

Logistic regression

This method focuses on $k=2$ assuming that **log odds** of class 1 vs. class 0 is linear in x :

$$\log \frac{P(Z=1 | X=x)}{P(Z=2 | X=x)} = \beta_0 + \beta_1^T x \text{ with } \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^p$$

- Note that LDA also has log odds linear in x with $\hat{\beta}_0, \hat{\beta}_1$ depend on $\hat{M}_1, \hat{M}_2, \hat{\Sigma}$.

$$\begin{aligned} \log \frac{P(Z=1 | X=x)}{P(Z=2 | X=x)} &= \log \frac{\pi_1 \cdot f(x; \mu_1, \Sigma)}{\pi_2 \cdot f(x; \mu_2, \Sigma)} \\ &= \log \frac{\frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)}}{\frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)}} + \log \frac{\pi_1}{\pi_2} = \\ &= -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1) + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2) + \log \frac{\pi_1}{\pi_2} = \\ &= (\mu_1 - \mu_2)^T \Sigma^{-1} x - \mu_1^T \Sigma^{-1} \mu_1 + \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{\pi_1}{\pi_2} \end{aligned}$$

To simplify the derivations we will relabel class 2 into 0, i.e. $y_i \in \{0, 1\}$, and include the intercept in features, i.e. $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \in \mathbb{R}^{p+1}$ and $X \in \mathbb{R}^{n \times (p+1)}$

Logistic regression directly estimates β .

$$\bullet P(Z=1|X=x) = \frac{1}{1+e^{-\beta^T x}}$$

Denote by $\pi(x) = P(Z=1|X=x) \Rightarrow 1-\pi(x) = P(Z=0|X=x)$

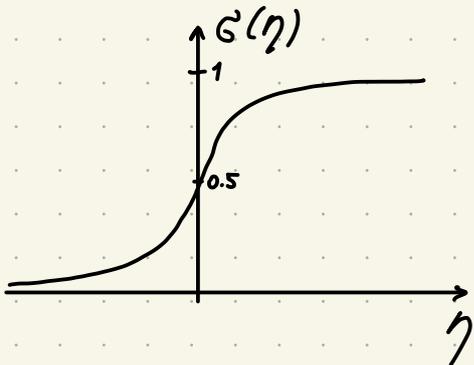
$$\log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \beta^T x \Rightarrow \frac{1-\pi(x)}{\pi(x)} = \frac{1}{\pi(x)} - 1 = e^{-\beta^T x} \Rightarrow$$

$$\pi(x) = \frac{1}{1+e^{-\beta^T x}}$$

$$P(Z=1|X=x) = \sigma(\beta^T x) \quad \text{where} \quad \sigma(\eta) = \frac{1}{1+e^{-\eta}} = \frac{e^\eta}{1+e^\eta}$$

This function is called sigmoid function.

Properties: $\sigma(\eta) = \frac{1}{1+e^{-\eta}}$



• $\sigma(\eta) \geq 0.5$ for $\eta \geq 0$

$\sigma(\eta) \leq 0.5$ for $\eta \leq 0$

$P(Z=1 | X=x) \geq 0.5$ if $\beta_0 + \beta^T x \geq 0$

$P(Z=2 | X=x) \leq 0.5$ if $\beta_0 + \beta^T x \leq 0$

Decision boundary $\beta_0 + \beta^T x = 0$.

• $\sigma'(\eta) = \sigma(\eta)(1-\sigma(\eta))$

$$\left| \sigma'(\eta) = \frac{e^{-\eta}}{(1+e^{-\eta})^2} = \frac{1}{1+e^{-\eta}} \cdot \frac{e^{-\eta}}{1+e^{-\eta}} = \sigma(\eta) \cdot (1-\sigma(\eta)) \right.$$

• $(\log \sigma(\eta))' = 1 - \sigma(\eta)$, $(\log(1-\sigma(\eta)))' = -\sigma(\eta)$

$$\left| (\log \sigma(\eta))' = \frac{\sigma'(\eta)}{\sigma(\eta)} = \frac{\cancel{\sigma(\eta)}(1-\sigma(\eta))}{\cancel{\sigma(\eta)}} \right.$$

Logistic regression finds β by maximizing log-likelihood:

$$e(x, y; \beta) = \sum_{i=1}^n \log P(z = y_i | X = x_i)$$

$$\begin{aligned} e(x, y; \beta) &= \sum_{i: y_i=1} \log \sigma(\beta^T x_i) + \sum_{i: y_i=0} \log (1 - \sigma(\beta^T x_i)) = \\ &= \sum_{i=1}^n y_i \log \sigma(\beta^T x_i) + (1 - y_i) \log (1 - \sigma(\beta^T x_i)) \end{aligned}$$

Maximizing log-likelihood has no explicit solution.

