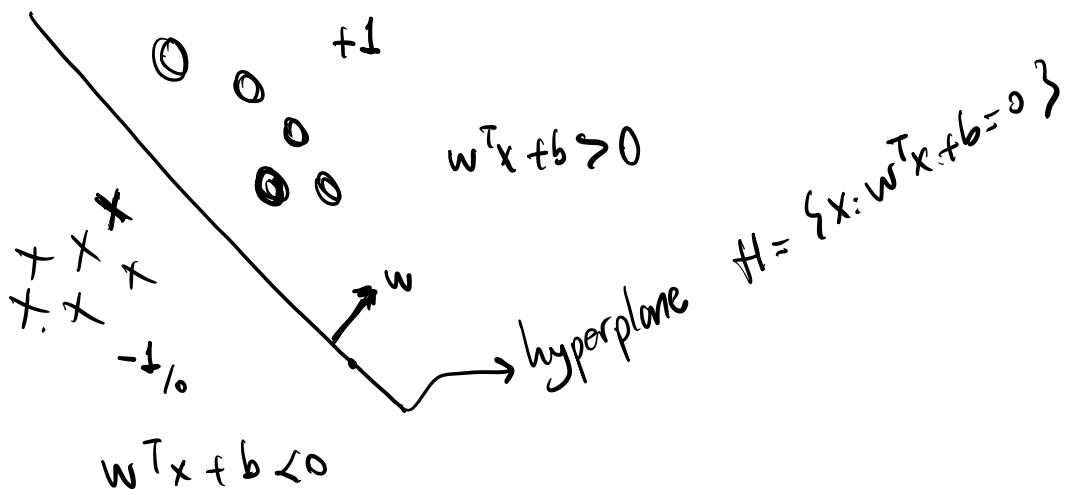


## Lecture - 4

Recap for Perceptron.



$P \leftarrow$  inputs with label  $\pm 1$

$N \leftarrow$   $"$        $"$        $"$        $"$        $0$

Initialize  $w \leftarrow$  randomly.

while ! convergence

Pick random  $x \in P \cup N$

if  $x \in P$  and  $\sum_{i=0}^n w_i x_i < 0$  then

$$w = w + x$$

if  $x \in N$  and  $\sum_{i=0}^n w_i x_i \geq 0$  then

$$w = w - x$$

$H$

MLE, MAP, Naive Bayes Classifier.

## Maximum Likelihood Estimation (MLE)

$P(x)$   
→ we don't know  $P$ .

$P(H) = ?$   
→ I can sample data.

Toss 10 times  $\Rightarrow n=10$

$$D = \{H, T, T, H, H, H, T, T, T, T\}$$

$$P(H) = ? \frac{4}{10} \quad n_H = 4 \Rightarrow P(H) \approx \frac{n_H}{n_H + n_T} = \boxed{\frac{4}{10}}$$

How can we derive this formula

MLE

$$P(D; \theta)$$

↳ probability of observing heads.  
 ↳ collected 10 flips.

$$\theta = \underbrace{\underset{\theta}{\operatorname{argmax}} \ P(D; \theta)}_{\text{---}}$$

$$\text{Binomial Distribution} = \Pr(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

- fixed number.
- independent Bernoulli trials.
- same probability of success.

$$P(D; \theta) = \binom{n_H + n_T}{n_H} \theta^{n_H} \cdot (1-\theta)^{n_T}$$

$$10 = n_T$$

MLE  $\Rightarrow$  frequentist approach.

$$\hat{\theta}_{\text{MLE}} = \operatorname{argmax} P(D; \theta)$$

$$= \arg \max \left( \frac{n_H + n_T}{n_H} \right) \theta^{n_H} (1-\theta)^{n_T}$$

$$= \underset{\theta}{\operatorname{argmax}} \log \left( \frac{n_H + n_T}{n_H} \right) + n_H \cdot \log(\theta) + n_T \cdot \log \frac{(1-\theta)}{(1-\theta)}$$

constant

$$= \underset{\theta}{\operatorname{argmax}} n_H \cdot \log(\theta) + n_T \cdot \log(1-\theta)$$

