

Lecture #3 - Classification (Ch 4)

$$Y = f(X) + \epsilon, E(\epsilon) = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{"the model"}$$

$\hookrightarrow Y \in \mathbb{R}$

In lecture #2, we "assumed" that $f(x) = \beta_0 + \beta_1 x$ or $f(x) = \beta_0 + \beta_1 \sin(x) + \beta_2 \cos(x)$

Today: (Ch 4)

Next week: (Ch 5, cross-validation; Ronan)

$Y \in \{A, B, O, AB\} \quad \left. \begin{array}{l} Y \text{ is qualitative.} \\ \text{If } X \text{ is qualitative} \end{array} \right\}$

How about: $y = \begin{cases} 1 & \text{if } A \\ 2 & \text{if } B \\ 3 & \text{if } O \\ 4 & \text{if } AB \end{cases}$

then $y = f(x) + \varepsilon!$

e.g. $y = \beta_0 + \beta_1 x + \varepsilon$

then a $\frac{1}{\beta_1}$ -unit increase in x
is assoc'd w/ blood type
 $A \Rightarrow \varepsilon$, or $B \Rightarrow 0$,
 or $O \Rightarrow AB$ - -

Switching to $y = \begin{cases} 1 & \text{if } B \\ 2 & \text{if } O \\ 3 & \text{if } AB \\ 4 & \text{if } A \end{cases}$

Note: if y is binary (e.g. default vs. non-default) then OK to code

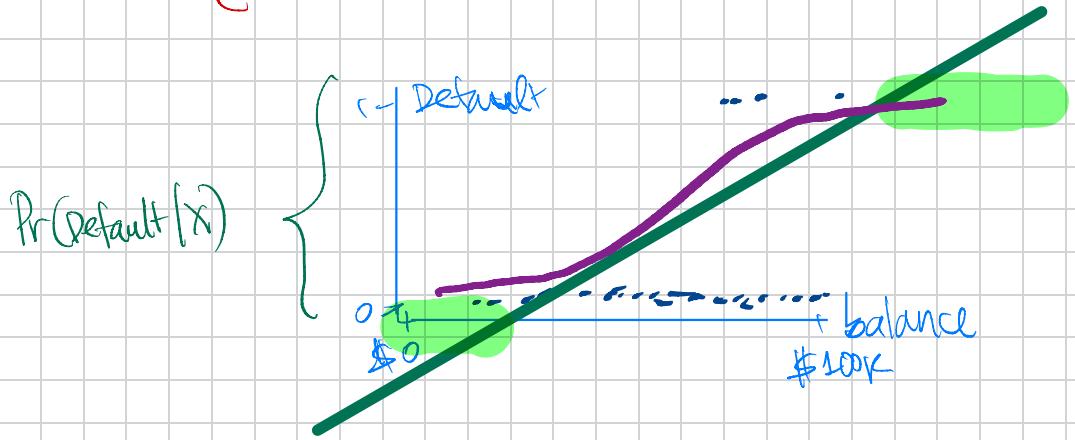
$$y = \begin{cases} 1 & \text{if default} \\ 0 & \text{if not default} \end{cases}$$

and do regression

if all you care about is prediction.

$$Y = \begin{cases} 1 & \text{if default} \\ 0 & \text{else} \end{cases}$$

$X = \text{balance}$



Logistic Regression.

