

2. Suppose we obtain a 'new' test $\langle \bar{x}_{n+1}, ? \rangle$

If $\underline{p(+ | \bar{x}_{n+1}) > p(- | \bar{x}_{n+1})}$ output '+'
otherwise output '-'

Qn:- what is $p(+ | \bar{x}_{n+1}) \rightarrow ?$

Bayes Rule to the rescue:-

$$P(+ | \bar{x}_{n+1}) = \frac{p(+, \bar{x}_{n+1})}{p(\bar{x}_{n+1})} = \frac{p(\bar{x}_{n+1} | +) \cdot p(+)}{p(\bar{x}_{n+1})}.$$

Similarly,

$$p(- | \bar{x}_{n+1}) = \frac{p(\bar{x}_{n+1} | -) \cdot p(-)}{p(\bar{x}_{n+1})}.$$

Comparison reduces to:- $\frac{p(\bar{x}_{n+1} | +) \cdot p(+)}{p(\bar{x}_{n+1} | -) \cdot p(-)} \geq 1$

In general, classes C_1, C_2, \dots, C_K (eg $K=10$ for digit recognition), $\{+, -\}, \{0, \dots, 9\}$.

Evaluate Find $\max_i P(C_i) \cdot P(\bar{x}_{n+1} | C_i) \rightarrow$ [output that].

Questions:- What is $P(C_i)$?

Prior belief about class C_i .

e.g., on 100 days, we played '30', maybe $P(\text{play}) \approx \frac{30}{100}$.

Say empirically 30% of our emails are spam,
then $\text{Prob}(\text{spam}) = 0.3$.

* What are the features?

| Who tells us to look @ weather, temp, ... before deciding to play, ... more philosophical.

Our features will be some high dimensional vector. we want nice representation.

"Represent as a (high dimensional) vector"

Examples :-

① MNIST dataset 20×20 , grayscale images.

each pixel $\in \{0, \dots, 255\} / \{0, 1\}$

Feature vector (\underline{x}) = (x_1, \dots, x_{400}) , $x_i \in \{0, \dots, 255\} / \{0, 1\}$

is the ~~fea~~ value of pixel 'i'.

② News document/email. each word is a ~~value~~ feature

\underline{x} - a length N document (has N words).
 $= (x_1, \dots, x_N)$, x_i - ith word.

Goal:- For each class estimate $P(\bar{x} | \text{class})$.

This is $P(\underline{x_1, \dots, x_k} | \text{class})$.

- What are the challenges in estimating these probabilities (from samples).

Suppose feature i has $|F_i|$ possible values.
e.g., MNIST, each pixel has $255/2$ values.

possible images feature vectors = $|F_1| \cdot |F_2| \cdots |F_k|$

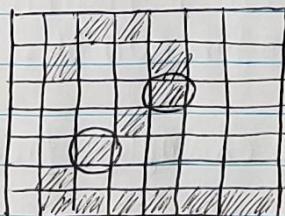
e.g., $255^{400} / 2^{400}$ → very very large.

Same with text data.

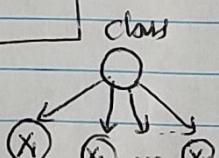
Result from stat/prob:- To estimate M parameters need M' samples. Impossible.

NAIVE BAYES assumption :- The features are "independent" given the class.

⟨very strong assumption⟩.



⟨next word is dependent on previous⟩.



$$P((x_1, \dots, x_k) | \text{class} = c_i) = P(x_1 | \text{class} = c_i) \cdot P(x_2 | \text{class} = c_i) \cdots P(x_k | \text{class} = c_i)$$

how many parameters now?

$$F_1 \dots F_k \rightarrow \boxed{(F_1 + F_2 + \dots + F_k)}.$$

= $F \cdot k$ if all F_i 's are the same.

~ 400 parameters for MNIST if on $\{0, 1\}$.

documents/text :- $F \approx$ all words, $k \rightarrow$ 1000's.

Can make one more assumption for text.

(BAG OF WORDS) model.

$$P(x_1 | \text{class}) = P(x_2 | \text{class}) = \dots = P(x_k | \text{class}).$$

Assumptions :- All positions have same conditional distribution.

(BAD IDEA FOR IMAGES, WHY?)

- Estimation of probabilities.

$$\arg \max_i P(c_i) \cdot P(\bar{x}_{n+1} | c_i)$$

$$= \max_i P(c_i) \cdot \prod_{j=1}^k P(x_{n+1,j} | c_i).$$

- Take a feature ' i ', with ' F ' possible values.
eg, $F=2$, for binary, $F=3$ for outlook.