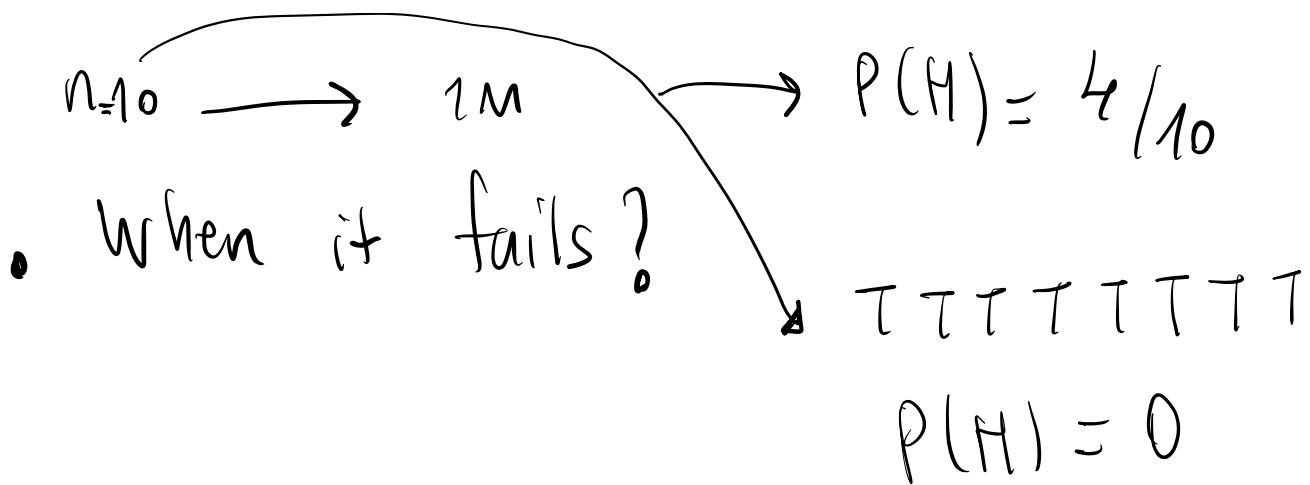
## Derivation

$$\frac{n_H}{\theta} - \frac{n_T}{1-\theta} = 0$$

$$\frac{n_H}{\theta} = \frac{n_T}{1-\theta} \Rightarrow n_H - \overbrace{\theta n_H}^{\rightarrow} = n_T \theta$$

$$\frac{n_H}{n_H + n_T} = \theta$$

~~n\_H = n\_T \theta + n\_H \theta~~



frequentist vs Bayesian

$P(D; \theta)$

parameter ↴

Random variable }

Simple fix: "m" imaginary throws

that would result in  $\theta'$  (e.g. 0.5). Add " $m'$ " heads and " $m''$ " tails to your data.

$$\hat{\theta} = \frac{n_H + m_H}{n_H + n_T + m_H + m_T} = \frac{h_H + m}{n_H + n_T + 2m}$$

$$= \frac{h_H + 1}{n_H + n_M + 2}$$

$\Rightarrow$  for small " $n$ ", it incorporates your "prior belief" about what  $\theta$  should be.



Laplace Smoothing:  $P(x) = \frac{\text{count}(x) + \alpha}{n + \alpha v}$

## Bayesian Approach

it combines prior belief, with observed data,  
 expressed as prior dist. likelihood.

to update our belief in the form

of posterior distribution.

$$P(D; \theta) \Rightarrow P(\tilde{P} | \theta)$$

parameter                          random variable.

### Maximum A Posterior (MAP) Estimation.

Prior distribution =  $P(\theta)$

likelihood function =  $P(D|\theta)$  for observed data D.

MAP estimate  $\hat{\theta}_{MAP}$  is.

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \underbrace{P(\theta|D)}$$

Using Bayes Theorem

prior belief

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(D|\theta) \cdot P(\theta)$$

$P(D)$

$P(\theta|D)$        $P(D|\theta) \cdot P(\theta)$       Bayes theorem.

$P(D)$       normalizer.

$$N = \{H_1, H_2, T, \dots, T\}$$

comes from  
 $H=4$   
 $T=6$

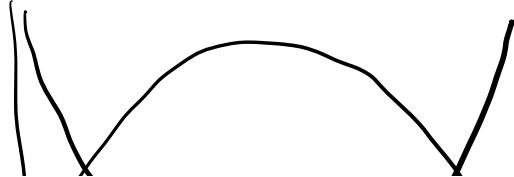
$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(D|\theta) \cdot P(\theta)$$

