

Faculty of Computer Science, Dalhousie University  
**CSCI 4152/6509 — Natural Language Processing**

16-Oct-2024

**Lecture 11: Naïve Bayes Model; P0 Topics Discussion (3)**

Location: Carleton Tupper Building Theatre C      Instructor: Vlado Keselj  
 Time: 16:05 – 17:25

**Previous Lectures**

- Joint distribution model
  - Spam example
- Fully independent model
- P0 Discussion (2): P-02

**3.b Conditioning in Fully Independent Model**

$$\begin{aligned}
 & P(V_{k+1} = y_{k+1}, \dots, V_n = y_n | V_1 = x_1, \dots, V_k = x_k) \\
 &= \frac{P(V_1 = x_1, \dots, V_k = x_k, V_{k+1} = y_{k+1}, \dots, V_n = y_n)}{P(V_1 = x_1, \dots, V_k = x_k)} \\
 &= \frac{P(V_1 = x_1) \cdots P(V_k = x_k) P(V_{k+1} = y_{k+1}) \cdots P(V_n = y_n)}{P(V_1 = x_1) \cdots P(V_k = x_k)} \\
 &= P(V_{k+1} = y_{k+1}) \cdots P(V_n = y_n)
 \end{aligned}$$

Only have to lookup and multiply  $n - k$  numbers.

**3.c Completion in Fully Independent Model**

$$\begin{aligned}
 y_{k+1}^*, \dots, y_n^* &= \arg \max_{y_{k+1}, \dots, y_n} P(V_{k+1} = y_{k+1}, \dots, V_n = y_n | V_1 = x_1, \dots, V_k = x_k) \\
 &= \arg \max_{y_{k+1}, \dots, y_n} P(V_{k+1} = y_{k+1}) \cdots P(V_n = y_n) \\
 &= \arg \max_{y_{k+1}} P(V_{k+1} = y_{k+1}) \left[ \arg \max_{y_{k+2}} \cdots \left[ \arg \max_{y_n} P(V_n = y_n) \right] \right]
 \end{aligned}$$

(Since max and arg max distributes over product just like sum.  
 That is,  $\max_i ax_i = a \max_i x_i$  (for  $a, x_i \geq 0$ )  
 just like  $\sum_i ax_i = a \sum_i x_i$ .)

$$\begin{aligned}
 &= \left[ \arg \max_{y_{k+1}} P(V_{k+1} = y_{k+1}) \right] \cdots \left[ \arg \max_{y_n} P(V_n = y_n) \right] \\
 &= \left[ \arg \max_{y_{k+1}} p_{k+1, y_{k+1}} \right] \cdots \left[ \arg \max_{y_n} p_{n, y_n} \right]
 \end{aligned}$$

Only have to search through  $m$  possible completions for each of the  $n - k$  variables separately.

### Joint Distribution Model vs. Fully Independent Model

The Fully Independent Model addresses the previous issues with the joint distribution model, but it suffers from a too strong assumption and too little structure, so it usually does not model accurately the real relationships among variables.

**Structured probability models** are a compromise solution between previous two models. Structured probability models are more efficient than the joint distribution model and they address the issue of the sparse training data, and in the same time they model important dependencies among random variables.

One of the simplest models of this kind is the Naïve Bayes Model.

## 12 Naïve Bayes Classification Model

Slide notes:

### Naïve Bayes Classification Model

- Fully independent model is not useful in classification: class variable should be dependent on other variables
- A solution: make class variable dependent, but everything else independent
- Let  $V_1$  be the class variable
- $V_2, V_3, \dots, V_n$  are input variables (features)
- Classification can be expressed as

$$\arg \max_{x_1} P(V_1 = x_1 | V_2 = x_2, V_3 = x_3, \dots, V_n = x_n)$$

In the Naïve Bayes model we assume that all variables are independent except one distinguished variable, which is usually called the *class variable* since the model is used for classification. The other variables are called *features* or *attributes*. Since in the classification task the features are used as input and the class variable produces the classification result or output, we also call the feature variables the *input variables* and the class variable the *output variable*.

