

N-gram language models

NLP Week3

Thanks to Dan Jurafsky for some of the slides and inspiration!

Plan for today


1. **Uni-gram or bag-of-words models**
2. **N-gram language models and conditional probabilities**
3. **Evaluating language models**
4. **Scoring vs Sampling for probabilistic generative models**
5. ***Group exercises***

What is a language model?

- They are probability models over contextualized language. They can assign probability to words in context and can thus do things like predict the next word in a sentence, or give you the probability of a whole sentence under the model.
- LLM = Large **language model**

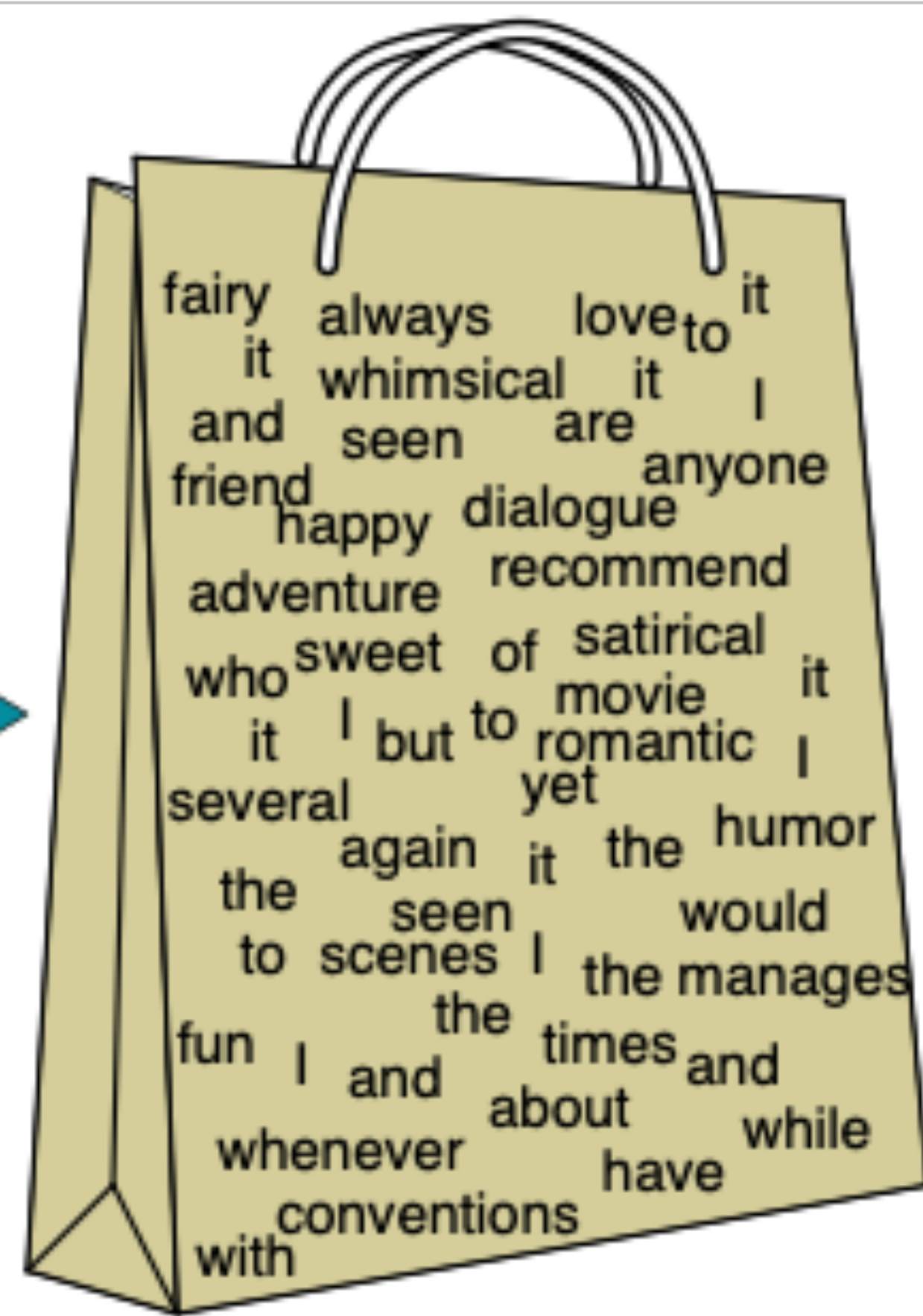
We will build language models adding to each layer of their complexity:

LANGUAGE MODELS

- 
- **N-gram language models** — Probabilistic models of language in context
 - **Embeddings** — vector semantics
 - **Recurrent neural language models, LSTMs** — neural networks for language
 - **Transformers** — large language model (LLM) architectures
 - **Encoder models** — scoring tasks
 - **Encoder-decoder models** — generative tasks
 - **Finetuning and in-context learning** — current foundation model paradigm

Bag-of-words (unigram) language model

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

Bag-of-words (unigram) language model

- Bag-of-words model (BoW) also known as a unigram model is a language model which assigns probability to words based on no prior context.
- So the probability of a word is simply its frequency, or normalized count in a document or corpus.

