# Lecture 13

$$\overrightarrow{net_1} = W_0^T \vec{x} + \vec{b_0}$$

$$\vec{z_1} = f^1(\overrightarrow{net_1})$$

$$\overrightarrow{net_2} = W_1^T \vec{z_1} + \vec{b_1}$$

$$\vec{z_2} = f^2(\overrightarrow{net_2})$$

$$\overrightarrow{net_o} = W_h^T \vec{z_u} + \vec{b_h}$$

$$\vec{o} = f^o(\overrightarrow{net_o})$$

$$\overrightarrow{net_h} = W_{h-1}^T \vec{z_{h-1}} + \vec{b_{h-1}}$$

$$\vec{z_u} = f^h(\overrightarrow{net_h})$$

1) feed-forward

$$\vec{x} \rightarrow \vec{o}$$

given input    compute the output

2)   Loss / Error

3)   backpropagation

Compute the gradients of all $W_i, \vec{b_i}$

$$\nabla_{W_i}, \nabla_{b_i}^s \qquad \forall\, i = 0 \dots h$$

$$\nabla W_1 = \frac{\partial \alpha}{\partial w_1}$$

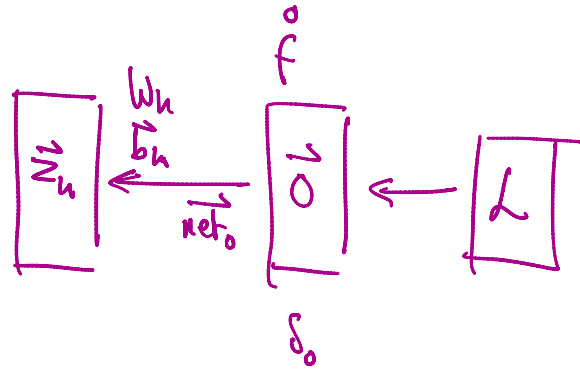$$\nabla_{W_i} = \frac{\partial \alpha}{\partial w_i}$$

weight gradient

$$\nabla_{b_i} = \frac{\partial \alpha}{\partial b_i}$$

bias gradient

# Computational graph approach

net gradient for each layer

$$\nabla_{net_i} = \boxed{\vec{\delta}_i} = \frac{\partial \mathcal{L}}{\partial net_i}$$



$$\vec{o} = \overset{\circ}{\vec{f}}(\vec{net_o})$$

$$\vec{net_o} = W_h^T \vec{z_h} + \vec{b_h}$$

$$\vec{\delta}_o = \frac{\partial \mathcal{L}}{\partial net_o}$$

$$= \frac{\partial \vec{o}}{\partial \vec{net_o}} \cdot \frac{\partial \mathcal{L}}{\partial \vec{o}}$$

$$= \boxed{\frac{\partial \overset{\circ}{f}(\vec{net_o})}{\partial \vec{net_o}} \cdot \frac{\partial \mathcal{L}}{\partial \vec{o}}}$$

derivative of the activation function $f^o$

$$1 \cdot (\vec{o} - \vec{y})$$

$$\boxed{\vec{\delta}_o = \vec{o} - \vec{y}}$$

e.g.

squared error $\quad \mathcal{L} = \frac{1}{2} \| \vec{y} - \vec{o} \|^2$

$$\frac{\partial L}{\partial \vec{o}} = \frac{1}{2} \left( 2 \cdot (\vec{y} - \vec{o}) \right) \cdot -1$$

$$= \vec{o} - \vec{y}$$

predicted value $\quad$ true value

$$\overset{\circ}{f} = linear$$

$$\frac{\partial \overset{\circ}{f}(net_o)}{\partial net_o} = 1$$

$$\frac{\partial \mathring{f}(net_o)}{\partial net_o} = \frac{\partial net_o}{\partial net_o} = \vec{1}$$