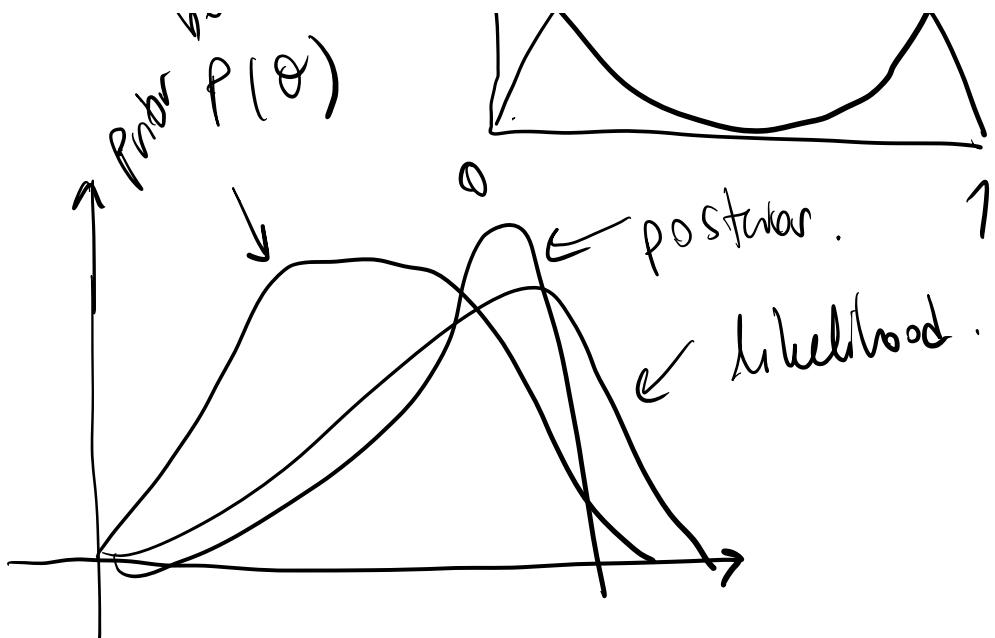
Beta distribution.

$$P(\theta) = \frac{\theta^{\alpha-1} \cdot (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

~~$B(\alpha, \beta)$~~

shape





$\hat{\theta}_{MLE}$

$$\hat{\theta}_{MAP} \underbrace{P(D|\theta)}_{\text{likelihood}} \cdot \underbrace{P(\theta)}_{\text{prior}}$$

$$P(D|\theta) = \left[ \theta^{n_H} (1-\theta)^{n_T} \cdot \theta^{\alpha-1} (1-\theta)^{B-1} \right]$$

$$= \theta^{n_H + \alpha - 1} \cdot (1-\theta)^{n_T + B - 1}$$

$$\hat{\theta}_{MAP} = \frac{n_H + \alpha - 1}{n_H + n_T + (\alpha + B - 2)}$$



- The MAP estimate is identical to MLE, with  $\alpha-1$  heads (belief) and  $\beta-1$  tails (belief)
- As  $n \rightarrow \infty$   $\hat{\theta}_{\text{MAP}} \xrightarrow{\text{P}} \theta_{\text{MLE}}$  as  $\underline{\alpha-1}$  and  $\underline{\beta-1}$  become irrelevant compared to very large  $n_H, n_T$ .

Ex/  $D = \{(x_1, y_1) \rightarrow \overset{d}{\sim} (x_n, y_n)\}$   
 drawn from some unknown  $P(x, y)$

for example  $(x, y)$

$$\hat{P}(x, y) = \frac{\sum_{i=1}^n I(x_i = x \wedge y_i = y)}{n}$$

I need one ---

How can we estimate  $\hat{P}(y|x) = ?$

spam email

$$\hat{P}(y|x) = ?$$

$$P(y|x) = \frac{P(y, x)}{P(x)} = \frac{\sum_{i=1}^n I(x_i = x \wedge y_i = y)}{\sum_{i=1}^n I(x_i = x)}$$

$$P(y=y | x_1 = x_1, x_2 = x_2, \dots, x_d = x_d)$$

MLE is only good if there are many examples (training vector) with the

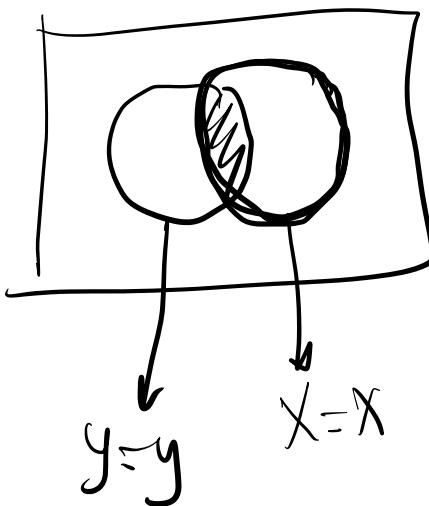
s are identical features of  $X$ !