If we assume that the variable  $V_1$  is the output variable, and the variables  $V_2, V_3, \dots, V_n$  are the input variables, then in the classification problem can be expressed as a conditional probability computation problem, or completion problem of the probability:

$$\arg \max_{x_1} P(V_1 = x_1 | V_2 = x_2, V_3 = x_3, \dots, V_n = x_n)$$

or

$$\arg \max_{V_1} P(V_1 | V_2, V_3, \dots, V_n)$$

for short. After applying Bayes theorem we obtain:

$$P(V_1 | V_2, V_3, \dots, V_n) = \frac{P(V_2, V_3, \dots, V_n | V_1) \cdot P(V_1)}{P(V_2, V_3, \dots, V_n)}$$

Now, we use the Naïve Bayes independence assumption, which is that the variables  $V_2, V_3, \dots, V_n$  are conditionally independent given  $V_1$ . Then, the above equation becomes:

$$\begin{aligned} P(V_1 | V_2, V_3, \dots, V_n) &= \frac{P(V_2, V_3, \dots, V_n | V_1) \cdot P(V_1)}{P(V_2, V_3, \dots, V_n)} \\ &= \frac{P(V_2 | V_1) \cdot P(V_3 | V_1) \cdot \dots \cdot P(V_n | V_1) \cdot P(V_1)}{P(V_2, V_3, \dots, V_n)} \end{aligned}$$

The conditional probabilities  $P(V_i|V_1)$  for  $i \in \{2 \dots n\}$  can be efficiently computed and stored, and they eliminate the sparse data problem. To be clear about the independence assumption that we made, let us repeat it here, the **Naïve Bayes Independence Assumption (1)** can be stated as follows:

$$P(V_2, V_3, \dots, V_n|V_1) = P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1)$$

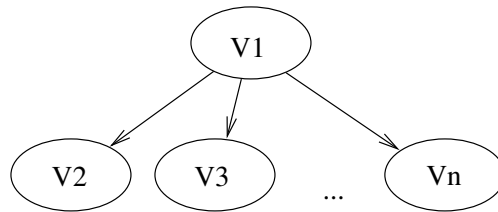
If we multiply both sides with  $P(V_1)$  and use definition of conditional probability, we the second way of expressing the **Naïve Bayes Independence Assumption (2)** is:

$$P(V_1, V_2, V_3, \dots, V_n) = P(V_1) \cdot P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1)$$

**Graphical Representation: Naïve Bayes Model**

Assumption:

$$P(V_1, V_2, V_3, \dots, V_n) = P(V_1) \cdot P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1)$$



**Naïve Bayes Classification**

– The classification formula becomes

$$\begin{aligned} \arg \max_{x_1} \frac{P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1) \cdot P(V_1)}{P(V_2, V_3, \dots, V_n)} = \\ \arg \max_{x_1} P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1) \cdot P(V_1) \end{aligned}$$

– To calculate marginal probability in the denominator we use

$$\begin{aligned} P(V_2, V_3, \dots, V_n) &= \sum_{V_1} P(V_1, V_2, V_3, \dots, V_n) = \\ &= \sum_{V_1} P(V_2|V_1) \cdot P(V_3|V_1) \cdot \dots \cdot P(V_n|V_1) \cdot P(V_1) \end{aligned}$$

Another way of deriving the Naïve Bayes assumption is the following:

$$\begin{aligned} P(V_1 = x_1, \dots, V_n = x_n) &= \tag{3} \\ &= P(V_1 = x_1)P(V_2 = x_2|V_1 = x_1)P(V_3 = x_3|V_1 = x_1, V_2 = x_2) \dots \tag{4} \\ &= P(V_n = x_n|V_1 = x_1, V_2 = x_2, \dots, V_{n-1} = x_{n-1}) \tag{5} \\ &\stackrel{NB}{\approx} P(V_1 = x_1)P(V_2 = x_2|V_1 = x_1)P(V_3 = x_3|V_1 = x_1) \dots \tag{6} \\ &= P(V_n = x_n|V_1 = x_1) \tag{7} \end{aligned}$$

Equality (3,4) holds always, and equality (5,6) is the Naïve Bayes assumption.

