

Author - Kanav Bansal
You can find me on **Linkedin** -
www.linkedin.com/in/kanavbansal
Or <http://www.thataiguy.com>

Building a Machine Learning Model

1. Training



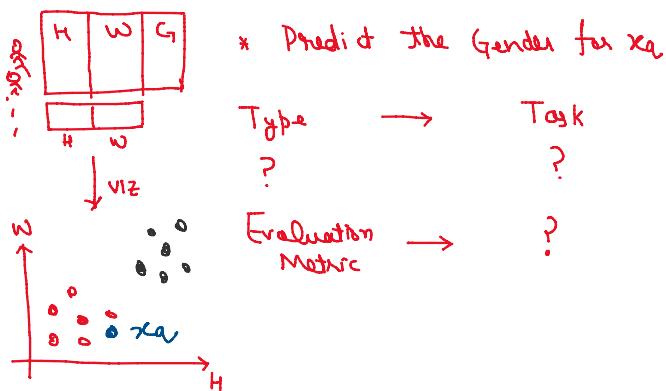
2. Predicting using unseen data



3. Evaluation



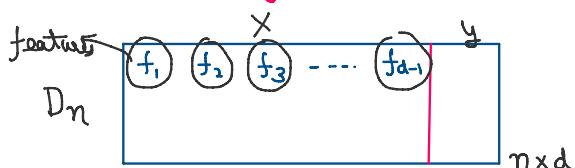
- Given Data containing Height, Weights & Gender of some individuals.



1. Target Variable → Gender (y)
 $\text{Gender} \in \{\text{M, F}\}$
 2. Task → Classification
 3. Evaluation Metrics → Accuracy
 4. Algorithms → Naive Bayes Algorithm
(How?) Probabilistic Approach

Naive Bayes Algorithm (AKA Maximum A Posteriori)

Type - Supervised Learning
Task - Classification.

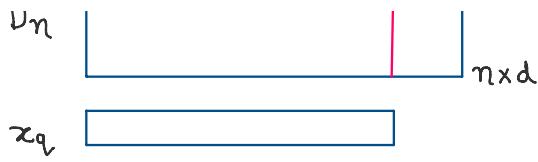


Given $\rightarrow D_n = \{ (x_i, y_i) \}_{i=1}^n \mid x_i \in \mathbb{R}^{d-1}, y_i \in \{c_1, c_2, \dots, c_k\}$

$x_q = f_1 n_{f_2} n_{f_3} \dots n_{f_{d-1}}$

Task \rightarrow Find y_q

How? \rightarrow Given x_q , compute the prob of x_q belonging to each of the classes.
Whichever is maximum, assign that class to x_q .



Whichever is maximum, assign that class to x_q .

i.e. $P(C_1|x_q) P(C_2|x_q) \dots P(C_k|x_q)$

$$x_q = \operatorname{arg\ max}_{C_1, C_2, \dots, C_k} \left\{ \frac{P(C_i|x_q)}{\text{Posterior Prob.}} \right\}$$

From Bayes Theorem \rightarrow

$$\underbrace{P(C_i|x_q)}_{\text{Posterior}} = \frac{\underbrace{P(x_q|C_i) * P(C_i)}_{\text{Likelyhood}}} {\underbrace{P(x_q)}_{\text{Marginal}} \text{ Prior}}$$

Key Concepts \rightarrow

\rightarrow Probability (Study of uncertainty)

Random Experiment

- It is a process for which outcome cannot be predicted with certainty.

Example \rightarrow

R.E \rightarrow Tossing a coin

Sample Space

- It is a set of all possible outcomes of a 'RANDOM EXPERIMENT'.

Example \rightarrow

S.S = {H, T}

Event

- It is a subset of a 'SAMPLE SPACE'.

