

# Probabilistic models

## MACHINE LEARNING

**Dr. Temitayo Olugbade**

# Learning outcome



After working through this mini-video, you'll see

- ☐ how ML could be framed in terms of estimation of likelihoods,
- ☐ how a Bayes classifier works, and
- ☐ how the logistic regression works.

# Mini-video outline

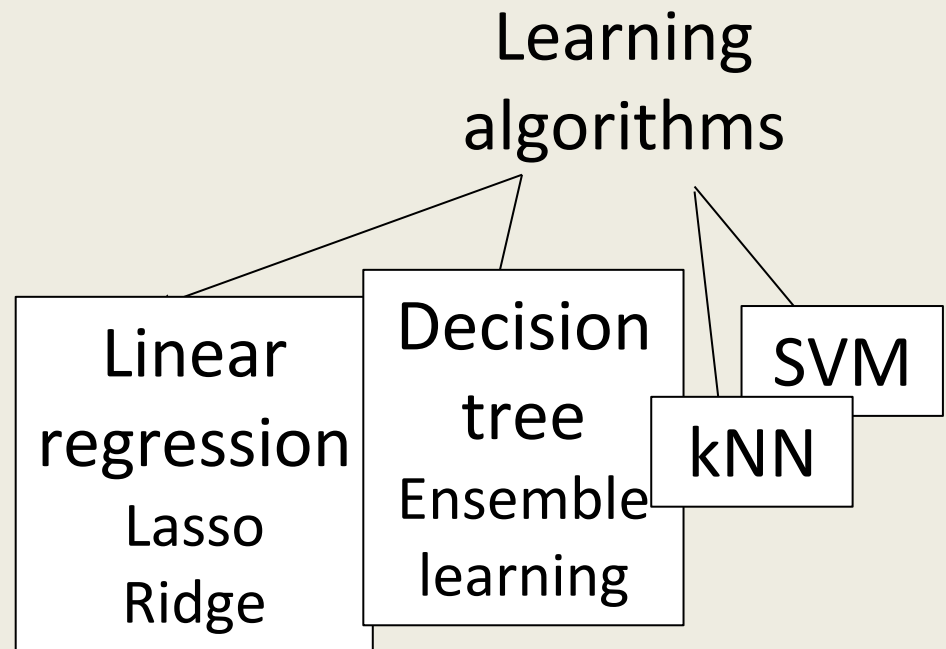
☐ Bayes classifier

☐ Logistic regression



# Recall from previous mini-videos

- The goal of ML is to learn a **model** that takes in input  $\mathbf{x}$  and gives label  $\hat{\mathbf{y}}$  close to  $\mathbf{y}$



- A model is characterised by its **learning algorithm**, **parameters**, **loss functions**, **hyperparameters**.

# Classification inference formulations

- Formulation so far
  - Given input data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $\hat{y}_i$ ?  $\hat{y}_i \subset \{c_k, \forall k, k = 1, 2, \dots K\}$
- A different (but equivalent) formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D}), \forall k, k = 1, 2, \dots K$ ?
  - $\hat{y}_i = c_k, \text{ for } \underset{k}{\operatorname{argmax}} p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D})$

# Bayes classifier

☐ Bayes classifier

☐ Logistic regression



# Bayes' Rule

For two events  $A$  and  $B$

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

where  $P(B|A)$  = conditional probability of B given observation of A,  
aka posterior probability of B

$P(B)$  = prior probability of B

$P(A|B)$  = the probability of A having occurred given observation of B  
(i.e. the probability of generating A from the associated  
probability distribution with parameters conditioned on B)  
aka likelihood of A

$P(A)$  = probability of A

# Probabilistic model formulation

*rewriting Bayes' rule for a set of possible classes  $\{c_k, \forall k, k = 1, 2, \dots, K\}$*

$$P(Y = c_k | \mathbf{X}) = \frac{P(\mathbf{X} | Y = c_k) P(Y = c_k)}{P(\mathbf{X})}$$

*let's ignore  $P(\mathbf{X})$  since it's common  $\forall k$  and so can be factored out*

$$P(Y = c_k | \mathbf{X}) = \frac{P(\mathbf{X} | Y = c_k) P(Y = c_k)}{P(\mathbf{X})}$$

*which is the same as*

$$P(Y = c_k | \mathbf{X}) = \frac{P(X_1 = x_{j_1}, X_2 = x_{j_2}, \dots, X_D = x_{j_D} | Y = c_k) P(Y = c_k)}{P(\mathbf{X})}$$



# Bayes classifier

$$P(c_k|\mathbf{X}) = \frac{P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) P(c_k)}{P(\mathbf{X})} \forall k$$

where  $c_k$  = the  $k$ th class

$K$  = total number of classes

$x_{j_1}, x_{j_2}, \dots, x_{j_D}$  = input data

# Bayes classifier – Probability term 1

$$P(c_k|\mathbf{X}) = \frac{P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) P(c_k)}{P(\mathbf{X})} \forall k$$

$P(c_k)$  can be estimated from the training data as

$$P(c_k) = \frac{N_{c_k}}{N}$$

where  $N_{c_k}$  = number of instances of class  $c_k$  in the training data

$N$  = total number of instances in the training data

# Bayes classifier – Probability term 2

$$P(c_k|\mathbf{X}) = \frac{P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) P(c_k)}{P(\mathbf{X})} \forall k$$

$P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k)$  can be estimated from the training data assuming a given distribution, e.g. Gaussian distribution