---

e.g. $\mathring{f} = $ Softmax

$\mathcal{L} = CE \text{ } \underline{loss} \text{ } (\text{cross entropy})$

$$\mathcal{L} = -\sum_{j=1}^{k} y_j \log(o_j)$$

$$\frac{\partial \mathcal{L}}{\partial \vec{o}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial o_1} \\ \frac{\partial \mathcal{L}}{\partial o_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial o_k} \end{bmatrix} = \begin{bmatrix} -y_1/o_1 \\ -y_2/o_2 \\ \vdots \\ -y_k/o_k \end{bmatrix}$$

$\vec{y} \equiv $ One-hot vector

$\vec{o} \equiv $ Vector of probabilities for $k$ classes

$$\frac{\partial \mathcal{L}}{\partial o_1} = -\frac{y_1}{o_1}$$

$$\vec{y} = \begin{pmatrix} 0 & 0 & -- & 1 & 0 & 0 & 0 \end{pmatrix} = (y_1, y_2 \dots y_k)^T$$

$$\vec{o} = (o_1, o_2, \dots o_k)^T$$

$$\frac{\partial \text{ Softmax}(\vec{net_o})}{\partial \vec{net_o}} = \begin{bmatrix} & & \\ & \text{matrix} & \\ & & \end{bmatrix}$$

$$\text{Softmax}(net_i) = \frac{e^{net_i}}{\sum_{j=1}^{k} e^{net_j}}$$

$$\underbrace{\frac{\partial f(net_o)}{\partial net_o} \cdot \frac{\partial \mathcal{L}}{\partial \vec{o}}}$$

$$\vec{o} - \vec{y}$$

$$w_h$$
$$\overrightarrow{z_h} \qquad \overleftarrow{net_o} \qquad \overrightarrow{o}$$
$$b_h \qquad \qquad \qquad \leftarrow \qquad \boxed{\mathcal{L}}$$

$$\uparrow \text{ loss}$$

$$\boxed{\delta_o}$$

$$\nabla w_h = \frac{\partial \mathcal{L}}{\partial w_h} = \frac{\partial net_o}{\partial w_h} \cdot \underbrace{\frac{\partial f(net_o)}{\partial net_o} \cdot \frac{\partial \mathcal{L}}{\partial \overrightarrow{o}}}_{\overrightarrow{\delta_o}}$$

$$= \frac{\partial w_h^T \overrightarrow{z_h} + b_h}{\partial w_h} \cdot \overrightarrow{\delta_o}$$

$$\boxed{\nabla w_h = \left( \overrightarrow{z_h} \cdot \overrightarrow{\delta_o}^T \right)}$$

$$\underbrace{\underbrace{m_h \cdot 1}_{\phantom{x}} \quad \underbrace{1 \cdot k}_{\phantom{x}}}_{m_h \cdot k}$$

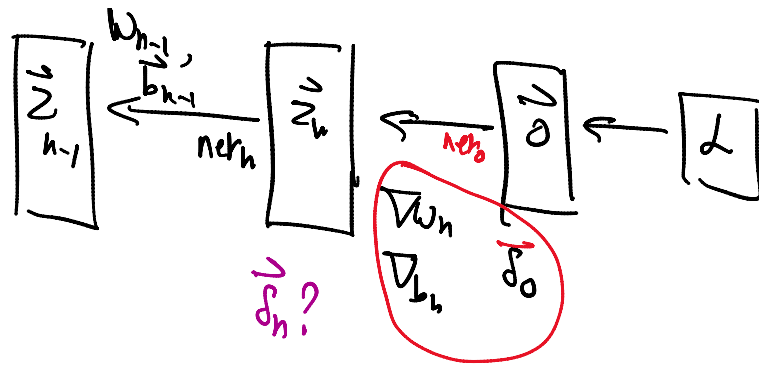$$W_h \in \mathbb{R}^{m_h \times k}$$

$$m_h = \text{size of } \overrightarrow{z_h}$$

$$k = \text{\# of output neurons}$$

$$\nabla b_h = \frac{\partial (w_h^T z_h) + b_h}{\partial b_h} \cdot \overrightarrow{\delta_o}$$

$$= \overrightarrow{1} \cdot \overrightarrow{\delta_b} = \overrightarrow{\delta_o}$$

