# Lecture 3

Cross Entropy

$$L = -\sum y_j \log S_j$$

estimation ⌐ ground truth

$$= -(0 + 0 + \log 0.978) = 0.02$$

GT $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$   S-outp $\begin{bmatrix} 0.002 \\ 0.02 \\ 0.978 \end{bmatrix}$
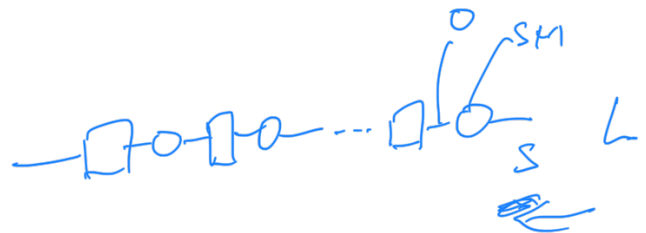
S-out $\begin{bmatrix} 0.2? \\ 0.?? \\ 0.?? \end{bmatrix}$   $L = -(0 + 0 + \log 0.33) = 1.1$



$$\frac{\sum L_i(S_i, Y_i)}{\sum (S-y)^2}$$

$$\frac{\partial L}{\partial O_i} = \frac{\partial L}{\partial S} \frac{\partial S}{\partial O_i}$$

$$L = -\sum_j y_j \log S_j$$

soft max

$$S_j = \frac{e^{O_j}}{\sum\limits_{i=1} e^{O_i}}$$

$$\frac{\partial s_j}{\partial O_i} = s_j(1-s$$

$$= -\sum_j y_j \frac{\partial \log s_j}{\partial s_*} \frac{\partial s_i}{\partial O_*}$$

$$\frac{\partial \log(x)}{\partial x} = \frac{1}{x}$$

$$= -\sum_j \left( y_j \frac{1}{s_j} \frac{\partial s_i}{\partial O_i} \right)$$

$$\frac{\partial \sum f(x_i)}{\partial x_i} = \sum \frac{\partial f(x_i)}{\partial x_i}$$

$i \ne j$

$i = j$

$$= -y_i \frac{s_i(1-s_i)}{s_i} - \sum_{i \ne j} y_j \frac{1}{s_j} \frac{\partial s_j}{\partial O_i}$$

$$-s_i s_j$$

$$= -y_i(1-s_i) - \sum_{i \ne j} y_j(-s_i)$$

$$=$$

$$\frac{\partial s_j}{\partial O_i} \quad \text{where } i \ne j \quad \sum e^{O_k} = C + e^{O_i}$$

$$= \partial \frac{\frac{e^{O_j}}{\sum e^{O_k}}}{\partial O_i} = \frac{\partial \frac{e^{O_j}}{C + e^{O_i}}}{\partial O_i}$$

$$= e^{O_j} \frac{\partial \frac{1}{C + e^{O_i}}}{\partial O_i}$$

$$= -e^{O_j} e^{O_i}/(C + e^{O_i})^2$$

$$= - \frac{e^{o_j}}{c + e^{o_i}} \cdot \frac{e^{o_i}}{c + e^{o_i}}$$

$$c + e^{o_i} = \sum_k e^{o_k}$$

$$= - \underbrace{\frac{e^{o_j}}{\sum_k e^{o_k}}}_{\parallel} \quad \underbrace{\frac{e^{o_i}}{\sum_k e^{o_k}}}_{\parallel}$$