We need to apply one of the optimization algorithms.

Detour: Newton's method

Given a function $f: \mathbb{R}^p \rightarrow \mathbb{R}$, solve

minimize $f(x)$

Let $x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ and $\nabla f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_p} \end{pmatrix} \in \mathbb{R}^p$ be the gradient

Let $\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_p} \\ \vdots & \vdots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_p} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_p} \dots & \frac{\partial^2 f(x)}{\partial x_p^2} \end{pmatrix} \in \mathbb{R}^{p \times p}$ be the Hessian.

Gradient descent update: $x^{(t+1)} = x^{(t)} - \nabla f(x^{(t)})$

Newton's method update: $x^{(t+1)} = x^{(t)} - (\nabla^2 f(x^{(t)}))^{-1} \nabla f(x^{(t)})$

Motivation for Newton's update:

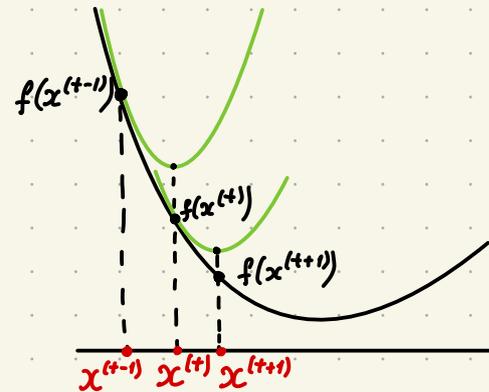
- Let's take the second order approximation of $f(x)$ at point $x \in \mathbb{R}^p$

$$f(y) \simeq f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x)$$

- value y that minimizes the approximation corresponds to Newton's update.

$$\nabla_y f(y) \simeq \nabla f(x) + \nabla^2 f(x) (y-x) = 0$$

$$\text{Then } y = x - (\nabla^2 f(x))^{-1} \nabla f(x)$$



Iteratively reweighted least squares

$$e(x, y; \beta) = \sum_{i=1}^n y_i \log \sigma(\beta^T x_i) + (1 - y_i) \log(1 - \sigma(\beta^T x_i))$$

$$\bullet \nabla_{\beta} e(x, y; \beta) = X^T (y - \pi) \quad \text{where} \quad \pi = \begin{pmatrix} P(z=1 | X=x_1) \\ \vdots \\ P(z=1 | X=x_n) \end{pmatrix}$$

Denote by $\pi_i = P(z=1 | X=x_i)$

$$\begin{aligned} \nabla_{\beta} &= \sum_{i=1}^n y_i (1 - \sigma(\beta^T x_i)) \cdot x_i - (1 - y_i) \sigma(\beta^T x_i) x_i = \\ &= \sum_{i=1}^n y_i x_i - \sigma(\beta^T x_i) x_i = \sum_{i=1}^n (y_i - \sigma(\beta^T x_i)) x_i = \\ &= \sum_{i=1}^n (y_i - \pi_i) x_i = X^T (y - \pi) \end{aligned}$$

- $\nabla_{\beta} \ell(x, y; \beta) = -X^T W X$ where $W = \text{diag}(\pi * (1 - \pi))$

$$\begin{aligned} \nabla_{\beta}^2 &= -\sum_{i=1}^n \sigma(\beta^T x_i) (1 - \sigma(\beta^T x_i)) x_i x_i^T = \\ &= -\sum_{i=1}^n \pi_i (1 - \pi_i) x_i x_i^T = -X^T W X \end{aligned}$$

- At iteration t , we obtain $\beta^{(t)}, \pi^{(t)}$ and $W^{(t)}$ then Newton's update implies

$$\beta^{(t+1)} = (X^T W^{(t)} X)^{-1} X^T W^{(t)} z^{(t)} \text{ for some } z^{(t)}$$

$$\begin{aligned} \beta^{(t+1)} &= \beta^{(t)} + (X^T W^{(t)} X)^{-1} X^T (y - \pi^{(t)}) = \\ &= (X^T W^{(t)} X)^{-1} X^T (W^{(t)} X \beta^{(t)} + y - \pi^{(t)}) = \\ &= (X^T W^{(t)} X)^{-1} X^T W^{(t)} \underbrace{(X \beta^{(t)} + (W^{(t)})^{-1} (y - \pi^{(t)}))}_{z^{(t)}} = \\ &= (X^T W^{(t)} X)^{-1} X^T W^{(t)} z^{(t)} \end{aligned}$$

This is regression $z^{(t)} \sim X$ with weights $W^{(t)}$.