$$\boxed{\overrightarrow{z}} \quad {w_{h-1}, \atop b_{h-1}} \quad \boxed{z} \qquad \boxed{z}$$

$$\vec{\delta_h} = \frac{\partial \mathcal{L}}{\partial net_h} = \frac{\partial \vec{z_h}}{\partial net_h} \frac{\partial net_o}{\partial \vec{z_h}} \underbrace{\frac{\partial \mathcal{L}}{\partial net_o}}_{\vec{\delta_o}}$$

$$\frac{\partial \overset{h}{f}(net_h)}{\partial net_h}$$

$\underbrace{\qquad}$ derivative of activation for hidden layer

$$\frac{\partial W_h^T \vec{z_h} + \vec{b_h}}{\partial \vec{z_h}} = W_h$$

$$= \frac{\partial \overset{h}{f}}{\partial net_h} \odot \left( W_h \cdot \delta_o \right)$$

element-wise product

## Backprop:

1) Compute $\vec{\delta_o} = \frac{\partial \overset{o}{f}}{\partial net_o} \cdot \frac{\partial \mathcal{L}}{\partial \vec{o}}$ $\longrightarrow \left( \vec{o} - \vec{y} \right)$

$\mathcal{L}$ = squared error
$f^o$ = linear

2) Compute **net** gradients

for $l = h, h-1, \ldots, 1$

2) for $\ell = h, h-1, \ldots, 1$

$$\vec{\delta}_\ell = \frac{\partial f^\ell}{\partial net_\ell} \odot \left( W_\ell \cdot \vec{\delta}_{\ell+1} \right)$$

$f$ = linear

$\mathcal{L}$ = CE
$f^o$ = softmax

$\mathcal{L}$ = BCE
$f^o$ = sigmid

3) Compute weight & bias gradient

$$\nabla_{W_\ell} = \underbrace{z^\ell \cdot \left( \delta^{\ell+1} \right)^T}_{\text{outerproduct} \leftarrow \text{matrix}}$$
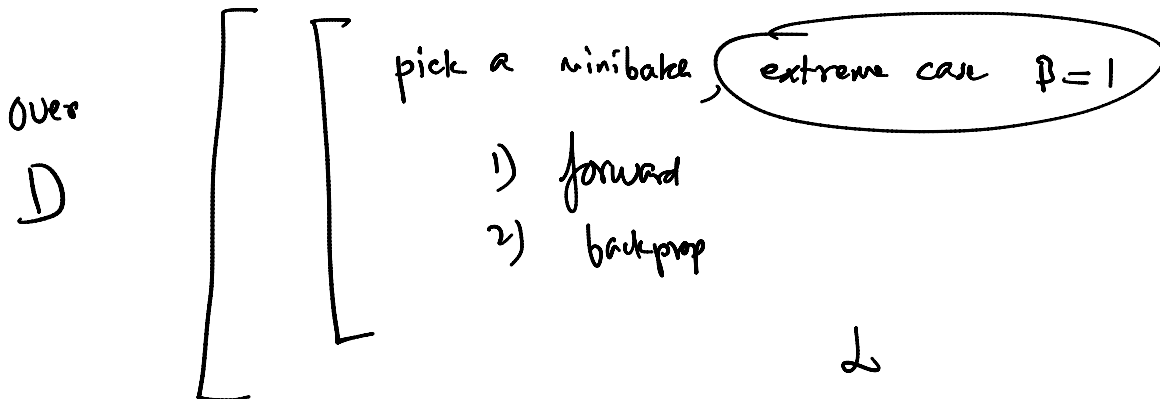
$$\nabla_{b_\ell} = \vec{\delta}^{\ell+1}$$

