

## Lecture 9: 5 February, 2026

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning  
January–April 2025

# Bayesian classifiers

- As before
  - Attributes  $\{A_1, A_2, \dots, A_k\}$  and
  - Classes  $C = \{c_1, c_2, \dots, c_\ell\}$
- Each class  $c_i$  defines a probabilistic model for attributes
  - $Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i)$
- Given a data item  $d = (a_1, a_2, \dots, a_k)$ , identify the best class  $c$  for  $d$
- Maximize  $Pr(C = c_i \mid A_1 = a_1, \dots, A_k = a_k)$

# Generative models

- To use probabilities, need to describe how data is randomly generated
  - Generative model
- Typically, assume a random instance is created as follows
  - Choose a class  $c_j$  with probability  $Pr(c_j)$
  - Choose attributes  $a_1, \dots, a_k$  with probability  $Pr(a_1, \dots, a_k | c_j)$
- Generative model has associated parameters  $\theta = (\theta_1, \dots, \theta_m)$ 
  - Each class probability  $Pr(c_j)$  is a parameter
  - Each conditional probability  $Pr(a_1, \dots, a_k | c_j)$  is a parameter
- We need to estimate these parameters

# Maximum Likelihood Estimators

- Our goal is to estimate parameters (probabilities)  $\theta = (\theta_1, \dots, \theta_m)$
- Law of large numbers allows us to estimate probabilities by counting frequencies
- Example: Tossing a biased coin, single parameter  $\theta = Pr(\text{heads})$ 
  - $N$  coin tosses,  $H$  heads and  $T$  tails
  - Why is  $\hat{\theta} = H/N$  the best estimate?
- Likelihood
  - Actual coin toss sequence is  $\tau = t_1 t_2 \dots t_N$
  - Given an estimate of  $\theta$ , compute  $Pr(\tau | \theta)$  — likelihood  $L(\theta)$
- $\hat{\theta} = H/N$  maximizes this likelihood —  $\arg \max_{\theta} L(\theta) = \hat{\theta} = H/N$ 
  - Maximum Likelihood Estimator (MLE)

# Bayesian classification

- Maximize  $Pr(C = c_i | A_1 = a_1, \dots, A_k = a_k)$
- By Bayes' rule,

$$\begin{aligned} & Pr(C = c_i | A_1 = a_1, \dots, A_k = a_k) \\ &= \frac{Pr(A_1 = a_1, \dots, A_k = a_k | C = c_i) \cdot Pr(C = c_i)}{Pr(A_1 = a_1, \dots, A_k = a_k)} \\ &= \frac{Pr(A_1 = a_1, \dots, A_k = a_k | C = c_i) \cdot Pr(C = c_i)}{\sum_{j=1}^{\ell} Pr(A_1 = a_1, \dots, A_k = a_k | C = c_j) \cdot Pr(C = c_j)} \end{aligned}$$

- Denominator is the same for all  $c_i$ , so sufficient to maximize

$$Pr(A_1 = a_1, \dots, A_k = a_k | C = c_i) \cdot Pr(C = c_i)$$

# Example

- To classify  $A = g, B = q$
- $Pr(C = t) = 5/10 = 1/2$
- $Pr(A = g, B = q | C = t) = 2/5$
- $Pr(A = g, B = q | C = t) \cdot Pr(C = t) = 1/5$
- $Pr(C = f) = 5/10 = 1/2$
- $Pr(A = g, B = q | C = f) = 1/5$
- $Pr(A = g, B = q | C = f) \cdot Pr(C = f) = 1/10$
- Hence, predict  $C = t$

$A$	$B$	$C$
$m$	$b$	$t$
$m$	$s$	$t$
$g$	$q$	$t$
$h$	$s$	$t$
$g$	$q$	$t$
$g$	$q$	$f$
$g$	$s$	$f$
$h$	$b$	$f$
$h$	$q$	$f$
$m$	$b$	$f$

## Example . . .

- What if we want to classify  $A = m, B = q$ ?
- $Pr(A = m, B = q | C = t) = 0$
- Also  $Pr(A = m, B = q | C = f) = 0!$
- To estimate joint probabilities across all combinations of attributes, we need a much larger set of training data

<i>A</i>	<i>B</i>	<i>C</i>
<i>m</i>	<i>b</i>	<i>t</i>
<i>m</i>	<i>s</i>	<i>t</i>
<i>g</i>	<i>q</i>	<i>t</i>
<i>h</i>	<i>s</i>	<i>t</i>
<i>g</i>	<i>q</i>	<i>t</i>
<i>g</i>	<i>q</i>	<i>f</i>
<i>g</i>	<i>s</i>	<i>f</i>
<i>h</i>	<i>b</i>	<i>f</i>
<i>h</i>	<i>q</i>	<i>f</i>
<i>m</i>	<i>b</i>	<i>f</i>

# Naïve Bayes classifier

- Strong simplifying assumption: attributes are pairwise independent

$$Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i) = \prod_{j=1}^k Pr(A_j = a_j \mid C = c_i)$$

- $Pr(C = c_i)$  is fraction of training data with class  $c_i$
  - $Pr(A_j = a_j \mid C = c_i)$  is fraction of training data labelled  $c_i$  for which  $A_j = a_j$
- Final classification is

$$\arg \max_{c_i} Pr(C = c_i) \prod_{j=1}^k Pr(A_j = a_j \mid C = c_i)$$

- Conditional independence is not theoretically justified
- For instance, text classification
  - Items are documents, attributes are words (absent or present)
  - Classes are topics
  - Conditional independence says that a document is a set of words: ignores sequence of words
  - Meaning of words is clearly affected by relative position, ordering
- However, naive Bayes classifiers work well in practice, even for text classification!
  - Many spam filters are built using this model

## Example revisited

- Want to classify  $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$
- $Pr(A = m \mid C = t) = 2/5$
- $Pr(B = q \mid C = t) = 2/5$
- $Pr(A = m \mid C = f) = 1/5$
- $Pr(B = q \mid C = f) = 2/5$
- $Pr(A = m \mid C = t) \cdot Pr(B = q \mid C = t) \cdot Pr(C = t) = 2/25$
- $Pr(A = m \mid C = f) \cdot Pr(B = q \mid C = f) \cdot Pr(C = f) = 1/25$
- Hence predict  $C = t$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f