

— INF4820 —
Algorithms for AI and NLP

Basic Probability Theory & Language Models

Murhaf Fares & Stephan Oepen

Language Technology Group (LTG)

October 20, 2016



- ▶ Vector space model
- ▶ Classification
 - ▶ Rocchio
 - ▶ kNN
- ▶ Clustering
 - ▶ K-means

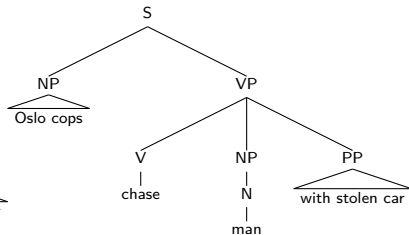
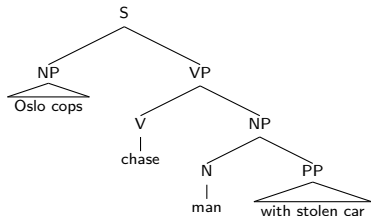
Point-wise prediction; geometric models.

Structured prediction; probabilistic models.

- ▶ Sequences
 - ▶ Language models
- ▶ Labelled sequences
 - ▶ Hidden Markov Models
- ▶ Trees
 - ▶ Statistical (Chart) Parsing

Probabilistic models to determine the *most likely interpretation*.

- ▶ Which string is **most likely**?
 - ▶ *She studies morphosyntax* vs. *She studies more faux syntax*
- ▶ Which category sequence is **most likely** for *Time flies like an arrow*?
 - ▶ $\text{Time}_N \text{flies}_N \text{like}_V \text{an}_D \text{arrow}_N$
 - ▶ $\text{Time}_N \text{flies}_V \text{like}_P \text{an}_D \text{arrow}_N$
- ▶ Which syntactic analysis is **most likely**?





- ▶ Experiment (or trial)
 - ▶ the process we are observing
- ▶ Sample space (Ω)
 - ▶ the set of all possible outcomes of a random experiment
- ▶ Event(s)
 - ▶ the subset of Ω we are interested in
- ▶ Our goal is to assign probability to certain events
 - ▶ $P(A)$ is the probability of event A , a real number $\in [0, 1]$



- ▶ Experiment (or trial)
 - ▶ rolling a die
- ▶ Sample space (Ω)
 - ▶ $\Omega = \{1, 2, 3, 4, 5, 6\}$
- ▶ Event(s)
 - ▶ $A =$ rolling a six: $\{6\}$
 - ▶ $B =$ getting an even number: $\{2, 4, 6\}$
- ▶ Our goal is to assign probability to certain events
 - ▶ $P(A) = ?$ $P(B) = ?$

- ▶ Experiment (or trial)
 - ▶ flipping two coins
- ▶ Sample space (Ω)
 - ▶ $\Omega = \{HH, HT, TH, TT\}$
- ▶ Event(s)
 - ▶ A = the same outcome both times: $\{HH, TT\}$
 - ▶ B = at least one head: $\{HH, HT, TH\}$
- ▶ Our goal is to assign probability to certain events
 - ▶ $P(A) = ?$ $P(B) = ?$

- ▶ Experiment (or trial)
 - ▶ rolling two dice
- ▶ Sample space (Ω)
 - ▶ $\Omega = \{11, 12, 13, 14, 15, 16, 21, 22, 23, 24, \dots, 63, 64, 65, 66\}$
- ▶ Event(s)
 - ▶ $A =$ results sum to 6: $\{15, 24, 33, 42, 51\}$
 - ▶ $B =$ both results are even: $\{22, 24, 26, 42, 44, 46, 62, 64, 66\}$
- ▶ Our goal is to assign probability to certain events
 - ▶ $P(A) = \frac{|A|}{|\Omega|}$ $P(B) = \frac{|B|}{|\Omega|}$

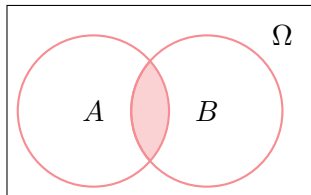
Probability Axioms

- ▶ $0 \leq P(A) \leq 1$
- ▶ $P(\Omega) = 1$
- ▶ $P(A \cup B) = P(A) + P(B)$ where A and B are mutually exclusive