**Summary of the Naïve Bayes Model**

Naïve Bayes assumption

$$P(V_2, V_3, \dots, V_n | V_1) = P(V_2 | V_1) P(V_3 | V_1) \dots P(V_n | V_1)$$

↓
class variable

↖
↙

text features
class variable

Second way of expression Naïve Bayes Assumption:

$$P(V_1, V_2, V_3, \dots, V_n) = P(V_1) P(V_2, V_3, \dots, V_n | V_1) = P(V_1) P(V_2 | V_1) P(V_3 | V_1) \dots P(V_n | V_1)$$

Naïve Bayes Model is a set of tables

V1	P(V1)

V1	V2	P(V2 V1)

V1	Vn	P(Vn V1)

(CPT -- Conditional Probability Tables)

**Example: A Naïve Bayes Model for Spam Detection**

In our spam detection example, the Naïve Bayes assumption is:

$$P(Free, Caps, Spam) = P(Spam) \cdot P(Free|Spam) \cdot P(Caps|Spam)$$

Hence, in order to create a Naïve Bayes model from our training data:

<i>Free</i>	<i>Caps</i>	<i>Spam</i>	Number of messages
Y	Y	Y	20
Y	Y	N	1
Y	N	Y	5
Y	N	N	0
N	Y	Y	20
N	Y	N	3
N	N	Y	2
N	N	N	49
Total:			100

we calculate the following tables:

<i>Spam</i>	$P(Spam)$
Y	$\frac{20+5+20+2}{100} = 0.47$
N	$\frac{1+0+3+49}{100} = 0.53$

<i>Caps</i>	<i>Spam</i>	$P(Caps Spam)$
Y	Y	$\frac{20+20}{20+5+20+2} \approx 0.8511$
Y	N	$\frac{1+3}{1+0+3+49} \approx 0.0755$
N	Y	$\frac{5+2}{20+5+20+2} \approx 0.1489$
N	N	$\frac{0+49}{1+0+3+49} \approx 0.9245$

<i>Free</i>	<i>Spam</i>	$P(Free Spam)$
Y	Y	$\frac{20+5}{20+5+20+2} \approx 0.5319$
Y	N	$\frac{1+0}{1+0+3+49} \approx 0.0189$
N	Y	$\frac{20+2}{20+5+20+2} \approx 0.4681$
N	N	$\frac{3+49}{1+0+3+49} \approx 0.9811$

The probability of a configuration in this model is calculated in the following way:

$$\begin{aligned}
 P(\text{Free} = Y, \text{Caps} = N, \text{Spam} = N) &= \\
 &= P(\text{Spam} = N) \cdot P(\text{Caps} = N | \text{Spam} = N) \cdot P(\text{Free} = Y | \text{Spam} = N) \\
 &\approx 0.53 \cdot 0.9245 \cdot 0.0189 \approx 0.0093
 \end{aligned} \tag{8}$$

## 12.1 Computational Tasks in the Naïve Bayes Model

We will cover the computational tasks in more details within the Bayesian Network in general.

### 1. Evaluation

The probability of a complete configuration is calculated using the Naïve Bayes assumption and table lookups. The formula (8) illustrates probability evaluation of a complete configuration:  $P(\text{Free} = Y, \text{Caps} = N, \text{Spam} = N)$

This example illustrates the fact that the Naïve Bayes model is less amenable to the sparse data problem than the joint distribution problem, in which the probability of this same configuration was estimated to be 0.

### 2. Simulation

Configurations are sampled by first sampling the output variable based on its table, and then the input variables using the corresponding conditional tables.