Logistic regression via IRLS:

Input: X, y , initialization $\beta^{(0)}$

At iteration $t=1, 2, \dots$

Step 1 Compute $z^{(t)} = X\beta^{(t)} + (W^{(t)})^{-1}(y - \pi^{(t)})$

$$W^{(t)} = \text{diag}(\pi^{(t)} * (1 - \pi^{(t)}))$$

$$\text{where } \pi^{(t)} \approx \begin{pmatrix} P(z=1 | x=x_1) \\ \vdots \\ P(z=1 | x=x_n) \end{pmatrix} = \begin{pmatrix} \sigma(\beta^{(t)\top} x_1) \\ \vdots \\ \sigma(\beta^{(t)\top} x_n) \end{pmatrix}$$

Step 2 Solve weighted regression problem

with response $z^{(t)}$, features X and weights $W^{(t)}$.

Output: coefficient vector β .

Extensions of logistic regression

① Multiple classes $k > 2$

$$P(Z=k | X=x) = \frac{e^{\beta_{0k} + \beta_{1k}^T x}}{\sum_{j=1}^k e^{\beta_{0j} + \beta_{1j}^T x}} \quad \text{for } k=1 \dots k$$

Softmax: $(\eta_1, \dots, \eta_k) \rightarrow \left(\frac{e^{\eta_1}}{\sum_{j=1}^k e^{\eta_j}}, \dots, \frac{e^{\eta_k}}{\sum_{j=1}^k e^{\eta_j}} \right)$

Fit $\{\beta_{0k}, \beta_{1k}\}_{k=1}^{k-1}$ by MLE.

Decision rule $h(x) = \operatorname{argmax}_{k=1 \dots k} P(Z=k | X=x)$

• Log-odds are linear functions in x

$$\log \frac{P(Z=k | X=x)}{P(Z=k | X=x)} = \log \left(\frac{e^{\beta_{0k} + \beta_{1k}^T x}}{e^{\beta_{0k} + \beta_{1k}^T x}} \right)$$

$$= \underbrace{(\beta_{0k} - \beta_{0k})}_{\tilde{\beta}_{0k}} + \underbrace{(\beta_{1k} - \beta_{1k})}_{\tilde{\beta}_{1k}}^T x = \tilde{\beta}_{0k} + \tilde{\beta}_{1k}^T x$$

② In high dimensions ($p > n$) add regularization:

minimize $-l(x, y; \beta) + \lambda \text{Pen}(\beta)$

Example: penalties

$\text{Pen}(\beta) = \|\beta\|_2^2$ "ridge"

$\text{Pen}(\beta) = \|\beta\|_1$ "lasso"

$\text{Pen}(\beta) = (1-d) \frac{\|\beta\|_2^2}{2} + d \|\beta\|_1$ "elastic net"

k-nearest neighbors (KNN)

This is a distribution-free and memory-based method.

Given data $X \in \mathbb{R}^{n \times p}$ labels $y \in \mathbb{R}^n$
classify a point $x \in \mathbb{R}^p$ by:

