**Naïve Bayes**

- **Maximum Likelihood Estimation (MLE)**
  - Consider a sequence of coin flips, for example ← *Coin*
    *prior* → HHTTTTTTHTTTTTHTTHTT  (5 times H, 15 times T)
  - Which Pr(H) and Pr(T) are the most likely?
  - Looks like Pr(H) = ¼ and Pr(T) = ¾ ...

$Pr(H) = Pr(T) = \frac{1}{2}$

$|flips| = 20$

$|H| = 5, |T| = 15$

$\theta = $

$P(T) = \frac{|T|}{|F|} = \frac{15}{20} = 3/4$

$P(H) = \frac{|H|}{|F|} = \frac{8}{20} = \frac{1}{4}$

- **Conditional probabilities**
  - Let A and B be events in a probability space $\Omega$
  - Denote by Pr(A | B) the probability of A ∩ B in the space B
  - **(1)** Pr(A | B) := Pr(A ∩ B) / Pr(B) ←
  - **(2)** Pr(A | B) · Pr(B) = Pr(B | A) · Pr(A)

$Pr(A|B) = Pr(B|A) \cdot P(A) / P(B)$

$\Omega = \{1, 2, 3, 4, 5, 6\}$

$A = \{2, 4, 6\} \Rightarrow P(A) = \frac{|A|}{|\Omega|} = \frac{3}{6} = \frac{1}{2}$

$B = \{1, 2, 3\} \Rightarrow P(B) = \frac{|B|}{|\Omega|} = \frac{3}{6} = \frac{1}{2}$

$P(A|B) = \frac{|A \cap B|}{|B|} = \frac{1}{3}$ ✓

- **Probabilistic assumptions**
  - Underlying probability distributions:
    A distribution $p_c$ over the classes ... where $\Sigma_c\, p_c = 1$
    For each c, a distr. $p_{wc}$ over the words ... where $\Sigma_w\, p_{wc} = 1$
  - Naïve Bayes assumes the following process for generating a document D with m words $W_1...W_m$ and class label C:

    Pick C=c with prob. $p_c$ , then pick each word $W_i$=w with probability $p_{wc}$ , independent of the other words

- **Learning phase**
  - For a **training set** T of objects, let:
    $T_c$ = the set of documents from class c
    $n_{wc}$ = #occurrences of word w in documents from $T_c$
    $n_c$ = #occurrences of all words in documents from $T_c$
  - We compute the following probabilities or likelihoods:
    → $p_c := |T_c| / |T|$      global likelihood of a class
    $p_{wc} := n_{wc} / n_c$      likelihood of a word for a class

$T = 90$

→ $T_{Horror} = 10$
→ $T_{Doc} = 25$
→ $T_{comedy} = 15$

$P_{Jupyter, Doc} = \frac{20}{80} = \frac{1}{4} = 0.25$

## Learning phase, example

- Consider Example 2 (artificial documents)

| | |
|---|---|
| aba | A |
| baabaaa | A |
| bbaabbab | B |
| abbaa | A |
| abbb | B |
| bbbaab | B |

*Handwritten annotations (right side):*

$$T_A = 3, \quad T_B = 3, \quad T = 6$$

$$\rightarrow P_A = P(T_A) = \frac{3}{6} = \frac{1}{2}, \quad \boxed{P_B = \frac{1}{2}}$$

$$n_{aA} = 10, \quad n_{bA} = 5, \quad n_A = 15$$

$$P_{aA} = \frac{10}{15} = \boxed{\frac{2}{3}}, \quad P_{bA} = \frac{5}{15} = \boxed{\frac{1}{3}}$$

$$n_{aB} = 6, \quad n_{bB} = 12, \quad n_B = 18$$

$$P_{aB} = \frac{6}{18} = \frac{1}{3}, \quad P_{bB} = \frac{12}{18} = \frac{2}{3}$$

## Prediction

- For a given document **d** we want to compute the probability of each class **c**, given document d:

  **Pr(C=c | D=d)**

- Using Bayes Theorem, we have:

  Pr(C=c | D=d) = Pr(D=d | C=c) · Pr(C=c) / Pr(D=d)

- Using our (naive) probabilistic assumptions, we have:

  $Pr(D=d \mid C=c) = Pr(W_1=w_1 \cap \dots \cap W_m=w_m \mid C=c)$
  $= \Pi_{i=1,\dots,m} Pr(W_i=w_i \mid C=c)$

*Handwritten: $\frac{16}{8} < \frac{24}{8}$, ②, ③, 2, 3, 0.5, A, a, b*

## Prediction ... continued

- We thus obtain that

  Pr(C=c | D=d)

  $= \Pi_{i=1,\dots,m} Pr(W_i=w_i \mid C=c) \cdot Pr(C=c) / Pr(D=d)$

  $= \Pi_{i=1,\dots,m} \mathbf{P_{w_i c}} \cdot \mathbf{P_c} / Pr(D=d)$

- Note 1: for $\Pi_{i=1,\dots,m} P_{w_i c}$ just take the $p_{wc}$ for all words w in the document and multiply them (if a word w occurs multiple times, also take the factor $p_{wc}$ multiple times)

- Note 2: the Pr(D=d) is the same for all classes c

## Prediction, example

- Consider Example 2 (artificial documents)

| | |
|---|---|
| aba | A |
| baabaaa | A |
| bbaabbab | B |
| abbaa | A |
| abbb | B |
| bbbaab | B |

- Let us predict the class for aab ... A or B ?

*Handwritten annotations (bottom right):*

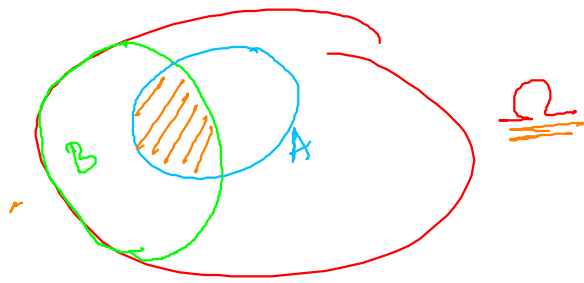$$P(A \mid aab) = P_{aA} \cdot P_{aA} \cdot P_{bA} \cdot P_A$$
$$= \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} = \frac{2}{27}$$

$$P(B \mid aab) = P_{aB} \cdot P_{aB} \cdot P_{bB} \cdot P_B$$
$$= \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{1}{2} = \frac{1}{27}$$

$$P(A \mid aab) > P(B \mid aab) \rightarrow \text{Predict } A$$

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

$$= \frac{|A \cap B| \;/\; |\Omega|}{|B| \;/\; |\Omega|} = \frac{Pr(A \cap B)}{Pr(B)}$$

$$P_{\text{Zombie, Horror}} = \frac{|\text{Zombie}|}{|\text{class, word}|} = \frac{5}{20} = \frac{1}{4} = 0.25$$