### 3. Inference

**3.a) Marginalization.** If the partial configuration includes the output variable, it can be shown that the marginal probability can be calculated using the following formula:

$$\begin{aligned}
 P(V_1 = x_1, \dots, V_k = x_k) &= \\
 &P(V_1 = x_1)P(V_2 = x_2 | V_1 = x_1)P(V_3 = x_3 | V_1 = x_1) \dots \\
 &P(V_k = x_k | V_1 = x_1)
 \end{aligned}$$

**3.b) Conditioning.** Example:

$$P(S = N | F = Y, C = N) = \frac{P(S = N, F = Y, C = N)}{P(F = Y, C = N)}$$

Using Naïve Bayes assumption:

$$\begin{aligned}
 P(S = N, F = Y, C = N) &= \\
 &= P(S = N)P(F = Y | S = N)P(C = N | S = N) \\
 &= 0.53 \cdot 0.9245 \cdot 0.0189 \approx 0.0093
 \end{aligned}$$

$$\begin{aligned}
P(F = Y, C = N) &= \text{(by definition)} \\
&= P(S = Y, F = Y, C = N) + P(S = N, F = Y, C = N) \\
&\approx P(S = Y)P(F = Y|S = Y)P(C = N|S = Y) + 0.0093 \\
&= 0.47 \cdot 0.5319 \cdot 0.1489 + 0.0093 \\
&\approx 0.0465
\end{aligned}$$

Finally,

$$P(S = N|F = Y, C = N) = \frac{0.0093}{0.0465} \approx 0.2$$

### 3.c) Completion in the Naïve Bayes Model

Slide notes:

#### 3.c) Completion in the NB Model

- Classification is the completion task:

$$\arg \max_{s \in \{Y, N\}} P(S = s|F = Y, C = N)$$

- It works out that we calculate:

$$P(S = Y, F = Y, C = N) = P(S) \cdot P(F|S) \cdot P(C|S)$$

and

$$P(S = N, F = Y, C = N) = P(S) \cdot P(F|S) \cdot P(C|S)$$

and choose the larger value.

Example:

$$\arg \max_{s \in \{Y, N\}} P(S = s|F = Y, C = N) \stackrel{\text{by definition}}{=} \arg \max_s \frac{P(S = s, F = Y, C = N)}{P(F = Y, C = N)}$$

$P(F = Y, C = N)$  does not depend on  $s$ , hence

$$= \arg \max_s P(S = s, F = Y, C = N)$$

and by using Naïve Bayes assumption)

$$= \arg \max_s \underbrace{P(S = s)P(F = Y|S = s)P(C = N|S = s)}_{A(s)}$$

For  $s = Y$   $A(s = Y) \approx 0.0465$ , and for  $s = N$   $A(s = N) \approx 0.0093$ ; hence

$$\arg \max_s A(s) = Y$$

**Learning**

Maximum Likelihood Estimation: The parameters are estimated using a corpus.

**12.2 Number of Parameters**

A Naïve Bayes model with  $n$  variables  $V_1, \dots, V_n$  is described with tables  $P(V_1), P(V_2|V_1), P(V_3|V_1), \dots, P(V_n|V_1)$ . These tables have constraints since each probability distribution must sum up to 1. If we assume that each variable can take one of  $m$  distinct values, then the number of parameters and constraints in required tables are:

	parameters	constraints
table $P(V_1)$	$m$	1
table $P(V_2 V_1)$	$m^2$	$m$
table $P(V_3 V_1)$	$m^2$	$m$
$\vdots$	$\vdots$	$\vdots$
table $P(V_n V_1)$	$m^2$	$m$
sum	$m + (n - 1)m^2$	$1 + (n - 1)m$

Hence, the number of free parameters is  $m + (n - 1)m^2 -$

$1 - (n - 1)m = O(m^2n)$ , which is not very large since the joint distribution model requires  $O(m^n)$  parameters.

**Pros and Cons of the Naïve Bayes Model**

Some advantages (pros) of the Naïve Bayes Model are:

**Efficiency:** It is a relatively efficient method, with good running-time complexity for inference and small memory size.

**No sparse data problem:** Since the number of parameters is relatively small, there is usually sufficient data to train all parameters, and smoothing is relatively easy.

**Performance:** Even though it has a very strong and unrealistic independence assumption, the model frequently shows surprisingly good classification performance.

