$$(x_i, y_i)_{i=1, 2, \cdots, n}$$

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \lambda \left( y_i, f(x_i, \theta) \right)$$

$f(x, \theta)$ is a 1) NN

feature      parameter.      $\lambda$: loss $\Big\{$ regression: $(y - \hat{y})^2$

classification:

$-y \log \hat{y} - (1-y) \log(1-\hat{y})$

$x_1$    $W x + b$    $\xi(W x + b)$

activation.

relu    $y=1$   $P_r(y=1|x)$

sigmoid

$x_2$    tanh

$y=0$   $P_r(y=0|x)$

$x_3$

Input layer      Hidden layer.    Output layer.

$$f(x, \theta) = \xi_2 \left( W \xi_1 \left( W x + b \right) + b \right)$$

① Loss? ✓.

② How to find $\hat{\theta} \in \arg\min L(\theta)$

Gradient Descent.

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \nabla_\theta L(\hat{\theta}_t)$$

③ How to estimate $\nabla_\theta L(\hat{\theta})$?

(i) BP / Chain-Rule.

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

$$\frac{dy}{dx} = \frac{df}{dz}\bigg|_{z=g(x)} \cdot \frac{dg}{dx}\bigg|_x.$$

(ii) SGD: $B \subseteq D$

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \frac{1}{B} \sum_{[i]}^{B} \nabla_\theta L(\hat{\theta}_t)$$

(iii) Adam