CLASS CONDITIONAL INDEPENDENCE

$2d+1$

---

PARAMETER ESTIMATION

$$p, \left\{ p_1^1, \cdots p_d^1 \right\}, \left\{ p_1^0, \cdots, p_d^0 \right\}$$

---

MAX. LIKELIHOOD ESTIMATES

$$1 \to \hat{p} = \frac{1}{n} \sum_{i=1}^{n} y_i \to \left\{ \text{Fraction of spam emails in the dataset} \right\}$$

$$\forall j \in \{1,\ldots,d\} \quad \hat{p}_j^y = \frac{\sum_{i=1}^{n} \mathbb{1}\left( f_j^i = 1, y_i = y \right)}{\sum_{i=1}^{n} \mathbb{1}\left( y_i = y \right)}$$
$$\forall y \in \{0,1\}$$

$\longleftarrow$ Number of emails with label $y$.

$\hookrightarrow$ Fraction of $y$-labelled emails that contain the $j^{th}$ word.

---

PREDICTION

Given $x^{test} \in \{0,1\}^d$, what is $\hat{y}^{test}$?

$$P\left( y^{test} = 1 \mid x^{test} \right) > P\left( y^{test} = 0 \mid x^{test} \right)$$

$$\Rightarrow \hat{y}^{test} = 1$$
$$= 0 \quad \text{otherwise.}$$

How to obtain $P(y/x)$ from $P(y)$ and $P(x/y)$ ?

**BAYES RULE !**

$$P\left(y^{test}=1 \mid x^{test}\right) = \frac{P\left(x^{test} \mid y^{test}=1\right) \cdot P\left(y^{test}=1\right)}{P\left(x^{test}\right)}$$

$$P\left(y^{test}=0 \mid x^{test}\right) = \frac{P\left(x^{test} \mid y^{test}=0\right) \cdot P\left(y^{test}=0\right)}{P\left(x^{test}\right)}$$

$$\underbrace{P\left(x^{test} \mid y^{test}=1\right)} \cdot P\left(y^{test}=1\right)$$

$$= P\left(x^{test}=[f_1 \ f_2 \cdots f_d] \mid y^{test}=1\right) \cdot P\left(y^{test}=1\right)$$

$$= \left(\prod_{j=1}^{d} \left(\hat{p}_j^1\right)^{f_j} \left(1-\hat{p}_j^1\right)^{(1-f_j)}\right) \cdot \hat{p}$$

$x_{test} = \begin{bmatrix} f_1 & f_2 & f_3 & f_4 \\ 0 & 0 & 1 & 0 \end{bmatrix}$

IF

$$\left(\prod_{j=1}^{d} \left(\hat{p}_j^1\right)^{f_j} \left(1-\hat{p}_j^1\right)^{(1-f_j)}\right) \cdot \hat{p} > \left(\prod_{j=1}^{d} \left(\hat{p}_j^0\right)^{f_j} \left(1-\hat{p}_j^0\right)^{(1-f_j)}\right) \cdot (1-\hat{p})$$

$\Rightarrow$ PREDICT $\hat{y}^{test} = 1$

else $\hat{y}^{test} = 0$.

MODEL USES 2 main Things

**NAIVE BAYES ALGORITHM.**

- CLASS CONDITIONAL INDEPENDENCE
- BAYES THEOREM

- may not hold in practice
- NAIVE ASSUMPTION
- still works well in practice.

# PITFALLS IN NAIVE BAYES TO WATCH OUT FOR.

- If a word does not appear in the train set but appears in a test datapoint,

$$\hat{p}_j^1 = 0 \qquad \hat{p}_j^0 = 0$$

$$P\left(y^{test} = 1 \mid x^{test} = [f_1, f_2 \cdots f_d]\right) \propto \left(\prod_{i=1}^{d} \underbrace{(\hat{p}_i^1)}_{0}^{\widehat{f_i}^1} \underbrace{(1-\hat{p}_i^1)}_{1}^{\overbrace{(1-f_i)}^{0}}\right) \hat{p}$$

$$\underset{\substack{\uparrow \\ Zebra \\ =1}}{}$$

$$\underbrace{\hspace{6cm}}_{0}$$

$$P\left(y^{test} = 0 \mid x^{test} = [f_1 \cdots f_d]\right) \propto \left(\prod_{i=1}^{d} \underbrace{(\hat{p}_i^0)}_{0}^{1 f_i} \underbrace{(1-\hat{p}_i^0)}_{1}^{(1-f_i)}\right)(1-\hat{p})$$

$$\underset{\substack{\uparrow \\ Zebra \\ =1}}{}$$

$$\underbrace{\hspace{6cm}}_{0}$$

---

## Possible Fix

- Can add two "pseudo" emails with all words present — one email has label 0 and another has label 1