# Bayes classifier – $P(X|c_k)$

$$\begin{aligned} P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) &= \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})} \end{aligned}$$

where

$$\boldsymbol{\mu}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})$$

*see math proof in last slide pages*

# Naïve Bayes classifier

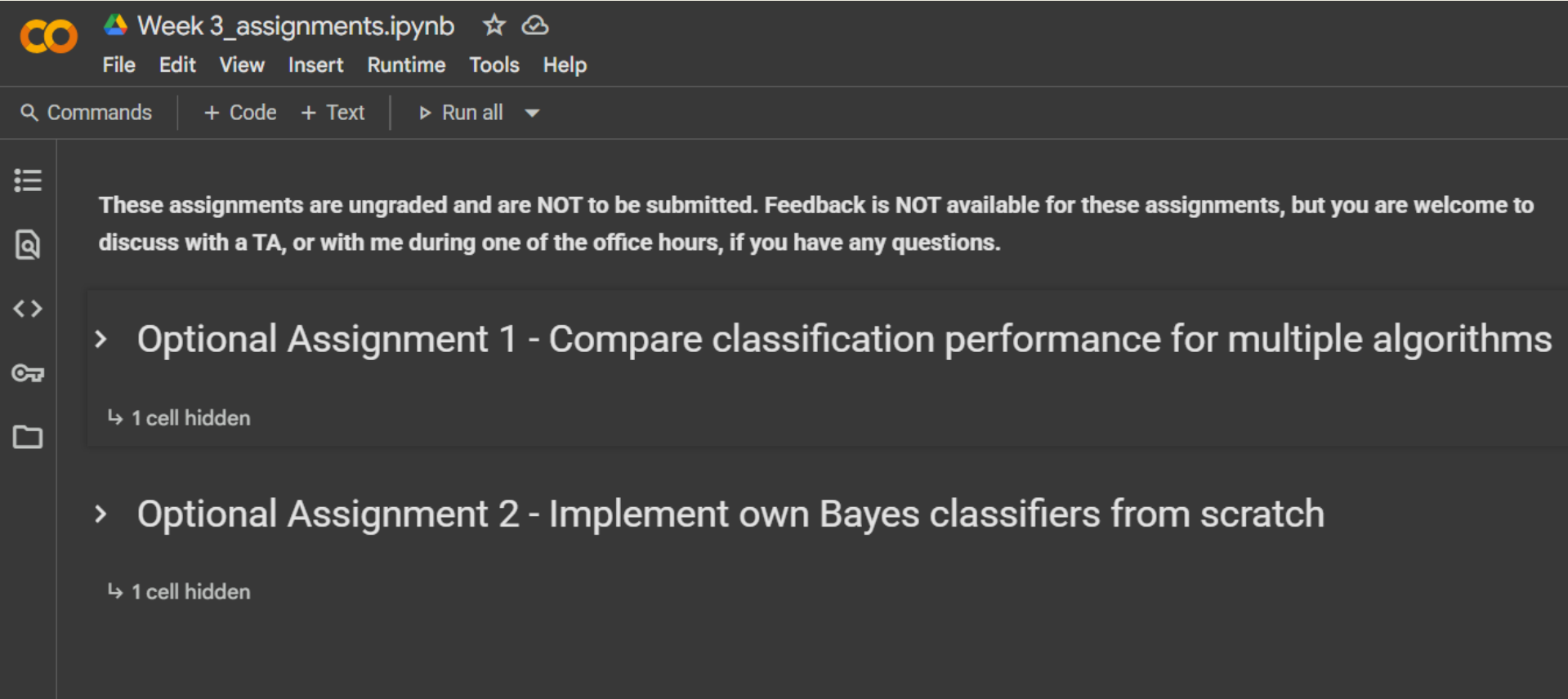
$$P(c_k|\mathbf{X}) = \frac{P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) P(c_k)}{P(\mathbf{X})} \forall k$$

Assuming that  $x_{j_1}, x_{j_2}, \dots, x_{j_D}$  are conditionally independent given class  $c_k$

$$\Rightarrow P(c_k|\mathbf{X}) = \frac{\prod_{d=1}^D P(x_{j_d} | c_k) P(c_k)}{P(\mathbf{X})}$$

- **Naïve!** – since this assumption is rarely true in the real world
- but the assumption simplifies  $P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k)$
- $P(x_{j_d} | c_k)$  can be estimated from the training data

# Implement a Bayes classifier yourself



Week 3\_assignments.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text ▶ Run all ▼

⋮

🔍

< >

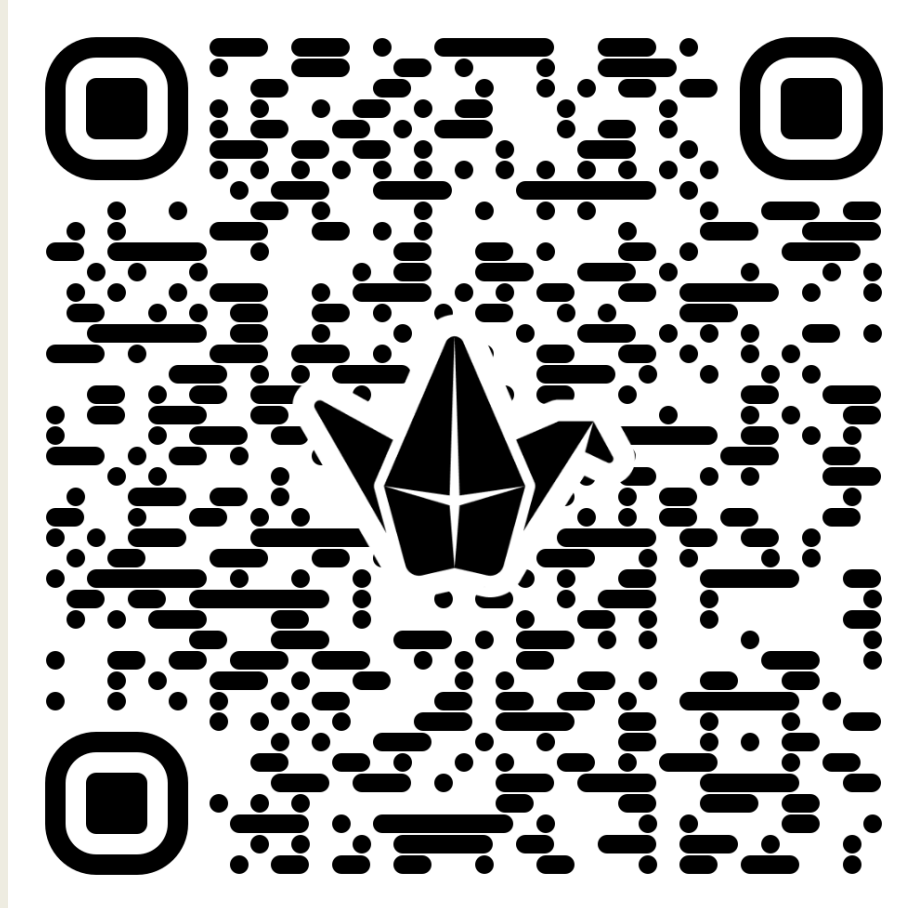
🔑

📁

**These assignments are ungraded and are NOT to be submitted. Feedback is NOT available for these assignments, but you are welcome to discuss with a TA, or with me during one of the office hours, if you have any questions.**