In high dimensional spaces it never happens

$P(X)$ ,  $D \sim P(Y|X) \rightarrow$  single feature.

$$P(Y=y | X=x)$$

$$\underline{X=x} / y=y$$



$$\frac{P(Y=y, X=x)}{P(X)}$$



$$\frac{\sum_{i=1}^n I(X_i = x \wedge Y_i = y)}{\sum_{i=1}^n I(X_i = x)}$$

- how many times do we see  $X$  and out of those how many times be also see  $Y=y$ .

In high dimension.

$$P(Y=y \mid X_1=x_1, X_2=x_2, \dots, X_d=x_d)$$

- how many times, i exactly see the same data point in my dataset?

## Naive Bayes

Features are independent given the label.]

n features ,  $X = (x_1, \dots, x_n)$

$$P(C_k | X) = \frac{\underbrace{P(X | C_k)}_{\text{likelihood}} \times P(C_k)}{P(X)}$$

feature independence assumption

n-dimensional space

$$\begin{aligned} P(X | C_k) &= \boxed{P(x_1, x_2, x_3, \dots, x_n | C_k)} \\ &\approx \underbrace{P(x_1 | C_k)}_{\text{.n.}} \times P(x_2 | C_k) \times \dots \times P(x_n | C_k) \end{aligned}$$

$\cap^P$

Plugging into Bayes' Theorem

$$P(C_k)$$

$$P(C_k | X) \propto \underbrace{P(x_1 | C_k) \times P(x_2 | C_k) \times \dots \times P(x_n | C_k)}_{\cancel{P(X)}} \times$$

Classification Decisions

- For a given obs.  $X$ , we compute

$$\underbrace{P(C_k | X)}_{\text{for each class } C_k}$$

and assign the obs. to the class that

max. value:

$$C(x) = \operatorname{argmax} P(C_k) \times P(x_1 | C_k) \times P(x_2 | C_k) \times \dots \times P(x_n | C_k)$$

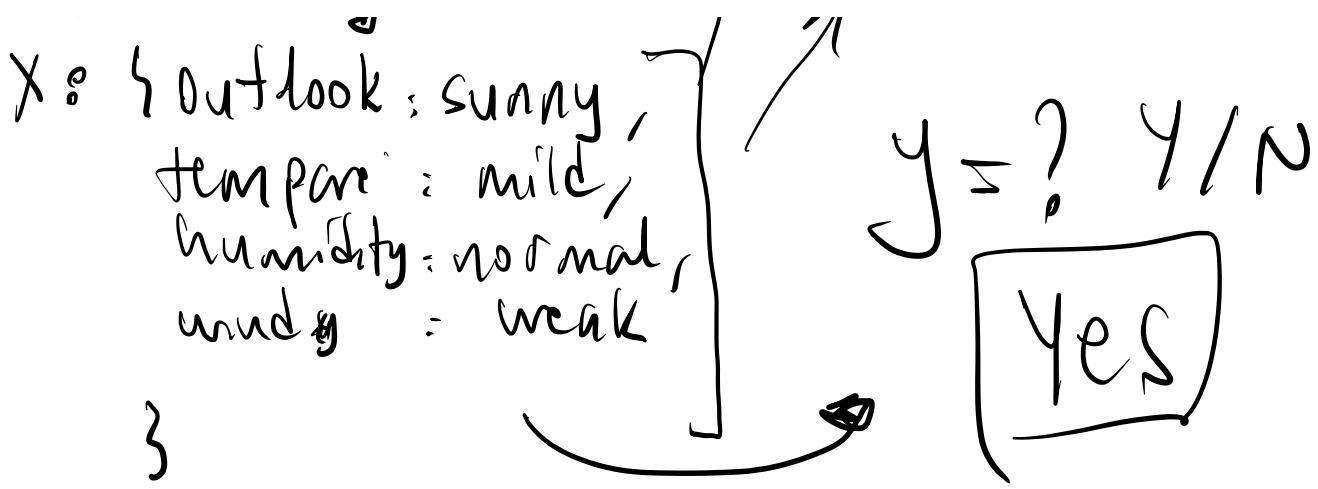
$$P(\text{spam} | \vec{x}) = \lambda$$

$$P(\neg \text{spam} | \vec{x}) = \mu$$

Example.