$$Y = \begin{cases} 1 & \text{if default} \\ 0 & \text{if not} \end{cases}$$

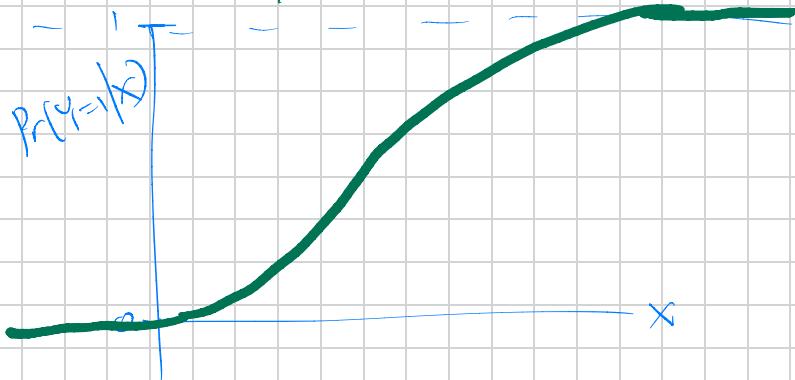
$X = \text{balance}$

Our goal: model

$$\Pr(Y=1|X=x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

$$\exp(a) = e^a = 2.71828^a$$

If $\beta_1 > 0$:



Define $p(x) := \Pr(Y=1|X)$

$$1 - p(x) = \Pr(Y=0|X)$$

$$= \frac{1}{1 + \exp(\beta_0 + \beta_1 x)}$$

$$\Pr(Y=1|X) = p(x) = \frac{\exp(f_0 + f_1 x)}{1 + \exp(f_0 + f_1 x)}$$

$$p(x)(1 + \exp(f_0 + f_1 x)) = \exp(f_0 + f_1 x)$$

$$p(x) = \exp(f_0 + f_1 x)[1 - p(x)]$$

$$\frac{p(x)}{1-p(x)} = \exp(f_0 + f_1 x)$$

$$f_0 + f_1 x = \log \left(\frac{p(x)}{1-p(x)} \right) = \underbrace{\log \left(\frac{\Pr(Y=1|X)}{\Pr(Y=0|X)} \right)}_{\text{logit}(x)}$$

Q: a one-unit increase in X is associated w/ how much increase in $\Pr(Y=1|X)$?

$$\underline{A}: \Pr(Y=1|X+1) - \Pr(Y=1|X)$$

$$= \frac{\exp(f_0 + f_1(x))}{1 + \exp(f_0 + f_1(x))} - \frac{\exp(f_0 + f_1(x))}{1 + \exp(f_0 + f_1(x))}$$

Logistic regression:

$$\Pr(Y=1|X) = \frac{\exp(f_0 + f_1(x))}{1 + \exp(f_0 + f_1(x))}$$

$$\Rightarrow Y \sim \text{Bernoulli} \left(\frac{\exp(f_0 + f_1(x))}{1 + \exp(f_0 + f_1(x))} \right)$$

Further assume that have n i.i.d.

draws of $\{(x_i, y_i)\}_{i=1}^n$ from this model

$$L(f_0, f_1) = \prod_{i=1}^n \left(\frac{\exp(f_0 + f_1 x_i)}{1 + \exp(f_0 + f_1 x_i)} \right)^{y_i} \left(\frac{1}{1 + \exp(f_0 + f_1 x_i)} \right)^{1-y_i}$$

$Z \sim \text{Bernoulli}(s)$

$$f(z) = s^z (1-s)^{1-z}$$

Select $\hat{\beta}_0, \hat{\beta}_1$ to maximize
 $\log([(\hat{\beta}_0, \hat{\beta}_1)])$

$$\hat{\Pr}(Y=1 | X=x) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x)}$$

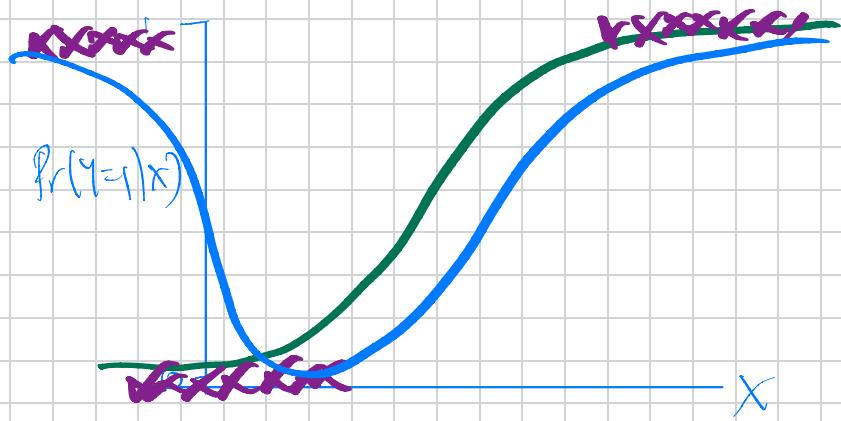
$$\hat{\Pr}(Y=0 | X=x) = \frac{1}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x)}$$