- **Optional Assignment 1 - Compare classification performance for multiple algorithms**  
↳ 1 cell hidden
- **Optional Assignment 2 - Implement own Bayes classifiers from scratch**  
↳ 1 cell hidden

# Any questions???



**scan the QR code to ask questions**

# Logistic regression

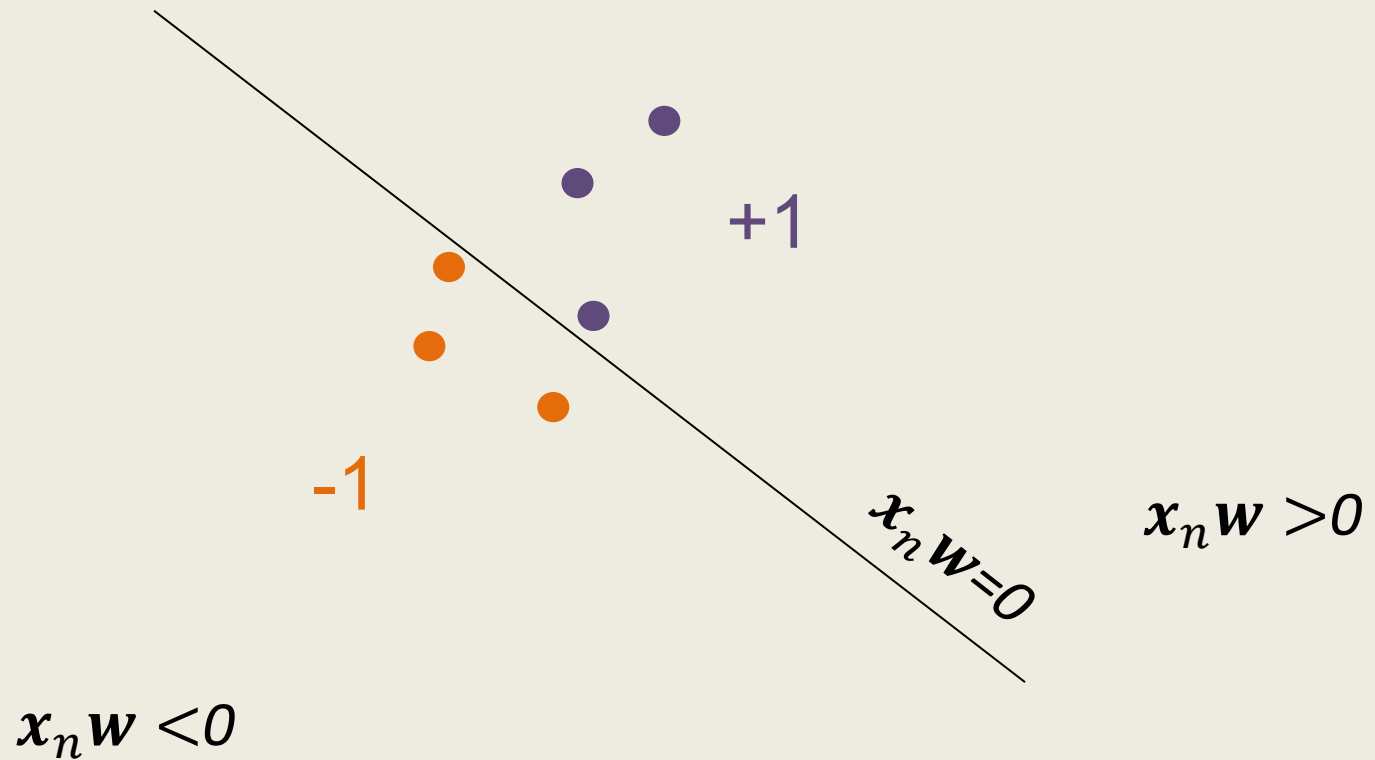
☐ Bayes classifier

☐ **Logistic regression**





# Re: Linear model (Classification)



The decision boundary is  $x_n w = 0$

# Re: Classification inference formulations

- Formulation so far
  - Given input data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $\hat{y}_i$ ?  $\hat{y}_i \subset \{c_k, \forall k, k = 1, 2, \dots K\}$
- A different (but equivalent) formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D}), \forall k, k = 1, 2, \dots K$ ?
  - $\hat{y}_i = c_k, \text{ for } \underset{k}{\operatorname{argmax}} p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D})$

# Decision boundaries (Binary classification)

- One formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $\hat{y}_i$ ?  $\hat{y}_i \in \{-1, 1\}$
  - $\mathbf{x}_i \mathbf{w} = 0$  is the decision boundary
- A different (but equivalent) formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D}), \forall k, k = 1, 2$ ?
  - $\hat{y}_i = c_k$ , for  $\operatorname{argmax}_k p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D})$
  - $\frac{p(Y=1|x_i)}{p(Y=-1|x_i)} = 1$  is the decision boundary

# Logistic regression (LR)

$$\frac{P(Y = 1 \mid \mathbf{x})}{P(Y = 0 \mid \mathbf{x})} = 1$$

$$\Rightarrow P(Y = 1 \mid \mathbf{x}) = \frac{1}{1 + e^{-xw}}$$

**the logistic regression (LR) model**

*see math proof on last slide pages*

# Decision boundaries (Binary classification)

- One formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $\hat{y}_i$ ?  $\hat{y}_i \in \{-1, 1\}$
  - $\mathbf{x}_i \mathbf{w} = 0$  is the decision boundary
- A different (but equivalent) formulation
  - Given data  $\{x_{i_1}, x_{i_2}, \dots x_{i_D}\}, \forall d, d = 1, 2, \dots D$
  - What is  $p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D}), \forall k, k = 1, 2$ ?
  - $\hat{y}_i = c_k$ , for  $\operatorname{argmax}_k p(c_k | x_{i_1}, x_{i_2}, \dots x_{i_D})$
  - $\frac{p(Y=1|x_i)}{p(Y=-1|x_i)} = 1$  is the decision boundary