① $\rightarrow \bar{x}_{(1)}, \dots, \bar{x}_{(N)} \rightarrow$ all ~~the~~ examples from class c_i

② In these examples,

$$n_j = \# \text{ examples with feature } i = j.$$

$$n_0, n_1 = \# \text{ examples with } i^{\text{th}} \text{ feature} = \{0/1\}$$

n_0, n_1 :- # examples with i^{th} feature = $\{0/1\}$. (37)

Pick a class, c_i , and a feature, say the j^{th} feature x_j .

We want to estimate $P(\bar{x}_j = x \mid \text{Class is } c_i)$.

Ex:- MNIST, $c_i = 2$, $j = 192$, what is the distribution of the 192nd pixel given the digit is 2.

Maximum Likelihood Estimates / Empirical estimators:- (Bishop Chap. 3)
(Murphy Chap 3).

$$P(f(\bar{x}) = c_i) = \frac{\# \text{ examples with } f(\bar{x}) = c_i}{\# \text{ total examples}}.$$

$$P(\bar{x}_j = x \mid f(\bar{x}) = c_i) = \frac{\# \text{ examples with } \bar{x}_j = x \wedge f(\bar{x}) = c_i}{\# \text{ examples with } f(\bar{x}) = c_i}.$$

Why is it called MLE estimator?

- Insights into NB classifiers:-

- ① NB assumption very strong, and does not hold in practice
 - digit continuity.

still surprisingly good, why?

- ② In sufficient data / zero probability assignment.

- Suppose $x_j = x$, never happens in the training set.

$X \parallel$ eg, say digit = 2, pixel 192, never had '101' ($\in \{0, \dots, 255\}$).
but had ...

Then, for the MLE estimator if the training example

has ~~or~~ j^{th} feature equal to 'x', $P(\bar{X}_{n+1} | c_j) = 0$.

Straight away gone, no matter what else happens.

- One erroneous pixel / one wrong feature in test can make the probability zero.
 - Much more likely in words example.

~~Laplace Smoothing~~ Laplace Smoothing. 'Smoothing of probabilities'.

- (-) Suppose there are D choices/possibilities of a random variable, $D=2$, coins, $D=6$, die.

Add- β smoothing,

Suppose we have N outcomes, $n_x = \# x's$ in the outcomes

$$\text{MLE-Prob}(x) = \frac{n_x}{N}.$$

$$\text{Laplace-Prob}(x) = \frac{n_x + 1}{N + D}. \quad (\text{Show this is indeed a probability}).$$

In general, 'add- β ' estimators:-

$$\text{Add-}\beta \text{ Prob}(x) = \left(\frac{n_x + \beta}{N + D \cdot \beta} \right).$$

- Continuous random variables, Gaussian Naive Bayes.

Each continuous feature \bar{X}_j is a Gaussian random variable, conditioned on the class.

e.g., pixels,

estimating distribution of a feature x_j under c_i ,

(\rightarrow) Take all the examples $\bar{x}_{j(1)}, \dots, \bar{x}_{j(N_i)}$ that are in class c_i .
estimate:

$$\text{Mean} = \frac{\bar{x}_{j(1)} + \dots + \bar{x}_{j(N_i)}}{N_i}, = \hat{\mu}_{j|c_i}$$

$$\text{Variance} = \frac{1}{N_i} \sum_{t=1}^{N_i} (\bar{x}_{j(t)} - \hat{\mu}_{j|c_i})^2$$

 Use logarithms for computation \gg

Naive Bayes Classifiers

- ① Assumptions? / generative model
- ② Probability Estimation.
- ③ NB-assumption.

④

Tennis example :-

$$P(\text{Class} = Y) = \frac{9}{14}, P(\text{No}) = \frac{5}{14}.$$

$$R(\text{out } \beta = '1') , \quad \boxed{\text{Footlock} = 3}$$

$$P(\text{outlook} = S | \text{Class} = \text{No}) = \frac{3+1}{5+3} = \frac{1}{2} = \frac{4}{8} \quad \left(\frac{3}{5} \right)$$

$$P(\text{" " } = R | \text{Class} = \text{No}) = \frac{2+1}{5+3} = \frac{3}{8} \quad \left(\frac{2}{5} \right)$$

$$P(\text{" " } = O | \text{Class} = \text{No}) = \frac{0+1}{5+3} = \frac{1}{8} \quad (0).$$

