Natural Language Processing CSCI 4152/6509 — Lecture 9 Introduction to Probabilistic Modeling

Instructors: Vlado Keselj

Time and date: 14:35 – 15:55, 23-Oct-2025 Location: Studley LSC-Psychology P5260

Previous Lecture

- Evaluation methods for Text Classification:
 - underfitting and overfitting
 - training error, train and test, n-fold cross-validation
- Text Clustering
- Discussion about evaluation methods for classifiers
- Similarity-based Text Classification
- CNG classification method
- Edit distance:
 - introduction, properties, dynamic programming approach, example, algorithm

Part III: Probabilistic Approach to NLP

Logical versus Plausible Reasoning

- As a part of AI (Artificial Intelligence), NLP follows two main approaches to *computer reasoning*, or *computer inference*:
- 1. logical reasoning
 - known also as classical, symbolic, knowledge-based Al
 - monotonic: once conclusion drawn, never retracted
 - certain: conclusions certain, given assumptions
- 2. plausible reasoning
 - examples: probabilistic, fuzzy logic, neural networks
 - non-monotonic
 - uncertain

Plausible Reasoning

- How to combine ambiguous, incomplete, and contradicting evidence to draw reasonable conclusions?
- Typical approach: make plausible inference of some hidden structure from observations
- Examples:

Observations (input)		Hidden Structure (output)
symptoms	\rightarrow	illness
pixel matrix	\rightarrow	object, relations
speech signal	\rightarrow	phonemes, words
word sequence	\rightarrow	meaning
sentence	\rightarrow	parse tree
word sequence	\rightarrow	POS tags, names, entities
words in e-mail Subject:	\rightarrow	Is message spam? Yes/No
text	\rightarrow	text category (class)

Probabilistic NLP as a Plausible Reasoning Approach

- Regular expressions and finite automata are example of logical or knowledge-based approach to NLP
- Plausible approaches to NLP:
 - 1. Probabilistic: use of Theory of Probability, also known as stochastic or statistical NLP
 - Alternative plausible approaches, examples:
 - 2. neural networks,
 - 3. kernel methods,
 - 4. fuzzy logic, fuzzy sets,
 - 5. Dempster-Shafer theory
 - 6. rough sets,
 - 7. default logic, ...

Review of Basics of Probability Theory

- You should have this background from previous courses; this is just a review,
 - discussed a bit in the textbook: [JM] 5.5, and [MS] 2.1
- Simple event or basic outcome
 - e.g., rolling a die, choosing a letter
- ullet Event space: the set of all outcomes, usually denoted Ω
- Event or outcome is a set of simple events or basic outcomes
- In other words event is any subset of Ω ; i.e., $A \subseteq \Omega$
- Each event is associated with a probability, which is a number between 0 and 1, inclusive: $0 \le P(A) \le 1$

Probability Examples

- P("rolling a 6 with a die") = 1/6
- Choosing a letter of English alphabet:
 - If we choose uniformly: $P('a') = 1/26 \approx 0.04$
 - Choosing from a text: $P('a') \approx 0.08$
 - Remember our output from "Tom Sawyer":

```
35697 0.1204 e
28897 0.0974 t
23528 0.0793 a
23264 0.0784 o
20200 0.0681 n
```

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ○○○

Probability Axioms

- (Nonnegativity) $P(A) \ge 0$, for any event A
- (Additivity) for disjoint events A and B, i.e., if $A, B \subset \Omega$ and $A \cap B = \emptyset$, then $P(A \cup B) = P(A) + P(B)$ or, more generally, $P(A_1 \cup A_2 \cup \ldots) = P(A_1) + P(A_2) + \ldots$
- (Normalization) $P(\Omega) = 1$, where Ω is the entire sample space.
- Some consequences of the above axioms are: $P(\emptyset) = 0$ and $P(\Omega A) = 1 P(A)$

Independent and Dependent Events

- Independent events A and B (definition): $P(A, B) = P(A) \cdot P(B)$
- Use of comma in: $P(A, B) = P(A \cap B)$
- Example: choosing two letters in text
 - Choosing independently: P('t') = 0.1, P('h') = 0.07, P('t', 'h') = 0.007
 - 2 Choosing two consecutive letters (dependent events): P('t', 'h') = 0.04

Conditional Probability

Conditional probability

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

ullet Expressing independency using conditional probability Two events A are B are independent if and only if:

$$P(A|B) = P(A)$$

This is an alternative definition of independent events.