# Re: Basic linear model for classification

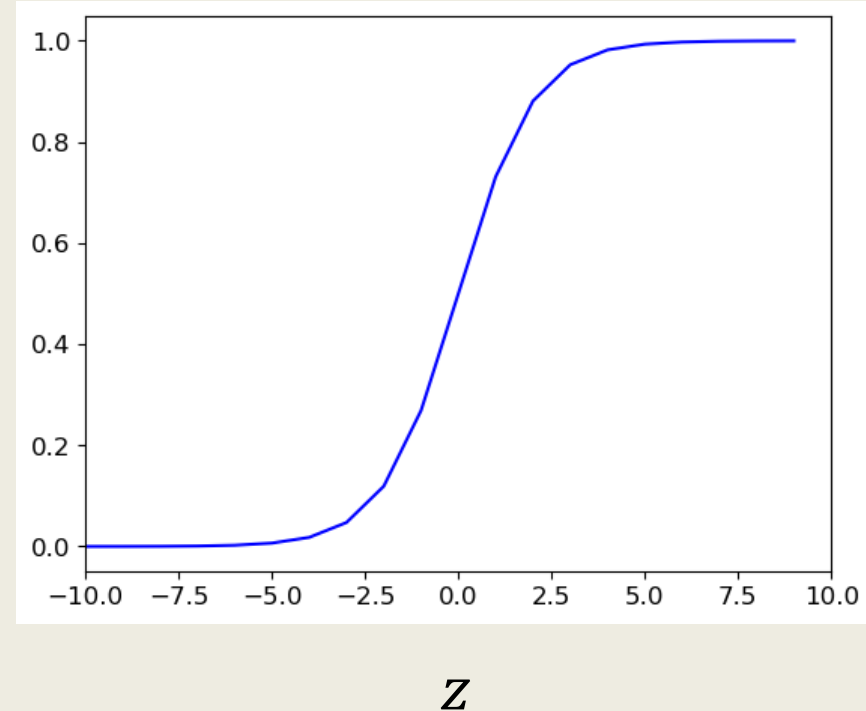
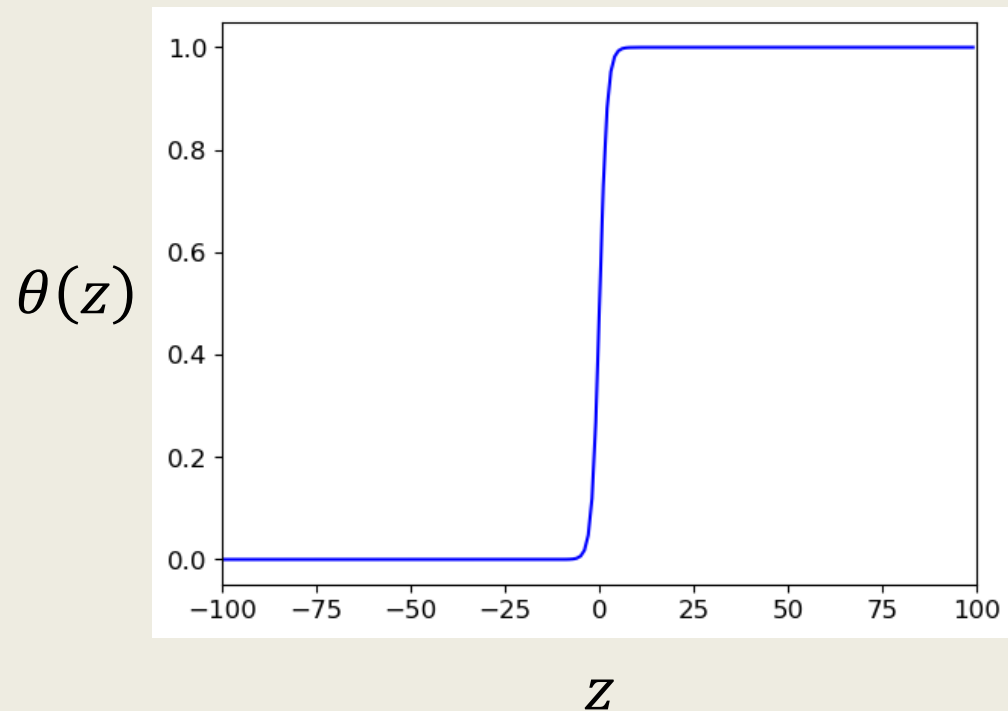
$$f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w} + b) = \hat{\mathbf{y}}$$

where

- $f(\cdot)$  – basic linear model
- $\mathbf{x}$  – features (or model input)
- $\hat{\mathbf{y}}$  – predicted labels/targets (or model output)
- $\mathbf{w}, b$  – weights, bias (or model parameters)
- $\sigma(\cdot)$  – activation function (for discretizing real values)

# (Logistic) sigmoid activation function

$$f(x) = \theta(xw + b) = \frac{e^{xw+b}}{e^{xw+b} + 1} = \frac{1}{1 + e^{-(xw+b)}}$$



# Optimal LR parameters (1)

- Likelihood based on a Bernoulli probability density function is **cross entropy loss**  $L_{CE}$

$$L_{CE} = \sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x}_n))$$

*see math proof on last slide pages*

- Model optimization involves maximizing likelihood (equivalent to finding the weights that minimize the loss)

$$\Rightarrow \sum_{n=1}^N \mathbf{x}_n (f(\mathbf{x}_n) - y_n) = 0$$

*see math proof on last slide pages*



# Optimal LR parameters (2)

- There is no closed form solution for

$$\sum_{n=1}^N \mathbf{x}_n (f(\mathbf{x}_n) - y_n) = 0$$

i.e. it cannot be solved analytically (to find the optimal  $\mathbf{w}$ )

