

Naïve Bayes

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



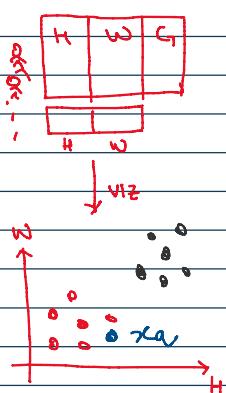
2. Predicting using Unseen data



3. Evaluation



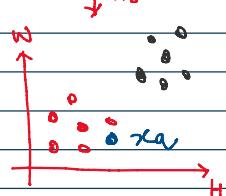
- Given Data containing Height, Weights & Gender of some individuals.



* Predict the Gender for x_q

Type \rightarrow Task

Evaluation Metric \rightarrow ?



- Target Variable \rightarrow Gender (y)

Gender $\in \{M, F\}$

- Task \rightarrow Classification

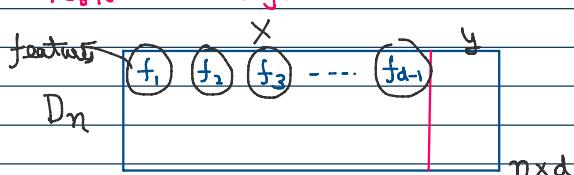
- Evaluation Metrics \rightarrow Accuracy

- (How?) Probabilistic Approach

Algorithms \rightarrow Naïve Bayes Algorithm

Naïve Bayes Algorithm (AKA Maximum A Posteriori)

Type - Supervised Learning
Task - Classification.



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

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

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 .

i.e. $D_n \sim \mathcal{N}(x_q, \Sigma)$

U_n | | | | | $n \times d$

x_q | | | |

Whichever is maximum, assign that class to x_q

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

$$y_q = \operatorname{argmax}_{C_1, C_2, \dots, C_k} \left\{ P(C_i|x_q) \right\}$$

Posterior Prob.

From Bayes Theorem \rightarrow

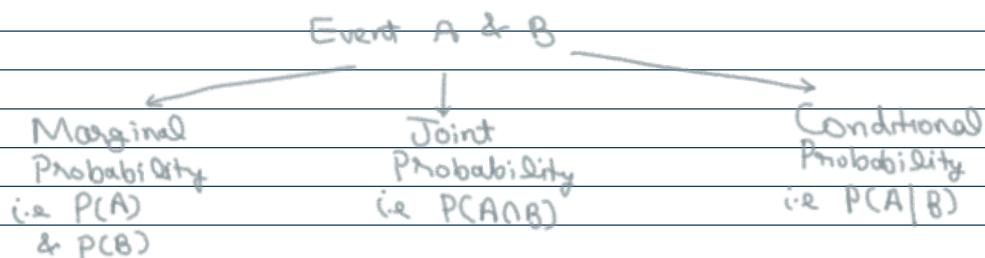
$$\frac{P(C_i|x_q)}{\text{Posterior}} = \frac{\underbrace{P(x_q|C_i) * P(C_i)}_{\text{Likelihood}}} {\underbrace{P(x_q)}_{\text{Marginal}} \underbrace{P(C_i)}_{\text{Prior}}}$$

Key Concepts \rightarrow

\rightarrow Probability

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

\rightarrow In case of two or more events \rightarrow



Bayes Theorem \rightarrow

$$\text{Acc}^o \text{ to Conditional Prob} \rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)}$$

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

\downarrow

Posterior

Likelihood

Prior

Marginal

Independent Events \rightarrow

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

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

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

Similarly:

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

Data →

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

f_1, f_2, \dots, f_{d-1}	y
$\leftarrow x_1 \rightarrow$	y_1
$\leftarrow x_2 \rightarrow$	y_2
!	:
$\leftarrow x_n \rightarrow$	y_n

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

Naive Bayes Algorithm

Task → 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.

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

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

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

⋮
⋮

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

Naive Bayes Derivation ↗

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

Using Bayes Theorem:

$$P(c_k|x_q) = \frac{P(x_q|c_k) * P(c_k)}{P(x_q)}$$

Compute → $P(c_k|x_q) \propto P(x_q|c_k) P(c_k)$

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

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

↳ 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 $f_1 \cap f_2 \cap \dots \cap f_{d-1}$

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

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

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

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

A B

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

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

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

A B

Expanding this expressing using
Chain Rule of conditional Probability

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

$$P(f_2 | f_3 \cap f_4 \cap \dots \cap f_{d-1} \cap C_k) *$$

:

$$P(f_{d-1} | C_k) * P(C_k)$$

∴ 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,

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

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

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

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

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

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

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)

* Coming Soon * → Log transformation

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.

$$\begin{array}{|c|c|c|c|} \hline f_1 & f_2 & f_3 & \rightarrow y \in \{0, 1\} \\ \hline \end{array}$$

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

$\uparrow \quad \uparrow \quad \uparrow$

$f_1 \in \{0, m\}$ $f_2 \in \{a, b\}$ $f_3 \in \{b, q\}$

Acc to NB ↴

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

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 → Apply regularization technique.

{ Here we will use → Laplace Smoothing }

i.e Whenever the $P(f_j | C_i) = 0$, apply a smoothing 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, $\alpha = \text{Hyperparameter}$
 $\alpha \in [0, +\infty]$

$$n = P(C_i)$$

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