More useful axioms

- ▶ $P(A) = 1 - P(\neg A)$
- ▶ $P(\emptyset) = 0$

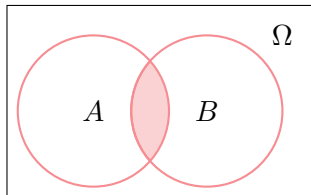
- ▶ $P(A, B)$: probability that both A and B happen
- ▶ also written: $P(A \cap B)$ or $P(A, B)$



What is the probability, when throwing two fair dice, that

- ▶ A : the results sum to 6 and
- ▶ B : at least one result is a 1?

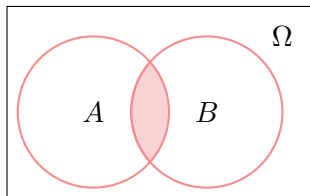
- ▶ $P(A, B)$: probability that both A and B happen
- ▶ also written: $P(A \cap B)$ or $P(A, B)$



What is the probability, when throwing two fair dice, that

- ▶ A : the results sum to 6 and
- ▶ B : at least one result is a 1?

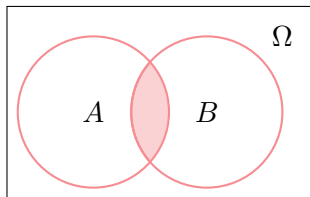
- ▶ $P(A, B)$: probability that both A and B happen
- ▶ also written: $P(A \cap B)$ or $P(A, B)$



What is the probability, when throwing two fair dice, that

- ▶ A : the results sum to 6 and $\frac{5}{36}$
- ▶ B : at least one result is a 1?

- ▶ $P(A, B)$: probability that both A and B happen
- ▶ also written: $P(A \cap B)$ or $P(A, B)$



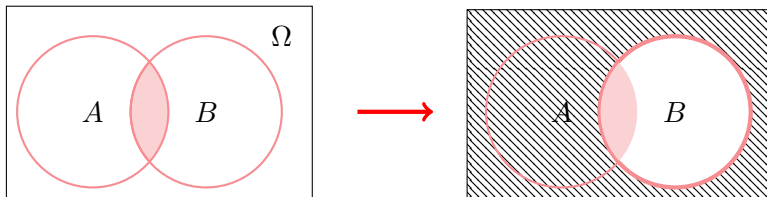
What is the probability, when throwing two fair dice, that

- ▶ A : the results sum to 6 and $\frac{5}{36}$
- ▶ B : at least one result is a 1? $\frac{11}{36}$

Often, we have partial knowledge about the outcome of an experiment.

What is the probability $P(A|B)$, when throwing two fair dice, that

- ▶ A : the results sum to 6 given
- ▶ B : at least one result is a 1?



$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (\text{where } P(B) > 0)$$

Joint probability is symmetric:

$$\begin{aligned} P(A \cap B) &= P(A) P(B|A) \\ &= P(B) P(A|B) \quad (\text{multiplication rule}) \end{aligned}$$

More generally, using the **chain rule**:

$$P(A_1 \cap \dots \cap A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \dots P(A_n | \cap_{i=1}^{n-1} A_i)$$

The chain rule will be very useful to us through the semester:

- ▶ it allows us to break a complicated situation into parts;
- ▶ we can choose the breakdown that suits our problem.

- ▶ Let A be the event that it rains tomorrow $P(A) = \frac{1}{3}$
- ▶ Let B be the event that flipping a coin results in heads $P(B) = \frac{1}{2}$
- ▶ What is $P(A|B)$?

If knowing event B is true has no effect on event A , we say

A and B are independent of each other.

If A and B are independent:

- ▶ $P(A \cap B) = P(A) P(B)$
- ▶ $P(A|B) = P(A)$
- ▶ $P(B|A) = P(B)$

- ▶ Your friend, Yoda, wakes up in the morning feeling sick.
- ▶ He uses a website to diagnose his disease by entering the symptoms
- ▶ The website returns that 99% of the people who had a disease D had the same symptoms Yoda has.
- ▶ Yoda freaks out, comes to your place and tells you the story.
- ▶ You are more relaxed, you continue reading the web page Yoda started reading, and you find the following information:
 - ▶ The prevalence of disease D : 1 in 1000 people
 - ▶ The reliability of the symptoms:
 - ▶ False negative rate: 1%
 - ▶ False positive rate: 2%

What is the probability that he has the disease?