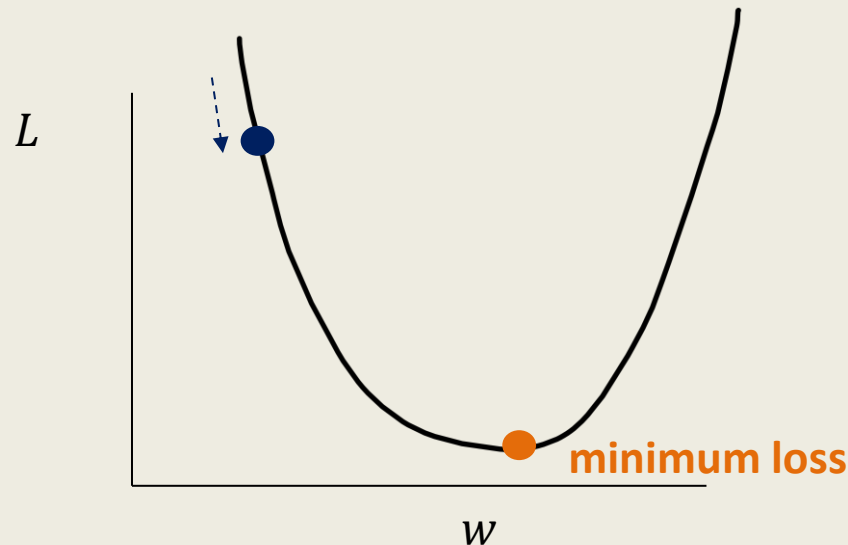
- The alternative is iterative solution using gradient descent algorithm

# Gradient descent algorithm

- The optimal model parameters are at

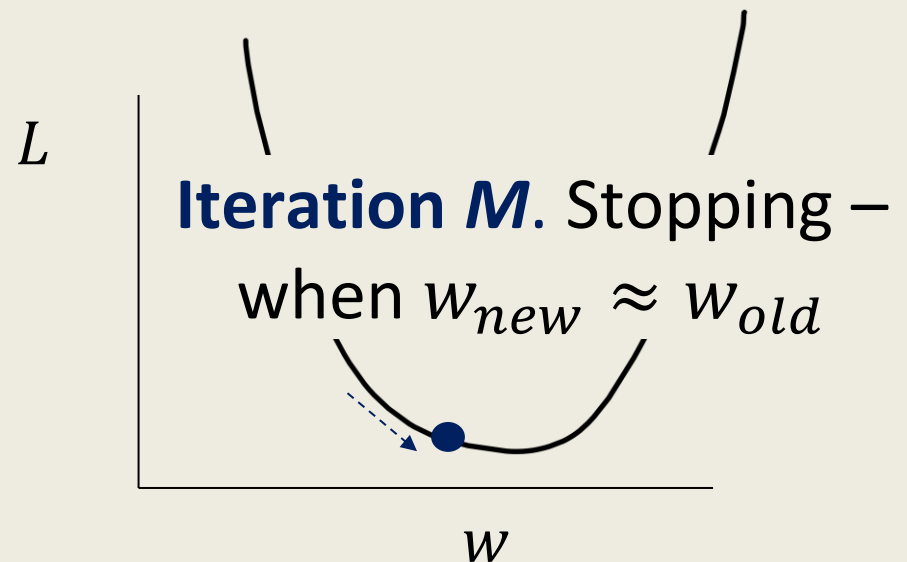
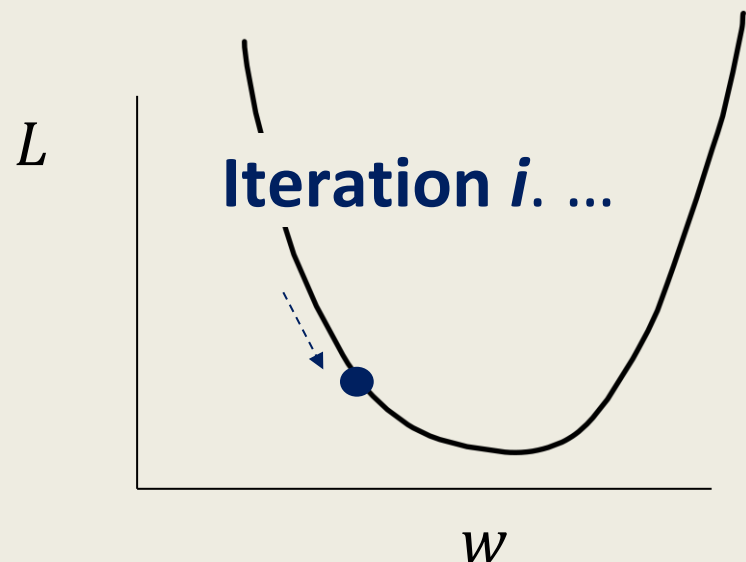
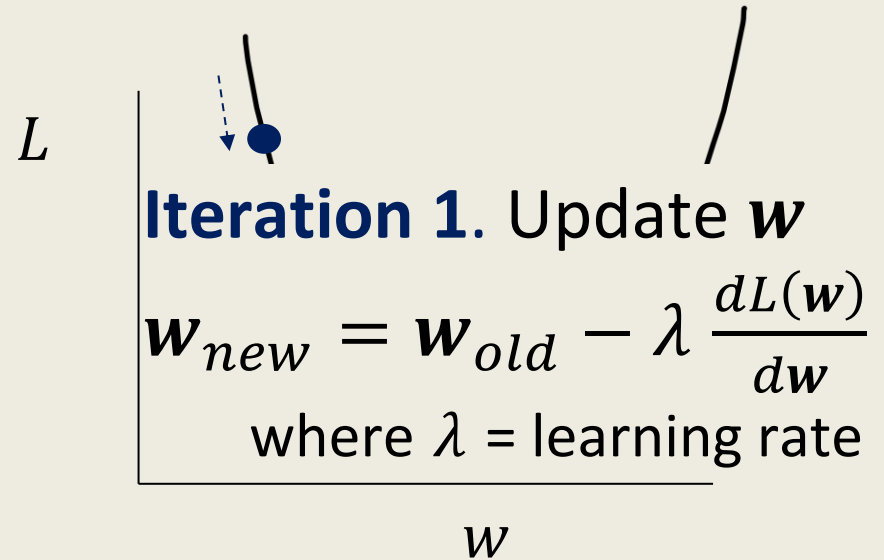
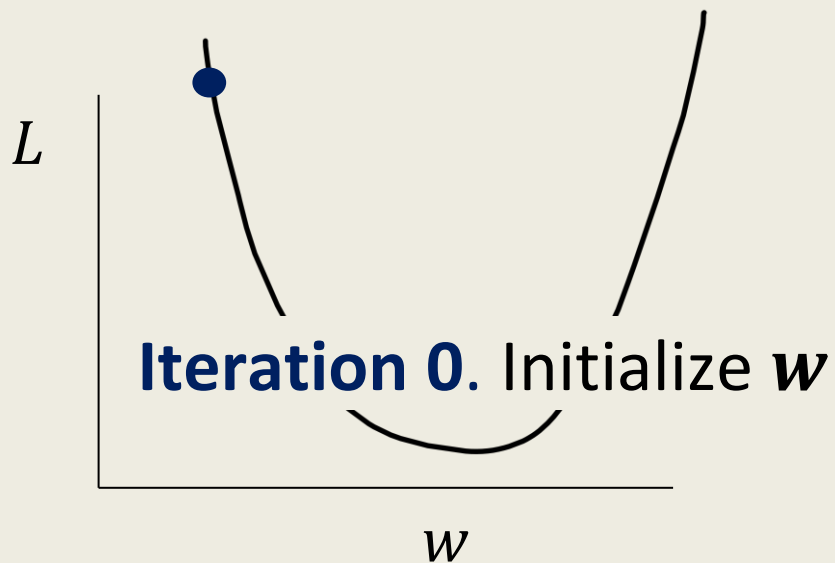
$$\frac{\partial L}{\partial \mathbf{w}} = 0$$