Example \rightarrow

Event - Getting a Head
{H}

\rightarrow Probability \rightarrow

$$P(\text{Event}) = \frac{\text{No. of favourable Outcomes}}{\text{Total no. of Outcomes}} = \frac{|\text{Event}|}{|\text{Sample Space}|}$$

\rightarrow Example-1 :

R.E \rightarrow Tossing 2 coins

S.S \rightarrow {HH, HT, TH, TT}

Event I,

A : Getting atleast 1 Head

A \rightarrow {HH, HT, TH}

$$P(A) = 3/4$$

\rightarrow Example-2 :

R.E \rightarrow Rolling a dice

S.S \rightarrow {1, 2, 3, 4, 5, 6}

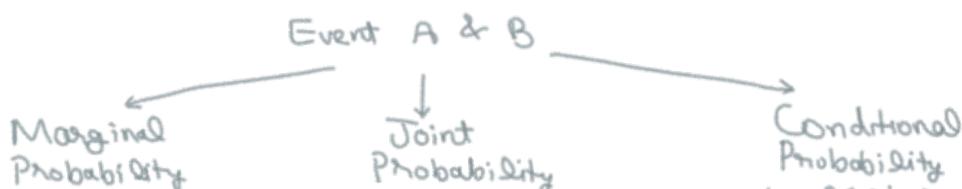
Event J,

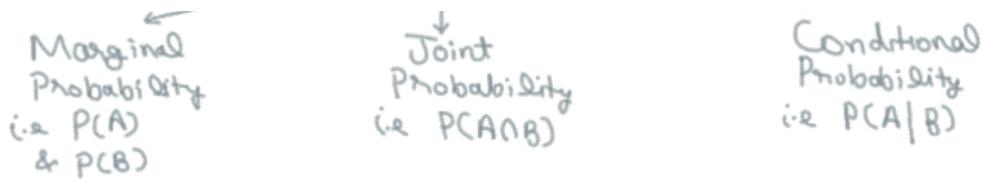
B : Getting a Prime No.

B \rightarrow {2, 3, 5}

$$P(B) = 3/6$$

\rightarrow In case of two or more events \rightarrow





- * Next we will see how to compute these probabilities with the help of 'FREQUENCY TABLE'.

Consider this binary classification data set:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Frequency Table

For Outlook vs PlayTennis

	Yes	No	Joint Freq.
Sunny	2	3	
	4	0	4
	3	2	5
	9	5	Marginal Frequencies
Overcast			
Rainy			Marginal Frequencies

Marginal Prob →

$$P(\text{PlayTennis} = \text{Yes}) = \frac{9}{14}$$

$$P(\text{PlayTennis} = \text{No}) = \frac{5}{14}$$

$$P(\text{Outlook} = \text{Sunny}) = \frac{5}{14}$$

Joint Prob →

$$P(\text{Outlook} = \text{Sunny} \cap \text{PlayTennis} = \text{Yes}) = \frac{2}{14}$$

$$P(\text{Outlook} = \text{Overcast} \cap \text{PlayTennis} = \text{No}) = \frac{0}{14}$$

$$P(\text{Outlook} = \text{Rainy} \cap \text{PlayTennis} = \text{No}) = \frac{2}{14}$$

Conditional Prob →

$$P(\text{Outlook} = \text{Sunny} \mid \text{PlayTennis} = \text{Yes}) = \frac{P(A \cap B)}{P(B)}$$

$$= \frac{2/14}{9/14} = \frac{2}{9}$$

Bayes Theorem →

$$\text{Acc}^n \text{ to Conditional Prob} \rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \quad -\textcircled{1}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$\therefore B \cap A = A \cap B$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A \cap B) = P(B|A)P(A) \quad -\textcircled{2}$$

Substituting $\textcircled{2}$ in eqⁿ $\textcircled{1}$: Likelyhood

Substituting ② in eqⁿ ①:

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

↑ ↓
Posteriori Marginal

Likelihood Prior

Independent Events \rightarrow

If given events A & B are known to be independent:

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

Similarly:

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

Data \rightarrow

$$D_n = \left\{ (x_i, y_i) \mid x_i \in \mathbb{R}^{d-1}, y_i \in \{0, 1\} \right\}$$

f_1 f_2 \vdots f_{d-1}	x_1 x_2 \vdots x_n	y_1 y_2 \vdots y_n
---	-------------------------------------	-------------------------------------

$$x_i = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{d-1} \end{bmatrix}$$

Naive Bayes Algorithm

Task \rightarrow Classification

How?