**LAPLACE SMOOTHING**

$$\overset{x_1^{pseudo}}{\left[1 \quad 1 \quad 1 \quad 1 \quad \cdots \cdots \quad 1\right]} \qquad \overset{y_1^{pseudo}}{1}$$

$$\underset{x_0^{pseudo}}{\left[1 \quad 1 \quad 1 \quad \cdots \cdots \quad 1\right]} \qquad 0$$

---

## DECISION FUNCTION OF NAIVE BAYES.

Given $x_{test}$; $\quad y_{test} = 1$ if $\quad \dfrac{P\left(y_{test} = 1 \mid x_{test}\right)}{P\left(y_{test} = 0 \mid x_{test}\right)} \geq 1$

$$\log\left(\frac{P(y_{test}=1 \mid x_{test})}{P(y_{test}=0 \mid x_{test})}\right) \geq 0$$

$$\log\left(\frac{P(x_{test} \mid y_{test}=1) \cdot P(y_{test}=1) \,/\, P(x_{test})}{P(x_{test} \mid y_{test}=0) \cdot P(y_{test}=0) \,/\, P(x_{test})}\right) \geq 0$$

$x_{test} = [f_1 \; f_2 \cdots f_d]$

$$\log\left(\prod_{i=1}^{d} \frac{(\hat{p}_i^1)^{f_i} (1-\hat{p}_i^1)^{(1-f_i)} \cdot \hat{p}}{(\hat{p}_i^0)^{f_i} (1-\hat{p}_i^0)^{(1-f_i)} (1-\hat{p})}\right) \geq 0$$

$$= \log\left(\prod_{i=1}^{d} \left(\frac{\hat{p}_i^1}{\hat{p}_i^0}\right)^{f_i} \left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0}\right)^{(1-f_i)} \cdot \frac{\hat{p}}{1-\hat{p}}\right) \geq 0$$

$$= \sum_{i=1}^{d}\left( f_i \log\left(\frac{\hat{p}_i^1}{\hat{p}_i^0}\right) + (1-f_i) \log\left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0}\right) + \log\left(\frac{\hat{p}}{1-\hat{p}}\right)\right) \geq 0$$
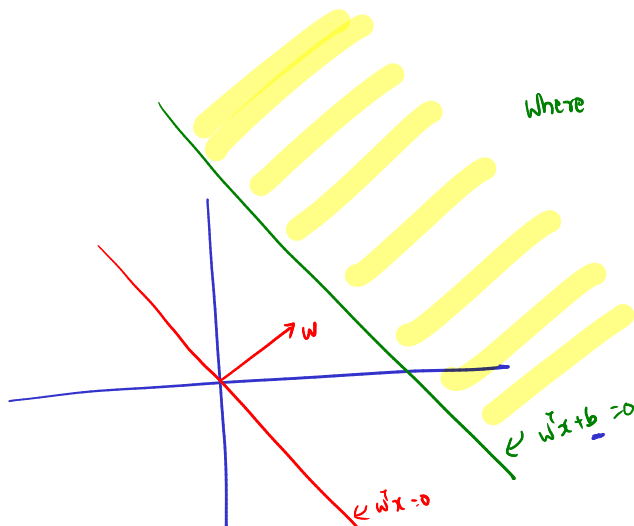
$x_{test} = [f_1 \cdots f_d]$

$$= \sum_{i=1}^{d} f_i\left( \log\left(\frac{\hat{p}_i^1 (1-\hat{p}_i^0)}{\hat{p}_i^0 (1-\hat{p}_i^1)}\right)\right) + \log\left(\frac{(1-\hat{p}_i^1)}{(1-\hat{p}_i^0)}\right) + \log\left(\frac{\hat{p}}{1-\hat{p}}\right) \geq 0$$

DECISION FUNCTION is of the form

Predict $y_{test} = 1$ if $\quad W^T x_{test} + b \geq 0$
$\quad \in \mathbb{R}^d$

Where $\quad W_i = \log\left(\dfrac{\hat{p}_i^1 (1-\hat{p}_i^0)}{\hat{p}_i^0 (1-\hat{p}_i^1)}\right) \quad , \quad b =$


$\leftarrow W^T x + b = 0$
$\leftarrow W^T x = 0$
$W$

CONCLUSION:

