



**Dr. Kamlesh Tiwari**

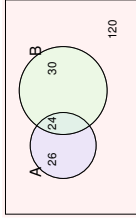
Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

March 08, 2021

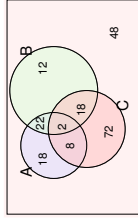
**ONLINE** (Campus @ BITS-Pilani Jan-May 2021)

<http://ktiwari.in/ml>

### Probability: Conditional independence



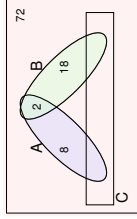
$$\begin{aligned} P(A) &= 50/200 = 0.25 \\ P(B) &= 54/200 = 0.27 \\ P(A) \times P(B) &= 0.0675 \\ P(A \cap B) &= 24/200 = 0.12 \neq P(A) \times P(B) \end{aligned}$$



$$\begin{aligned} P(A|C) &= 10/100 = 0.1 \\ P(B|C) &= 20/100 = 0.2 \\ P(A) \times P(B|C) &= 2/100 = 0.02 = 0.02 = P(A) \times P(B|C) \end{aligned}$$

Independence is obtained by C

Independence does not mean conditional independence



- A, and B are independent due to the construction.
- But, in C. Happening of one completely rules out the other.

### Bayes Optimal Classifier

$$\operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D)$$

$$V = \{\oplus, \ominus\}$$

$$\begin{aligned} P(h_1|D) &= 0.4 & P(\ominus|h_1) &= 0 & P(\oplus|h_1) &= 1 \\ P(h_2|D) &= 0.3 & P(\ominus|h_2) &= 1 & P(\oplus|h_2) &= 0 \\ P(h_3|D) &= 0.3 & P(\ominus|h_3) &= 1 & P(\oplus|h_3) &= 0 \end{aligned}$$

Therefore,

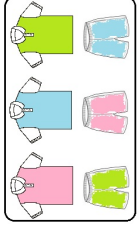
$$\sum_{h_i \in H} P(\oplus|h_i)P(h_i|D) = 0.4 \quad \sum_{h_i \in H} P(\ominus|h_i)P(h_i|D) = 0.6$$

and

$$\operatorname{argmax}_{v_j \in \{\oplus, \ominus\}} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D) = \ominus$$

This type of classifier is called a **Bayes optimal classifier**, or Bayes optimal learner.

### Probability: Independence



Assume there are three T-shirts  $T_r, T_g, T_b$  and three shorts  $S_r, S_g, S_b$  in a shop

What is the probability that someone would purchase  $(T_r, S_r)$

#### Independence

Selling any T-shirt or shorts

$$\begin{aligned} Pr(T_r, S_r) &= Pr(T_r) \times Pr(S_r) \\ &= 1/3 \times 1/3 = 1/9 \end{aligned}$$

#### Dependence

Selling T-shirt and shorts in set of same color

$$Pr(T_r, S_r) = Pr(T_r \& S_r) = 1/3$$

There is something called **conditional independence**

### Bayes Optimal Classifier

**Switching the question**, from "which is most probable hypothesis?" to "what is the most probable classification of the new instance?"

Is it possible to do better than MAP?

**Example:** Let posterior probabilities of three hypotheses  $h_1, h_2, h_3$  given the training data are 0.4, 0.3, and 0.3 (obviously  $h_1$  is MAP)

- Let classification of a new instance  $x$  is **positive** by  $h_1$  and **negative** by  $h_2$  and  $h_3$
- By taking all hypotheses into account, the probability that  $x$  is positive is 0.4, and negative is 0.6
  - ▶ Most probable classification is **negative** and it differs from MAP

#### Bayes optimal classification:

$$\operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D)$$

where classification  $v_j \in V$  and  $P(v_j|D)$  is the correct classification

### Bayes Optimal Classifier

- **Best Classifier:** No other classification method using the same hypothesis space and same prior knowledge can outperform this method on an average.
- This method maximizes the probability that the new instance is classified correctly, given the available data, hypothesis space, and prior probabilities over the hypotheses
- **Note** that the predictions made by Bayes optimal classifier may not be contained in  $H$ . (It is possible that there is no hypothesis in  $H$  giving same classification as Bayes optimal classifier.)

**Issue:** Although the Bayes optimal classifier obtains the best performance that can be achieved from the given training data, it can be quite costly to apply!

