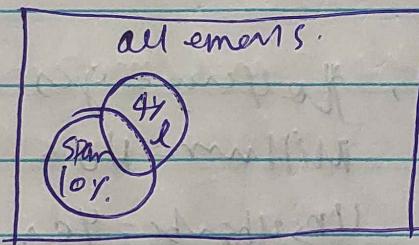


SMC

U-3

① Naive Bayes Classification.

- It is a classification algo based on bayes theorem
- It describes the prob. of an event based on prior knowledge
- It calc. the prob. of an event based on prior prob. of the event ~~and~~ and cond. prob. of the event.
- Using Joint prob., we can predict the event outcome
ex → The venn diag. showing the joint prob. of spam with lottery



→ lottery
→ not all spam msg contain word lottery and not every mail with lottery is spam

- Using cond. prob.; we can calc. the relationship b/w dependent ents. for exmp., A & B are two ents, calc. $P(A|B)$, i.e. ~~occurred~~ the prob. of event A given that fact that event B is already occurred.

$$\begin{aligned} P(A|B) &= \frac{P(B|A) * P(A)}{P(B)} \\ &= \frac{P(A \cap B)}{P(B)} \end{aligned}$$

Classification:

- construct a likelihood table for three words (w_1, w_2, w_3) for 100 emails

	<u>Catagry (w_1)</u>	<u>Milum (w_2)</u>	<u>Unsubscrib (w_3)</u>				
<u>labeled</u>	yes no	yes no	yes no	<u>total</u>			
spam	3/22	19/22	11/22	11/22	9/22	22	
ham	7/78	76/78	15/78	63/78	7/78	78	
<u>total</u>	5/100	95/100	26/100	14/100	34/100	66/100	100

Using Bayes Theorem,
 Category = yes
 $w_1/w_2/w_3 = \text{No}$
 Unsubscribe = Yes.

$$P(\text{Spam} | w_1 \cap w_2 \cap w_3) = \frac{P(w_1 | \text{spam}) * P(w_2 | \text{spam}) * P(w_3 | \text{spam})}{P(w_1) * P(w_2) * P(w_3)}$$

$$P(\text{Ham} | w_1 \cap w_2 \cap w_3) = \frac{P(w_1 | \text{ham}) * P(w_2 | \text{ham}) * P(w_3 | \text{ham})}{P(w_1) * P(w_2) * P(w_3)}$$

$$P(\text{Spam} | w_1 \cap w_2 \cap w_3) =$$

$$P(\text{Spam}) = \frac{3}{22} + \frac{11}{22} \times \frac{13}{22} + \frac{22}{100} = 0.00864.$$

$$P(\text{Ham} | w_1 \cap w_2 \cap w_3) =$$

$$\frac{2}{78} \times \frac{15}{78} \times \frac{21}{78} \times \frac{78}{100} = 0.004349$$

$$P(\text{Spam}) = \frac{0.00864}{(0.00864 + 0.004349)} = 0.67$$

$$P(\text{Ham}) = \frac{0.004349}{(0.00864 + 0.004349)} = 0.33$$

② Curse dimensionality with 1D, 2D, 3D example:

- Curse of dimensionality depends on KNN
- Let's select 60 Random points

1-D plot.

import numpy as np

import pandas as pd

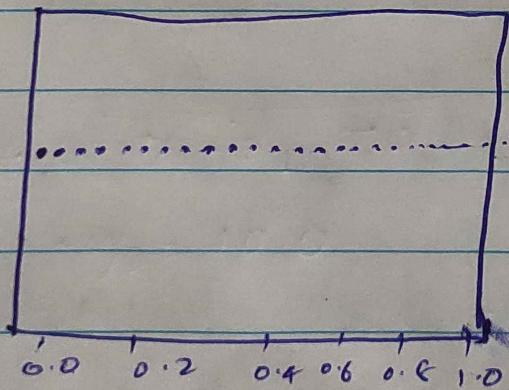
import matplotlib.pyplot as plt

```
one_d_data = np.random.rand(100, 1)
```

```
one_d_data_df = pd.DataFrame(one_d_data)
```

```
one_d_data_df.columns = ["1D_data"]
```

```
one_d_data_df["height"] = 1
```



60 data points reported in 1-D

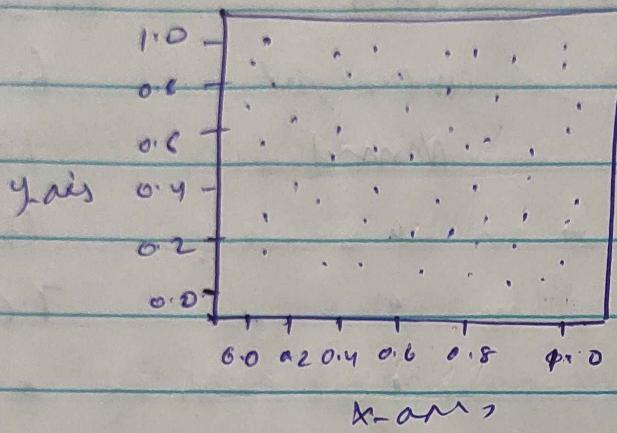
2-D plot

using x and y, plot them

```
two_d_data = np.random.rand(60, 2)
```

```
two_d_data_df = pd.DataFrame(two_d_data)
```

```
two_d_data_df.columns = ["x-axis", "y-axis"]
```