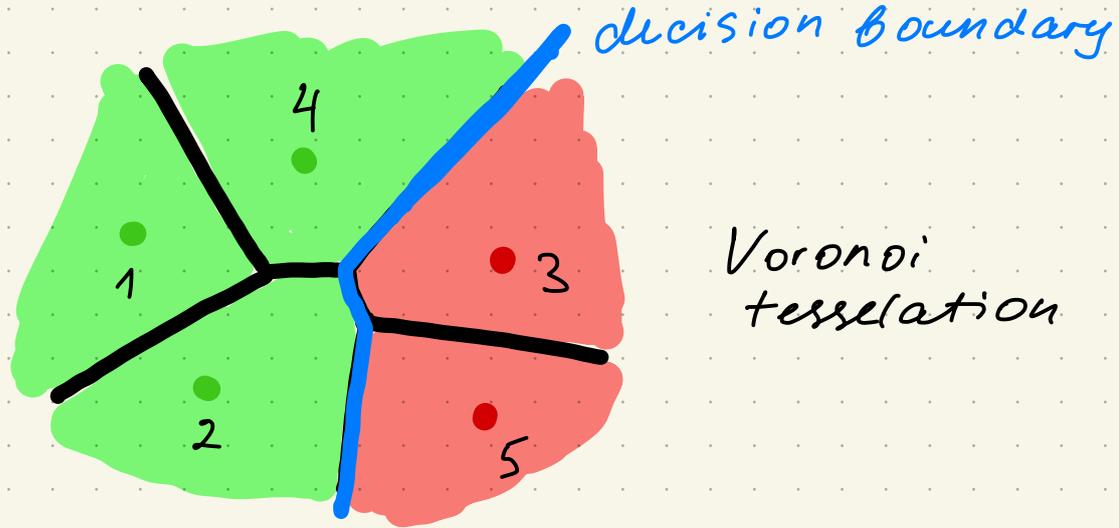
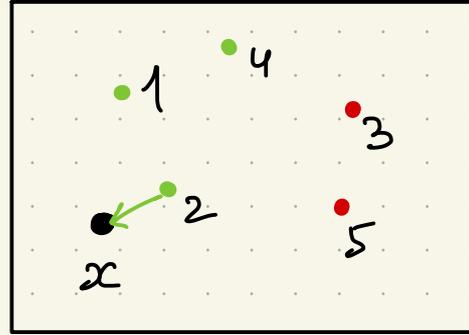
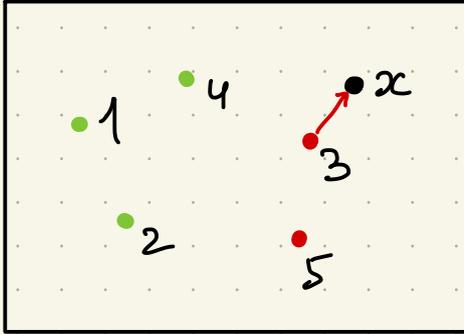
- find $N_k(x) \subseteq \{1, \dots, n\}$ the set of indices for k-nearest neighbors of x in the training data
- Classify x according to a majority vote

$$h(x) = \operatorname{argmax}_{k=1 \dots K}$$

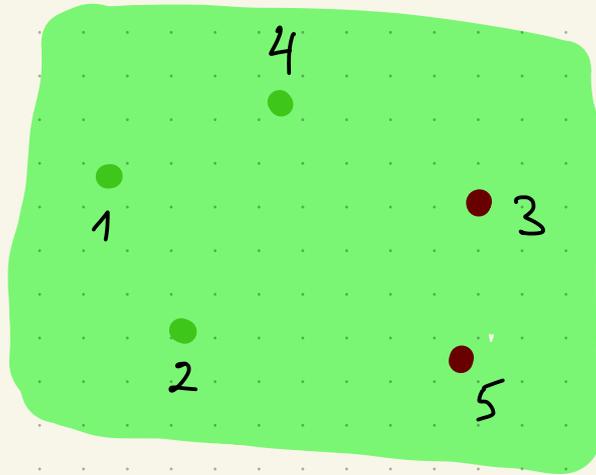
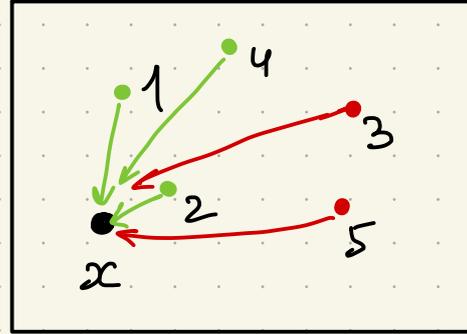
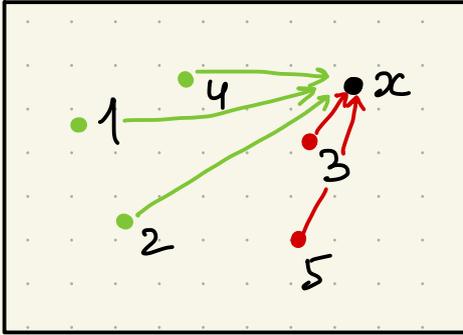
$$\frac{1}{K} \sum_{i \in N_k(x)} I(y_i = k)$$

