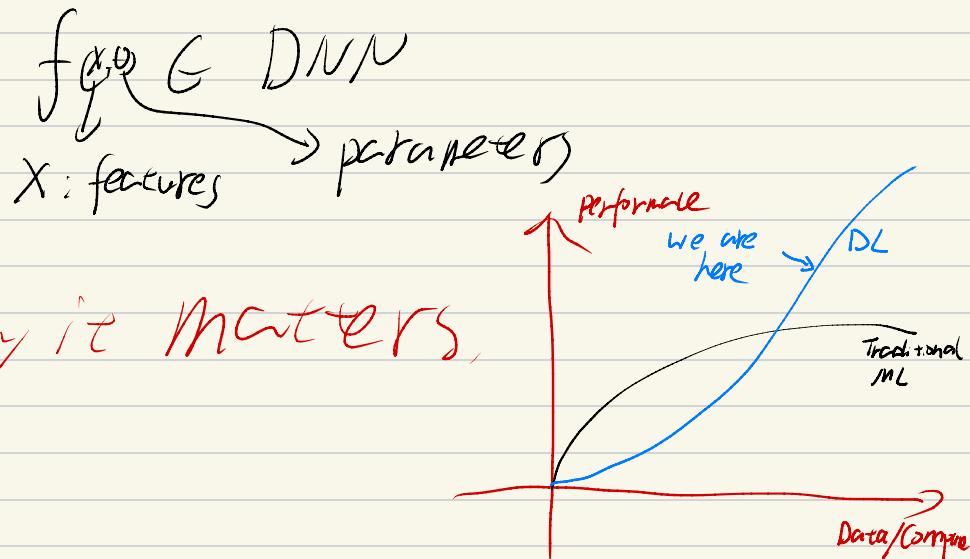
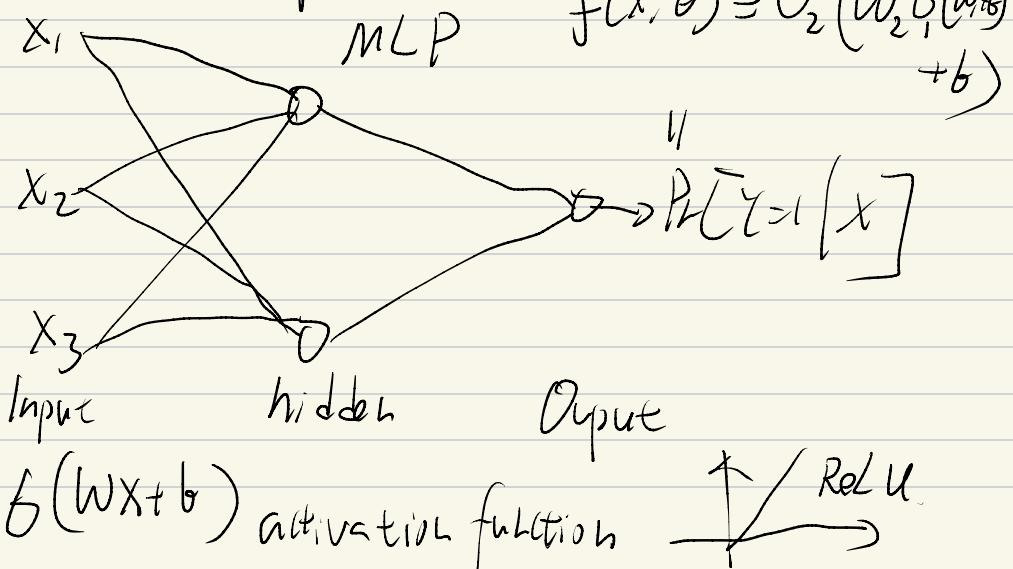


26.01.20

$$\text{DL: } \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i; \theta))$$



(1) What is $f(x, \theta)$



(2) What is the loss $\mathcal{L}(\cdot)$

Regression: $\mathcal{L}(y, \hat{y}) = \frac{1}{2}(\hat{y} - y)^2$

Classification: $\mathcal{L}(\hat{y}, y) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$

$$\hat{y} = f(x, \theta)$$

(3) How do we find $\hat{\theta}$ (argmin $_{\theta} \mathcal{L}(y, \hat{y})$)

Gradient Descent: $\hat{\theta}_{t+1} = \hat{\theta}_t - \nabla_{\theta} \mathcal{L}(\hat{\theta}_t)$

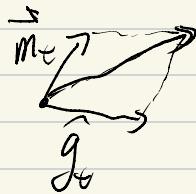
(4) How do we find $\nabla_{\theta} \mathcal{L}(\hat{\theta}_t)$

B.P. $y = f(z), z = g(x)$

$$\frac{dy}{dx} = \left. \frac{dy}{dz} \right|_{z=g(x)} \cdot \left. \frac{dz}{dx} \right|_x$$

(iii) SGD $\{l, \dots, B\} \subset \{1, \dots, n\}$
 random

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \frac{1}{B} \sum_{i=1}^B \nabla f_i(\hat{\theta}_t)$$

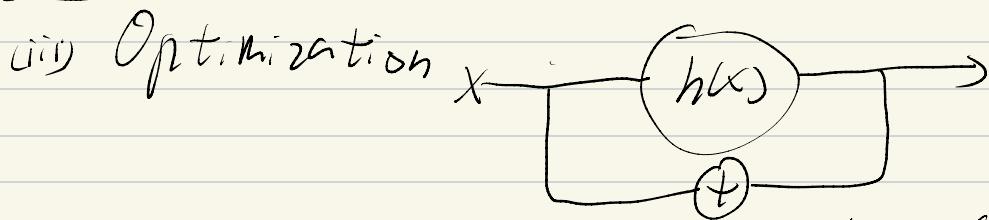


$$\vec{m}_{t+1} = \mu \vec{m}_t + (1-\mu) \vec{g}_t$$

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \vec{m}_t$$

learning rate

Adam = Momentum + Adaptive Learning Rate



(a) Regu [w/ize]

Optimization landscape.

(b) Gradient Vanishing / Explosion

$$g_1, g_2, \dots, g_k; g_i = \nabla f_i \quad f_i = x + h_i, \alpha$$

$$\nabla f_i = I + \nabla h_i$$

for Initialization of the Initialization

(a) Batch Normalization

$$2 \text{ arms / Ads} \quad \text{ad}_1 \quad \text{ad}_2 \\ \theta_1 \quad \theta_2$$

Prior $\theta_1 \sim \text{Beta}(\alpha_1, \beta_1)$; $\theta_2 \sim \text{Beta}(\alpha_2, \beta_2)$

$$\text{pdf } x^{\alpha_1}(1-x)^{\beta_1} \quad x \in [0, 1]$$

$\alpha_1 = \beta_1 = \alpha_2 = \beta_2 \Rightarrow \theta_1, \theta_2 \sim \text{Uniform}[0, 1]$

TS: Period $\alpha_1(t), \alpha_2(t),$
 $\beta_1(t), \beta_2(t)$

Sample $\theta_1(t) \sim \text{Beta}(\alpha_1(t), \beta_1(t))$

$\theta_2(t) \sim \text{Beta}(\alpha_2(t), \beta_2(t))$

$$\alpha_t^* = \begin{cases} 1 & \text{if } \theta_1(t) > \theta_2(t) \\ 2 & \text{if o.w.} \end{cases}$$

$$\alpha_{at}^*(t+1) = \alpha_{at}^*(t) + R_t$$

$$\beta_{at}^*(t+1) = \beta_{at}^*(t) + 1 - R_t$$