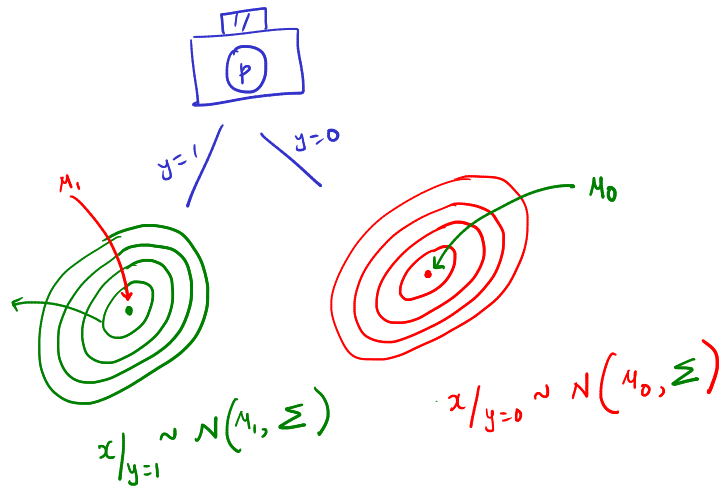$\rightarrow$ DECISION FUNCTION OF NAIVE BAYES is LINEAR!

DATA: $\left\{ (x_1, y_1), \ldots\ldots (x_n, y_n) \right\}$

$$x_i \in \mathbb{R}^d \qquad y_i \in \{0, 1\}$$

## A GENERATIVE STORY



PARAMETERS

- $p$
- $\mu_0, \mu_1$
- $\Sigma$

$$x/_{y=1} \sim N(\mu_1, \Sigma) \qquad x/_{y=0} \sim N(\mu_0, \Sigma)$$

NOTE: In this model, covariances are assumed to be same

## MAXIMUM LIKELIHOOD ESTIMATES

$$\hat{p} = \frac{\sum_{i=1}^{n} y_i}{n} \qquad \longleftarrow \text{FRACTION OF points labelled 1.}$$

$$\hat{\mu}_1 = \frac{\sum_{i=1}^{n} \mathbb{1}(y_i = 1) \cdot x_i}{\sum_{i=1}^{n} \mathbb{1}(y_i = 1)} \qquad \longleftarrow \text{Sample mean of data points labelled 1.}$$

$$\hat{\mu}_0 = \frac{\sum_{i=1}^{n} \mathbb{1}(y_i = 0) \cdot x_i}{\sum_{i=1}^{n} \mathbb{1}(y_i = 0)} \qquad \longleftarrow \text{Sample mean of data points labelled 1.}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} \left( x_i - \hat{M}_{y_i} \right) \left( x_i - \hat{M}_{y_i} \right)^T$$

---

PREDICTION ?    Bayes rule.

$$P\left( y_{test} \mid x_{test} \right) \quad \propto \quad \underbrace{P\left( x_{test} \mid y_{test} \right)}_{f\left( x_{test}; \hat{M}_{y_{test}}, \hat{\Sigma} \right)} \cdot \underbrace{P\left( y_{test} \right)}_{\hat{p}}$$

Predict $y_{test} = 1$ if

$$f\left( x_{test}; \hat{M}_1, \hat{\Sigma} \right) \cdot \hat{p} \quad \geq \quad f\left( x_{test}; \hat{M}_0, \hat{\Sigma} \right) \cdot (1-\hat{p})$$

$$e^{-\left( x_{test} - \hat{M}_1 \right)^T \cdot \hat{\Sigma}^{-1} \left( x_{test} - \hat{M}_1 \right)} \cdot \hat{p} \quad \geq \quad e^{-\left( x_{test} - \hat{M}_0 \right)^T \hat{\Sigma}^{-1} \left( x_{test} - \hat{M}_0 \right)} \quad (1-\hat{p})$$

on  Simplification    [Take log]