4) gradient descent

$$W_\ell^{(t+1)} = W_\ell^{(t)} - \eta \cdot \nabla_{W_\ell}$$

$$b_\ell^{(t+1)} = b_\ell^{(t)} - \eta \cdot \nabla_{b_\ell}$$

---

## NN training

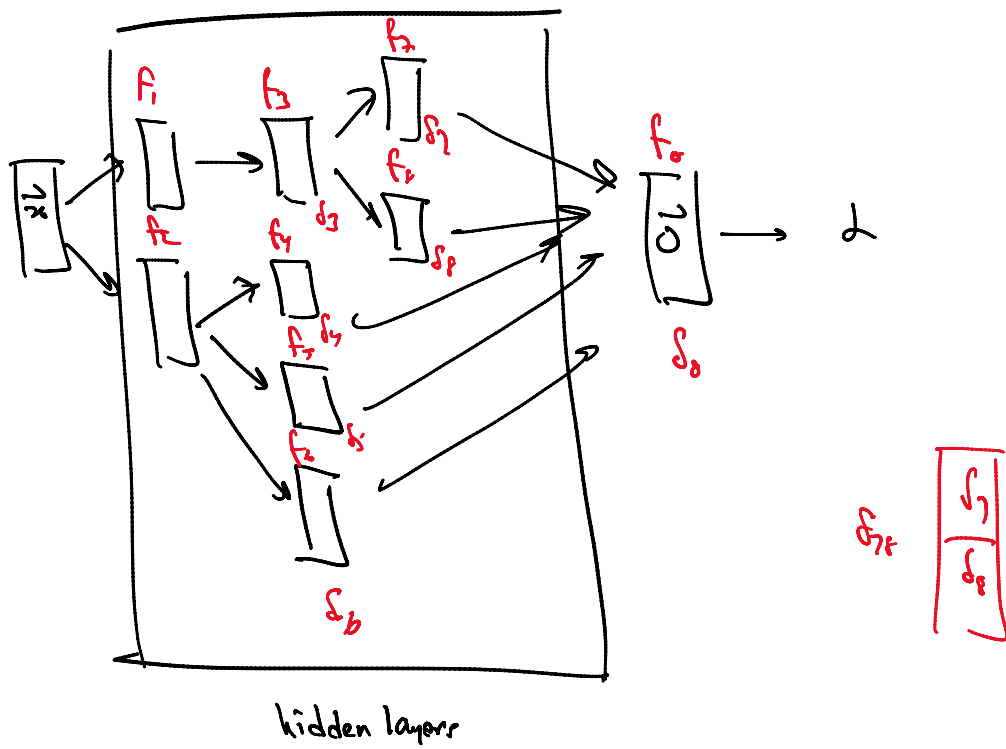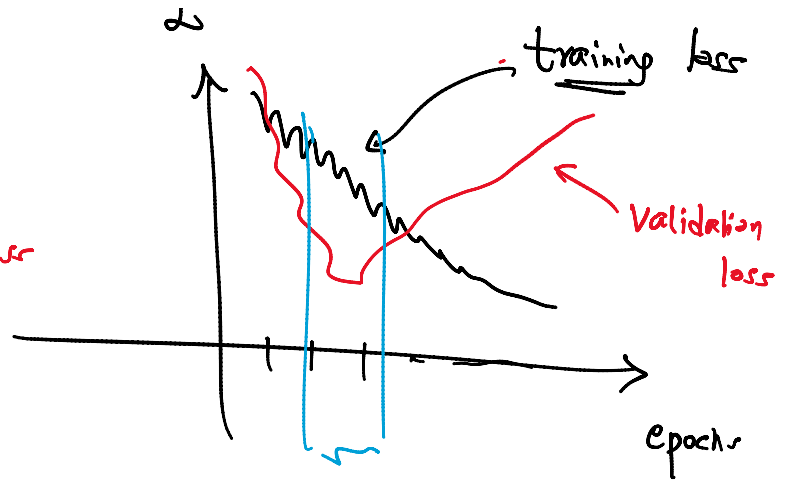for $e = 1 \ldots$ max_epochs $\quad\leftarrow$ pass over the data

over
D

pick a minibatch ⟨ extreme case $B = 1$ ⟩

1) forward
2) backprop

$\downarrow$

$\uparrow$ ⟩

training loss

L

when to stop?