Multiple logistic regression:

$Y \in \{0, 1\}$; $X_1, \dots, X_p \in \mathbb{R}$

$$\Pr(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}{1 + \exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}$$

$$\log\left(\frac{\Pr(Y=1|X)}{1 - \Pr(Y=1|X)}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$



$$\Pr(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3)}{1 + \exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3)}$$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\underline{Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \varepsilon}$$

$$\Pr(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 \sin X)}{1 + \exp(\beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 \sin X)}$$

$$\text{logit} \left(\frac{\Pr(Y=1|X)}{1 - \Pr(Y=1|X)} \right) = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 \sin X$$

Logistic, more generally:

$$\Pr(Y=1|X) = \frac{\exp(f(X))}{1 + \exp(f(X))}$$

where $f(\cdot)$ is unknown

Compare:

$$Y = f(X) + \varepsilon$$

Multinomial logistic regression.

Now, instead of Y taking on 2 possible values, it can take $K \geq 2$ values.

e.g. $Y \in \{A, B, AB, O\}$.

$$Y = \begin{cases} 1 & \text{if } A \\ 2 & \text{if } B \\ 3 & \text{if } AB \\ 4 & \text{if } O \end{cases}$$

$(K-1)(p+1)$



For $k=1, \dots, K-1$:

$$\Pr(Y=k | X=x) = \frac{\exp(\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p)}{1 + \sum_{l=1}^{K-1} \exp(\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p)}$$

and

$$\Pr(Y=K | X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p)}$$

Check: $\sum_{k=1}^K \Pr(Y=k | X=x) = 1$

Check: if $K=2$,

$$\Pr(Y=1 | X=x) = \frac{\exp(\beta_{10} + \beta_{11}x_1 + \dots + \beta_{1p}x_p)}{1 + \exp(\beta_{10} + \beta_{11}x_1 + \dots + \beta_{1p}x_p)}$$

$$\Pr(Y=2 | X=x) = \frac{1}{1 + \exp(\beta_{20} + \beta_{21}x_1 + \dots + \beta_{2p}x_p)}$$

In general,

[aka Soft max]

For $k=1, \dots, K-1$,

$$\Pr(Y=k | X=x) = \frac{\exp(f_k(x))}{1 + \sum_{l=1}^{K-1} \exp(f_l(x))}$$

and

$$\Pr(Y=K | X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(f_l(x))}$$

Another way to write (multinomial) logistic regression:

$$K \cdot (p+1)$$

For $k=1, \dots, K$:

$$\Pr(Y=k|X=x) = \frac{\exp(\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p)}{\sum_{l=1}^K \exp(\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p)}$$

and

~~$$\Pr(Y=1|X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p)}$$~~

LINEAR REG:

$$Y = f(X) + \varepsilon \quad E(\varepsilon) = 0$$

and I think it's

reasonable to suppose

$$f(X) \approx \beta_0 + \beta_1 X,$$

and $\varepsilon_1, \dots, \varepsilon_n$ indept

LOGISTIC REG

$[4 \text{ to } 13]$

$$\Pr(Y=1 | X) = \frac{\exp(f(X))}{1 + \exp(f(X))}$$

$$1 + \exp(f(X))$$

and reasonable $f(X) \approx \beta_0 + \beta_1 X$

and y_1, \dots, y_n indept.

