

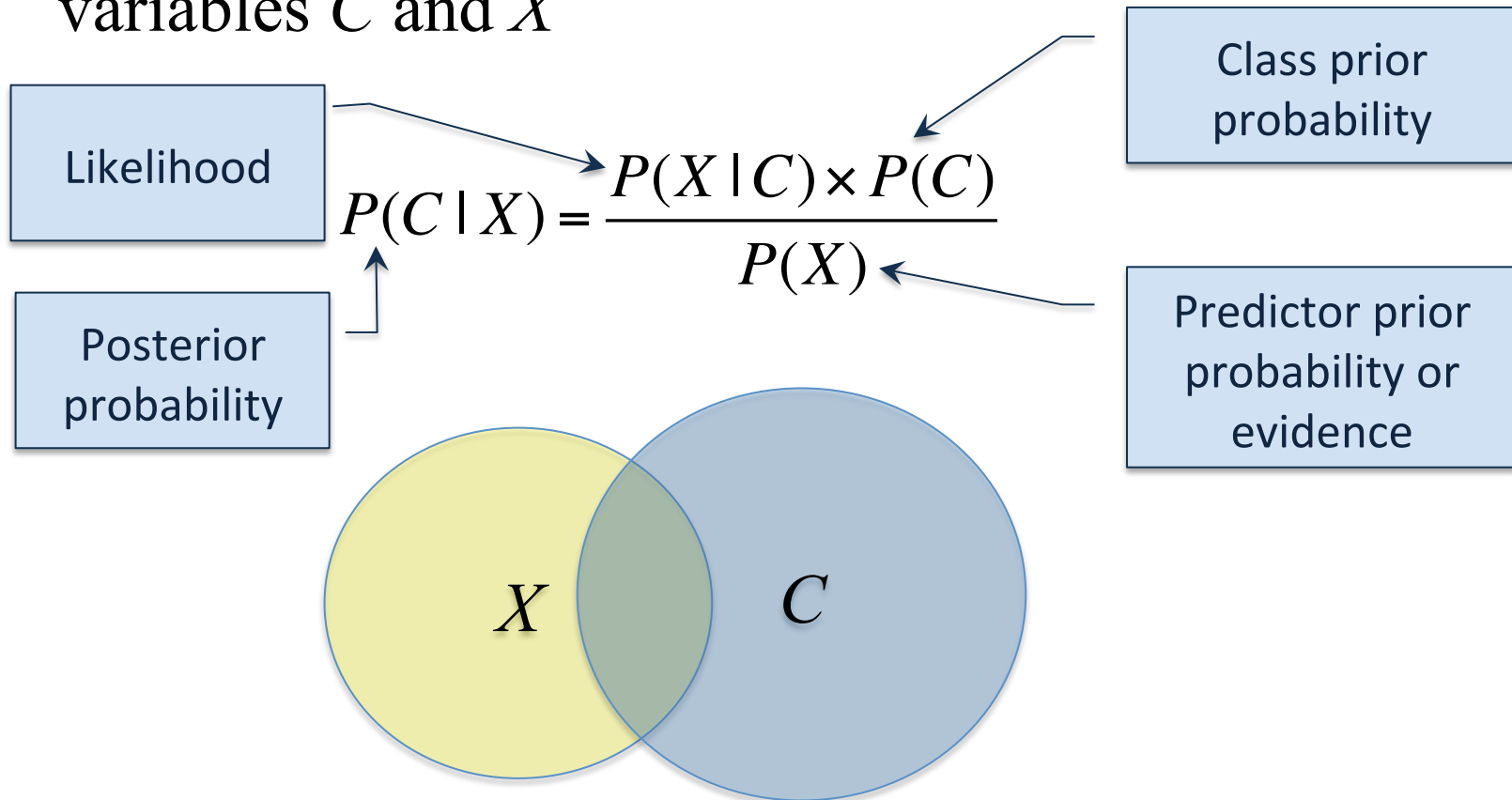
Naïve Bayes Classification

Debapriyo Majumdar
Data Mining – Fall 2014
Indian Statistical Institute Kolkata

August 14, 2014

Bayes' Theorem

- Thomas Bayes (1701-1761)
- Simple form of Bayes' Theorem, for two random variables C and X