Annotation with More Events

- There is a bit of flexibility in using notation; e.g.,
- $P(A, B, C) = P(A \cap B \cap C)$
- $P(A|B,C) = P(A|B \cap C)$
- $P(A, B, C|D, E, F) = P(A \cap B \cap C|D \cap E \cap F)$
- and so on.
- Three independent events: P(A, B, C) = P(A)P(B)P(C)
- Conditionally independent events

$$P(A, B|C) = P(A|C) \cdot P(B|C)$$

Bayes' Theorem

Bayes' theorem (one form):

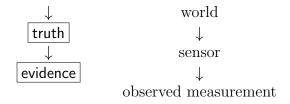
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

• The second form is based on breaking the set B into disjoint sets $B = A_1 \cup A_2 \cup \ldots$

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_i P(B|A_i) P(A_i)}$$

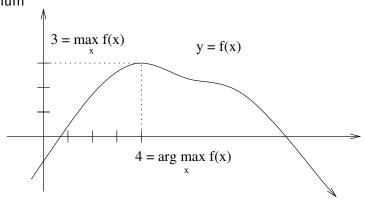
Bayesian Inference and Generative Models

- We will use Bayesian Inference on Generative Models
- Generative Models, also known as Forward Generative Models
- One way of representing knowledge with a probabilistic model



Notation Remark: max and argmax

- max is the maximum value of a function
- $\arg\max$ is an argument value for which function achieves the maximum



Bayesian Inference: Using Bayes' Theorem

Bayesian inference is a principle of combining evidence

$$\begin{array}{lll} \text{conclusion} &=& \underset{\text{possible truth}}{\operatorname{arg\ max}} \ P(\text{possible truth}|\text{evidence}) \\ &=& \underset{\text{possible truth}}{\operatorname{arg\ max}} \ \frac{P(\text{evidence}|\text{possible truth})P(\text{possible truth})}{P(\text{evidence})} \\ &=& \underset{\text{possible truth}}{\operatorname{arg\ max}} \ P(\text{evidence}|\text{possible truth})P(\text{possible truth}) \end{array}$$

application to speech recognition: acoustic model and language model

Bayesian Inference: Speech Recognition Example

- evidence → sound
- lacktriangle possible truth o utterance (words spoken)
- our best guess about utterance → utterance*

Probabilistic Modeling

- How do we create and use a probabilistic model?
- Model elements:
 - Random variables
 - Model configuration (Random configuration)
 - Variable dependencies
 - Model parameters
- Computational tasks

Random Variables

- Random variable V, defining an event as V=x for some value x from a domain of values D; i.e., $x \in D$
- ullet V=x is usually not a **basic** event due to having more variables
- An event with two random variables: $V_1 = x_1, V_2 = x_2$
- Multiple random variables: $\mathbf{V} = (V_1, V_2, ..., V_n)$

Model Configuration (Random Configuration)

• Full Configuration: If a model has n random variables, then a Full Model Configuration is an assignment of all the variables:

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

 Partial configuration: only some variables are assigned, e.g.:

$$V_1 = x_1, V_2 = x_2, \dots, V_k = x_k \quad (k < n)$$

Probabilistic Modeling in NLP

Probabilistic Modeling in NLP is a general framework for modeling NLP problems using random variables, random configurations, and an effective ways to reason about probabilities of these configurations.