Some disadvantages (cons) of the Naïve Bayes Model are:

**Too strong independence assumption:** The strong independence assumption often affects performance for many domains. In other words, the model is too simplistic.

**Only one “output” variable:** The model is designed as a classification problem; i.e., it contains only one hidden, or output, variable; which value can be inferred. Many problems require that we infer the value of multiple variables, and the only way to apply Naïve Bayes model to those problems is to build separate models for all hidden variables. In that case we would not capture any inter-dependencies among those variables.

**Additional Notes on Naïve Bayes Model**

- Text classification: how do we choose features?
- Two options:
  - Bernoulli Naïve Bayes — binary variables for each word
  - Multinomial Naïve Bayes — variable for each word position
- Zero-probability problem
  - Smoothing using +1 or similar addition (Laplace smoothing)

The *Bernoulli Naïve Bayes* model uses a variable for each distinct word in the vocabulary, with values 1 if the word is present, or 0 if not. Training is done on per-document basis. The name comes from the Bernoulli distribution as defined in the probability theory, which is distribution of a random variable having value 1 with a probability  $p$  and 0 with the probability  $q = 1 - p$ . This is the distribution we use to model probability that a word is in a document of a given class.

The *Multinomial Naïve Bayes* model uses a variable for each word position, and the value of the variable is the actual word. All conditional probabilities for these variables are the same, but they are collected in one large table. The model is trained on one ‘mega-document’; i.e., a document with concatenated all documents of a class. The model is named after the Multinomial distribution in the probability theory, which models the outcome of  $n$  repeated trials, where each trial can have one of  $k$  different results, with probabilities  $p_1, p_2, \dots, p_k$ . In the Multinomial Naïve Bayes model,  $n$  is the length of a text, and individual trials are word positions, where words are taken from a vocabulary of size  $k$ .

### 12.3 Spam Example Summary

Let us take a look at a summary of the Spam Example for the three discussed models: Joint Distribution, Fully Independent, and Naïve Bayes model. In all three models, the initial training data was the same, represented in the following table:

<i>Free</i>	<i>Caps</i>	<i>Spam</i>	Number of messages
Y	Y	Y	20
Y	Y	N	1
Y	N	Y	5
Y	N	N	0
N	Y	Y	20
N	Y	N	3
N	N	Y	2
N	N	N	49
Total:			100

The **Joint Distribution Model** is represented using a joint probability distribution table, learned from the training data as:

<i>Free</i>	<i>Caps</i>	<i>Spam</i>	Number of messages	$p$
Y	Y	Y	20	0.20
Y	Y	N	1	0.01
Y	N	Y	5	0.05
Y	N	N	0	0.00
N	Y	Y	20	0.20
N	Y	N	3	0.03
N	N	Y	2	0.02
N	N	N	49	0.49
Total:			100	1.00

As an example, the conditional probability  $P(\text{Spam} = Y | \text{Free} = Y, \text{Caps} = N)$  would be evaluated as:

$$\begin{aligned}
 P(\text{Spam} = Y | \text{Free} = Y, \text{Caps} = N) &= \\
 &= \frac{P(\text{Spam} = Y, \text{Free} = Y, \text{Caps} = N)}{P(\text{Free} = Y, \text{Caps} = N)} \\
 &= \frac{P(\text{Spam} = Y, \text{Free} = Y, \text{Caps} = N)}{P(\text{Spam} = Y, \text{Free} = Y, \text{Caps} = N) + P(\text{Spam} = N, \text{Free} = Y, \text{Caps} = N)} \\
 &= \frac{0.05}{0.05 + 0.00} = 1.00
 \end{aligned}$$

The **Fully Independent Model** is represented using a set of independent probability tables for all variables, learned from the training data as:



<i>Free</i>	$P(\textit{Free})$	<i>Caps</i>	$P(\textit{Caps})$	and	<i>Spam</i>	$P(\textit{Spam})$
Y	$\frac{20+1+5+0}{100} = 0.26$	Y	$\frac{20+1+20+3}{100} = 0.44$		Y	$\frac{20+5+20+2}{100} = 0.47$
N	$\frac{20+3+2+49}{100} = 0.74$	N	$\frac{5+0+2+49}{100} = 0.56$		N	$\frac{1+0+3+49}{100} = 0.53$