- i.e. in the direction of negative gradient (derivative)  
aka 'gradient descent'



1D loss landscape for a simplified illustration

# Descending iteratively




# Addressing LR overfitting

$$L_{LR} =$$

$$\sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x})) + \alpha \frac{\|w\|^2}{2}$$

cross-entropy  
loss



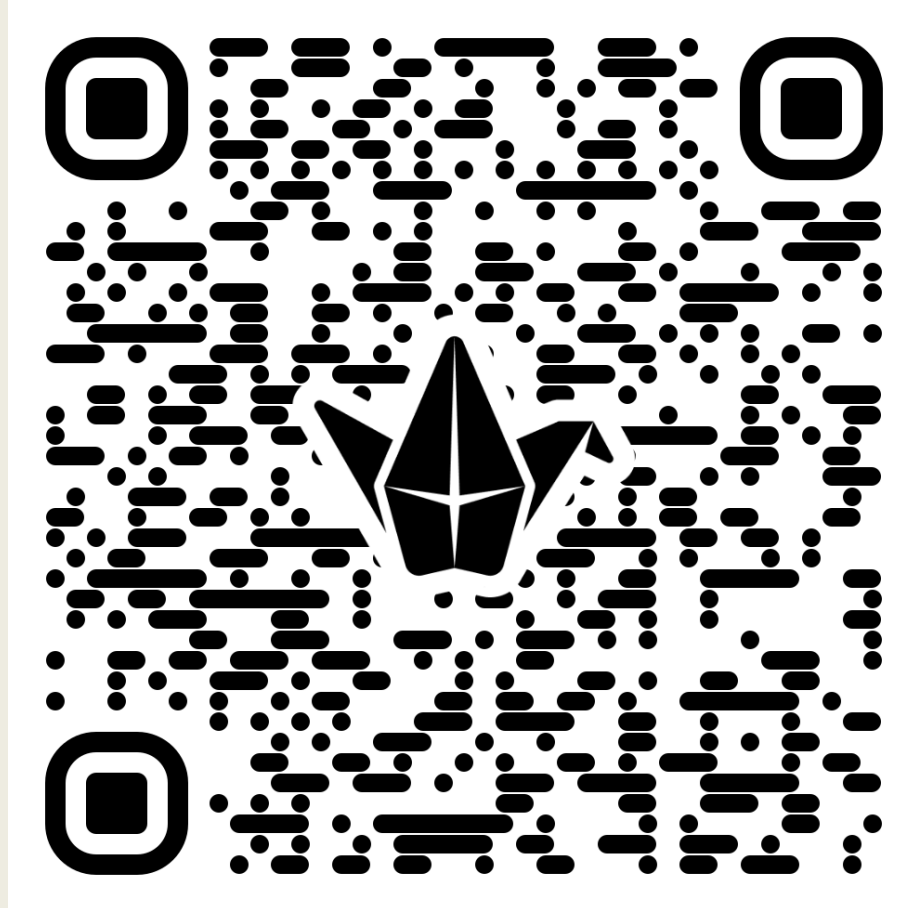
regularization  
term



# Summary

1. The **Bayes classifier** gives the probability of each class by updating the prior probability of the class with the likelihood of the observed data.
2. (Naïve) assumption of conditional independence in the Bayes model results in the **Naïve Bayes classifier**.
3. A **logistic regression** model is a classifier.
4. Its loss function is **cross-entropy loss**, and its parameters of classifier are optimized using **gradient descent algorithm**.