Discriminant Methods. ($K \geq 2$ classes)

Bayes Rule: $\Pr(A|B) = \frac{\Pr(B|A)\Pr(A)}{\Pr(B)}$

$$\Pr(Y=k | X=x) = \frac{\Pr(X=x | Y=k) \Pr(Y=k)}{\Pr(X=x)}$$
$$= \frac{\Pr(X=x | Y=k) \Pr(Y=k)}{\sum_{l=1}^K \Pr(X=x | Y=l) \Pr(Y=l)}$$

Goal: Use this \uparrow to develop a classification method.

Notation: $\Pi_k := \Pr(Y=k)$

$$f_k(x) = \Pr(X=x | Y=k)$$

$$\Pr(Y=k|X=x) = \frac{f_k(x) \pi_k}{\sum_{l=1}^K f_l(x) \pi_l}$$

① est. π_k
 $f_k(\cdot)$

② plug in those
 ests to get
 an est of
 $\Pr(Y=k|X=x)$

↗
discriminant method

All that remains is to est. π_k and $f_k(\cdot)$.

To estimate π_k , use $\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{Y_i=k\}}$

How to estimate $f_k(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$?

TODAY

{ Option 1: linear discriminant analysis

Option 2: Naive Bayes

Linear discriminant analysis.

We will assume that within the k^{th} class ($k=1, \dots, K$),

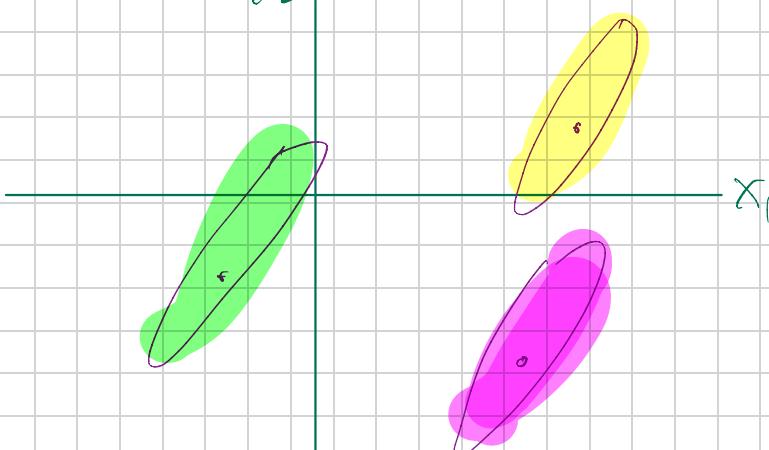
$$X \sim N_p(\mu_k, \Sigma) \quad \begin{matrix} \text{p-vector} \\ \text{p-vector} \end{matrix}$$

X is MVN w/ class-specific mean, and a common variance

Then,

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \right)$$

(where $|\Sigma| = \det(\Sigma)$)



So,

$$\Pr(Y=k|X=x) = \frac{f_k(x) \hat{\pi}_k}{\sum_{l=1}^K f_l(x) \hat{\pi}_l}$$

$$= \frac{\hat{f}_k(x) \hat{\pi}_k}{\text{denom}(x)}$$

$$= \frac{c \exp\left\{-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right\} \cdot \hat{\pi}_k}{\text{denom}(x)}$$

same for all $k=1, \dots, K$

$$\propto \exp\left\{-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right\} \cdot \hat{\pi}_k$$

- will assign x to the class for which $\Pr(Y=k|X=x)$ largest

To est. μ_k :

$\hat{\mu}_k$ = sample mean of the
obs. in the k^{th} class.

$\hat{\Sigma}$ = "pooled" est. of Σ .