Variable Independence and Dependence

- Random variables V_1 and V_2 are independent if $P(V_1\!=\!x_1,V_2\!=\!x_2)=P(V_1\!=\!x_1)P(V_2\!=\!x_2)$ for all x_1,x_2
- or expressed in a different way: $P(V_1 = x_1 | V_2 = x_2) = P(V_1 = x_1)$ for all x_1, x_2, x_3 .
- Random variables V_1 and V_2 are conditionally independent given V_3 if, for all x_1, x_2, x_3 : $P(V_1 = x_1, V_2 = x_2 | V_3 = x_3) = P(V_1 = x_1 | V_3 = x_3) P(V_2 = x_2 | V_3 = x_3)$
- or $P(V_1 = x_1 | V_2 = x_2, V_3 = x_3) = P(V_1 = x_1 | V_3 = x_3)$

Computational Tasks in Probabilistic Modeling

- 1. Evaluation: compute probability of a complete configuration
- 2. Simulation: generate random configurations
- 3. Inference: has the following sub-tasks:
 - 3.a Marginalization: computing probability of a partial configuration,
 - 3.b Conditioning: computing conditional probability of a completion given an observation,
 - 3.c Completion: finding the most probable completion, given an observation
- 4. Learning: learning parameters of a model from data.

Illustrative Example: Spam Detection

- the problem of spam detection
- a probabilistic model for spam detection; random variables:
 - Caps = 'Y' if the message subject line does not contain lowercase letter, 'N' otherwise,
 - Free = 'Y' if the word 'free' appears in the message subject line (letter case is ignored), 'N' otherwise, and
 - Spam = 'Y' if the message is spam, and 'N' otherwise.
- one random configuration represents one e-mail message

Random Sample

• Data based on sample of 100 email messages

Free	Caps	Spam	Number of messages
Y	Y	Y	20
Υ	Y	N	1
Υ	N	Y	5
Υ	N	N	0
N	Y	Y	20
N	Y	N	3
N	N	Y	2
N	N	N	49
Total:			100

What are examples of computational tasks in this example?

Joint Distribution Model

 Probability of each complete configuration is specified; i.e., the joint probability distribution:

$$P(V_1 = x_1, ..., V_n = x_n)$$

- If each variable can have m possible values, the model has m^n parameters
- The model is a large lookup table: For each full configuration $\mathbf{x}=(V_1\!=\!x_1,...,V_n\!=\!x_n)$, a parameter $p_{\mathbf{x}}$ is specified such that

$$0 \le p_{\mathbf{x}} \le 1$$
 and $\sum_{\mathbf{x}} p_{\mathbf{x}} = 1$

Example: Spam Detection (Joint Distribution Model)

MLE — Maximum Likelihood Estimation of probabilities:

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

Computational Tasks in Joint Distribution Model:

1. Evaluation

- Evaluate the probability of a complete configuration $\mathbf{x} = (x_1, ..., x_n)$.
- Use a table lookup:

$$P(V_1 = x_1, ..., V_n = x_n) = p_{(x_1, x_2, ..., x_n)}$$

For example:

$$P(Free = Y, Caps = N, Spam = N) = 0.00$$

- This example illustrates the sparse data problem
- Inferred that the probability is zero since the configuration was not seen before.

2. Simulation (Joint Distribution Model)

2. Simulation (Joint Distribution Model)

- Simulation is performed by randomly selecting a configuration according to the probability distribution in the table
- Known as the "roulette wheel" method
- 1. Divide the interval [0,1] into subintervals of the lengths: p_1 , p_2 , ..., p_{m^n} : $I_1 = [0,p_1)$, $I_2 = [p_1,p_1+p_2)$, $I_3 = [p_1+p_2,p_1+p_2+p_3)$, ... $I_{m^n} = [p_1+p_2+\ldots+p_{m^n-1},1)$
- 2. Generate a random number r from the interval [0,1)
- 3. r will fall exactly into one of the above intervals, e.g.: $I_i = [p_1 + \ldots + p_{i-1}, p_1 + \ldots + p_{i-1} + p_i)$
- 4. Generate the configuration number i from the table
- 5. Repeat steps 2–4 for as many times as the number of configurations we need to generate

Joint Distribution Model: 3. Inference

3.a Marginalization

• Compute the probability of an *incomplete* configuration $P(V_1 = x_1, ..., V_k = x_k)$, where k < n:

$$P(V_1 = x_1, ..., V_k = x_k)$$

$$= \sum_{y_{k+1}} ... \sum_{y_n} P(V_1 = x_1, ..., V_k = x_k, V_{k+1} = y_{k+1}, ..., V_n = y_n)$$

$$= \sum_{y_{k+1}} ... \sum_{y_n} p_{(x_1, ..., x_k, y_{k+1}, ..., y_n)}$$

 Implementation: iterate through the lookup table and accumulate probabilities for matching configurations