Using the same example, the conditional probability  $P(\textit{Spam} = Y | \textit{Free} = Y, \textit{Caps} = N)$  would be evaluated as:

$$\begin{aligned}
 P(\textit{Spam} = Y | \textit{Free} = Y, \textit{Caps} = N) &= \\
 &= \frac{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N)}{P(\textit{Free} = Y, \textit{Caps} = N)} \\
 &= \frac{P(\textit{Spam} = Y) \cdot P(\textit{Free} = Y) \cdot P(\textit{Caps} = N)}{P(\textit{Free} = Y) \cdot P(\textit{Caps} = N)} \\
 &= P(\textit{Spam} = Y) = 0.47
 \end{aligned}$$

The **Naïve Bayes Model** is represented using a set of conditional probability tables, learned from the training data as:

<i>Spam</i>	$P(\textit{Spam})$	<i>Caps</i>	<i>Spam</i>	$P(\textit{Caps}   \textit{Spam})$	<i>Free</i>	<i>Spam</i>	$P(\textit{Free}   \textit{Spam})$
Y	$\frac{20+5+20+2}{100} = 0.47$	Y	Y	$\frac{20+20}{20+5+20+2} \approx 0.8511$	Y	Y	$\frac{20+5}{20+5+20+2} \approx 0.5319$
N	$\frac{1+0+3+49}{100} = 0.53$	Y	N	$\frac{1+3}{1+0+3+49} \approx 0.0755$	Y	N	$\frac{1+0}{1+0+3+49} \approx 0.0189$
		N	Y	$\frac{5+2}{20+5+20+2} \approx 0.1489$	N	Y	$\frac{20+2}{20+5+20+2} \approx 0.4681$
		N	N	$\frac{0+49}{1+0+3+49} \approx 0.9245$	N	N	$\frac{3+49}{1+0+3+49} \approx 0.9811$

Using the same example, the conditional probability  $P(\textit{Spam} = Y | \textit{Free} = Y, \textit{Caps} = N)$  would be evaluated as:

$$\begin{aligned}
 P(\textit{Spam} = Y | \textit{Free} = Y, \textit{Caps} = N) &= \\
 &= \frac{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N)}{P(\textit{Free} = Y, \textit{Caps} = N)} \\
 &= \frac{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N)}{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N) + P(\textit{Spam} = N, \textit{Free} = Y, \textit{Caps} = N)}
 \end{aligned}$$

We first calculate:

$$\begin{aligned}
 P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N) &= \\
 &= P(\textit{Spam} = Y) \cdot P(\textit{Free} = Y | \textit{Spam} = Y) \cdot P(\textit{Caps} = N | \textit{Spam} = Y) \\
 &= 0.47 \cdot 0.5319 \cdot 0.1489 \approx 0.047248677
 \end{aligned}$$

and

$$\begin{aligned}
 P(\textit{Spam} = N, \textit{Free} = Y, \textit{Caps} = N) &= \\
 &= P(\textit{Spam} = N) \cdot P(\textit{Free} = Y | \textit{Spam} = N) \cdot P(\textit{Caps} = N | \textit{Spam} = N) \\
 &= 0.53 \cdot 0.0189 \cdot 0.9245 \approx 0.009260717
 \end{aligned}$$

so finally, based on the above equation,

$$\begin{aligned}
 P(\textit{Spam} = Y | \textit{Free} = Y, \textit{Caps} = N) &= \\
 &= \frac{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N)}{P(\textit{Spam} = Y, \textit{Free} = Y, \textit{Caps} = N) + P(\textit{Spam} = N, \textit{Free} = Y, \textit{Caps} = N)} \\
 &= \frac{0.047248677}{0.047248677 + 0.009260717} \approx 0.836120752
 \end{aligned}$$

**P0 Topics Discussion (3)**

- Discussion of individual projects as proposed in P0 submissions
- Projects discussed: P-07, P-08, P-09, P-11, P-12, P-13