



CSE6390E Introduction to Computational Linguistics

Fully Independent Model and Naïve Bayes Model

Presented by Nikolay Yakovets



Agenda

1

Preliminaries/Last Day

2

Fully Independent Model

3

Naïve Bayes Model

Slides adapted from lectures by
Vlado Keselj
Dalhousie University



Probabilistic Model

- ❖ Set of n **random** variables that capture the **outcome** in a model:

$$V = (V_1, V_2, \dots, V_n)$$

- ❖ Each variable can be assigned a value from a **finite set** of different values:

$$\{x_1, x_2, \dots, x_m\}$$

- ❖ **Random configuration**

- A tuple of values where each value is assigned to a variable

$$\vec{x} = (x_1, x_2, \dots, x_n)$$

$$V_1 = x_1, V_2 = x_2, \dots, V_n = x_n$$



Probabilistic Model

- ❖ In modeling our problem we assume that a **sequence of configurations** is drawn from some **random source**:

$$\begin{aligned}x^{(1)} &= (x_{11}, x_{12}, \dots, x_{1n}) \\x^{(2)} &= (x_{21}, x_{22}, \dots, x_{2n}) \\&\vdots \\x^{(t)} &= (x_{t1}, x_{t2}, \dots, x_{tn})\end{aligned}$$

- ❖ **Probabilistic Modeling in NLP** is a general framework for modeling NLP problems using random variables, random configurations, and **finding effective ways of reasoning** about probabilities of these configurations.



Computational Tasks

❖ Evaluation

- Compute probability of a **complete** configuration

❖ Simulation (aka Generation or Sampling)

- Generate **random** configurations

❖ Inference

▪ Marginalization

- Computing probability of a **partial** configuration

▪ Conditioning

- Computing **conditional probability** of a completion given an observation

▪ Completion

- Finding the **most probable completion**, given an observation

❖ Learning

- Learning **parameters** of a model from **data**



Example: Spam Detection

- ❖ **GOAL:** Automatically detect whether an arbitrary email message is spam or not
- ❖ **Have 3 random variables in the model:**
 - *Caps* is 'Y' if the message subject line does not contain lowercase letter, 'N' otherwise
 - *Free* is 'Y' if the word 'free' appears in the message subject line, 'N' otherwise
 - *Spam* is 'Y' if message is spam, and 'N' otherwise
- ❖ **Mailbox is a random source**
 - Randomly select 100 messages
 - Count how many times each configuration (each email) appears

<i>Free</i>	<i>Caps</i>	<i>Spam</i>	Number of messages
Y	Y	Y	20
Y	Y	N	1
Y	N	Y	5



❖ Last day

- Joint Distribution Model
 - Specify **complete** joint probability distribution, i.e. **the probability of each complete configuration**
- Drawbacks
 - **Memory cost** to store table
 - **Running-time** to do summations
 - The **sparse** data problem in learning

❖ Today

- Fully Independent Model



Fully Independent Model

- ❖ Assume all variables are **independent**:

$$P(V_1 = x_1, \dots, V_n = x_n) = P(V_1 = x_1) \cdots P(V_n = x_n)$$

- ❖ This is an **efficient** model

Evaluation

- **Small** number of parameters: $O(nm)$
- Represent each component of distribution separately
 - Fetch $P(V_j = x)$ from lookup table with m parameters
- ❖ **BUT usually a too strong assumption!**
 - Very **restricted** form of joint distribution
 - **Silly** model as far as real applications go, not very useful
- ❖ **Translated into SPAM example:**

$$P(\text{Free}, \text{Caps}, \text{Spam}) = P(\text{Free}) \cdot P(\text{Caps}) \cdot P(\text{Spam})$$



Example: Spam Detection

❖ Assume

- *Caps*, *Free* and *Spam* are **independent**

❖ Say we have the following data:

<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

Any message is **Spam** with **P=0.47** no matter what **Free** or **Caps** is!