look @ validation loss



$\alpha$ — training loss / Validation loss vs epochs

---



hidden layers

---

# Naive Bayes classifier

↳ Bayes classifier

↳ probabilistic classifier

$y_i \in \{ c_1, c_2 \dots c_k \}$

$$D \overset{\overrightarrow{x_i}}{\Big|} \Big| y_i \Big|$$

$$y_i \in \{ c_1 \; c_2 \; .. \; c_k \}$$

$$\underbrace{\qquad\qquad\qquad}$$

k-way classification

Posterior probability

$$P( c_i | \overrightarrow{x} )$$

$$\begin{array}{l} P(c_1|\overrightarrow{x}) \\ P(c_2|\overrightarrow{x}) \\ \vdots \\ P(c_k|\overrightarrow{x}) \end{array}$$

Optimal decision rule :

$$c_* = \underset{i=1}{\overset{k}{\arg\max}} \left\{ P(c_i|\overrightarrow{x}) \right\}$$

Unknown

a) logistic regression

b) Bayes classifier

$$P(c_i|\overrightarrow{x}) = \frac{\overset{\text{likelihood}}{P(\overrightarrow{x}|c_i)} \cdot \overset{\text{Prior}}{P(c_i)}}{\underset{\text{evidence}}{P(\overrightarrow{x})}} = \frac{P(\overrightarrow{x}|c_i) \cdot P(c_i)}{\sum\limits_{j=1}^{k} P(\overrightarrow{x}|c_j) \cdot P(c_j)}$$

Posterior

$$D \quad (\vec{x}_i, y_i)$$

$$D_1 = \{ \vec{x}_i \mid y_i = c_1 \}$$

$$D_2 = \{ \vec{x}_i \mid y_i = c_2 \}$$

$$\vdots$$

$$D_k = \{ \vec{x}_i \mid y_i = c_k \}$$

$$n_1 = |D_1| \qquad \hat{P}(c_1) = n_1/n$$

$$n_2 = |D_2| \qquad \hat{P}(c_2) = n_2/n$$

$$\vdots$$

$$n_k = |D_k| \qquad \hat{P}(c_k) = n_k/n$$

$$P(c_i) = \frac{n_i}{n} \quad \leftarrow \quad \frac{\text{\# in class } c_i}{\text{\# of total points}}$$

$$P(\vec{x} \mid c_i) = \quad \begin{array}{l} \text{a) } \textcolor{red}{\text{parametric ( model assumptions)}} \\ \text{b) } \textcolor{red}{\text{non-parametric ( no prior assumption)}} \end{array}$$

Parametric :  Assume class $c_i$ is normally distributed

$$P(\vec{x} \mid c_i) \equiv N(\vec{x} \mid \vec{\mu}_i, \Sigma_i)$$

mean    cov

for $c_i$

Unknown

$$\vec{x}_i \in \mathbb{R}^d$$

$$\vec{\mu}_i = \text{mean}(D_i) \qquad d$$

$$\Sigma_i = \text{cov}(D_i) \qquad d \times d$$

$$D \begin{array}{l} \rightarrow D_1 \quad n_1 \\ \searrow D_2 \quad n_2 \\ \qquad \vdots \\ \searrow \quad n \end{array}$$

$$\Sigma_i = cov(\mu_i) \quad d \times d$$

$n \searrow \begin{matrix} \vdots \\ D_k \quad n_k \end{matrix}$

1) Class imbalance

$$n_i \ll n$$

2) dimensionality $\underline{d}$ 

e.g. $d = 1000$

$$\Sigma_i \simeq 10^6 \text{ parameters}$$

$$\boxed{O(d^2) \overset{vs}{vs} n}$$

$$P(c_i|x) \cong \frac{N(\vec{x} \mid \mu_i, \Sigma_i) \cdot \frac{n_i}{n}}{P(\vec{x})}$$

$$c^* = \text{argmax} \left\{ \frac{N(\vec{x} \mid \mu_i, \Sigma_i) \, n_i}{\boxed{P(\vec{x}) \cdot n}} \right\}$$

$$\hat{y} = c^* = \text{argmax} \left\{ N(\vec{x} \mid \underline{\mu_i}, \underline{\Sigma_i}) \cdot \underline{n_i} \right\}$$