We assign obs. x to the class

for which

$$\Pr(Y=k | X=x)$$

$$\log[C^A \cdot B] = A + \log(B)$$

$$\propto \exp\left\{-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right\} \cdot \hat{\pi}_k$$

is largest.

i.e. we will assign x to class for which

$$-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\hat{\pi}_k)$$

is largest

$$= -\frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\hat{\pi}_k)$$

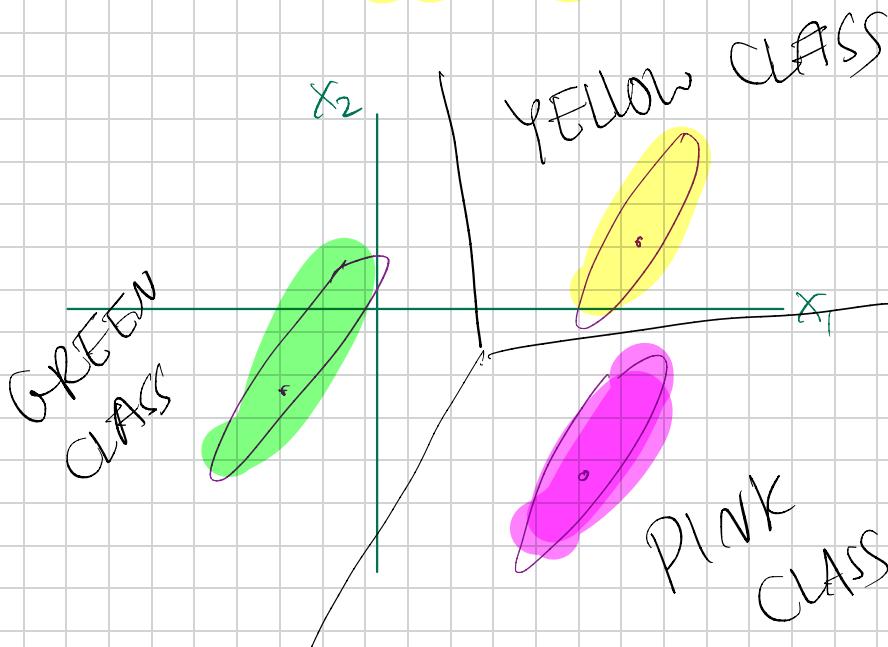
$$= \underbrace{\left(-\frac{1}{2} x^T \Sigma^{-1} x\right)}_{c(x)} - \underbrace{\left(\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\hat{\pi}_k)\right)}_{d_k} + x^T (\Sigma^{-1} \mu_k)$$

$$= c(x) + d_k + x^\top (\Sigma^{-1} \mu_k)$$

⇒ So, I'm assigning x to the class for which

$$\cancel{c(x)} + d_k + x^\top (\Sigma^{-1} \mu_k)$$

is largest.



Linear disc. analysis:

① recalling Bayes:

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

② applied it to classifn:

$$\Pr(Y=k | X=x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l}$$

③ assume that if $Y=k$, then

$$X \sim N_p(\mu_k, \Sigma)$$

④ we est. $\hat{\pi}_k, \hat{\mu}_k, \hat{\Sigma}$

⑤ we got $\hat{\Pr}(Y=k|X)$, and a linear decision boundary

NAIVE BAYES

① recalling Bayes:

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

② applied it to classifn:

