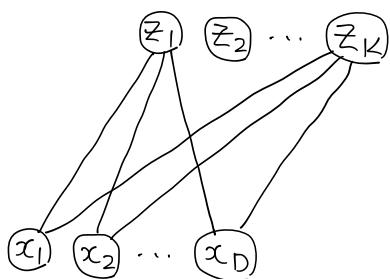


1. What weight are used in multiclass classification?



$$z_1 = \vec{w}_1^T \vec{x} + b_1$$

$$\begin{matrix} \uparrow & \uparrow \\ \begin{pmatrix} w_{11} \\ w_{12} \\ \vdots \\ w_{1D} \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix} \end{matrix}$$

$$z_k = \vec{w}_k^T \vec{x} + b_k$$

$$\begin{matrix} \uparrow & \uparrow \\ \begin{pmatrix} w_{k1} \\ w_{k2} \\ \vdots \\ w_{kD} \end{pmatrix} & \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix} \end{matrix}$$

same K linear models

$$\begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{pmatrix} \leftarrow \vec{z} = \vec{w} \vec{x} + \vec{b} \leftarrow \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{pmatrix}$$

$$\vec{w} = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1D} \\ \vdots & \vdots & & \vdots \\ w_{K1} & \dots & \dots & w_{KD} \end{pmatrix}$$

how do we obtain \vec{w} & \vec{b} ? vid 22

2. provide examples of softmax

$$\vec{y} = \text{softmax}(\vec{z}) = \left[\frac{e^{z_k}}{\sum_{k'=1}^K e^{z_{k'}}} \right]_{k \in 1 \dots K}$$

$$\vec{z} = [3, -1, 4]^T \rightarrow \vec{y} = [0.3, 0, 0.7] \quad \frac{e^4}{e^3 + e^4}$$

$$\vec{z} = [200, 10, -4]^T \rightarrow \vec{y} = [1, 3 \cdot 10^{-83}, 2.5 \cdot 10^{-89}]$$

$$\approx [1, 0, 0]$$

3. More intuition behind connection between softmax & logistic activation?

$$\text{recall } y = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$(\text{binary}) \quad \mathcal{E} = -t \log(y) - (1-t) \log(1-y)$$

$$(\text{multi-class}) \quad \mathcal{E} = -\sum_{k=1}^K t_k \log(y_k)$$

$$\vec{t} = [0, 0, \dots, 0, \underset{\uparrow}{1}, 0, \dots, 0]$$

index $j=1$ if sample is in class j

take $K=2$ (2 classes)

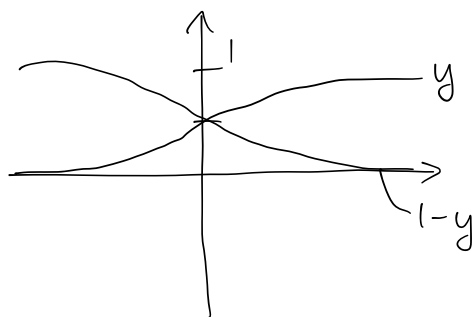
$$\mathcal{E} = -\sum_{k=1}^2 t_k \log(y_k)$$

$$= -t_1 \log(y_1) - t_2 \log(y_2)$$

$$\left. \begin{array}{l} \text{let } t_2 = 1 - t_1 \\ y_2 = 1 - y_1 \end{array} \right\} \mathcal{E} = -t_1 \log(y_1) - (1 - t_1) \log(y_2)$$

what about softmax?

$$1-y = 1 - \sigma(z) = 1 - \frac{1}{1+e^{-z}} = \frac{e^{-z}}{1+e^{-z}} = \frac{1}{1+e^z}$$



• sum always 1 for softmax

observe

$$y = \frac{1}{1+e^{-z}} = \frac{e^z}{e^z + e^0}$$

$$1-y = \frac{1}{1+e^z} = \frac{e^0}{e^0 + e^z}$$

$$\begin{pmatrix} y \\ 1-y \end{pmatrix} = \begin{pmatrix} \frac{e^z}{e^z + e^0} \\ \frac{e^0}{e^0 + e^z} \end{pmatrix} = \text{softmax} \begin{pmatrix} z \\ 0 \end{pmatrix}$$

$z \gg 0$ output 1

$z \ll 0$ output 0

4. How does cross-entropy encourage correct preds?

$$\mathcal{E} = - \sum_{k=1}^K t_k \log \left(\frac{e^{z_k}}{\sum_{k'=1}^K e^{z_{k'}}} \right)$$

$$= - \sum_{k=1}^K t_k \left(z_k - \log \left(\sum_{k'=1}^K e^{z_{k'}} \right) \right)$$

= 1 only if $k=j$ should be huge for z_j

$$- \log \left(\sum_{k'=1}^K e^{z_{k'}} \right)$$

Let $m = \operatorname{argmax}_k z_k$

if $m=j$ and $z_m \gg z_k$ for $k \neq j$

$$\log \left(\sum_{k'=1}^K e^{z_{k'}} \right) \approx \log(e^{z_m}) \approx z_m = z_j$$

$$\hookrightarrow \varepsilon \approx -(z_j - z_j) = 0$$

$$\text{if } m \neq j \Rightarrow \varepsilon \approx -(z_j - z_m)$$

encourage correct prediction

penalize incorrect prediction

5. How to decide on # neurons/layers & activation function

• model zoo

\hookrightarrow CNN: images, videos

\hookrightarrow language models (attention)

- hyperparameter training
 - ↳ grid/random search
- neural architecture search