Sequential dependencies

‘Not only was it the greatest cast’

-> positive

‘Only it was not the greatest cast’

-> negative

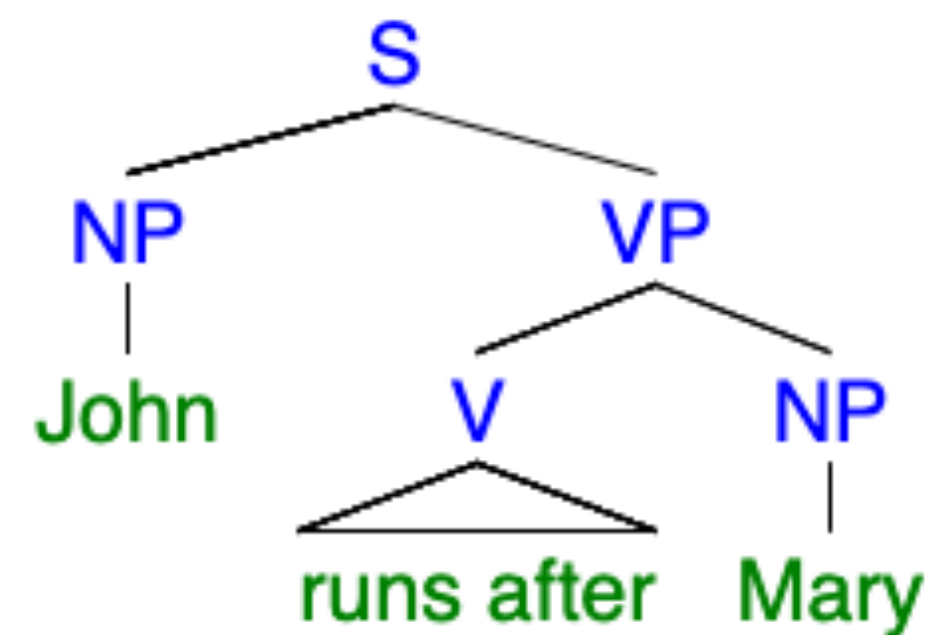
Unigram features are the same!

[cast, greatest, it, only, not, the, was]

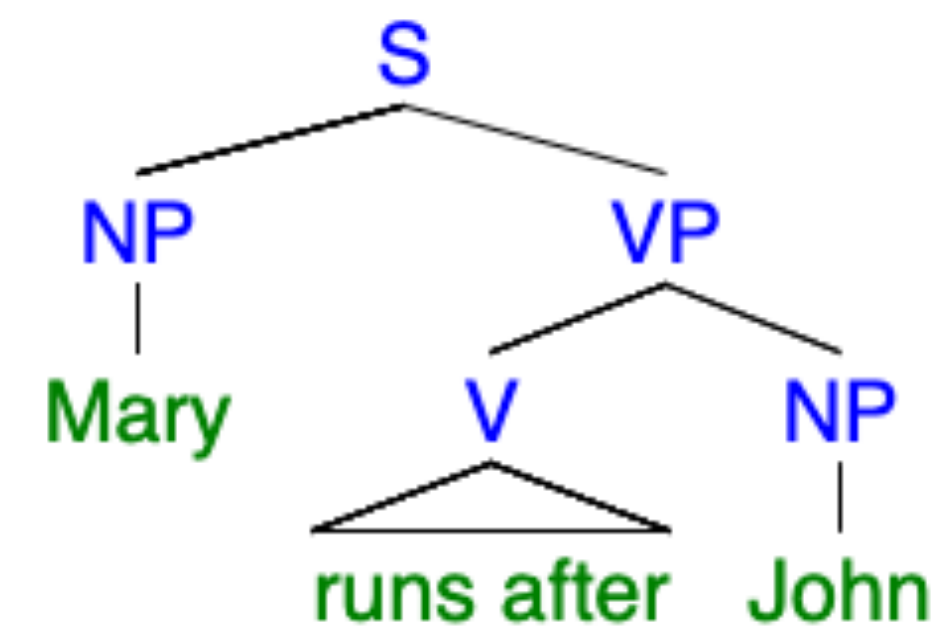
The order of words matters, there is meaning in sequential dependencies.

Linguistic data... why so special?

There exists a mapping between structure and meaning.



VS



N-gram models

- BoW models do not represent sequential dependencies in their distributions.
- **The solution: N-gram models, a generalization of BoWs to larger context sizes**
- The 'N' stands for the dependency context size:
 - Uni-gram models consider 1-word contexts a.k.a uni-gram = BoW
 - Bi-gram models consider 2-word dependencies
 - Tri-gram models consider 3-word dependencies
 - ...