Given:

- ▶ event A: has disease
- ▶ event B: has the symptoms

We know:

- ▶ $P(A) = 0.001$
- ▶ $P(B|A) = 0.99$
- ▶ $P(B|\neg A) = 0.02$

We want

- ▶ $P(A|B) = ?$

	A	$\neg A$	
B	0.00099	0.01998	0.02097
$\neg B$	0.00001	0.97902	0.97903
	0.001	0.999	1

$$P(A) = 0.001; \quad P(B|A) = 0.99; \quad P(B|\neg A) = 0.02$$

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

$$P(\neg A \cap B) = P(B|\neg A)P(\neg A)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.00099}{0.02097} = 0.0472$$



- From the two 'symmetric' sides of the joint probability equation:

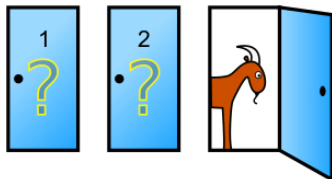
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- reverses the order of dependence (which can be useful)
- in conjunction with the chain rule, allows us to determine the probabilities we want from the probabilities we know

Bonus: The Monty Hall Problem



- ▶ On a gameshow, there are three doors.
- ▶ Behind 2 doors, there is a goat.
- ▶ Behind the 3rd door, there is a car.



- ▶ The contestant selects a door that she hopes has the car behind it.
- ▶ Before she opens that door, the gameshow host opens one of the other doors to reveal a goat.
- ▶ The contestant now has the choice of opening the door she originally chose, or switching to the other unopened door.

What should she do?

- ▶ Now that we have ____ the basics in probability theory we can move to a new topic.
- ▶ Det var en ____?
- ▶ Je ne parle pas ____?

Natural language contains redundancy, hence can be predictable.

Previous context can constrain the next word

- ▶ semantically;
 - ▶ syntactically;
- by frequency – **language models**.

Language model: a probabilistic model that assigns *approximate* probability to an arbitrary sequence of words.

- ▶ Machine translation
 - ▶ She is going home vs. She is going house
- ▶ Speech recognition
 - ▶ She studies morphosyntax vs. She studies more faux syntax
- ▶ Spell checkers
 - ▶ *Their* are many NLP applications that use language models.
- ▶ Input prediction (predictive keyboards)

- ▶ A probabilistic **language model** M assigns probabilities $P_M(x)$ to all strings x in language L .
 - ▶ L is the sample space
 - ▶ $0 \leq P_M(x) \leq 1$
 - ▶ $\sum_{x \in L} P_M(x) = 1$

- ▶ Given a sentence $S = w_1 \dots w_n$, we want to estimate $P(S)$
- ▶ $P(S)$ is the joint probability over the words in S : $P(w_1, w_2 \dots, w_n)$
- ▶ We can calculate the probability of S using the chain rule:

$$P(w_1 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1 \cap w_2) \dots P(w_n | \cap_{i=1}^{n-1} w_i)$$

- ▶ Example:

$$\begin{aligned} P(I \text{ want to go to the beach}) = \\ P(I) P(\text{want}|I) P(\text{to}|I \text{ want}) P(\text{go}|I \text{ want to}) P(\text{to}|I \text{ want to go}) \dots \end{aligned}$$

- ▶ Given a sentence $S = w_1 \dots w_n$, we want to estimate $P(S)$
- ▶ $P(S)$ is the joint probability over the words in S : $P(w_1, w_2 \dots, w_n)$
- ▶ We can calculate the probability of S using the chain rule:

$$P(w_1 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1 \cap w_2) \dots P(w_n | \cap_{i=1}^{n-1} w_i)$$

- ▶ Example:

$$\begin{aligned} P(I \text{ want to go to the beach}) = \\ P(I) P(\text{want}|I) P(\text{to}|I \text{ want}) P(\text{go}|I \text{ want to}) P(\text{to}|I \text{ want to go}) \dots \end{aligned}$$