$$\Pr(Y=k|X=x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l}$$

③ assume that

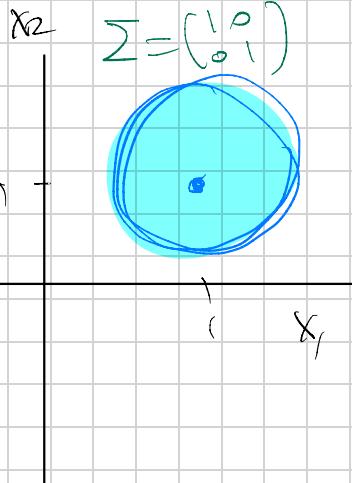
$$f_k(x) = \prod_{j=1}^p f_{kj}(x_j)$$

} each of the p features is independent

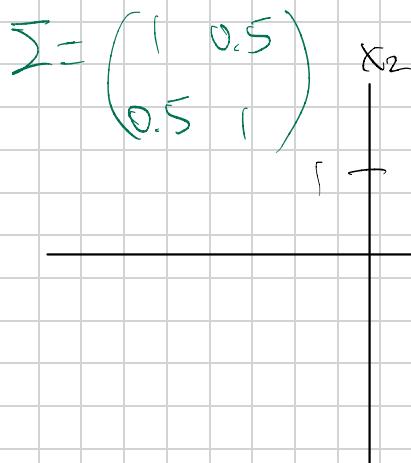
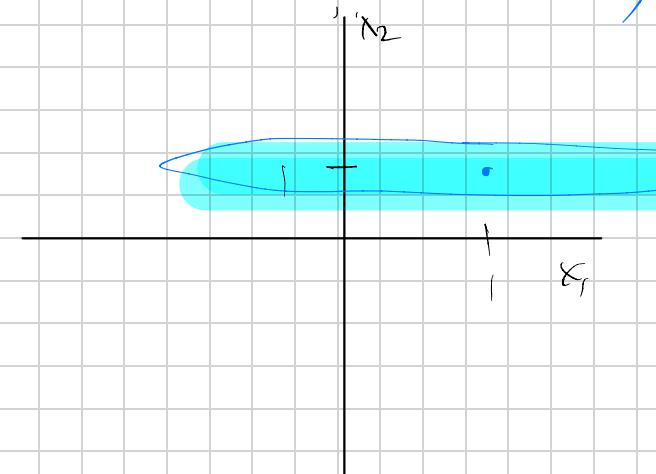
④ we est. $\{\hat{f}_{kj}(\cdot)\}_{j=1, \dots, p; k=1, \dots, K}$

⑤ we got $\hat{f}_k(x)$

$$X \sim N_2(\mu, \Sigma) \quad ; \quad \mu = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 0.2 \end{pmatrix}$$



Wrapping up Classification

3 options for classification:

① logistic regression;

$$\Pr(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X)}{1 + \exp(\beta_0 + \beta_1 X)}$$



[or more generally]



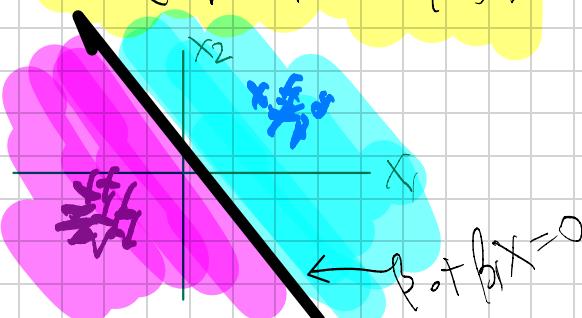
$$\Pr(Y=1|X) = \frac{\exp(f(X))}{1 + \exp(f(X))}$$