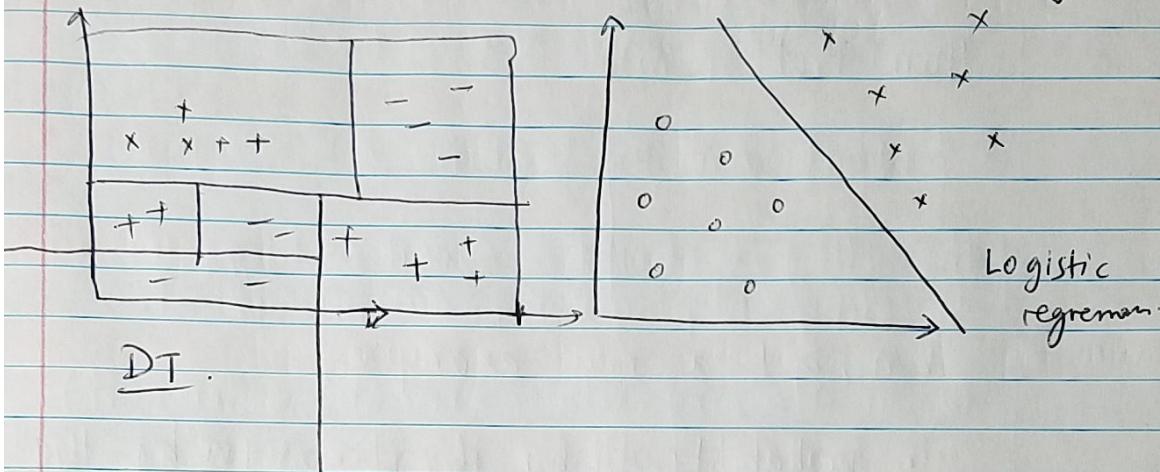
$$P(\text{Humidity} = H | \text{Class} = \text{No}) = \frac{4+1}{5+2} = \frac{5}{7} \quad \left(\frac{4}{5} \right)$$

$$P(\text{Humidity} = N | \text{Class} = \text{No}) = \frac{1+1}{5+2} = \frac{2}{7} \quad \left(\frac{1}{5} \right)$$

$$P(\text{@ Sunny, High temp, Normal Humidity, } \cancel{\text{High Strong Wind}} | \text{Class} = \text{No})$$

NB $= (\frac{4}{8} \cdot \frac{3}{8} \cdot \frac{2}{7} \cdot \frac{4}{7}) \cdot$

Discriminative vs Generative Learning



Generative models

(posi

y

breaking the "feature space" into regions.

(set of all possible features).

$$\text{Ex:- } N \rightarrow \{\text{Age}\} \times \{\text{Boy}\}$$

$$\in \mathbb{Z}_+ \quad \{0, 1\}$$

Depending on where the data lies, we assign appropriate values.

There is no 'model' for generating the examples.

Generative model :- algorithm.

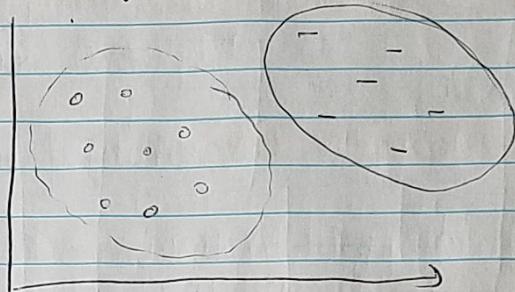
Look at one class at a time.

① say all the labels.

② Fit a probabilistic model to this part of the data.

* Detour a bit into estimating probabilities from data.
($P \rightarrow \underline{\quad}$).

Then do the same for -ve examples.



Then, get a new example,

- probability higher ~~if predicted by the~~ by the '+'ve' model / '-ve' model?

Discriminative :- \bar{x} $\xrightarrow[\text{Learn}]{\substack{\text{try to} \\ \text{find}}} f(\bar{x})$ (label).

Generative :- Learn $P(\bar{x}|f(\bar{x}))$ (generative model).
↓ ↓
features labels.

Generative classification steps

< Θ Find a generative model >

1. Using training examples "learn":-

$$P(\bar{x}|f(\bar{x})) \text{ and } p(f(\bar{x}))$$
$$\downarrow$$
$$P(\bar{x}|+) \& P(\bar{x}|-) , p(+)&p(-).$$