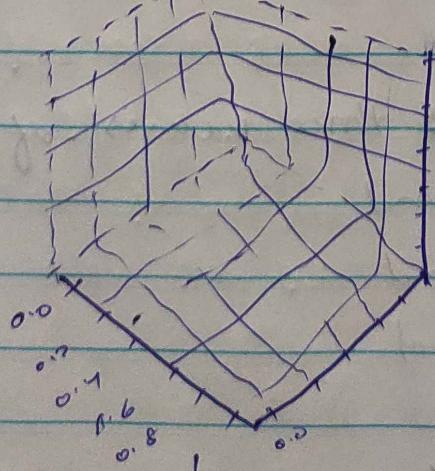
3-D plot

plotting 3D Space of 60 data pts

```
three_d_data = np.random.rand(60, 3)
```

```
three_d_data_df = pd.DataFrame(three_d_data)
```

```
three_d_data_df.columns = ["x axis", "y axis", "z axis"]
```



60 random pts
on 3D

③ KNN classifies with breast cancer WISconsin death example:

115kg 170cm

<u>wt</u> (x_1)	<u>ht</u> (y_1)	<u>class</u>	<u>Eucleandist</u>
51	167	underweight	6.7
62	182	Normal	13
69	176	"	13.41
64	173	"	7.615
65	172	"	8.24
56	174	underweight	4.12
58	169	Normal	(1.41) N ₁
57	173	"	(3) N ₃
55	170	"	(2) N ₂

$$\text{Eucleandist} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$x_2 = 57, y_2 = 100$$

Select the least ~~is~~ three values of K .

$$N_1 = 1.41 = \text{normal}$$

$$N_2 = 2 = \text{normal}$$

$$N_3 = 3 = \text{normal}$$

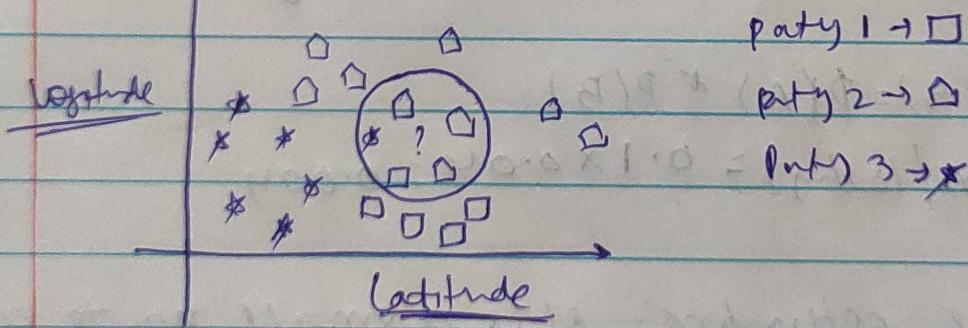
Person is normal.

⑤ KNN

- Non-parametric ml model
- Instance-based learning
- known as lazy learning as it does not learn during the training phase, it starts working only during test/evaluation phase for comparison.

KNN voter example:

- Objective is to predict the party
- Voters will vote based on their neighbourhood
- Voters will vote based on majority voters did for that particular party.
- Tuning the K value:



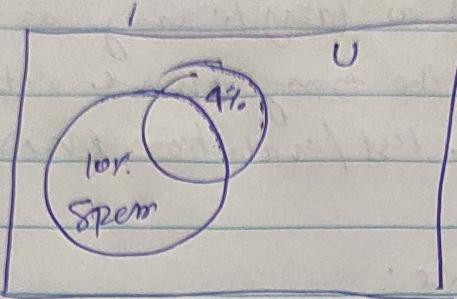
Here, the one voter will vote party 3

one ————— 1
3 ————— party 2

so, the ? can be vote for party 2 as they have majority

⑥ Joint probability.

- two or more events occurring together
- It is the prob. of intersection of two or more events
-



4% \rightarrow lottery

10% \rightarrow spam

\cap \rightarrow elements

- If email msg present in the world, which is likely to be lottery or spam.
- This venn-diag. indicates the joint prob. of spam with lottery.
- It shows all spam msg contains word lottery and not every email with the word lottery is spam.

$$P(A \cap B) = P(A) * P(B)$$

$$P(\text{spam} \cap \text{lottery}) = 0.1 \times 0.04 = 0.004$$

- ## ⑦ Laplace estimator:
- adds a small no. to each of the counts in the freq. table
 this ensures that each feature has non-zero occurring prob.

: Usually Laplace estimator set to 1

ex $P(\text{Spam} | w_1, w_2, w_3) = \frac{3}{22} * \frac{11}{22} * \frac{0}{22} * \frac{22}{100} = 0$
 using Laplace estimator ~~for~~ for $w_3 = 1$

$$P(\text{Spam} | w_1, w_2, w_3) = \frac{3}{25} * \frac{11}{28} * \frac{1}{25} * \frac{22}{100} = 0$$

③ Curse of Dimensionality

- ↳ phenomenon where the performance of the ML algs degrades as no. of features in data increases
- ↳ As features ↑, amount of data req. to maintain the same level grows exponentially.
- ↳ It is because, the no. of dimensions increasing, the volume of feature grows exponentially and data points become ~~more~~ increasingly sparse.
- ↳ It affects many ML algs like clustering, classification and regression.
- ↳ It can lead to overfitting also.
- ↳ ~~To understand these processes~~
- ↳ To check these curse of dimensionality, we can use 1D, 2D or 3D space also.