- ▶ Given a sentence $S = w_1 \dots w_n$, we want to estimate $P(S)$
- ▶ $P(S)$ is the joint probability over the words in S : $P(w_1, w_2 \dots, w_n)$

Recall The Chain Rule

$$P(A_1 \cap \dots \cap A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \dots P(A_n | \cap_{i=1}^{n-1} A_i)$$

- ▶ We can calculate the probability of S using the chain rule:

$$P(w_1 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1 \cap w_2) \dots P(w_n | \cap_{i=1}^{n-1} w_i)$$

- ▶ Example:

$$\begin{aligned} P(I \text{ want to go to the beach}) = \\ P(I) P(want|I) P(to|I \text{ want}) P(go|I \text{ want to}) P(to|I \text{ want to go}) \dots \end{aligned}$$

$P(I \text{ want to go to the beach to read about Markov assumption}) =$
 $P(I) P(\text{want}|I) P(\text{to}|I \text{ want}) P(\text{go}|I \text{ want to})$
... $P(\text{assumption}|I \text{ want to go to the beach to read about Markov})$



- Simplifying assumption (limited history):

$$P(\text{assumption}|I \text{ want to go to the beach to read about Markov}) \approx P(\text{assumption}|\text{Markov})$$

- Or:

$$P(\text{assumption}|I \text{ want to go to the beach to read about Markov}) \approx P(\text{assumption}|\text{about Markov})$$

We simplify using the **Markov assumption** (limited history):

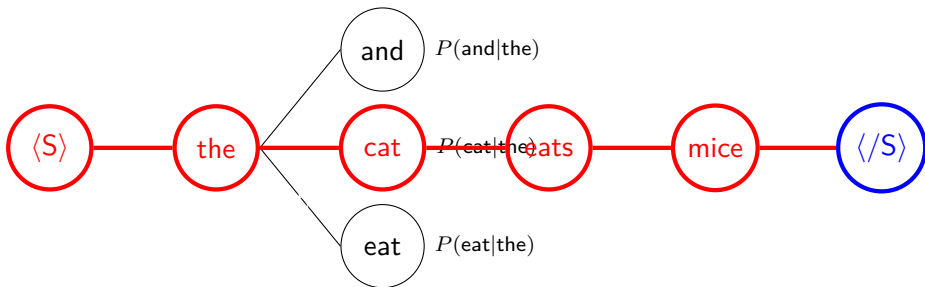
the last $n - 1$ elements can approximate the effect of the full sequence.

$$P(w_1 \dots w_n) \approx \prod_i P(w_i | w_{i-n+1} \dots w_{i-1})$$

We call these short sequences of words n -grams:

- ▶ bigrams ($n = 2$): *I want, want to, to go, go to, to the, the beach*
- ▶ trigrams ($n = 3$): *I want to, want to go, to go to, go to the*
- ▶ 4-grams ($n = 4$): *I want to go, want to go to, to go to the*

- We can think of an n -gram model as a probabilistic automaton that generates sentences.



$$P(S) = P(\text{the}|\langle S \rangle) P(\text{cat}|\text{the}) P(\text{eats}|\text{cat}) P(\text{mice}|\text{eats}) P(\langle /S \rangle|\text{mice})$$

An n -gram language model records the n -gram conditional probabilities:

$$\begin{array}{ll} P(I|\langle S \rangle) &= 0.0429 & P(to|go) &= 0.1540 \\ P(want|I) &= 0.0111 & P(the|to) &= 0.1219 \\ P(to|want) &= 0.4810 & P(beach|the) &= 0.0006 \\ P(go|to) &= 0.0131 & & \end{array}$$

We calculate the probability of a sentence as (assuming bi-grams):

$$\begin{aligned} P(w_1^n) &\approx \prod_{k=1}^n P(w_k|w_{k-1}) \\ &\approx P(I|\langle S \rangle) \times P(want|I) \times P(to|want) \times P(go|to) \times P(to|go) \times \\ &\quad P(the|to) \times P(beach|the) \\ &\approx 0.0429 \times 0.0111 \times 0.4810 \times 0.0131 \times 0.1540 \times \\ &\quad 0.1219 \times 0.0006 = 3.38 \times 10^{-11} \end{aligned}$$

How to estimate the probabilities of n -grams?