N-Gram

N=1:

This is a sentence

Uni-grams

this,
is,
a,
sentence

N=2:

This is a sentence

Bi-grams

this is,
is a,
a sentence

N=3:

This is a sentence

Tri-grams

this is a,
is a sentence

Language modeling

Language models are systems that can predict upcoming words:

- They can assign a probability to each potential next word
- They can assign a probability to a whole sentence

Language modeling

Why word prediction?

It's how **large language models (LLMs)** work!

LLMs are **trained** to predict words

- Left-to-right (autoregressive) LMs learn to predict next word

LLMs **generate** text by predicting words

- By predicting the next word over and over again

Language modeling

Goal: compute the probability of a sentence or sequence of words W :

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

Related task: probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4) \text{ or } P(w_n | w_1, w_2 \dots w_{n-1})$$

An LM computes either of these:

$$P(W) \quad \text{or} \quad P(w_n | w_1, w_2 \dots w_{n-1})$$

Joint distributions

- A distribution over a series of random variables or events
- The support of a joint distribution is the cartesian product of the individual random variables' supports.
- The probability mass function will depend on whether or not there are independence assumptions.

$$P(X_1, \dots, X_n)$$

Marginal distributions

Given some joint distribution: $P(X, Y)$

- The support of a marginal distribution $P(Y)$ is simply the set of outcomes for the variable Y
- The marginal probability of some outcome $y \in Y$ is :

$$P(y) = \sum_{x \in X} P(x, y)$$

Conditional distributions

- The support of a conditional distribution is the subset of the joint distributions support where the conditioner Y is true.
- The probability mass function will return the renormalized probabilities for the support such that they sum to one (we do so by dividing by the marginal probability of Y)

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Marginalization versus conditioning:

- In marginalization we are asking: what is the probability of the last word in our sentence being “movie”
- In conditioning we are asking: given that the last word in the sentence is “movie”, what is the probability of any of the other words in the sentence.

These are complementary operations in relation to the joint distribution

Variable independence

- **Independence** can be defined in two ways
 - Variables are independent if their joint distribution is the product of the marginal distributions

$$P(X_1, \dots, X_n) = P(X_1) \cdot \dots \cdot P(X_n)$$

- Variables are independent if their conditional distributions are equal to their marginal distributions

$$P(X | Y) = P(X)$$

The Chain Rule

- A decomposition of joint probability using conditional and marginal probabilities
- Works for both independent and dependent variables
- Is order invariant!

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

The Chain Rule

The Chain Rule applied to compute joint probability of words in sentence

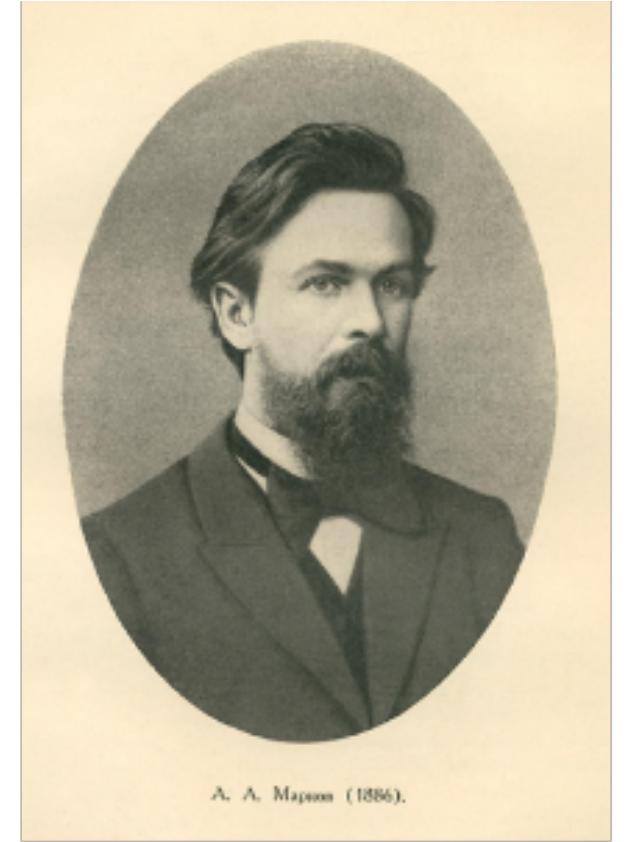
$$\begin{aligned} P(w_{1:n}) &= P(w_1)P(w_2|w_1)P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k-1}) \end{aligned}$$

$P(\text{“The water of Walden Pond”}) =$

$P(\text{The}) \times P(\text{water}|\text{The}) \times P(\text{of}|\text{The water})$

$\times P(\text{Walden}|\text{The water of}) \times P(\text{Pond}|\text{The water of Walden})$

Markov assumption



Andrei Markov

- N-grams models are also known as Markov models, because they follow the Markov assumption.
- **Markov assumption:** The probability of future events in a joint distribution is dependant on the current event only and conditionally independent from all past events. It can be decomposed as such.