<sup>1</sup>we could have infinitely many hypothesis

## GIBBS Algorithm 2

A less optimal method is the Gibbs algorithm

- For a new instance  $x$ 
  - 1 Choose a hypothesis  $h \in H$  at random, according to the posterior probability distribution over  $H$
  - 2 Use  $h$  to predict the classification of the instance  $x$

### Importance

Under certain conditions the expected misclassification error for the Gibbs algorithm is at most twice the expected error of the Bayes optimal classifier

<sup>2</sup> Oppor, Manfred and Haussler, David. "Generalization performance of Bayes optimal classification algorithm for learning a percepton." in Physical Review Letters. 68(20), pp-2677. AFS-1991

## Naive Bayes Classifier

Attribute values are conditionally independent given the target value

- Under this assumption,
- Given a target value, the probability of observing the conjunction  $\langle a_1, a_2, \dots, a_n \rangle$  is just the product of the probabilities.

$$P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$$

### Naive Bayes classifier

is the one which

$$\operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

**Caution:** "conditionally independent" assumption may NOT be true

## Example: Naive Bayes Classification

Day	Outlook	Temperature	Humidity	Wind	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Rainy	Cool	Normal	Strong	No
D4	Rainy	Mild	High	Strong	No
D14	Rainy	Mild	High	Strong	No

Day	Outlook	Temperature	Humidity	Wind	Play
D3	Overcast	Hot	High	Weak	Yes
D4	Rainy	Mild	High	Weak	Yes
D5	Sunny	Cool	Normal	Weak	Yes
D6	Overcast	Hot	Normal	Strong	Yes
D8	Sunny	Cool	Normal	Weak	Yes
D9	Rainy	Mild	Normal	Weak	Yes
D10	Rainy	Mild	Normal	Strong	Yes
D11	Sunny	Mild	High	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes

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

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

### Outlook

	Yes	No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rainy	3/9	2/5

## Naive Bayes Classifier

Bayes classifier is a highly practical Bayesian learning method

- In some domains, its performance found to be comparable to neural network and decision tree
- The Bayesian approach to classify a new instance is to assign the most probable target value describing the instance
- We can use Bayes theorem to rewrite this expression as

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n) \quad (1)$$

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j) P(v_j) \quad (1)$$

Naive Bayes has assumption is that the attribute values are conditionally independent given the target value

## Example: Naive Bayes Classification

Given the data

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

Determine classification for  $\langle \text{Rainy}, \text{Hot}, \text{High}, \text{Strong} \rangle$

## Example: Naive Bayes Classification

$$P(\text{Yes}) = 9/14 \quad P(\text{No}) = 5/14$$

### Outlook

	Yes	No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rainy	3/9	2/5

### Humidity

	Yes	No
High	3/9	4/5
Low	6/9	1/5

### Temperature

	Yes	No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

### Wind

	Yes	No
Strong	3/9	3/5
Weak	6/9	2/5

## Example: Naive Bayes Classification

For  $x = \langle \text{Rainy, Hot, High, Strong} \rangle$

### P(Yes)

- $P(\text{Yes}) \times P(x|\text{Yes})$
- $P(\text{Yes}) \times [P(\text{Rainy}|\text{Yes}) \times P(\text{Hot}|\text{Yes}) \times P(\text{High}|\text{Yes}) \times P(\text{Strong}|\text{Yes})]$
- $9/14 \times [3/9 \times 2/9 \times 3/9 \times 3/9]$
- 0.005291...

### P(No)

- $P(\text{No}) \times P(x|\text{No})$
- $P(\text{No}) \times [P(\text{Rainy}|\text{No}) \times P(\text{Hot}|\text{No}) \times P(\text{High}|\text{No}) \times P(\text{Strong}|\text{No})]$
- $5/14 \times [2/5 \times 2/5 \times 4/5 \times 3/5]$
- 0.027428...

So the classification of the new data item  $x$  is

Thank You!

## Laplace smoothing

Consider data =

[apple,mango,mango,apple,banana,apple,apple,mango,apple,mango,apple,mango]

- Here there are three class in the data {mango,apple,banana}
- So  $p(\text{apple})=5/10$ ,  $p(\text{mango})=4/10$ ,  $p(\text{banana})=1/10$
- Can we add some hypothetical points in data? say  $k = 1$  we are adding one instance of all the items in the data. then the probabilities are  $p(\text{apple})=6/13$ ,  $p(\text{mango})=5/13$ ,  $p(\text{banana})=2/13$
- This method is called **Laplace smoothing**
- Very useful for the cases when number of classes is large and data size is small. Due to this we may miss some of the classes appearing in the data. See the example below

	Yes	No	Outlook	Yes	No
Sunny	2/9	3/5	Sunny	3/12	4/8
Overcast	4/9	0/5	Overcast	5/12	1/8
Rainy	3/9	2/5	Rainy	4/12	3/8

Thank you very much for your attention! (Reference<sup>3</sup>)

Queries ?

<sup>3</sup> [1] Book - Machine Learning, ch-6, Tom M. Mitchell.