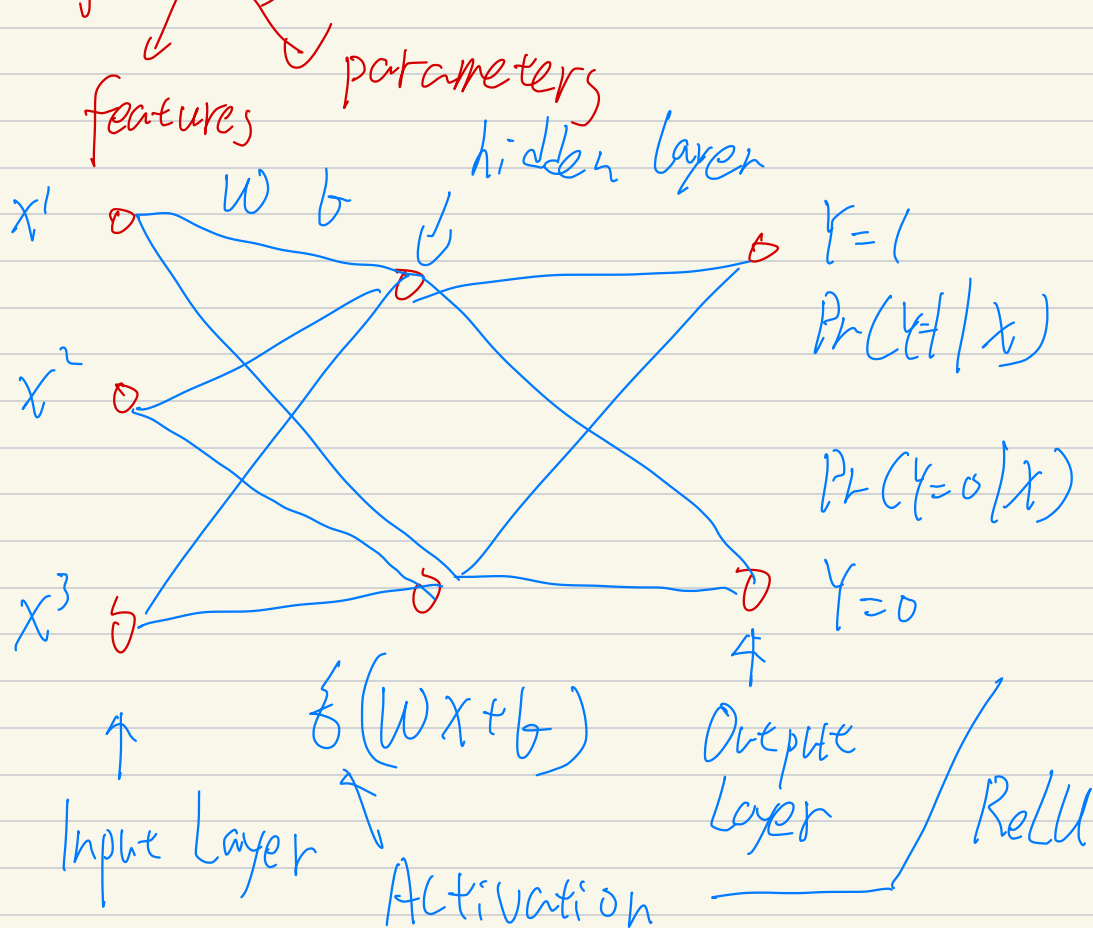


$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(x_i, f(x_i; \theta)) \equiv \mathcal{L}(\theta)$$

$f(x, \theta)$ is a DNN,



$$f(x; \theta) = b_2 \left(\sigma_1(W_1 x + b_1) \right) + w_2 + b_2$$

Sigmoid / tanh

① What is the loss $L(y, \hat{y})$

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

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

② How do we find $\hat{\theta}_{\text{Gradient Descent}}$

$$\hat{\theta}_{\text{GD}} = \hat{\theta}_0 - 2 \cdot \nabla_{\theta} L(\theta)$$

③ How do we estimate $\nabla_{\theta} L(\theta)$

(i): Chain Rule / BP

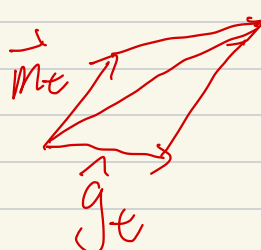
$$y = f(z) \quad z = g(x)$$

$$\frac{dy}{dx} = \frac{df}{dz} \bigg|_{y=f(z)} \cdot \frac{dz}{dx} \bigg|_x,$$

(i) $SGID \{1, 2, \dots, B\} \subseteq \{1, \dots, n\}$

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \mathcal{L}(y_i, f(x_i, \theta))$$

Momentum:



$$\vec{m}_{t+1} = \mu \cdot \vec{m}_t + \vec{g}_t$$

Adam = Momentum + RMSprop

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

(i) Initialization $(w, b) \sim N(0, \frac{1}{n_0})$
 (ii) B/Normalization $\frac{x_i - \bar{x}_i}{SD(x_i)}$
 (iii) Skip Connection g_1, g_2, \dots, g_k

\downarrow
 Dist. of $\hat{\theta}_0$

$g_1, g_2, g_3, \dots, g_k$: Gradients
Vanishing

$$\underbrace{-x + f_i(x) + x}_{\uparrow} \\ I.$$