Given x_q , compute $P(c_1|x_q), P(c_2|x_q), \dots, P(c_k|x_q)$.

Whichever probability is maximum, prediction would be the corresponding class.

$$\hat{y}_q = \arg \max_{i \in \{1, 2, \dots, k\}} \{P(c_i|x_q)\} \quad \text{ie maximum a posteriori}$$

$$\rightarrow P(c_1|x_q) = \frac{P(x_q|c_1) P(c_1)}{P(x_q)} \Rightarrow P(c_1|x_q) \propto P(x_q|c_1) P(c_1)$$

$$\rightarrow P(c_2|x_q) = \frac{P(x_q|c_2) P(c_2)}{P(x_q)} \Rightarrow P(c_2|x_q) \propto P(x_q|c_2) P(c_2)$$

⋮

⋮

$$\rightarrow P(c_k|x_q) = \frac{P(x_q|c_k) P(c_k)}{P(x_q)} \Rightarrow P(c_k|x_q) \propto P(x_q|c_k) P(c_k)$$

$$\rightarrow P(C_k | X_q) = \frac{P(X_q | C_k) P(C_k)}{P(X_q)} \Rightarrow P(C_k | X_q) \propto P(X_q | C_k) P(C_k)$$

Naive Bayes Derivation

Given $X_q = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{d-1} \end{bmatrix}$ i.e. X_q is $f_1 \cap f_2 \cap f_3 \cap \dots \cap f_{d-1}$

Using Bayes Theorem:

$$P(C_k | X_q) = \frac{P(X_q | C_k) * P(C_k)}{P(X_q)}$$

Compute $\rightarrow P(C_k | X_q) \propto P(X_q | C_k) P(C_k)$

$$\hat{y}_q = \arg \max_{k \in \{1, 2, \dots, K\}} \left\{ P(C_k | X_q) \right\}$$

$$\hat{y}_q = \arg \max_{k \in \{1, 2, \dots, K\}} \left\{ P(X_q | C_k) P(C_k) \right\}$$

↳ Let's zoom inside this term and understand how to compute it step by step.

According to Conditional Probability ↴

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$\Rightarrow P(A \cap B) = P(A | B) P(B)$$

$P(X_q | C_k) P(C_k)$ can be written as $P(X_q \cap C_k)$

We have already seen that X_q is $f_1 \cap f_2 \cap \dots \cap f_{d-1}$

$$\rightarrow P(X_q | C_k) P(C_k)$$

$$\rightarrow P(X_q \cap C_k)$$

$$\rightarrow P(f_1 \cap f_2 \cap f_3 \cap \dots \cap f_{d-1} \cap C_k)$$

$$P(f_1 \cap f_2 \cap f_3 \dots \cap f_{d-1} \cap c_k)$$

A

B

$$\text{Using } P(A \cap B) = P(A|B) P(B)$$

$$\rightarrow P(f_1 | f_2 \cap f_3 \dots \cap f_{d-1} \cap c_k) * \\ P(f_2 \cap f_3 \dots \cap f_{d-1} \cap c_k)$$

Expanding this expressing using
chain Rule of conditional Probability

$$\rightarrow P(f_1 | f_2 \cap f_3 \dots \cap f_{d-1} \cap c_k) * \\ P(f_2 | f_3 \cap f_4 \dots \cap f_{d-1} \cap c_k) * \\ \vdots \\ P(f_{d-1} | c_k) * P(c_k)$$

\therefore Above expression is extremely hard to compute.

