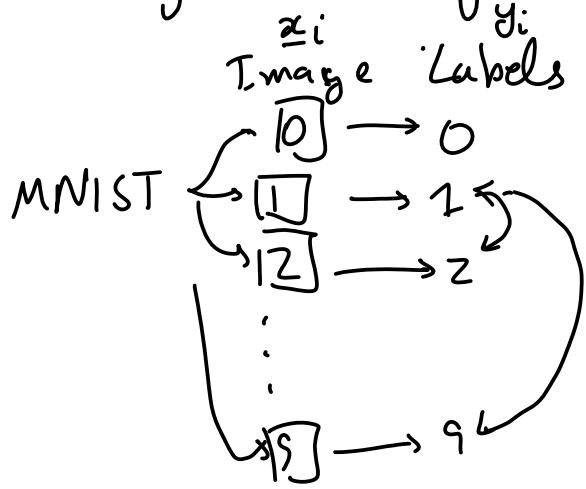


① Cross Entropy loss : Multi class - classification

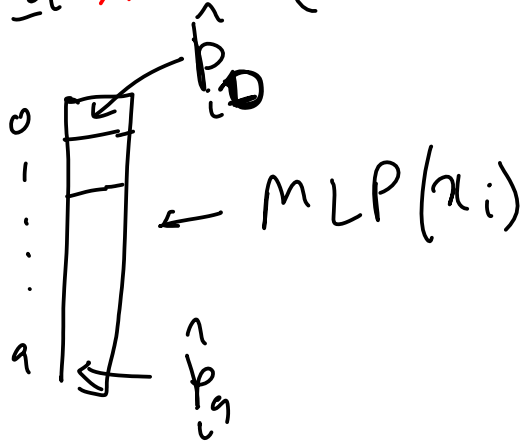
② Weight Decay as Regularizer



label 1 is no more
closer to 2 than 9

$$\hat{y}_i \neq \hat{MLP}(x_i)$$

$$\|y_i - \hat{y}_i\|_2^2 \times$$



$$\hat{p}_i = \text{MLP}(\underline{x}_i)$$

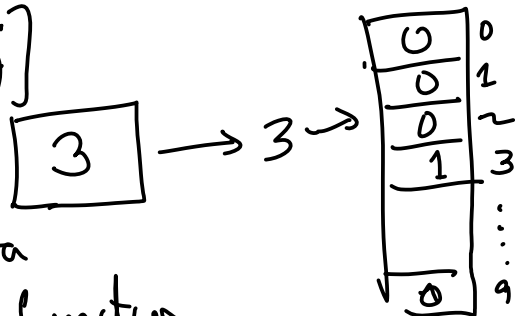
How can we ensure that $\sum_{j=0}^q \hat{p}_{ij} = 1$?
 $0 \leq \hat{p}_{ij} \leq 1$?

True probability

$$p_{ij} = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{otherwise} \end{cases}$$

$$p_{ij} = \mathbb{1}_{\{y_i = j\}}$$

Dirac delta
indicator function



Cross entropy loss

$$l(p_{ij}, \hat{p}_{ij}) = \sum_{j=0}^q -p_{ij} \log_e \hat{p}_{ij}$$

$$= -\log_e \hat{p}_{i[y_i]}$$

$$\hat{p}_{ij} < 1 \\ \Rightarrow \log \hat{p}_{ij} < 0$$

$$\hat{p}_{i[y_i]} = 1 \\ \Rightarrow \log \hat{p}_{i[y_i]} = 0$$

Entropy

$$H(x) = \sum_{x \in \Omega} -P[X=x] \log_e(P[X=x])$$

$$= \int_{x \in \Omega} -f(x) \log_e f(x) dx$$

$$= \mathbb{E}_x [\log_e(P[x])]$$

How can we ensure that $\sum_{j=0}^q \hat{p}_{ij} = 1$?
 $0 \leq \hat{p}_{ij} \leq 1$?