By counting:

$$P(w_i | w_{i-n+1}, \dots, w_{i-1}) = \frac{\text{Count}(w_{i-n+1}, \dots, w_i)}{\text{Count}(w_{i-n+1}, \dots, w_{i-1})}$$

E.g. for trigrams:

$$P(\text{go} | \text{want to}) = \frac{C(\text{want to go})}{C(\text{want to})}$$

The probabilities are estimated using the **relative frequencies** of observed outcomes. This process is called **Maximum Likelihood Estimation (MLE)**.

- ▶ Suppose “Chinese” occurs 400 times in a corpus of 1M words
- ▶ MLE of “Chinese” is $\frac{400}{1000000} = 0.0004$
- ▶ Is this a good estimate for all corpora?

w_1	w_2	$C(w_1w_2)$	$C(w_1)$	$P(w_2 w_1)$
$\langle S \rangle$	I	1039	24243	0.0429
I	want	46	4131	0.0111
want	to	101	210	0.4810
to	go	128	9778	0.0131
go	to	59	383	0.1540
to	the	1192	9778	0.1219
the	beach	14	22244	0.0006

What's the probability of *Others want to go to the beach* ?

$$P(\text{Others} | \langle S \rangle) = ?$$



- ▶ Like many statistical models, n -grams is dependent on the training corpus
- ▶ Data sparseness: many perfectly acceptable n -grams will not be observed in the training data
- ▶ Zero counts will result in an estimated probability of 0
- ▶ But we still want to have good intuition about more likely sentences:
 - ▶ Others want to go to the beach.
 - ▶ Others the beach go to want to.



- ▶ Reassign some of the probability mass of frequent events to less frequent (or unseen) events.
- ▶ Known as **smoothing** or **discounting**
- ▶ The simplest approach is **Laplace** ('add-one') smoothing:

$$P_L(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V}$$

“Others want to go to the beach”

w_1	w_2	$C(w_1 w_2)$	$C(w_1)$	$P(w_2 w_1)$	$P_L(w_2 w_1)$
$\langle S \rangle$	I	1039	24243	0.0429	0.01934
$\langle S \rangle$	Others	17	24243	0.0007	0.00033
I	want	46	4131	0.0111	0.00140
Others	want	0	4131	0	0.00003
want	to	101	210	0.4810	0.00343
to	go	128	9778	0.0131	0.00328
go	to	59	383	0.1540	0.00201
to	the	1192	9778	0.1219	0.03035
the	beach	14	22244	0.0006	0.00029

$$P_L(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + 29534}$$

$S =$ I want to go to the beach

$$\begin{aligned}P(S) &= P(I | \langle S \rangle) \times P(\text{want} | I) \times P(\text{to} | \text{want}) \times P(\text{go} | \text{to}) \times P(\text{to} | \text{go}) \times \\&\quad P(\text{the} | \text{to}) \times P(\text{beach} | \text{the}) \\&= 0.0429 \times 0.0111 \times 0.4810 \times 0.0131 \times 0.1540 \times \\&\quad 0.1219 \times 0.0006 = 3.38 \times 10^{-11}\end{aligned}$$

- ▶ Multiplying many small probabilities \rightarrow Risk underflow
- ▶ Solution: work in $\log(\text{arithmetic})$ space:
 - ▶ $\log(AB) = \log(A) + \log(B)$
 - ▶ hence $P(A)P(B) = \exp(\log(A) + \log(B))$
 - ▶ $\log(P(S)) = -1.368 + -1.954 + -0.317 + -1.882 \dots$



- ▶ The likelihood of the next word depends on its context.
- ▶ We can calculate this using the chain rule:

$$P(w_1^N) = \prod_{i=1}^N P(w_i | w_1^{i-1})$$

- ▶ In an n -gram model, we approximate this with a Markov chain:

$$P(w_1^N) \approx \prod_{i=1}^N P(w_i | w_{i-n+1}^{i-1})$$

- ▶ We use Maximum Likelihood Estimation to estimate the conditional probabilities.
- ▶ Smoothing techniques are used to avoid zero probabilities.