Markov assumption in bi-gram model

$P(\text{blue} | \text{The water of Walden Pond is so beautifully})$

$\approx P(\text{blue} | \text{beautifully})$

$$P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-1})$$

N-gram models

- **Uni-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

- **Bi-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

- **Tri-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_{k=1}^n P(w_k | w_{k-2} w_{k-1})$$

N-gram models

P('John chased Mary') =

- **Uni-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

P('John') x P('chased') x P('Mary')

- **Bi-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

P('\PAD John') x P('John chased') x P('chased Mary') x P('Mary \PAD')

- **Tri-gram LM:**

$$P(w_1 w_2 \dots w_n) \approx \prod_{k=1}^n P(w_k | w_{k-2} w_{k-1})$$

P('\PAD \PAD John') x P('\PAD John chased') x P('John chased Mary') x P('chased Mary \PAD') x P('Mary \PAD \PAD')

Evaluations

How can we evaluate the ‘goodness’ or ‘accuracy’ of a language model?

- **Extrinsic evaluations**: Evaluating models on their external behaviour, such as their performance on downstream tasks (eg. If we use an LM to generate features for classification)
- **Intrinsic evaluations**: Evaluating the models internal behaviour, a.k.a. the goodness of its contextual probability distributions. But how?

Perplexity

- Is the common intrinsic metric used to evaluate LMs.
- Based on the intuition that LMs should mimic our grammaticality judgements, such that grammatical sentences should be more likely than ungrammatical ones... or sentences we think are more likely should also be more likely for the model.
- **Perplexity:** A measure of how likely *or surprising* a test set of sentences is under a model.

Perplexity

Perplexity = the inverse probability of the test set (one over the probability of the test set), normalized by the number of words (or tokens)

$$\begin{aligned}\text{perplexity}(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}\end{aligned}$$

Or we can use the chain rule to expand the probability of W :