❖ Estimate probability tables for independent variables:

<i>Free</i>	P(<i>Free</i>)
Y	$\frac{20+1+5+0}{100} = 0.26$
N	$\frac{20+3+2+49}{100} = 0.74$

<i>Caps</i>	P(<i>Caps</i>)
Y	$\frac{20+1+20+3}{100} = 0.44$
N	$\frac{5+0+2+49}{100} = 0.56$

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



Example: Spam Detection

- ❖ Say want to know the probability of configuration:

$$(Caps = Y, Free = N, Spam = N)$$

- ❖ Then the probability in fully independent model evaluates to the following:

$$P(Free = Y, Caps = N, Spam = N) =$$

$$= P(Free = Y) \cdot P(Caps = N) \cdot P(Spam = N) = 0.26 \cdot 0.56 \cdot 0.53$$

$$= 0.077168 \approx 0.08$$

Fetch from
lookup table

Evaluation



Computational Tasks

❖ Evaluation

❖ Simulation (Sampling)

- For every $j = 1, \dots, n$ **independently sample** x_j according to lookup table value of $P(V_j = x_j)$
- **Conjoin** (x_1, \dots, x_n) to form a complete configuration

❖ Inference – Marginalization

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

Sum-product computation:
 y_k is constant for summation over $y_{\{k+1\}}$!

❖ Only have to lookup and multiply n numbers!



Computational Tasks

❖ Inference – Conditioning

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

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



Example: Spam Detection

❖ Inference – Completion

$$y_{k+1}^*, \dots, y_n^* = \arg \max_{y_{k+1}, \dots, y_n} P(V_{k+1} = y_{k+1}, \dots, V_n = y_n | V_1 = x_1, \dots, V_k = x_k)$$

Conditioning

$$= \arg \max_{y_{k+1}, \dots, y_n} P(V_{k+1} = y_{k+1}) \cdots P(V_n = y_n)$$

$$= \left[\arg \max_{y_{k+1}} P(V_{k+1} = y_{k+1}) \right] \cdots \left[\arg \max_{y_n} P(V_n = y_n) \right]$$

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

Sum-product computation:
 y_k is constant for argmax over $y_{\{k+1\}}$!



Pros and Cons

❖ Joint Distribution Model vs. Fully Independent Model

Advantages

○ Efficient

$$O(nm) \text{ vs. } O(m^n)$$

Lot less to compute!

○ No sparse data problem

Disadvantages

▪ Too strong assumption!

▪ Too little structure

▪ Usually does not model accurately

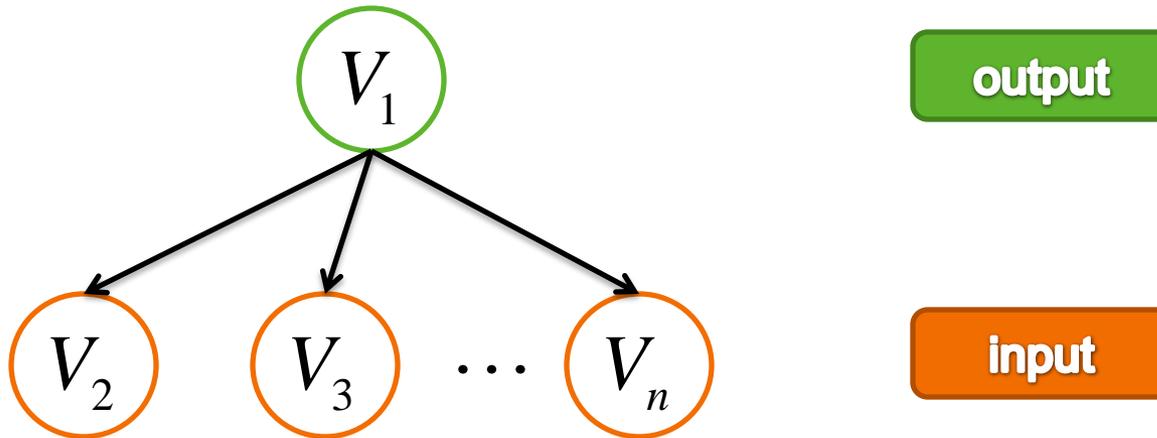
❖ **Structured probability models** are a compromise solution between these two models

- Address the issue of **sparse data**
- Model **important dependencies among** random variables



Naïve Bayes Model

- ❖ **Structured Probability Model**
- ❖ **Assume that all variables are **independent** except one distinguished variable – **class variable****
- ❖ **Graphical representation:**





Naïve Bayes Model

- ❖ Assume V_1 is the **output** variable and V_2, \dots, V_n are **input** variables
- ❖ Then classification problem is a **conditional probability computation problem**:

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

- ❖ After applying Bayes theorem we obtain:

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

Assume V_2, V_3, \dots, V_n are conditionally independent given V_1

The **conditional probabilities** can be efficiently computed and stored, and they eliminate **sparse data** problem



Computational Tasks

❖ Evaluation

Holds always

$$P(V_1 = x_1, \dots, V_n = x_n) =$$

$$\begin{aligned} \Rightarrow & P(V_1 = x_1)P(V_2 = x_2|V_1 = x_1)P(V_3 = x_3|V_1 = x_1, V_2 = x_2) \dots \\ & P(V_n = x_n|V_1 = x_1, V_2 = x_2, \dots, V_{n-1} = x_{n-1}) \end{aligned}$$

NB
≈

$$P(V_1 = x_1)P(V_2 = x_2|V_1 = x_1)P(V_3 = x_3|V_1 = x_1) \dots$$

$$P(V_n = x_n|V_1 = x_1)$$

By Naïve Bayes
assumption!