Softmax layer (temperature)

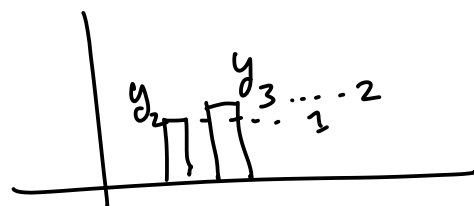
$$s(\underline{y}) = \frac{\exp(\underline{y})}{\sum_{i=1}^n \exp(y_i)}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$= \frac{\exp(\underline{y})}{\mathbf{1}^T \exp(\underline{y})}$$

$$-\infty < y_i < +\infty \xrightarrow{+ve} 0 < \exp(y_i) < \infty \\ \downarrow \text{normalize} \\ 0 < \exp(y_i) / \sum \exp(y_i) < 1$$

Why name softmax?



$$y_2 = 10$$

$$y_3 = 20$$

$$\exp(1) < \exp(2)$$

$$1 < 2$$

$$e^1 < e^2$$

2.71 times

$$e^{10} < e^{20}$$

$\sim e^{10}$ times

$$\sum \frac{\exp(y_i)}{\sum \exp(y_i)} = 1$$

$\underbrace{\hspace{1cm}}_{p_i}$

softmax exaggerates differences so much that the highest item goes close to 1 while others go close to zero.

Increasing temperature

$$S(y_i) = \frac{\exp(T y_i)}{\sum \exp(T y_i)}$$

loss cross entropy \leftarrow softmax \leftarrow MLP

$$\text{Cross entropy with softmax} = \sum_{j=0}^a p_{ij} \log \left[\frac{\exp(\hat{y}_{ij})}{\sum_{k=0}^a \exp(\hat{y}_{ik})} \right] = \sum_{j=0}^a p_{ij} \hat{y}_{ij} - \sum_{j=0}^a p_{ij} \log \left(\sum_{k=0}^a \exp(\hat{y}_{ik}) \right)$$

minimize \sum loss
Dataset

$$= \sum_{j=0}^a p_{ij} \hat{y}_{ij} - \log \left(\sum_{k=0}^a \exp(\hat{y}_{ik}) \right) \underbrace{\sum_{j=0}^a p_{ij}}_1$$

loss $\left\{ \begin{array}{l} \rightarrow \text{least square loss for regression} \\ \rightarrow \text{Hinge loss for two class classification} \\ \rightarrow \text{Cross entropy loss for multi-class classification} \end{array} \right.$

$$\text{cross entropy loss with softmax} (p_{ij}, \hat{y}_{ij}) = \sum_{j=0}^a p_{ij} \hat{y}_{ij} - \log \left(\sum_{k=0}^a \exp(\hat{y}_{ik}) \right)$$

logits (log+bits)
 $\log(\hat{p}_{ij})$

Weight Decay

= L_2 -regularization

minimize w

$$\sum_{\text{Data}} l(y_i, \hat{y}_i) + \lambda \|w\|_2^2$$

Prob. Prior Occam razor smoother functions

$$1 > (1 - 2\alpha^+ \lambda) > 0$$

$$\underline{w}_{t+1} = \underline{w}_t - \alpha^t \sum_B \nabla_{\underline{w}} l(y_i, \hat{y}_i) - \underbrace{2\alpha^t \lambda \underline{w}_t}_{\text{weight decay}} \quad \left. \begin{array}{l} \text{weight decay parameter} \approx 0.1 \\ \|\underline{w}\|^2 = \underline{w}^T \underline{w} \end{array} \right\} = \underline{w}^T \underline{I} \underline{w}$$

$$\underline{w}_{t+1} = \underbrace{(1 - 2\alpha^+ \lambda)}_{\text{weight decay}} \underline{w}_t - \alpha^t \sum_B \nabla_{\underline{w}} l(y_i, \hat{y}_i)$$

$$\frac{\partial \underline{w}^T \underline{I} \underline{w}}{\partial \underline{w}} = 2 \underline{w}^T \underline{I} = 2 \underline{w}^T$$

$$\frac{\partial}{\partial \underline{w}} \underbrace{\underline{w}^T}_{1 \times n} \underbrace{\underline{I}}_{n \times n} \underline{w} = 2 \underbrace{\underline{w}^T}_{1 \times n} \underline{I} = 2 \underline{w}^T_{1 \times n}$$

$$\frac{\partial}{\partial \underline{w}} \underline{w}^T A \underline{w} = \underline{w}^T (A + A^T) = 2 \underline{w}^T A$$