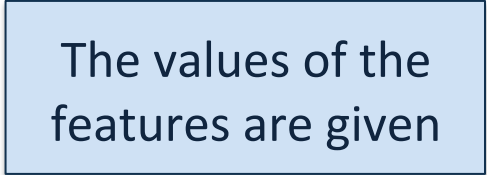


Probability Model

- Probability model: for a *target class variable* C which is dependent over features X_1, \dots, X_n

$$P(C | X_1, \dots, X_n) = \frac{P(C) \times P(X_1, \dots, X_n | C)}{P(X_1, \dots, X_n)}$$

The values of the features are given



- So the denominator is effectively constant
- Goal: calculating probabilities for the possible values of C
- We are interested in the numerator:

$$P(C) \times P(X_1, \dots, X_n | C)$$

Probability Model

- The conditional probability is equivalent to the joint probability

$$P(C) \times P(X_1, \dots, X_n | C) = P(C, X_1, \dots, X_n)$$

- Applying the chain rule for joint probability

$$P(A, B) = P(A) \times P(B | A)$$

$$\begin{aligned} & P(C) \times P(X_1, \dots, X_n | C) \\ &= P(C) \times P(X_1 | C) \times P(X_2, \dots, X_n | C, X_1) \\ &= P(C) \times P(X_1 | C) \times P(X_2 | C, X_1) P(X_3, \dots, X_n | C, X_1, X_2) \\ & \quad \dots \\ & \quad \dots \\ &= P(C) \times P(X_1 | C) \times P(X_2 | C, X_1) \dots P(X_n | C, X_1, \dots, X_{n-1}) \end{aligned}$$

Strong Independence Assumption (Naïve)

- Assume the features X_1, \dots, X_n are conditionally independent given C
 - Given C , occurrence of X_i does not influence the occurrence of X_j , for $i \neq j$.

$$P(X_i | C, X_j) = P(X_i | C)$$

- Similarly,

$$P(X_i | C, X_j, \dots, X_k) = P(X_i | C)$$


- Hence:

$$\begin{aligned} & P(C) \times P(X_1 | C) \times P(X_2 | C, X_1) \dots P(X_n | C, X_1, \dots, X_{n-1}) \\ &= P(C) \times P(X_1 | C) \times P(X_2 | C) \dots P(X_n | C) \end{aligned}$$

Naïve Bayes Probability Model

$$P(C | X_1, \dots, X_n) = \frac{P(C) \times P(X_1 | C) \times P(X_2 | C) \times \dots \times P(X_n | C)}{P(X_1, \dots, X_n)}$$

Class posterior
probability



Known values:
constant

Classifier based on Naïve Bayes

- Decision rule: pick the hypothesis (value c of C) that has highest probability
 - Maximum A-Posteriori (MAP) decision rule

$$\operatorname{argmax}_c \left\{ P(C = c) \prod_{i=1}^n P(X_i = x_i | C = c) \right\}$$

Approximated
from frequency in
the training set

Approximated
from relative
frequencies in
the training set

The values of
features are
known for the
new observation

Example of Naïve Bayes

Reference: The IR Book by Raghavan et al, Chapter 6

Text Classification with Naïve Bayes

The Text Classification Problem

- Set of class labels / tags / categories: C
- Training set: set D of documents with labels
 $\langle d, c \rangle \in D \times C$
- *Example*: a document, and a class label
 $\langle \text{Beijing joins the World Trade Organization, } \textit{China} \rangle$
- Set of all terms: V
- Given $d \in D'$, a set of new documents, the goal is to find a class of the document $c(d)$

Multinomial Naïve Bayes Model

- Probability of a document d being in class c

$$P(c | d) \propto P(c) \times \prod_{k=1}^{n_d} P(t_k | c)$$

where $P(t_k | c)$ = probability of a term t_k occurring in a document of class c

Intuitively:

- $P(t_k | c)$ ~ how much evidence t_k contributes to the class c
- $P(c)$ ~ prior probability of a document being labeled by class c

Multinomial Naïve Bayes Model

- The expression

$$P(c) \times \prod_{k=1}^{n_d} P(t_k | c)$$

has many probabilities.