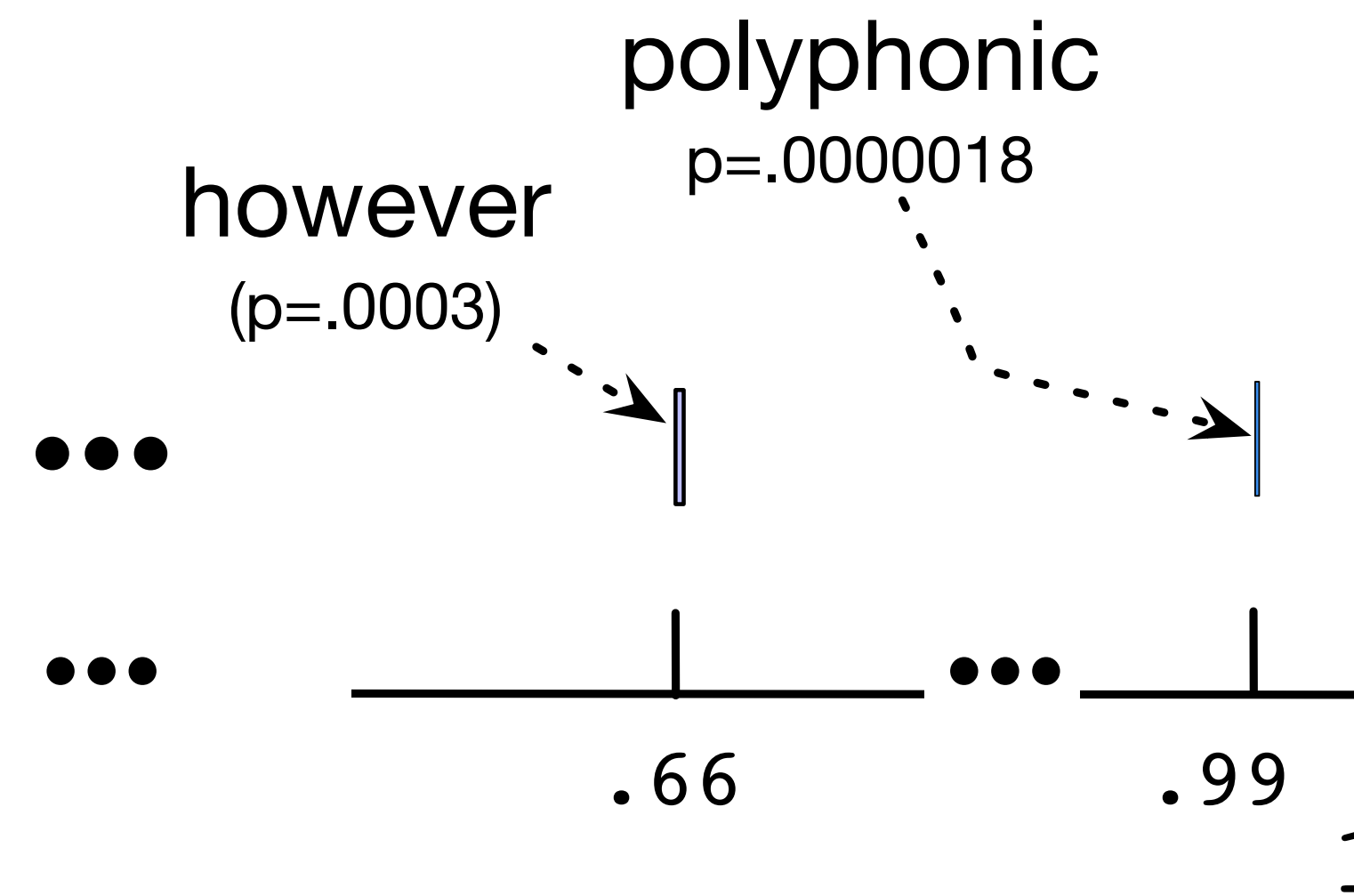
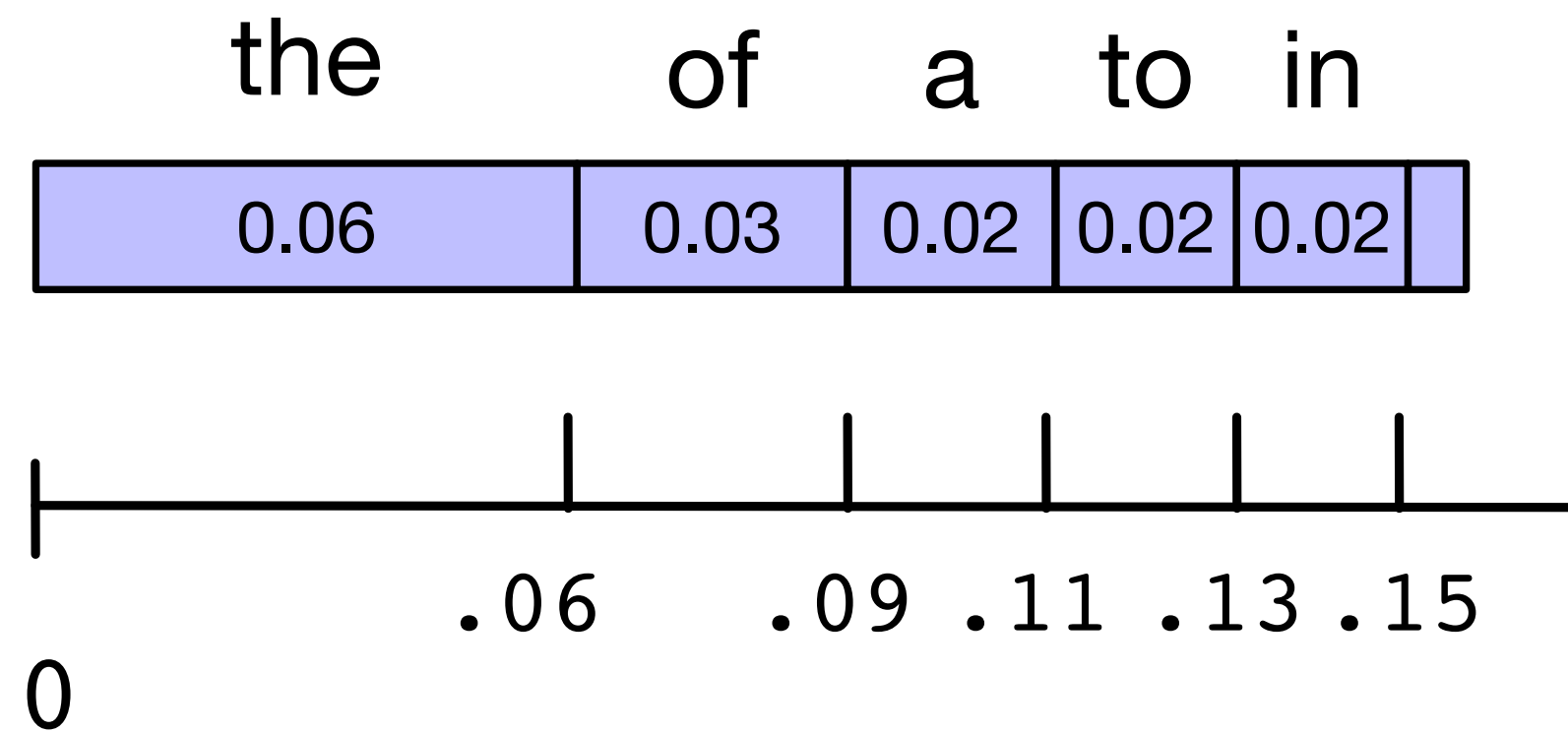
$$\text{perplexity}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

Sampling/Generating

We have seen how to **score** a sentence under our language model — i.e. return its probability — but how can we use a language model to **generate** new likely sentences under our model?