$$\log \left(\frac{\Pr(Y=1|X)}{1 - \Pr(Y=1|X)} \right)$$

$$\beta_0 + \beta_1 X$$



$$f(X)$$

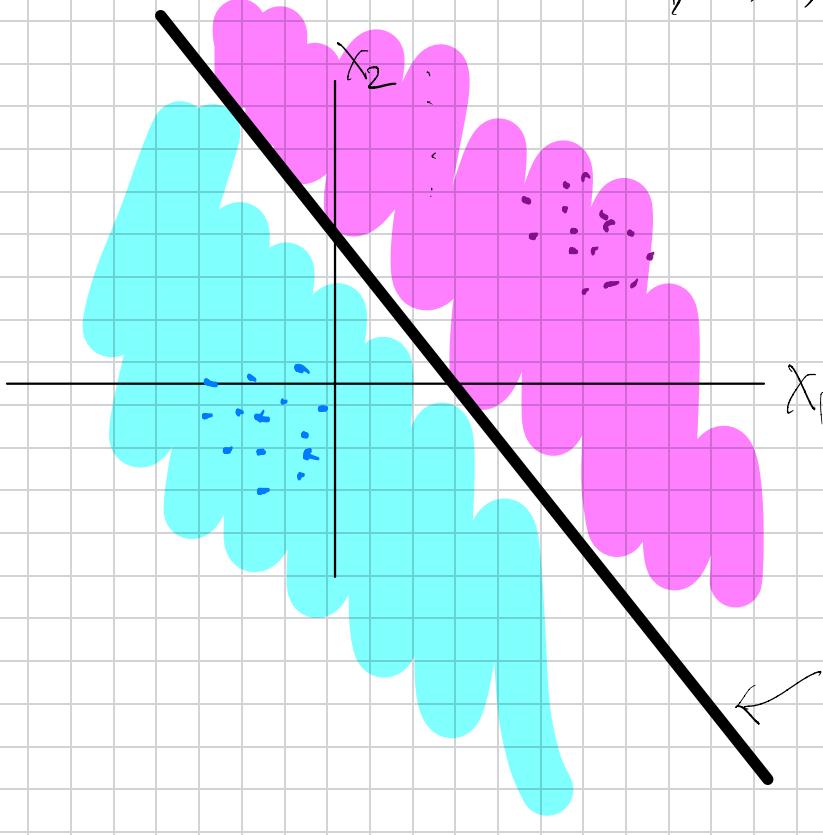


② linear disc. analysis;

for k^{th} class,

$$X \sim N(\mu_k, \Sigma)$$

} different assumption gives naive Bayes

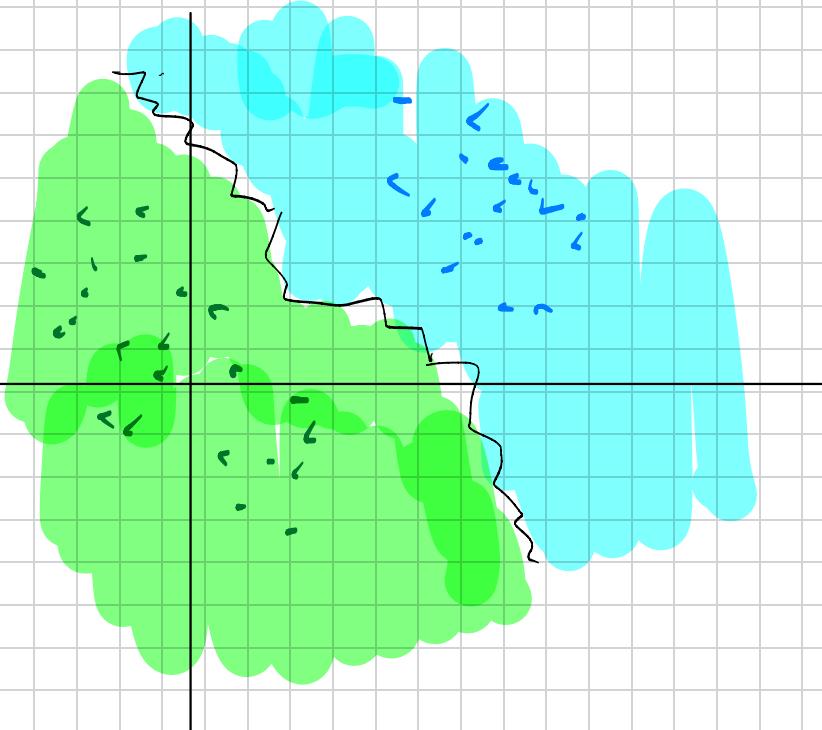


Use it to est

$$\Pr(Y=k | X=x)$$

③ K-nearest neighbors.

$$\hat{\Pr}(Y=l \mid X=x) = \left\{ \begin{array}{l} \text{among the } k \\ \text{nearest neighbors of} \\ x, \text{ what fraction} \\ \text{are in } l^{\text{th}} \text{ class?} \end{array} \right.$$



Reminder: Bayes classifier is the unattainable ideal -- it assigns the obs x to the class k for which

$$\Pr(Y=k|X=x) = \frac{\Pr_{k|k}(x)}{\sum_{l=1}^K \Pr_{l|k}(x)}$$

~~is larger!~~

THE POINT IS: all of these classifiers are doing their best to approx. the Bayes classifier!