Joint Distribution Model: 3.b Conditioning

 Compute a conditional probability of assignments of some variables given the assignments of other variables; for example,

$$P(V_1 = x_1, ..., V_k = x_k | V_{k+1} = y_1, ..., V_{k+l} = y_l)$$

$$= \frac{P(V_1 = x_1, ..., V_k = x_k, V_{k+1} = y_1, ..., V_{k+l} = y_l)}{P(V_{k+1} = y_1, ..., V_{k+l} = y_l)}$$

- This task can be reduced to two marginalization tasks
- If the configuration in the numerator happens to be a full configuration, that the task is even easier and reduces to one evaluation and one marginalization.

Joint Distribution Model: 3.c Completion

• Find the most probable completion $(y_{k+1}^*,...,y_n^*)$ given a partial configuration $(x_1,...,x_k)$.

$$\begin{split} y_{k+1}^*,...,y_n^* &= \underset{y_{k+1},...,y_n}{\arg\max} & \mathbf{P}(V_{k+1}\!=\!y_{k+1},...,V_n\!=\!y_n|V_1\!=\!x_1,...,V_k\!=\!x_k) \\ &= \underset{y_{k+1},...,y_n}{\arg\max} & \frac{\mathbf{P}(V_1\!=\!x_1,...,V_k\!=\!x_k,V_{k+1}\!=\!y_{k+1},...,V_n\!=\!y_n)}{\mathbf{P}(V_1\!=\!x_1,...,V_k\!=\!x_k,V_{k+1}\!=\!y_{k+1},...,V_n\!=\!y_n)} \\ &= \underset{y_{k+1},...,y_n}{\arg\max} & \mathbf{P}(V_1\!=\!x_1,...,V_k\!=\!x_k,V_{k+1}\!=\!y_{k+1},...,V_n\!=\!y_n) \\ &= \underset{y_{k+1},...,y_n}{\arg\max} & p_{(x_1,...,x_k,y_{k+1},...,y_n)} \end{split}$$

• Implementation: search through the model table, and from all configurations that satisfy assignments in the partial configuration, chose the one with maximal probability.

Joint Distribution Model: 4. Learning

- Estimate the parameters in the model based on given data
- Use Maximum Likelihood Estimation (MLE)
- Count all full configurations, divide the count by the total number of configurations, and fill the table:

$$p_{(x_1,\ldots,x_n)} = \frac{\#(V_1 = x_1,\ldots,V_n = x_n)}{\#(*,\ldots,*)}$$

 With a large number of variables the data size easily becomes insufficient and we get many zero probabilities — sparse data problem

Drawbacks of Joint Distribution Model

- memory cost to store table,
- running-time cost to do summations, and
- the sparse data problem in learning (i.e., training).

Other probability models are found by specifying specialized joint distributions, which satisfy certain independence assumptions.

The goal is to impose structure on joint distribution $P(V_1 = x_1, ..., V_n = x_n)$. One key tool for imposing structure is variable independence.

Fully Independent Model

Assumption: all variables are independent

$$P(V_1 = x_1, ..., V_n = x_n) = P(V_1 = x_1) \cdot \cdot \cdot P(V_n = x_n).$$

- \bullet Efficient model with a small number of parameters: O(nm)
- Drawback: usually a too strong assumption
- Fully independent model for the Spam example:

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

Fully Independent Model: [4.] Learning

Spam example:

Free	$P(\mathit{Free})$			
Υ	$\frac{20+1+5+0}{100} = 0.26$	and similarly,		
N	$\frac{20+3+2+49}{100} = 0.74$			
Caps	P(Caps)		Spam	P(Spam)
Y	100	and	Υ	$\frac{20+5+20+2}{100} = 0.47$
N	$\frac{5+0+2+49}{100} = 0.56$	_	N	$\frac{1+0+3+49}{100} = 0.53$

Hence, in this model any message is a spam with probability 0.47, no matter what the values of *Caps* and *Free* are.