$\underset{\uparrow}{\text{predicted class}}$

$\underset{\uparrow}{d \times d}$

---

## Naive Bayes

$$P(\vec{x}|c_i) \equiv \text{joint prob} \quad P((x_1, x_2 \ldots x_d)^T \mid c_i)$$

$$P(\vec{x} \mid c_i) \equiv \text{joint prob} \quad P\left(\left(x_1\, x_2 \ldots x_d\right)^{\top} \mid c_i\right)$$

$$\Downarrow$$

Assume all attributes are independent

$$= P(x_1 \mid c_i) \times P(x_2 \mid c_i) \times \ldots \times P(x_d \mid c_i)$$

$$P(\vec{x} \mid c_i) = \prod_{j=1}^{k} P(x_j \mid c_i)$$

$$\Downarrow$$

approx via univariate normal

$$= \prod_{j=1}^{k} N_j\left(x_j \mid \underbrace{\mu_{ij}}_{\text{mean}},\ \underbrace{\sigma_{ij}^2}_{\text{var}}\right) \qquad \text{normal per attribute} \atop \text{per class}$$

$\{$ diagonal cov!

$$\Sigma_i = \begin{bmatrix} \sigma_{i1}^2 & & 0 \\ & \sigma_{i2}^2 & \\ & & \ddots \\ 0 & & \sigma_{id}^2 \end{bmatrix}$$
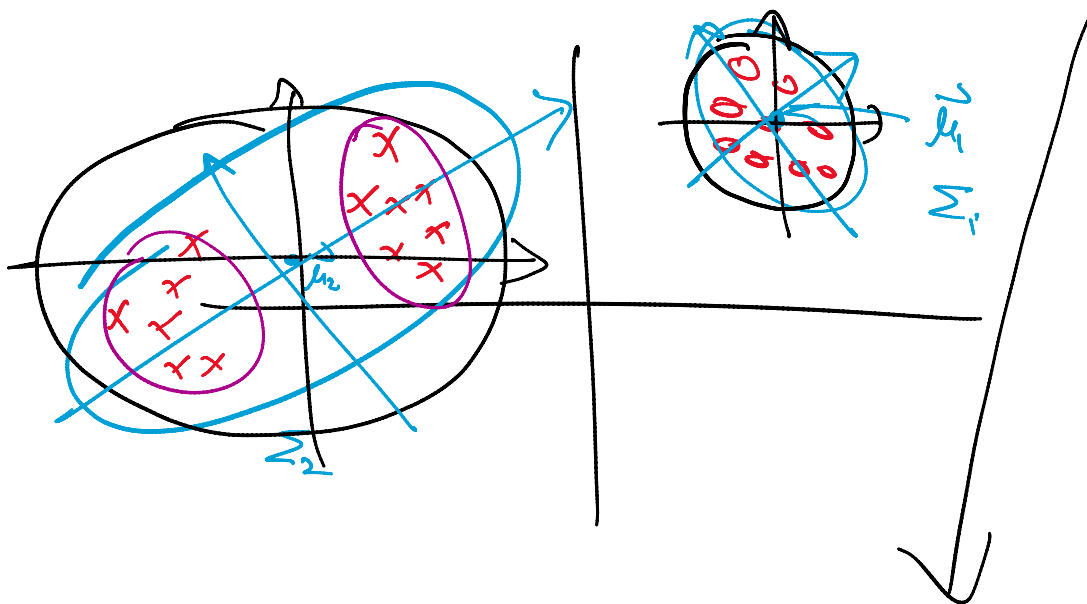
$d$ value to be estimated

$j \equiv$ attribute

$i = $ class

$\mu_{ij} = $ mean

$\sigma_{ij}^2 = $ variance

---

$$P(c_i \mid \vec{x}) = \frac{\boxed{P(\vec{x} \mid c_i) \cdot P(c_i)}}{P(\vec{x})}$$

likelihood    a)    $p(\vec{x}|c_i) = N(\vec{x}|\mu_i, \Sigma_i)$



Bayes approach
$(full\ \Sigma_i)$

Naive Bayes
$(diagonal\ \Sigma_i)$
axis aligned

mixture of Gaussians