

# Naive Bayes Classifier

# Bayesian Methods

- Learning and classification methods based on probability theory.
- Bayes theorem plays a critical role in probabilistic learning and classification.
- Uses *prior* probability of each category given no information about an item.
- Categorization produces a *posterior* probability distribution over the possible categories given a description of an item.

# Basic Probability Formulas

- Product rule

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

- Sum rule

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

- Bayes theorem

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- Theorem of total probability, if event  $A_i$  is mutually exclusive and probability sum to 1

$$P(B) = \sum_{i=1}^n P(B | A_i)P(A_i)$$

# Bayes Theorem

- Given a hypothesis  $h$  and data  $D$  which bears on the hypothesis:

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- $P(h)$ : independent probability of  $h$ : *prior probability*
- $P(D)$ : independent probability of  $D$
- $P(D|h)$ : conditional probability of  $D$  given  $h$ : *likelihood*
- $P(h|D)$ : conditional probability of  $h$  given  $D$ : *posterior probability*

# Does patient have cancer or not?

- A patient takes a lab test and the result comes back positive. It is known that the test returns a correct positive result in only 99% of the cases and a correct negative result in only 95% of the cases. Furthermore, only 0.03 of the entire population has this disease.
  1. What is the probability that this patient has cancer?
  2. What is the probability that he does not have cancer?
  3. What is the diagnosis?

# Maximum A Posterior

- Based on Bayes Theorem, we can compute the *Maximum A Posterior* (MAP) hypothesis for the data
- We are interested in the best hypothesis for some space  $H$  given observed training data  $D$ .

$$\begin{aligned}h_{MAP} &\equiv \operatorname{argmax}_{h \in H} P(h | D) \\ &= \operatorname{argmax}_{h \in H} \frac{P(D | h)P(h)}{P(D)} \\ &= \operatorname{argmax}_{h \in H} P(D | h)P(h)\end{aligned}$$

$H$ : set of all hypothesis.

Note that we can drop  $P(D)$  as the probability of the data is constant (and independent of the hypothesis).

# Maximum Likelihood

- Now assume that all hypotheses are equally probable a priori, i.e.,  $P(h_i) = P(h_j)$  for all  $h_i, h_j$  belong to  $H$ .
- This is called assuming a *uniform prior*. It simplifies computing the posterior:

$$h_{ML} = \arg \max_{h \in H} P(D | h)$$

- This hypothesis is called the *maximum likelihood hypothesis*.

## Desirable Properties of Bayes Classifier

- *Incrementality*: with each training example, the prior and the likelihood can be updated dynamically: flexible and robust to errors.
- *Combines prior knowledge and observed data*: prior probability of a hypothesis multiplied with probability of the hypothesis given the training data
- *Probabilistic hypothesis*: outputs not only a classification, but a probability distribution over all classes

# Bayes Classifiers

**Assumption:** training set consists of instances of different classes described  $c_j$  as conjunctions of attributes values

**Task:** Classify a new instance  $d$  based on a tuple of attribute values into one of the classes  $c_j \in C$

**Key idea:** assign the most probable class  $c_{MAP}$  using Bayes Theorem.

$$\begin{aligned}c_{MAP} &= \operatorname{argmax}_{c_j \in C} P(c_j | x_1, x_2, \dots, x_n) \\ &= \operatorname{argmax}_{c_j \in C} \frac{P(x_1, x_2, \dots, x_n | c_j)P(c_j)}{P(x_1, x_2, \dots, x_n)} \\ &= \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j)P(c_j)\end{aligned}$$

# Parameters estimation

- $P(c_j)$ 
  - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_n | c_j)$ 
  - $O(|X|^n \cdot |C|)$  parameters
  - Could only be estimated if a very, very large number of training examples was available.
- **Independence Assumption**: attribute values are conditionally independent given the target value: **naïve Bayes**.

$$P(x_1, x_2, \dots, x_n | c_j) = \prod_i P(x_i | c_j)$$

$$c_{NB} = \arg \max_{c_j \in C} P(c_j) \prod_i P(x_i | c_j)$$

# Properties

- Estimating  $P(x_i | c_j)$  instead of  $P(x_1, x_2, \dots, x_n | c_j)$  greatly reduces the number of parameters (and the data sparseness).
- The learning step in Naïve Bayes consists of estimating and based on the frequencies in the training data
- An unseen instance is classified by computing the class that maximizes the posterior  $P(c_j | x)$
- When conditioned independence is satisfied, Naïve Bayes corresponds to MAP classification.

# Example. 'Play Tennis' data

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Question: For the day <sunny, cool, high, strong>, what's the play prediction?

# Naive Bayes solution

Classify any new datum instance  $\mathbf{x}=(a_1, \dots, a_T)$  as:

$$h_{Naive\ Bayes} = \arg \max_h P(h)P(\mathbf{x} | h) = \arg \max_h P(h) \prod_t P(a_t | h)$$

- To do this based on training examples, we need to estimate the parameters from the training examples:
  - For each target value (hypothesis)  $h$

$$\hat{P}(h) := \text{estimate } P(h)$$

- For each attribute value  $a_t$  of each datum instance

$$\hat{P}(a_t | h) := \text{estimate } P(a_t | h)$$

Based on the examples in the table, classify the following datum  $\mathbf{x}$ :

$\mathbf{x}=(\text{Outl}=\text{Sunny}, \text{Temp}=\text{Cool}, \text{Hum}=\text{High}, \text{Wind}=\text{strong})$

- That means: Play tennis or not?

$$h_{NB} = \arg \max_{h \in [\text{yes}, \text{no}]} P(h)P(\mathbf{x} | h) = \arg \max_{h \in [\text{yes}, \text{no}]} P(h) \prod_t P(a_t | h)$$

$$= \arg \max_{h \in [\text{yes}, \text{no}]} P(h)P(\text{Outlook} = \text{sunny} | h)P(\text{Temp} = \text{cool} | h)P(\text{Humidity} = \text{high} | h)P(\text{Wind} = \text{strong} | h)$$

- Working:

$$P(\text{PlayTennis} = \text{yes}) = 9/14 = 0.64$$

$$P(\text{PlayTennis} = \text{no}) = 5/14 = 0.36$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{yes}) = 3/9 = 0.33$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{no}) = 3/5 = 0.60$$

*etc.*

$$P(\text{yes})P(\text{sunny} | \text{yes})P(\text{cool} | \text{yes})P(\text{high} | \text{yes})P(\text{strong} | \text{yes}) = 0.0053$$

$$P(\text{no})P(\text{sunny} | \text{no})P(\text{cool} | \text{no})P(\text{high} | \text{no})P(\text{strong} | \text{no}) = \mathbf{0.0206}$$

$\Rightarrow \text{answer} : \text{PlayTennis}(x) = \text{no}$

# Underflow Prevention

- Multiplying lots of probabilities, which are between 0 and 1 by definition, can result in floating-point underflow.
- Since  $\log(xy) = \log(x) + \log(y)$ , it is better to perform all computations by summing logs of probabilities rather than multiplying probabilities.
- Class with highest final un-normalized log probability score is still the most probable.

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j)$$