Fetch from lookup
table



Computational Tasks

❖ Assume

- *Free* and *Caps* are **input** variables
- *Spam* is **output** variable

❖ Naïve Bayes Assumption:

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

❖ Say we have the following 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



Computational Tasks

❖ Compute the following tables:

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

Maximum Likelihood Estimation (MLE)

The parameters of the model are estimated using a corpus

<i>Free</i>	<i>Spam</i>	$P(\textit{Free} \textit{Spam})$
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	$\frac{3+49}{1+0+3+49} \approx 0.9811$

<i>Caps</i>	<i>Spam</i>	$P(\textit{Caps} \textit{Spam})$
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	$\frac{0+49}{1+0+3+49} \approx 0.9245$



Computational Tasks

- ❖ Say want to **evaluate** the following configuration:

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

Fetch from lookup table

- ❖ Observe that in the Joint Distribution Model we had:

$$P(\text{Free} = Y, \text{Caps} = N, \text{Spam} = N) = 0.00$$

- ❖ This illustrates the fact that that the Naïve Bayes model is **less amenable** to the **sparse data problem!**



Computational Tasks

❖ Simulation:

- Configurations are sampled by:
 - Sample the **output** variable based on its **table**
 - Sample the **input** variables using corresponding **conditional tables**

❖ Inference – Marginalization

$$P(V_1 = x_1, \dots, V_k = x_k) =$$

$$P(V_1 = x_1)P(V_2 = x_2|V_1 = x_1)P(V_3 = x_3|V_1 = x_1) \dots$$

$$P(V_k = x_k|V_1 = x_1)$$

Fetch from lookup table

If the partial configuration includes the output variable, then compute marginal configuration as shown



Computational Tasks

❖ Inference – Conditioning – Example

Want to find the probability of message being spam given that header contains word “Free” and not all letters are uppercase

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

Fetch from lookup table

By Naïve Bayes assumption!

$$\begin{aligned} P(S = N, F = Y, C = N) &= \\ &= P(S = N)P(F = Y | S = N)P(C = N | S = N) \\ &= 0.53 \cdot 0.9245 \cdot 0.0189 \approx 0.093 \end{aligned}$$



Computational Tasks

❖ Inference – Conditioning – Example Continued

$$\begin{aligned} P(F = Y, C = N) & \stackrel{\text{By definition}}{=} \\ &= P(S = Y, F = Y, C = N) + P(S = N, F = Y, C = N) \\ &\approx P(S = Y)P(F = Y|S = Y)P(C = N|S = Y) + 0.093 \\ &= 0.47 \cdot 0.5319 \cdot 0.1489 + 0.093 \\ &\approx 0.0465 \end{aligned}$$

Fetch from lookup table

- Finally,

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



Computational Tasks

❖ Inference – Completion - Example

By definition

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

Does not depend on s

By Naïve Bayes Assumption

$$.. \Leftrightarrow \arg \max_s P(S = s)P(F = Y | S = s)P(C = N | S = s)$$

$$s = Y$$

$$A(S) = 0.0465$$

$$s = N$$

$$A(S) = 0.0093$$



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



Computational Tasks

❖ Learning

- **Maximum Likelihood Estimation:** The parameters are estimated using a corpus

❖ Number of Parameters

- A Naïve Bayes model with n variables V_1, \dots, V_n is described with tables:

	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
⋮	⋮	⋮
table $P(V_n V_1)$	m^2	m
sum	$m + (n - 1)m^2$	$1 + (n - 1)m$

Total: $O(m^2n)$



Pros and Cons

❖ Joint Distribution Model vs. Fully Independent Model

Advantages

○ Efficient

$$O(m^2n) \text{ vs. } O(nm) \text{ vs. } O(m^n)$$

Lot less to compute than Joint Distribution Model!

○ No sparse data problem

○ Surprisingly good performance (accuracy), e.g. in text classification

Disadvantages

▪ Can be over-simplifying!

▪ Practically, dependencies exist among attributes

▪ Cannot model more than one “output” variable

Thank You!