- May result a floating point *underflow*.
 - Add logarithms of the probabilities instead of multiplying the original probability values

$$c_{map} = \operatorname{argmax}_{c \in C} \left[\log P(c) + \sum_{k=1}^{n_d} \log P(t_k | c) \right]$$

Todo: estimating these probabilities

Maximum Likelihood Estimate

- Based on relative frequencies
- Class prior probability

$$P(c) = \frac{N_c}{N}$$

#of documents labeled as class c in training data

total #of documents in training data

- Estimate $P(t|c)$ as the relative frequency of t occurring in documents labeled as class c

$$P(t | c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

total #of occurrences of t in documents $d \in c$

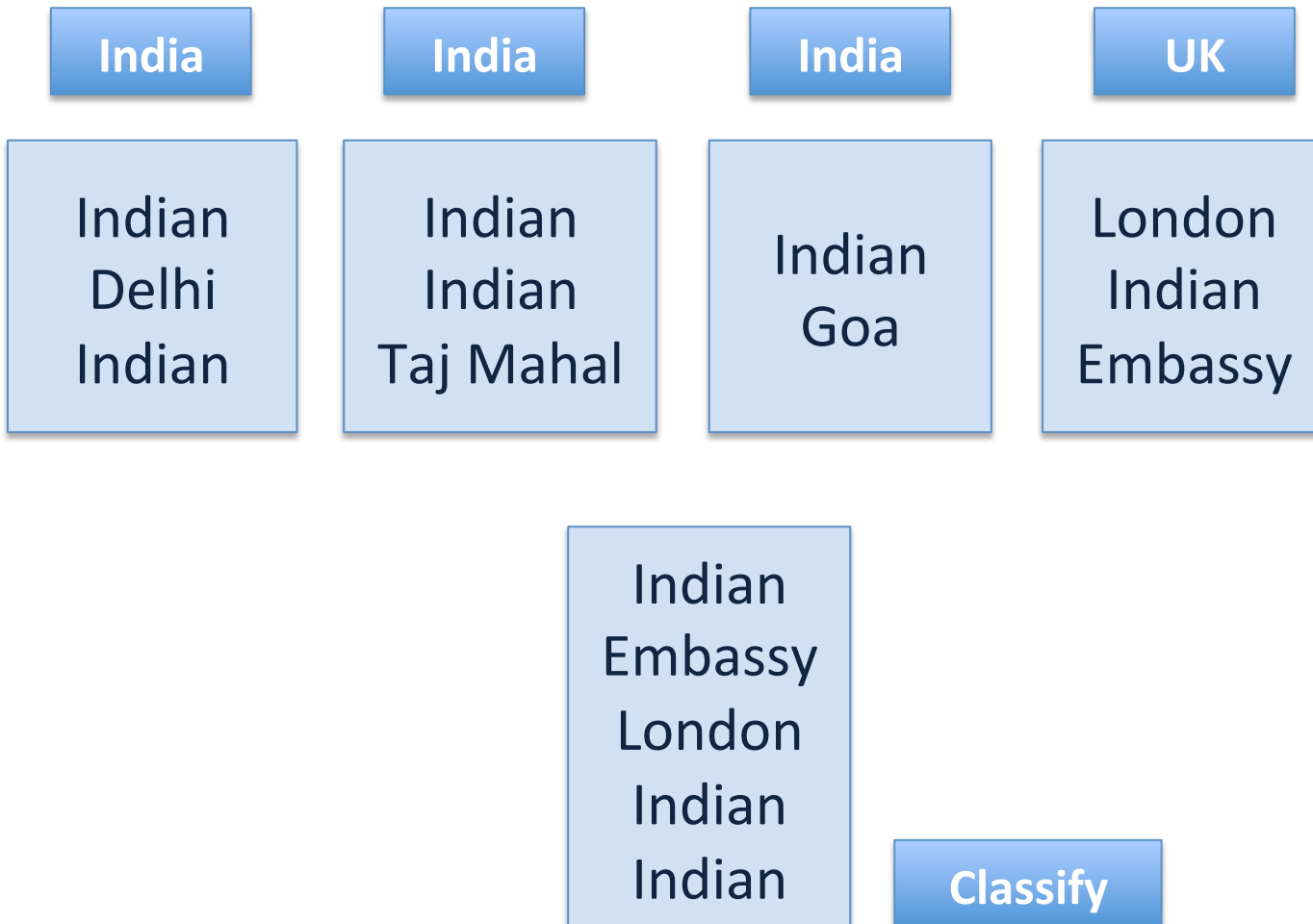
total #of occurrences of all terms in documents $d \in c$

Handling Rare Events

- What if: a term t did not occur in documents belonging to class c in the training data?
 - Quite common. Terms are sparse in documents.
- Problem: $P(t|c)$ becomes zero, the whole expression becomes zero
- Use *add-one* or *Laplace smoothing*

$$P(t | c) = \frac{T_{ct} + 1}{\sum_{t' \in V} (T_{ct'} + 1)} = \frac{T_{ct} + 1}{\left(\sum_{t' \in V} T_{ct'} \right) + |V|}$$