# Any questions???



**scan the QR code to ask questions**

# Math details and proofs

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 0)

$$P(c_k|\mathbf{X}) = \frac{P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) P(c_k)}{P(\mathbf{X})} \forall k$$

$P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k)$  can be estimated from the training data assuming a given distribution, e.g. Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})$

$$\Rightarrow P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) = \mathbb{E}[P(\mathbf{X} | \boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})]$$

where  $\boldsymbol{\mu}_{c_k}$  =  $D$ -dimensional mean vector

$\boldsymbol{\Sigma}_{c_k}$  =  $D \times D$  covariance matrix

$\mathbb{E}[\cdot]$  = expected value



# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 1)

- Assuming a Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})$ , the probability distribution parameters are  $\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}$
- The optimal distribution parameters  $\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}$  maximize the likelihood of the training data, equivalent to  $E[P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})]$
- Recall that optimum value of parameter  $\theta$  for a function  $L$  is at

$$\frac{\partial L}{\partial \theta} = 0$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 2)

The likelihood is Gaussian probability density function

$$P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})}$$

*taking log of both sides*

$$\begin{aligned} & \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) \\ &= \log \left( \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})} \right) \end{aligned}$$

*log of products is equivalent to sum of logs*

$$\begin{aligned} & \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) \\ &= \log \left( \frac{1}{(2\pi)^{\frac{D}{2}}} \right) + \log \left( \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{\frac{1}{2}}} \right) + \log \left( e^{-\frac{(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})}{2}} \right) \end{aligned}$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 3)

$$\begin{aligned} & \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) \\ = & \log \left( \frac{1}{(2\pi)^{\frac{D}{2}}} \right) + \log \left( \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{\frac{1}{2}}} \right) + \log \left( e^{-\frac{(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})}{2}} \right) \end{aligned}$$

*applying log rules*

$$\begin{aligned} & \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) \\ = & -N_{c_k} \frac{D}{2} \log(2\pi) - N_{c_k} \frac{1}{2} \log |\boldsymbol{\Sigma}_{c_k}| \\ & - \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \end{aligned}$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 4)

$$\Rightarrow \frac{\partial}{\partial \boldsymbol{\mu}} \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) = 0$$

*substituting for  $P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})$*

$$\Rightarrow \frac{\partial}{\partial \boldsymbol{\mu}} \left( \frac{-N_{c_k} D}{2} \log(2\pi) - \frac{N_{c_k}}{2} \log |\boldsymbol{\Sigma}_{c_k}| \right. \\ \left. - \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$

*removing terms without  $\boldsymbol{\mu}$  based on differentiation rules*

$$\Rightarrow \frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 5)

$$\frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$

*expanding the product*

$$\frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} \sum_{n=1}^{N_{c_k}} \mathbf{x}_n^T \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n - 2\boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n + \boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \boldsymbol{\mu}_{c_k} \right) = 0$$

*further removing terms without  $\boldsymbol{\mu}$  based on differentiation rules*

$$\frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} \sum_{n=1}^{N_{c_k}} -2\boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n + \boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \boldsymbol{\mu}_{c_k} \right) = 0$$

# Bayes classifier – $P(X|c_k)$ (MATH 6)

$$\frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{1}{2} \sum_{n=1}^{N_{c_k}} -2\boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n + \boldsymbol{\mu}_{c_k}^T \boldsymbol{\Sigma}_{c_k}^{-1} \boldsymbol{\mu}_{c_k} \right) = 0$$

*applying differentiation*

$$\sum_{n=1}^{N_{c_k}} \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n - \boldsymbol{\Sigma}_{c_k}^{-1} \boldsymbol{\mu}_{c_k} = 0$$

*collecting terms with  $\boldsymbol{\mu}_{c_k}$  to the right hand side*

$$\sum_{n=1}^{N_{c_k}} \boldsymbol{\Sigma}_{c_k}^{-1} \mathbf{x}_n = N_{c_k} \boldsymbol{\Sigma}_{c_k}^{-1} \boldsymbol{\mu}_{c_k}$$

# Bayes classifier – $P(X|c_k)$ (MATH 7)

$$\sum_{n=1}^{N_{c_k}} \Sigma_{c_k}^{-1} \mathbf{x}_n = N_{c_k} \Sigma_{c_k}^{-1} \boldsymbol{\mu}_{c_k}$$

*dividing through by  $\Sigma_{c_k}^{-1}$*

$$\sum_{n=1}^{N_{c_k}} \mathbf{x}_n = N_{c_k} \boldsymbol{\mu}_{c_k}$$

*making  $\boldsymbol{\mu}_{c_k}$  the subject of the equation*

$$\Rightarrow \boldsymbol{\mu}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} \mathbf{x}_n$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 8)

$$\Rightarrow \frac{\partial}{\partial \boldsymbol{\Sigma}} \log P(\mathbf{X}|\boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) = 0$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \left( \frac{-N_{c_k} D}{2} \log(2\pi) - \frac{N_{c_k}}{2} \log |\boldsymbol{\Sigma}_{c_k}| \right. \\ \left. - \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \left( -\frac{N_{c_k}}{2} \log |\boldsymbol{\Sigma}_{c_k}| - \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$



# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 9)

$$\frac{\partial}{\partial \mathbf{\Sigma}} \left( -\frac{N_{c_k}}{2} \log |\mathbf{\Sigma}_{c_k}| \right) + \frac{\partial}{\partial \mathbf{\Sigma}} \left( \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \mathbf{\Sigma}_{c_k}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) \right) = 0$$

$$-\frac{N_{c_k}}{2} \mathbf{\Sigma}_{c_k}^{-1} + \frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \mathbf{\Sigma}_{c_k}^{-2} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k}) = 0$$

$$-\frac{N_{c_k}}{2} \mathbf{\Sigma}_{c_k}^{-1} = -\frac{1}{2} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T \mathbf{\Sigma}_{c_k}^{-2} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})$$

$$\Rightarrow \mathbf{\Sigma}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})$$

# Bayes classifier – $P(\mathbf{X}|c_k)$ (MATH 10)

$$P(x_{j_1}, x_{j_2}, \dots, x_{j_D} | c_k) = P(\mathbf{X} | \boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k})$$

where

$$\begin{aligned} P(\mathbf{X} | \boldsymbol{\mu}_{c_k}, \boldsymbol{\Sigma}_{c_k}) \\ = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_{c_k}|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})^T \boldsymbol{\Sigma}_{c_k}^{-1} (\mathbf{X}_{c_k} - \boldsymbol{\mu}_{c_k})} \end{aligned}$$

$$\boldsymbol{\mu}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_{c_k} = \frac{1}{N_{c_k}} \sum_{n=1}^{N_{c_k}} (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})^T (\mathbf{x}_n - \boldsymbol{\mu}_{c_k})$$