" π_k "

Example: 1-NN, Euclidean distance



Example: n -NN, Euclidean distance



Dominant
class
prediction

1NN vs Bayes classifier

If we know $P(z|x)$, Bayes optimal classifier

$$\text{is } h(x) = \operatorname{argmax}_{k=1 \dots K} P(z=k | x=x) = \operatorname{argmax}_{k=1 \dots K} \pi_k(x) = k^*$$

Then the error for x is $\epsilon_{BO} = 1 - \pi_{k^*}(x)$

For 1NN denote by x_{NN} the nearest neighbor of x

Then the error for x is $\epsilon_{NN} = \sum_{k=1}^K \pi_k(x) (1 - \pi_k(x))$

If $n \rightarrow \infty$ then $\pi_k(x) = \pi_k(x_{NN})$ for $k=1 \dots K$

x is in class k with probability $\pi_k(x)$

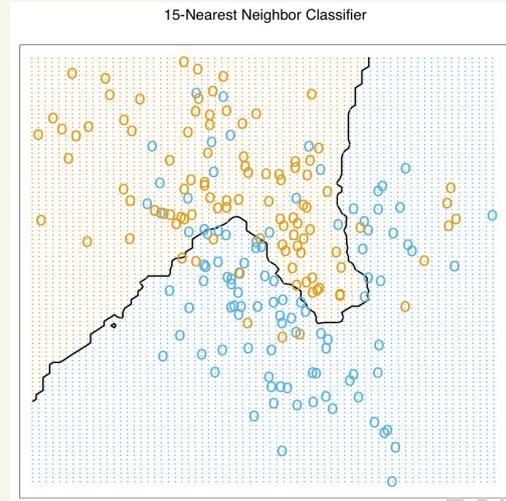
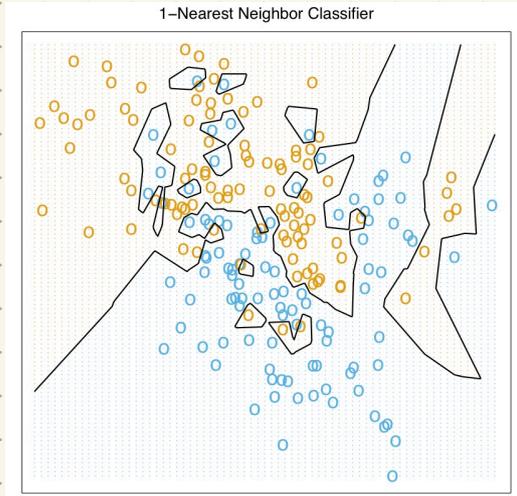
x_{NN} is not in class k with probability $1 - \pi_k(x)$

One can show that $\epsilon_{BO} \leq \epsilon_{NN} \leq 2 \epsilon_{BO}$

| $K=2$, $\epsilon_{NN} = 2 p_{k^*}(x) (1 - p_{k^*}(x))$

Practical aspects of KNN

① k has significant impact on the result

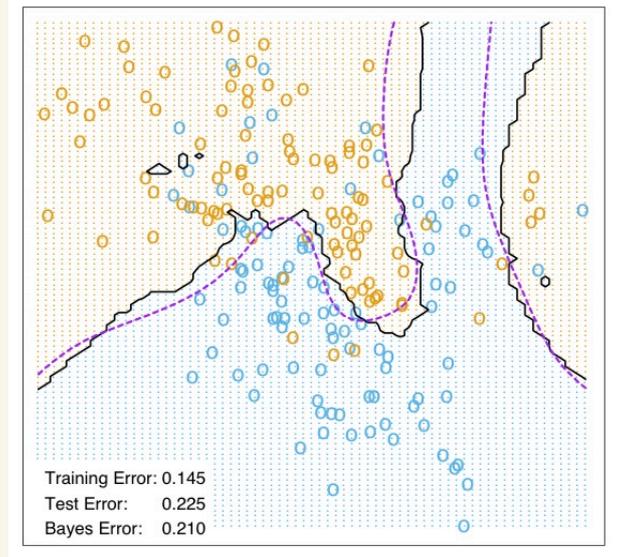
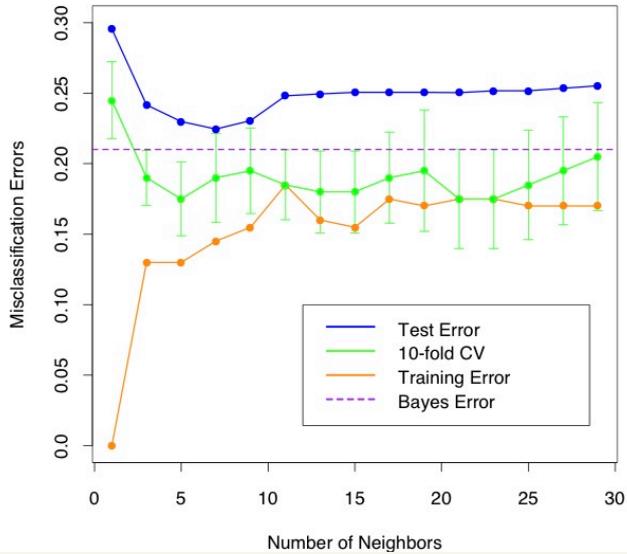