	features				↓ label
	outlook	temperature	humidity	wind	play football
1	Sunny	hot	high	weak	N
2	S	h	H	strong	Y
3	overcast	h	H	(w)	Y
4	rainy	mild	H	(w)	Y
5	R	cool	normal	(w)	Y
6	R	C	N	S	Y
7	O	C	N	S	Y
8	S	M	High	W	Y
9	S	C	N	(w)	Y
10	R	M	N	S	Y
11	R	M	N	S	Y
12	O	M	N	S	Y
13	O	T	H	W	N
14	R	M	H	S	N

test data ↴ ↑ ↳



$$P(\text{yes} | X) \propto P(X|y) \cdot P(y) \quad \cancel{\text{P}(y)}$$

$$\propto \underbrace{P(x_1|y) \cdot P(x_2|y) \cdot P(x_3|y) \cdot P(x_4|y)}_{\text{P}(X|y)}.$$

P(y)

$$\propto P(\text{sunny} | y) \cdot P(\text{mild} | y) \cdot P(\text{normal} | y),$$

$$P(\text{weak} | y) \cdot P(y)$$

$$\frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \cdot \frac{9}{44}$$

$$= 0.02$$

or 1 in 51 times

$$P(\text{No} | X) = \frac{P(\text{Sunny} | N) \cdot P(\text{Wet} | N)}{P(\text{Normal} | N) \cdot P(\text{Wet} | N)} \cdot P(N)$$

$$\Rightarrow P(\text{No}|X) \approx 0.006$$

$$\overbrace{P(\text{Yes}|X)}^{\geq} > P(\text{No}|X)$$

$X = (x_1, \dots, x_d)^T$  <sup>d features</sup>,  $Y$  is the label  $0, 1$

$$\operatorname{argmax}_y P(Y=y | X=(x_1, \dots, x_d))$$

Bayes theorem  $\Rightarrow P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$

$X = (\text{Outlook}, \text{Temp}, \text{humidity}, \text{windy})$  and  $Y = \text{play}$

$$P(Y = \text{play}) = \frac{\# \text{ play}}{\# \text{ play} + \# \text{ not play}} = \frac{9}{14}$$

↓

$$P(X = \text{sunny} | Y = \text{play}) = \frac{\# \text{ play} \& \# \text{ sunny}}{\# \text{ play}}$$

$$= 2/9$$

$$\Rightarrow P(X_1 = 1, X_2 = 2, \dots | Y = \text{play})$$

→ think all features are independent!

$$P(X_1, X_2 | Y) = P(X_1 | Y) \cdot P(X_2 | Y)$$

conditional independence.

$$P(Y | X) = \frac{\prod_{i=1}^d P(X_i | Y) \cdot P(Y)}{P(X)}$$

$$P(x|y) = P(x_1|y) \times P(x_2|y) \times \dots \times P(x_n|y)$$

$\propto P(y|x) \propto \prod_{i=1}^d P(x_i|y) \times P(y)$

Types of NB classifiers.

- Bernoulli NB  $\rightarrow X = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$
- Multinomial NB  $\rightarrow \begin{bmatrix} 5 \\ 2 \\ 3 \\ 2 \\ 1 \\ \vdots \end{bmatrix}$
- Gaussian Naive Bayes.

$t_1$	label
125	2
100	2
70	2
120	2
95	2
60	2
220	2
85	2
75	2
90	2

Test Data

$$X = \frac{120}{?}$$

$$\mu = 110$$

$$\sigma^2 = 2975$$

$$P(x|y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(x_i - \mu_i)^2}{2\sigma^2} \right)$$

$$P(\text{Yes} | 120)_2$$

$$P(Mo(120)) \propto$$

$$\frac{1}{\sqrt{2\pi 2975}} \exp \left( -\frac{(120-110)^2}{2 \times 2975} \right)$$

Adv / Dis . NB

- Adv
- Simplicity
  - Scalability

• Full functioning ...

- Fast computation
- Online updates possible
- Works with high dimensions.
  - ↳ text classification problems.
- Low training work well.
- probabilistic output.

### Disadvantages

1. Naive assumption

2 - Data Distribution Ass.

Gaussian, Bernoulli, ...

3 - Zero frequency problem.