Evaluation Example

As an example of evaluation, the probability of configuration (Caps = Y, Free = N, Spam = N) in the fully independent model is:

$$\begin{split} & \text{P}(\textit{Free} = Y, \textit{Caps} = N, \textit{Spam} = N) = \\ & = & \text{P}(\textit{Free} = Y) \cdot \text{P}(\textit{Caps} = N) \cdot \text{P}(\textit{Spam} = N) = \\ & = & 0.26 \cdot 0.56 \cdot 0.53 \\ & = & 0.077168 \approx 0.08 \end{split}$$

Fully Independent Model: 2. Simulation

- For j=1,...,n, independently draw x_j according to $P(V_j\!=\!x_j)$ using "roulette wheel" for one variable
- Conjoin $(x_1,...,x_n)$ to form a complete configuration.

3. Inference in Fully Independent Model

3.a Marginalization in Fully Independent Model

The probability of a partial configuration $(V_1 = x_1, \dots, V_k = x_k)$ is

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

This formula can be obvious, but it can also be derived.

Derivation of Marginalization Formula

$$P(V_{1} = x_{1}, ..., V_{k} = x_{k}) = \sum_{y_{k+1}} \cdots \sum_{y_{n}} P(V_{1} = x_{1}, ..., V_{k} = x_{k}, V_{k+1} = y_{k+1}, ..., V_{n} = y_{n})$$

$$= \sum_{y_{k+1}} \cdots \sum_{y_{n}} P(V_{1} = x_{1}) \cdots P(V_{k} = x_{k}) P(V_{k+1} = y_{k+1}) \cdots P(V_{n} = y_{n})$$

$$= P(V_{1} = x_{1}) \cdots P(V_{k} = x_{k}) \left[\sum_{y_{k+1}} P(V_{k+1} = y_{k+1}) \left[\sum_{y_{k+2}} \cdots \left[\sum_{y_{n}} P(V_{n} = y_{n}) \right] \right] \right]$$

$$= P(V_{1} = x_{1}) \cdots P(V_{k} = x_{k}) \left[\sum_{y_{k+1}} P(V_{k+1} = y_{k+1}) \right] \cdots \left[\sum_{y_{n}} P(V_{n} = y_{n}) \right]$$

$$= P(V_{1} = x_{1}) \cdots P(V_{k} = x_{k})$$

A Note on Sum-Product Computation

$$\sum_{a} \sum_{b} f(a)g(b) = \sum_{a} f(a) \left(\sum_{b} g(b)\right)$$
(because $f(a)$ is a constant for summation over b)
$$= \left(\sum_{b} g(b)\right) \cdot \left(\sum_{a} f(a)\right)$$
(because $\sum_{b} g(b)$ is a constant for sumation over a)
$$= \left(\sum_{a} f(a)\right) \cdot \left(\sum_{b} g(b)\right)$$

Similar Note for Max-Product Computation

If we assume that $f(a) \ge 0$ and $g(b) \ge 0$, the same rule applies for \max_a and \max_b :

$$\max_{a} \max_{b} f(a)g(b) =$$

$$= \max_{a} f(a) \left(\max_{b} g(b)\right)$$
(because $f(a)$ is a constant for maximization over b)
$$= \left(\max_{b} g(b)\right) \cdot \left(\max_{a} f(a)\right)$$
(because $\max_{b} g(b)$ is a constant for maximization over a)
$$= \left(\max_{a} f(a)\right) \cdot \left(\max_{b} g(b)\right)$$

3.c Completion in Fully Independent Model

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

$$= \underset{y_{k+1}, ..., y_n}{\operatorname{arg max}} P(V_{k+1} = y_{k+1}) \cdots P(V_n = y_n)$$

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

Joint Distribution Model vs. Fully Independent Model

- Fully Independent Model addresses some issues of the Joint Distribution Model
- Efficient and small number of parameters
- However: too strong assumption, no structure
- Too trivial to be usable
- Better method: Structured probability models
 - compromise between no dependence and too much dependence