ESL11

Small k : capture local information,
may overfit

Large k : more stable, may underfit

Use cross-validation to select k .



② Scaling is important, standardize X .

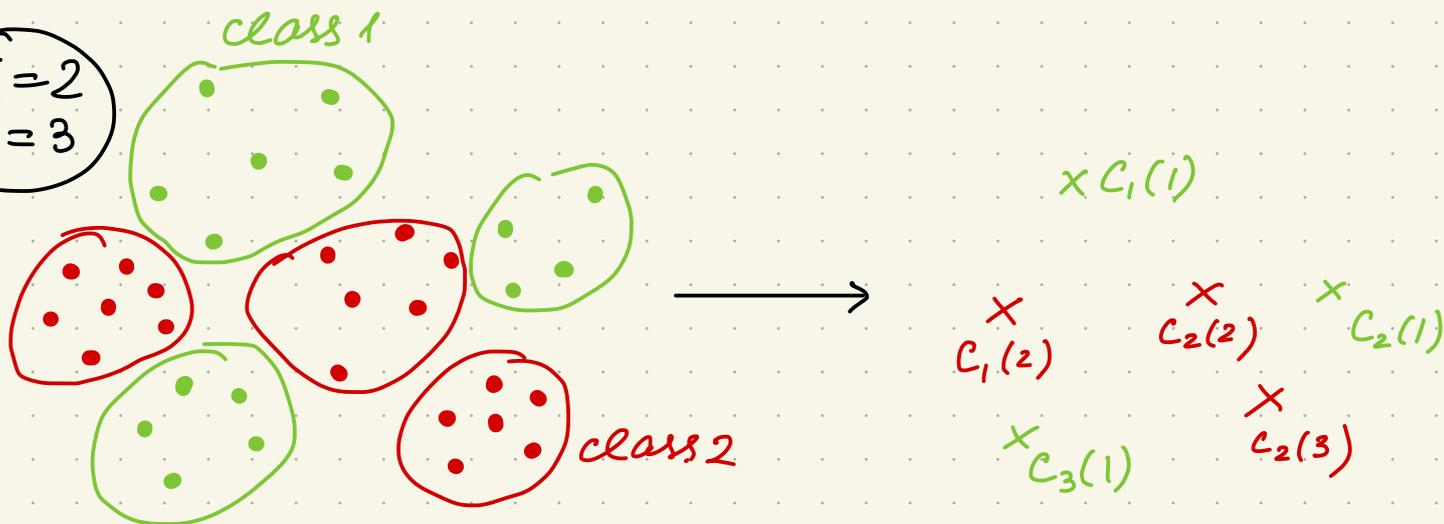
③ KNN requires storing X and Y

For class k , take $\{x_i : i \in C_k\}$ and perform R -Means clustering.

Denote the centroids by $C_1(k) \dots C_r(k)$

Replace $\{x_i : i \in C_k\} \rightarrow \{C_r(k) : r = 1 \dots R\}$

$K=2$
 $R=3$



④ KNN has equal vote weight

$$h(x) = \operatorname{argmax}_{k=1 \dots K} \sum_{i \in N_k(x)} \frac{1}{K} I(y_i = k)$$

Ties can be broken at random, alternatively:

$$h(x) = \operatorname{argmax}_{k=1 \dots K} \sum_{i \in N_k(x)} w(x, x_i) I(y_i = k)$$

e.g. $w(x, x_i) = \frac{1}{\|x_i - x\|^2}$