$$\qquad S_j \qquad S_i$$

$$= - S_j \, S_i$$

$$= - y_i (1 - S_i) + \sum_{i \neq j} y_j \, S_i$$

$$= - y_i + y_i S_i + \sum_{i \neq j} y_j \, S_i$$

$$= - y_i + \sum_j y_j \, S_i$$

$$= - y_i + S_i \sum y_j$$

$$= S_i - y_i$$

input

$$\boxed{\phantom{x}} - \bigcirc - \cdots \cdots - \bigcirc - \quad L$$

# Training

- preprocessing

- overfitting, Generalization

- Large data set

input $X_i$ $D-$ dimension

$i = 1 \sim N$, $N$ samples

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} X_{ij} \qquad \leftarrow \text{mean}$$
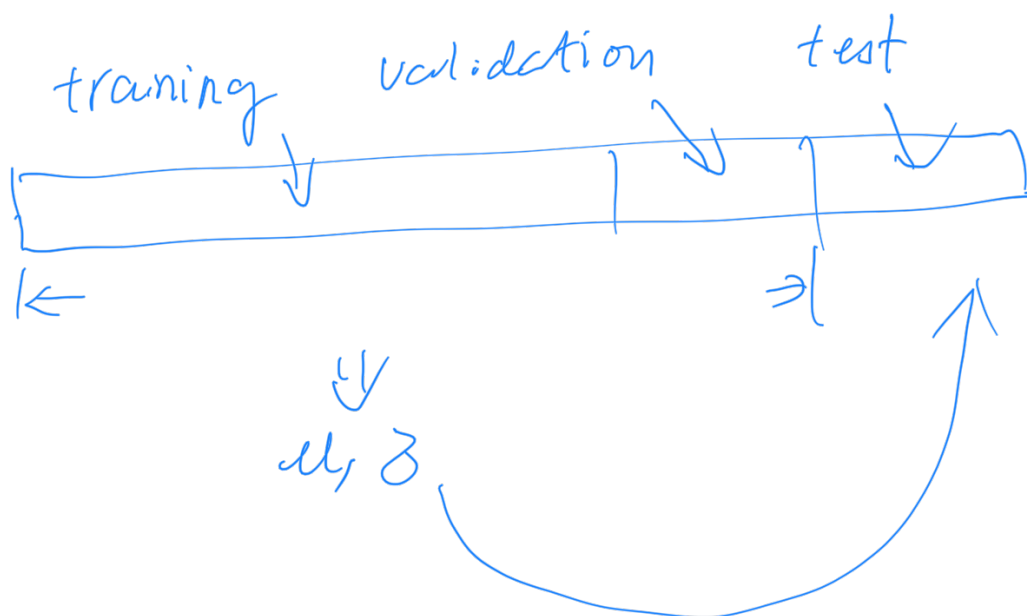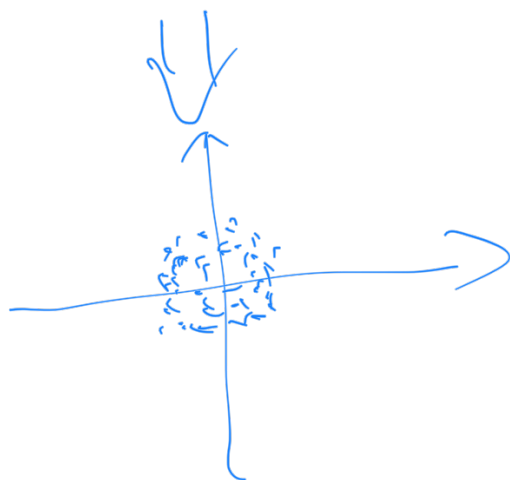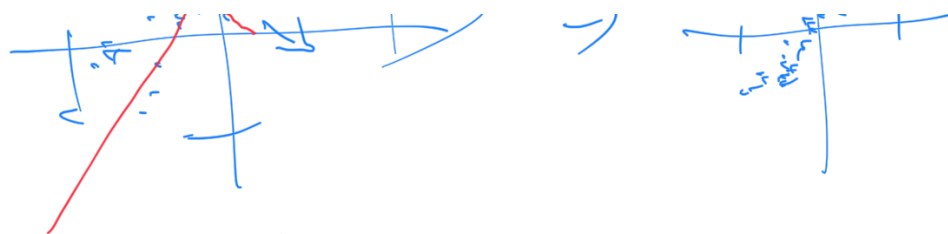
$j = 1 \rightarrow D$

$$\partial_j = \frac{1}{N} \sum_{i=1}^{N} (X_{ij} - \mu_j)^2 \quad \leftarrow \text{variance}$$