To simplify the computation lets take a
NAIVE assumption that:

Features are conditionally independent
of each other.

i.e if $f_1, f_2, f_3, \dots, f_{d-1}$ are independent
of each other,

$$P(f_1 | f_2 \cap f_3 \dots \cap f_{d-1} \cap c_k) = P(f_1 | c_k)$$

$$\rightarrow P(f_1 | c_k) * P(f_2 | c_k) * \dots * P(f_{d-1} | c_k) * P(c_k)$$

$$\rightarrow P(c_k) * \prod_{i=1}^{d-1} P(f_i | c_k)$$

$$\hat{y}_a = \arg \max_{j \in \{1, 2, \dots, k\}} \left\{ P(x_a | c_j) P(c_j) \right\}$$

↓

$$P(c_j) * \prod_{i=1}^{d-1} P(f_i | c_j)$$

$$\text{i.e. } \hat{y}_a = \arg \max_{j \in \{1, 2, \dots, k\}} \left\{ P(c_j) * \prod_{i=1}^{d-1} P(f_i | c_j) \right\}$$

Flavours of Naive Bayes

1. Gaussian Naive Bayes

If data contains numerical features use Gaussian NB.

Assumptions :

- (a) X contains features which are conditionally independent
- (b) All the numerical features follows Normal Distribution

$$P(f_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} * \exp \left\{ -\frac{(f_i - \mu_y)^2}{2\sigma_y^2} \right\}$$

2. Multinomial Naive Bayes

If data contains text use Multinomial NB.

3. Bernoulli Naive Bayes

If data contains Categorical features use Bernoulli NB.

Issues with Naive Bayes

1. Numerical Instability (For all flavours of NB)

2. Zero probability issue. (For MultinomialNB & BernoulliNB)

In x_q , if a feature contains a value that was never observed in the historical data, we will encounter a zero probability issue.

Diagram illustrating the zero probability issue:

f_1	f_2	f_3	$y \in \{0, 1\}$
\uparrow	\uparrow	\uparrow	$\rightarrow f_3 \in \{b, q\}$
\uparrow	\uparrow	\uparrow	$f_2 \in \{a, b\}$
\uparrow	\uparrow	\uparrow	$f_1 \in \{l, m\}$

$x_q \rightarrow f_1 = m, f_2 = c, f_3 = b$

Acc to NB \uparrow

$$\hat{y}_q = \arg \max_{i=1, 2, \dots, k} \left\{ P(C_i) * \prod_{j=1}^{d-1} P(f_j | C_i) \right\}$$

i.e Find $P(C_1 = 0 | x_q) \propto P(C_1 = 0) * P(f_1 | C_1) * P(f_2 | C_1) * P(f_3 | C_1)$

& $P(C_2 = 1 | x_q) \propto P(C_2 = 1) * P(f_1 | C_2) * P(f_2 | C_2) * P(f_3 | C_2)$

Observe that $\rightarrow P(f_2 | C_1) = 0$
& $P(f_2 | C_2) = 0$

which means $\rightarrow P(C_1=0 | x_q) = 0$
 $\& P(C_2=1 | x_q) = 0$

Now, how do we decide the class of x_q ?

Solution \rightarrow Apply regularization technique.

{ Here we will use \rightarrow Laplace Smoothing }

i.e Whenever the $P(f_j | C_i) = 0$, apply a smoothening technique so that the value is not exactly equal to 0.

$$\text{Laplace Smoothing} \Rightarrow \frac{P(f_j | C_i) + \alpha}{n + k + \alpha}$$

Here, $\alpha = \text{Hyperparameter}$
 $\alpha \in [0, +\infty]$

$$n = P(C_i)$$

$$k = \# \text{ of unique values in } f_j$$