Example



Bernoulli Naïve Bayes Model

- Binary indicator of occurrence instead of frequency
- Estimate $P(t|c)$ as the *fraction* of documents in c containing the term t
- Models absence of terms explicitly:

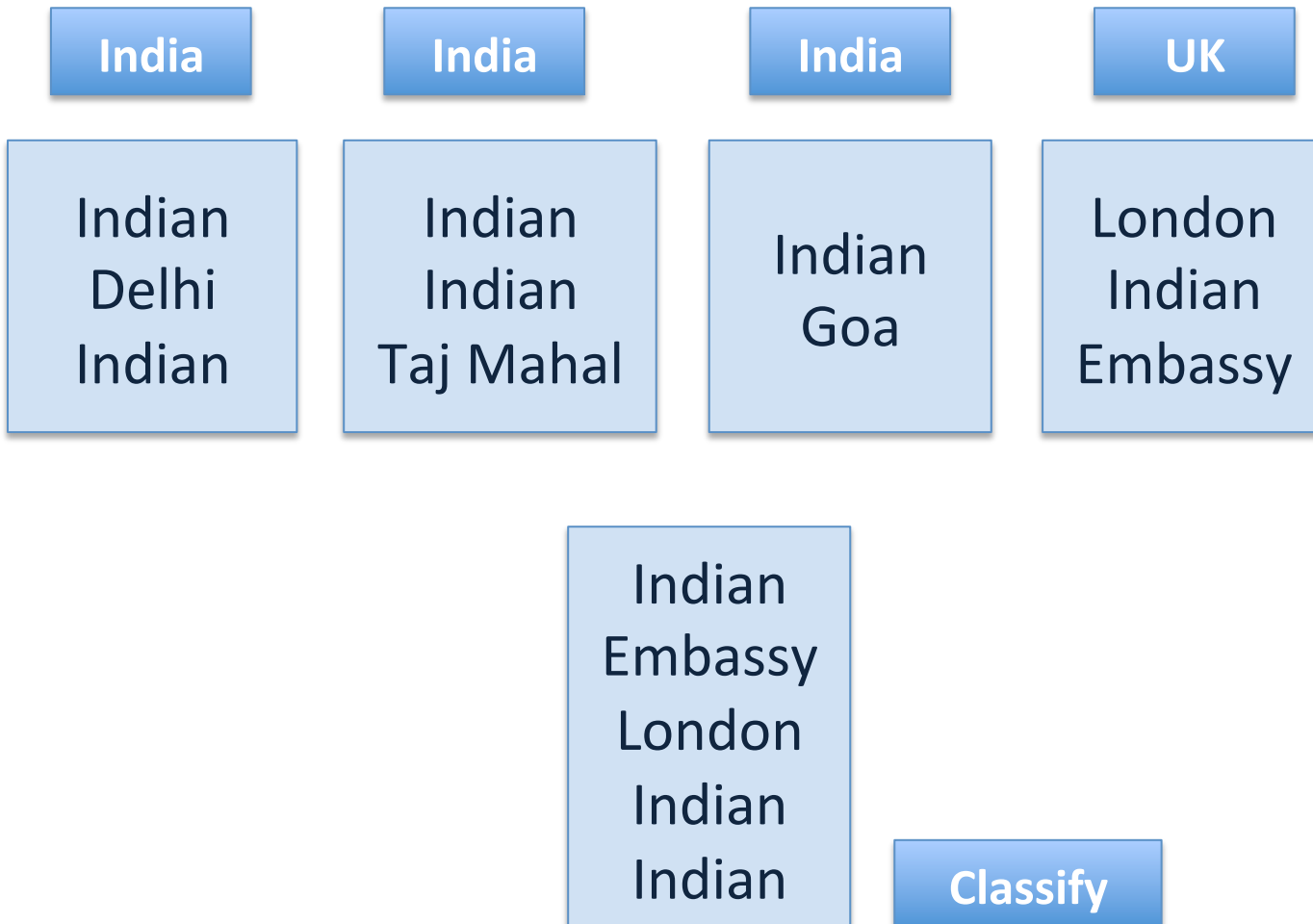
$$P(C | X_1, \dots, X_n) = \prod_{i=1}^n [X_i P(t_i | C) + (1 - X_i)(1 - P(t_i | C))]$$

$X_i = 1$ if t_i is present
0 otherwise

Absence of terms

Difference between Multinomial with frequencies truncated to 1, and Bernoulli Naïve Bayes?