# Logistic regression (LR) – MATH PROOF (1)

$$\frac{P(Y = 1 \mid \mathbf{x})}{P(Y = 0 \mid \mathbf{x})} = 1$$

*taking log of both sides*

$$\log \left( \frac{P(Y = 1 \mid \mathbf{x})}{P(Y = 0 \mid \mathbf{x})} \right) = 0$$

*substituting for  $\mathbf{xw} = 0$  on the right hand side*

$$\log \left( \frac{P(Y = 1 \mid \mathbf{x})}{P(Y = 0 \mid \mathbf{x})} \right) = \mathbf{xw}$$

*applying probability rule*

$$\log \left( \frac{P(Y = 1 \mid \mathbf{x})}{1 - P(Y = 1 \mid \mathbf{x})} \right) = \mathbf{xw}$$

# Logistic regression (LR) – MATH PROOF (2)

$$\log \left( \frac{P(Y = 1 \mid \mathbf{x})}{1 - P(Y = 1 \mid \mathbf{x})} \right) = \mathbf{x}\mathbf{w}$$

*taking exponent function of both sides*

$$\frac{P(Y = 1 \mid \mathbf{x})}{1 - P(Y = 1 \mid \mathbf{x})} = e^{\mathbf{x}\mathbf{w}}$$

*cross-multiplying*

$$P(Y = 1 \mid \mathbf{x}) = e^{\mathbf{x}\mathbf{w}} - e^{\mathbf{x}\mathbf{w}}(P(Y = 1 \mid \mathbf{x}))$$

*making  $P(Y = 1 \mid \mathbf{x})$  the subject of the formula*

$$P(Y = 1 \mid \mathbf{x}) = \frac{e^{\mathbf{x}\mathbf{w}}}{1 + e^{\mathbf{x}\mathbf{w}}}$$

# Logistic regression (LR) – MATH PROOF (3)

$$P(Y = 1 \mid \mathbf{x}) = \frac{e^{xw}}{1 + e^{xw}}$$

*dividing numerator and denominator by  $e^{xw}$*

$$P(Y = 1 \mid \mathbf{x}) = \frac{1}{1 + e^{-xw}}$$

**the logistic regression (LR) model**

# Optimal LR parameters (1) – PROOF (a)

- Consider binary classification and so, assuming Bernoulli distribution  $Bern(\theta)$
- Likelihood based on a Bernoulli probability density function is

$$\begin{aligned} P(Y|\mathbf{X}, \mathbf{w}) &= \prod_{n=1}^N P(Y|\mathbf{x}_n)^{y_n} (1 - P(Y|\mathbf{x}_n))^{1-y_n} \\ &= \prod_{n=1}^N f(\mathbf{x}_n)^{y_n} (1 - f(\mathbf{x}_n))^{1-y_n} \end{aligned}$$

- Optimization involves maximizing likelihood  $P(Y|\mathbf{X}, \mathbf{w})$  equivalent to minimizing negative log likelihood

# Optimal LR parameters (1) – PROOF (b)

$$\frac{\partial}{\partial \mathbf{w}} (-\log(P(Y|\mathbf{X}, \mathbf{w}))) = 0$$

*substituting for the formula of the likelihood*

$$\frac{\partial}{\partial \mathbf{w}} \left( -\log \left( \prod_{n=1}^N f(\mathbf{x}_n)^{y_n} (1 - f(\mathbf{x}_n))^{1-y_n} \right) \right) = 0$$

*applying logarithm rules*

$$\frac{\partial}{\partial \mathbf{w}} \left( \sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x}_n)) \right) = 0$$

# Optimal LR parameters (1) – PROOF (c)

$$\frac{\partial}{\partial \mathbf{w}} \left( \sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x}_n)) \right)$$

where  $f(\mathbf{x}_n, \mathbf{w}) = \frac{1}{1 + e^{-\theta}}$  and  $\theta = \mathbf{x}_n \mathbf{w}$

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial L}{\partial f(\theta)} \cdot \frac{\partial f(\theta)}{\partial \theta} \cdot \frac{\partial \theta}{\partial \mathbf{w}}$$

where  $L = \sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x}_n))$

**cross entropy loss**





# Optimal LR parameters (1) – PROOF (d)

$$\frac{\partial}{\partial \mathbf{w}} \left( \sum_{n=1}^N -y_n \log f(\mathbf{x}_n) - (1 - y_n) \log(1 - f(\mathbf{x}_n)) \right) = 0$$

*applying rule of derivative of a function of a function*

$$\sum_{n=1}^N - \left( \frac{y_n}{f(\mathbf{x}_n)} - \frac{1 - y_n}{1 - f(\mathbf{x}_n)} \right) \left( \frac{e^{-\mathbf{x}_n \mathbf{w}}}{(1 + e^{-\mathbf{x}_n \mathbf{w}})^2} \right) \mathbf{x}_n = 0$$


$$\frac{\partial L}{\partial f(\boldsymbol{\theta})}$$


$$\frac{\partial f(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$


$$\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{w}}$$

# Optimal LR parameters (1) – PROOF (e)

$$\sum_{n=1}^N - \left( \frac{y_n}{f(\mathbf{x}_n)} - \frac{1 - y_n}{1 - f(\mathbf{x}_n)} \right) \left( \frac{e^{-\mathbf{x}_n \mathbf{w}}}{(1 + e^{-\mathbf{x}_n \mathbf{w}})^2} \right) \mathbf{x}_n = 0$$

*expanding each of the two terms and substituting for  $f(\mathbf{x}_n) = \frac{1}{1+e^{-\mathbf{x}_n \mathbf{w}}}$*

$$\sum_{n=1}^N - \left( \frac{y_n - f(\mathbf{x}_n)}{f(\mathbf{x}_n)(1 - f(\mathbf{x}_n))} \right) \left( \mathbf{x}_n f(\mathbf{x}_n)(1 - f(\mathbf{x}_n)) \right) = 0$$

*cancelling terms common to both numerator and denominator*

$$\sum_{n=1}^N - \mathbf{x}_n (y_n - f(\mathbf{x}_n)) = 0$$

*dividing through by -1*

$$\sum_{n=1}^N \mathbf{x}_n (f(\mathbf{x}_n) - y_n) = 0$$