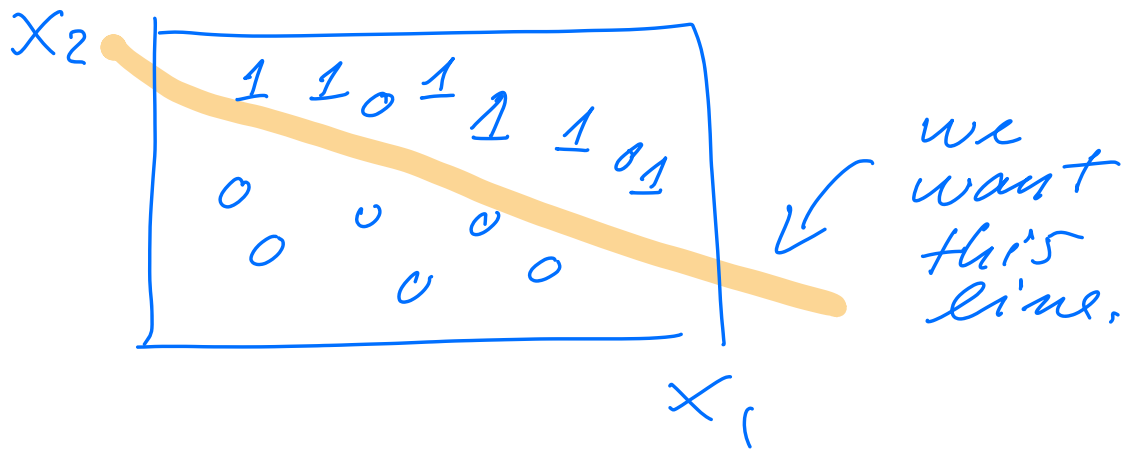
$$f(x_1 w_1 + x_2 w_2 + b) = y_i$$

Example data set:

$X_1$  = # hours slept

$X_2$  = # hours studied

$y_i = \begin{cases} 0 & \text{grade below average} \\ 1 & \text{grade above average} \end{cases}$



Universal approximation theorem: with one or more hidden layers we can have a network which can fit non-linear functions.

Single perceptron model can only fit a linear model.