Predict 1 if $\left( \left( \hat{M}_1 - \hat{M}_0 \right)^T \hat{\Sigma}^{-1} \right) x_{test} + \hat{M}_0^T \hat{\Sigma}^{-1} \hat{M}_0 - \hat{M}_1^T \hat{\Sigma}^{-1} \hat{M}_1 + \log\left( \frac{1-\hat{p}}{\hat{p}} \right) \geq 0$
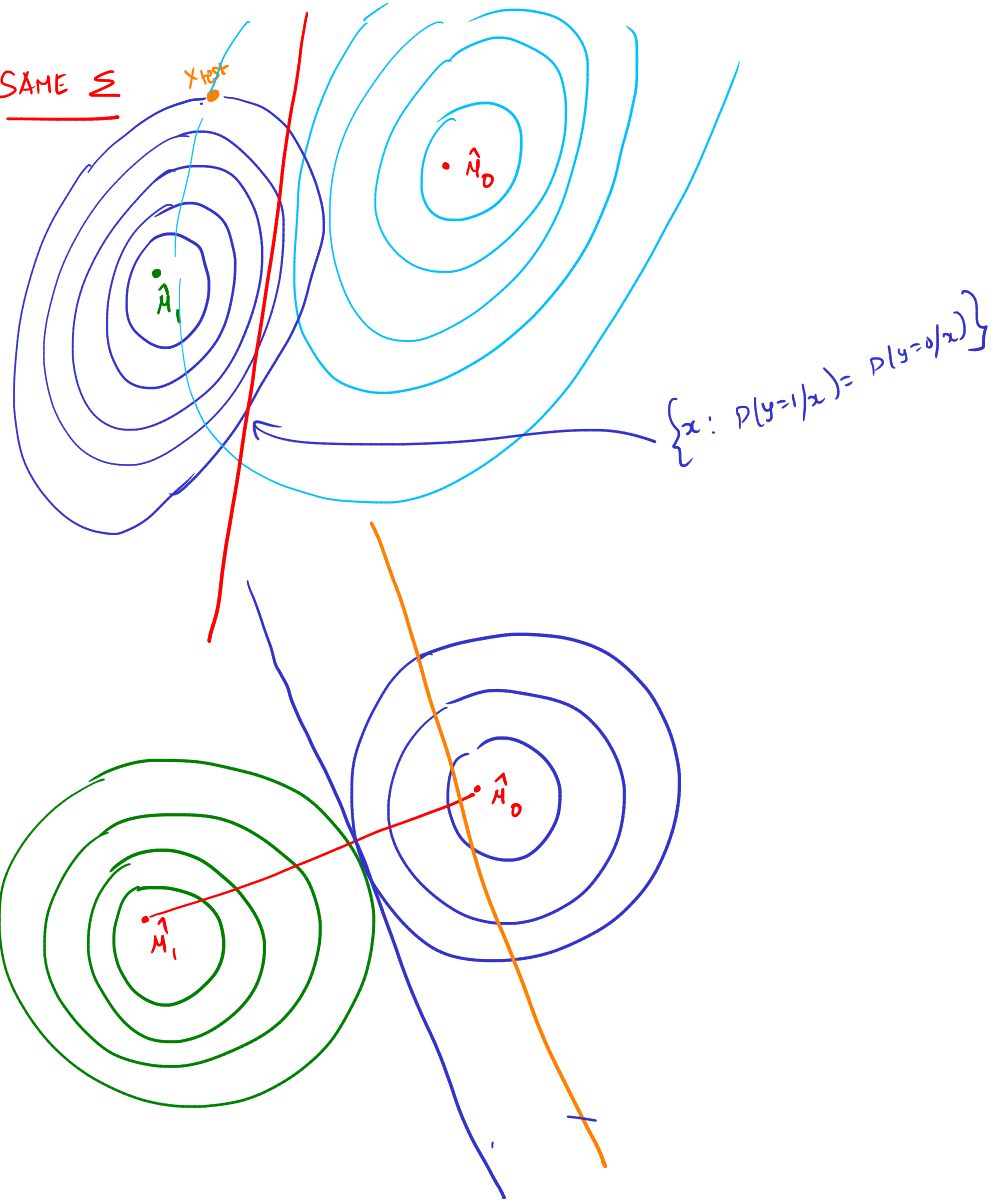
$$w^T x_{test} + b \geq 0$$



Predict 1

Predict 0    $\leftarrow w^T x + b$

DECISION FUNCTION is LINEAR!

$\rightarrow \Sigma$ is same for both classes.

SAME $\Sigma$

$x_{test}$

$\cdot \hat{M}_0$

$\hat{M}_1$

$\left\{ x: P(y=1/x) = P(y=0/x) \right\}$

$\cdot \hat{M}_0$

$\hat{M}_1$

DIFFERENT $\Sigma$

In general
Quadratic
decision
boundary.

$\cdot \hat{M}_1$

$\cdot \hat{M}_0$

$\leftarrow \left\{ x: P(y=1/x) = P(y=0/x) \right\}$

GAUSSIAN NAIVE BAYES