Example

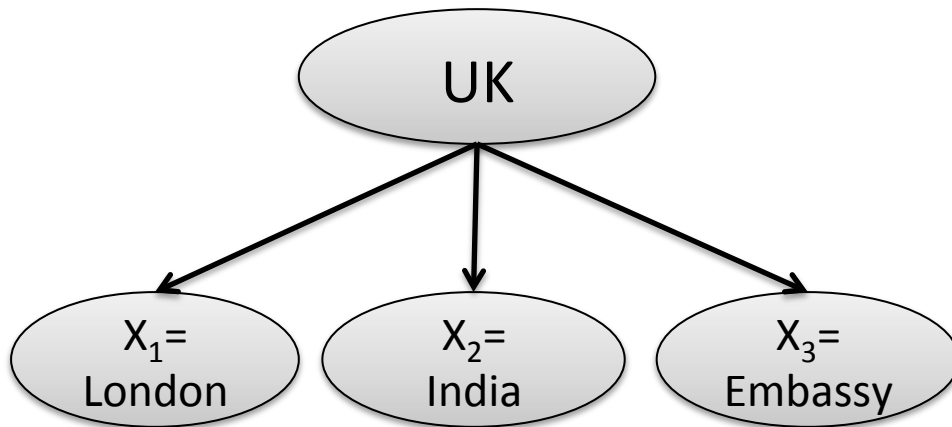


Naïve Bayes as a Generative Model

- The probability model:

$$P(c | d) = \frac{P(c) \times P(d | c)}{P(d)}$$

Multinomial model



$$P(d | c) = P(\langle t_1, \dots, t_{n_d} \rangle | c)$$

Terms as they occur in d , exclude other terms

where X_i is the random variable for position i in the document

- Takes values as terms of the vocabulary

- Positional independence assumption \rightarrow Bag of words model:

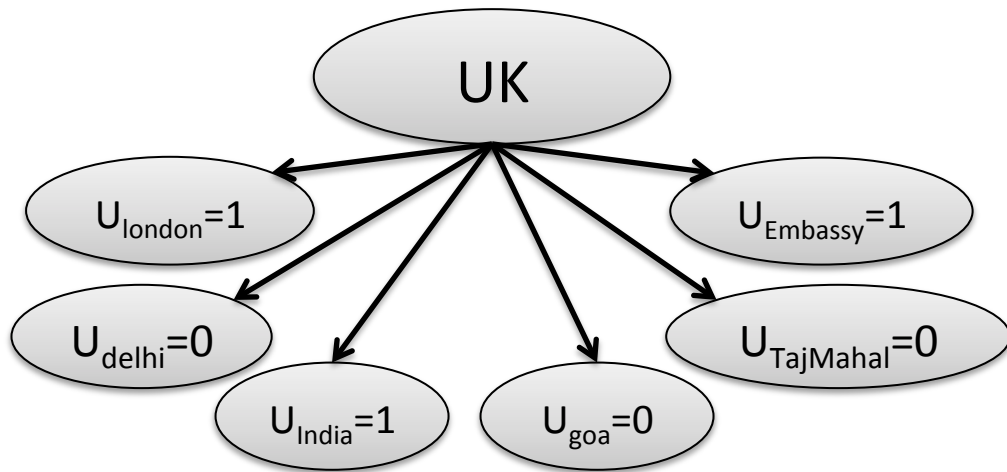
$$P(X_{k_1} = t | c) = P(X_{k_2} = t | c)$$

Naïve Bayes as a Generative Model

- The probability model:

$$P(c | d) = \frac{P(c) \times P(d | c)}{P(d)}$$

Bernoulli model



$$P(d | c) = P(\langle e_1, \dots, e_{|V|} \rangle | c)$$

All terms in the vocabulary

$P(U_i=1|c)$ is the probability that term t_i will occur in any position in a document of class c

Multinomial vs Bernoulli

	Multinomial	Bernoulli
Event Model	Generation of token	Generation of document
Multiple occurrences	Matters	Does not matter
Length of documents	Better for larger documents	Better for shorter documents
#Features	Can handle more	Works best with fewer

On Naïve Bayes

- Text classification
 - Spam filtering (email classification) [Sahami et al. 1998]
 - Adapted by many commercial spam filters
 - SpamAssassin, SpamBayes, CRM114, ...
- Simple: the conditional independence assumption is very strong (naïve)
 - Naïve Bayes may not estimate right in many cases, but ends up classifying correctly quite often
 - It is difficult to understand the dependencies between features in real problems