Sampling/Generating

Sampling a word from a distribution



Sampling/Generating

To sample a sentence, you recursively sample from the conditional N-gram distribution. Here's an example with a bi-gram model.

Choose a random bigram ($\langle s \rangle$, w)

according to its probability $p(w | \langle s \rangle)$

Now choose a random bigram (w , x)
according to its probability $p(x | w)$

And so on until we choose $\langle /s \rangle$

Then string the words together

$\langle s \rangle$ I
I want
want to
to eat
eat Chinese
Chinese food
food $\langle /s \rangle$
I want to eat Chinese food

Sampling/Generating

Samples from different N-gram models fit to the Wall Street Journal corpus

1
gram

Months the my and issue of year foreign new exchange's september
were recession exchange new endorsed a acquire to six executives

2
gram

Last December through the way to preserve the Hudson corporation N.
B. E. C. Taylor would seem to complete the major central planners one
point five percent of U. S. E. has already old M. X. corporation of living
on information such as more frequently fishing to keep her

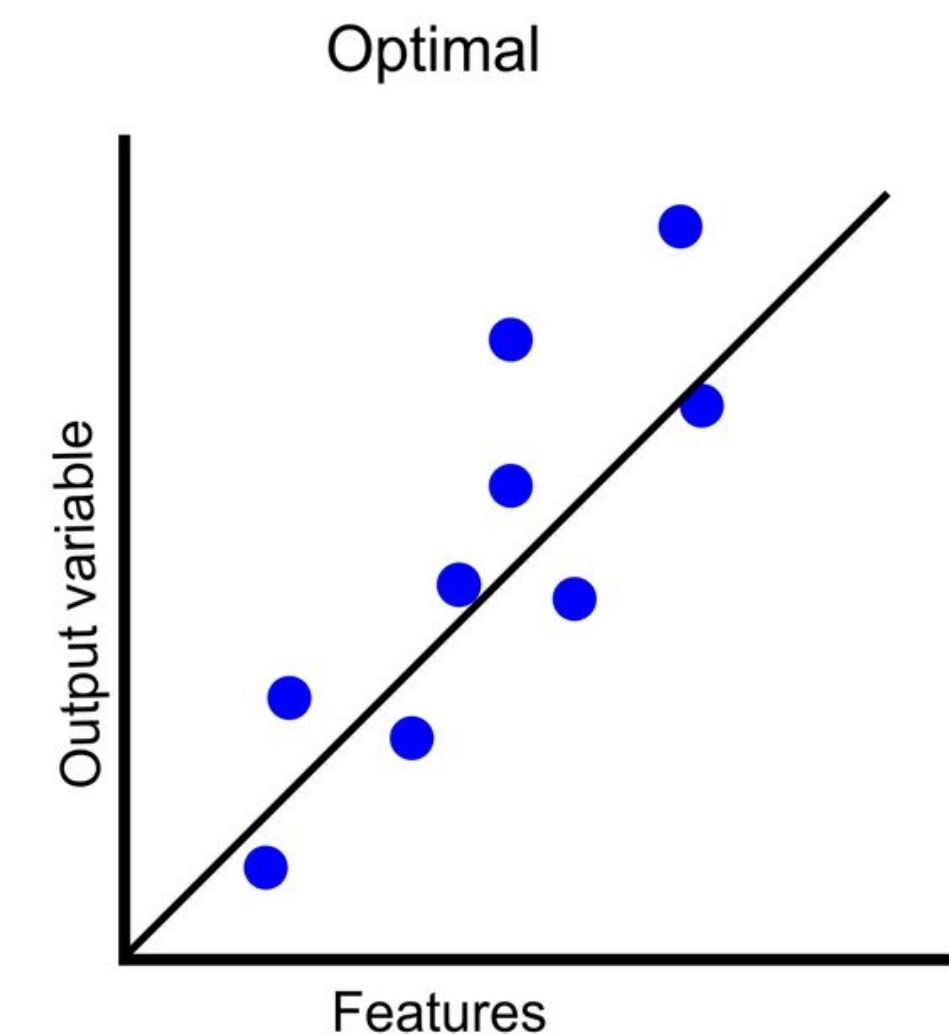
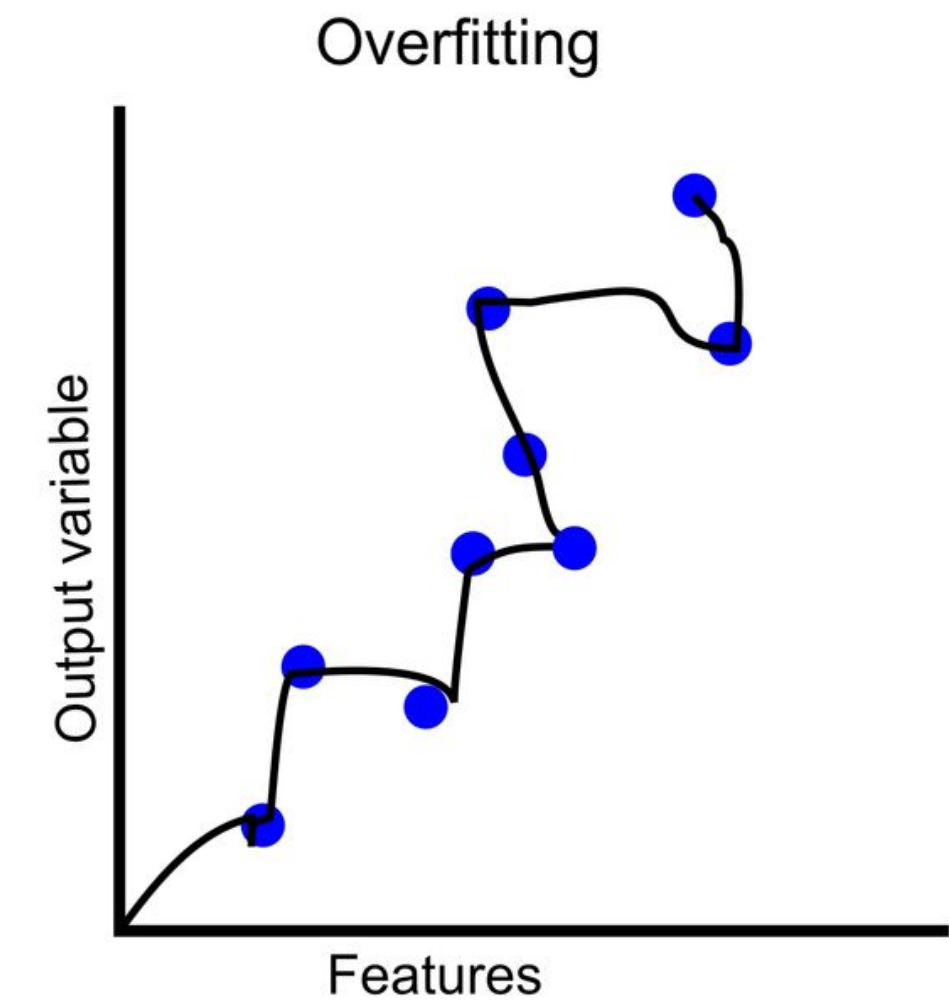
3
gram