$j = 1 \rightarrow D$

$$\hat{X}_{ij} = \frac{X_{ij} - \mu_j}{\sqrt{\partial_j + \varepsilon}}$$

$\phantom{xxxxxxxxxxxxxxxxxx}\llcorner$ small #

training    validation    test

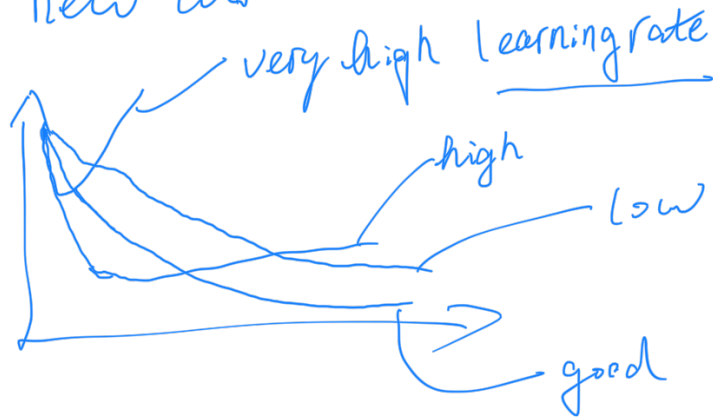$$\mu, \partial$$

## Early Stopping

loss

validation

training

$t$

$$w_{t+1} = w_t - \lambda \frac{\partial L}{\partial w}$$

keep a copy when validation reach a
new low



Regularization

  $w^T w$ to be small

$$\sum w_i^2 = w^T w$$

$$L = L + \gamma L(w)$$
$$\quad\quad\quad\quad L\ w^T w \text{ or } \sum w_i^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial w} + r \frac{\partial L(w)}{\partial w}$$
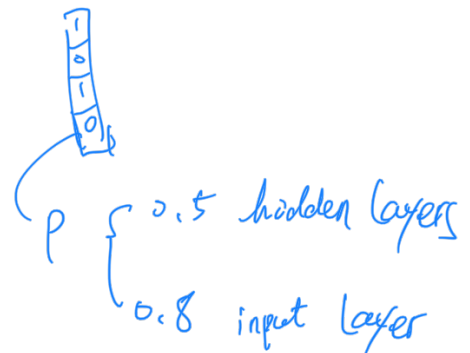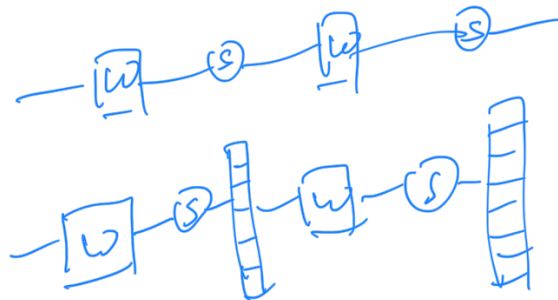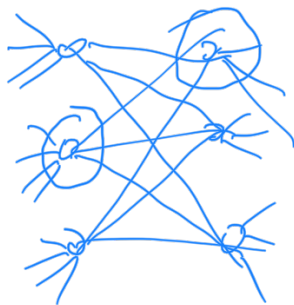
$$\frac{\partial L(w)}{\partial w_j} = \frac{\sum w_i^2}{\partial w_j} = 2w_j$$

$$w_j^{t+1} = w_j^t - \lambda \left( \frac{\partial L}{\partial w_j} + r 2 w_j^t \right)$$

$$= W_j^t - \lambda \frac{\partial L}{\partial \omega_j} - 2\lambda \gamma W_j^t$$

$$= W_j^t (1 - 2\lambda \gamma) - \lambda \frac{\partial L}{\partial \omega_j}$$

Drop out



$$P \begin{cases} 0.5 \text{ hidden layers} \\ 0.8 \text{ input layer} \end{cases}$$

Bagging / Ensemble

short for

bootstrap aggregation

- Different dataset

  { common data , different data }

- Different initializations

- Different batches

# Stochastic Gradient Descent (SGD)

Loop

Sample a batch of data from whole

Use the batch to compute feed forward

- - - - - - - backpropagation to get

$\Delta w$

update $w_{t+1} = w_t - \lambda \delta w$

Data size = 100,000 , batch size 1,000

- randomize the whole data set in terms of order

- segment it based on batch size

- run SGD on batch i

- till to the end of all batches

epoch