They also point to ninety nine point six billion dollars from two hundred
four oh six three percent of the rates of interest stores as Mexico and
Brazil on market conditions

Overfitting versus generalisation

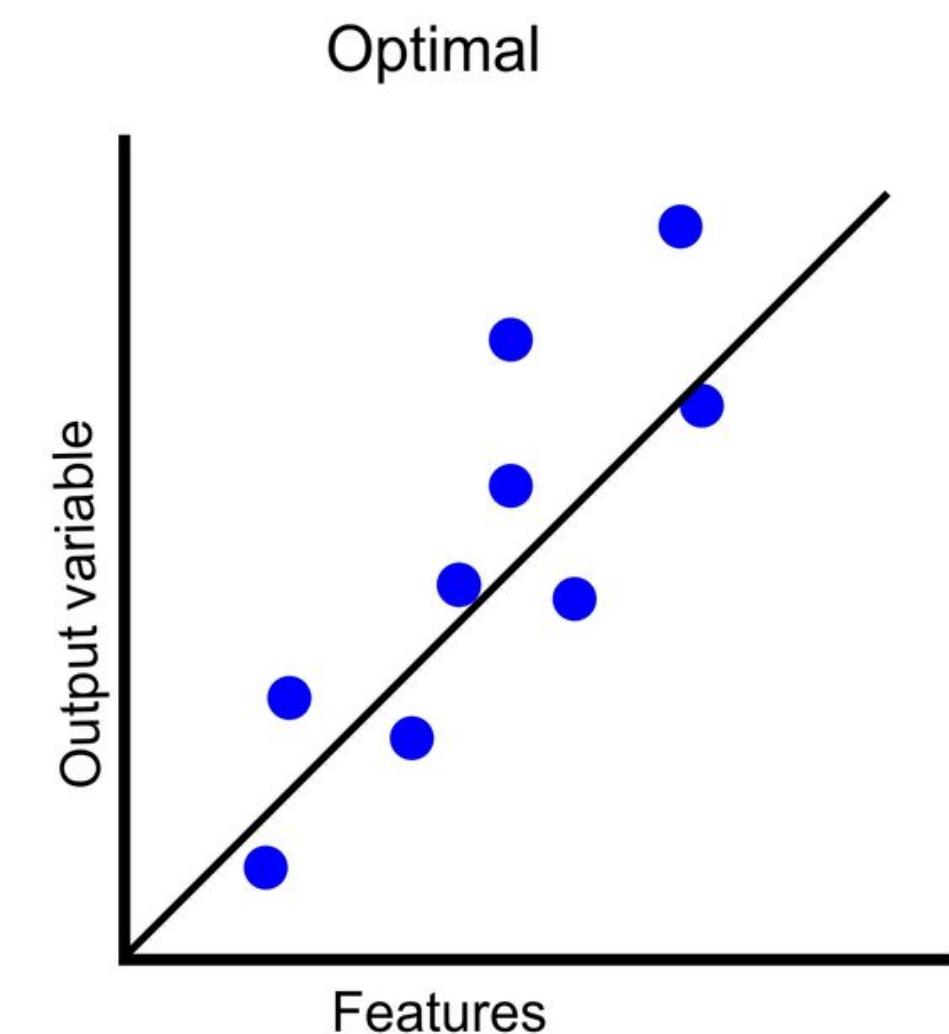
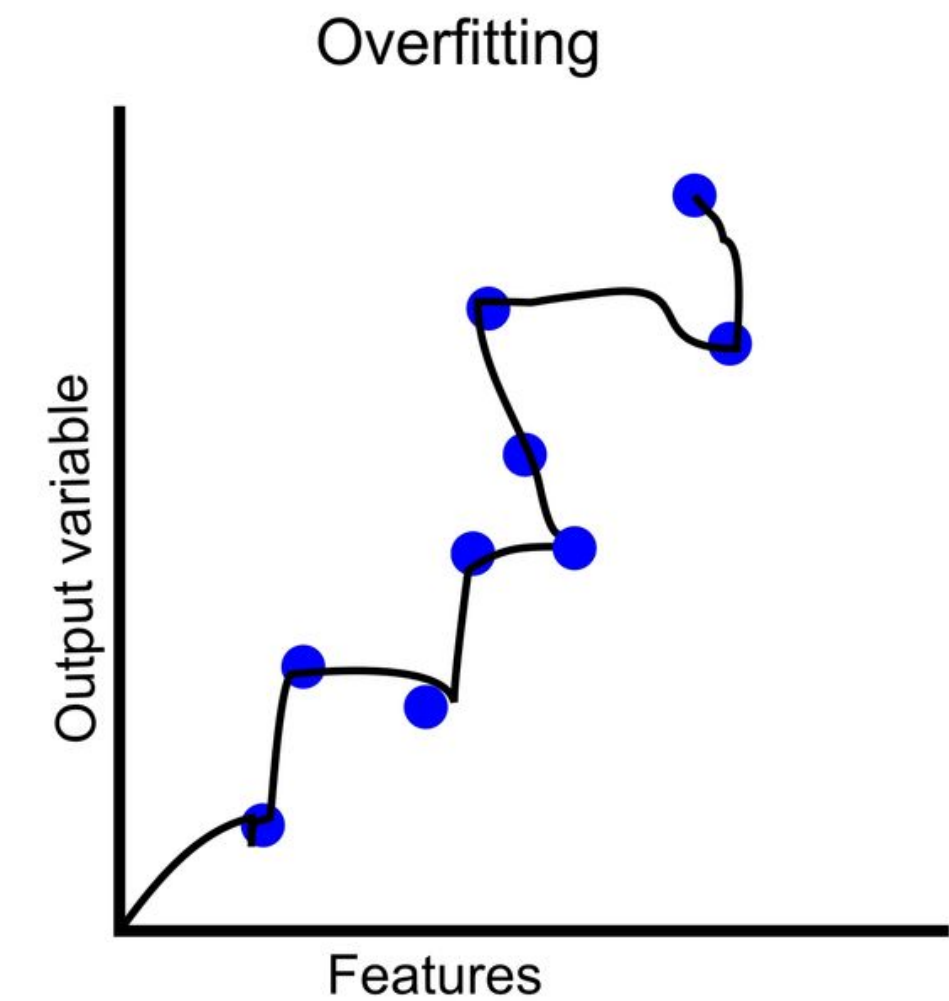
- Language models are trained/fit to a training dataset — i.e. they learn to represent the contextual distribution of the training data.
- If the true train data distribution diverges from the true test data distribution:

➡ The model can be *overfit* to the train data distribution and be unable to *generalise* to unseen test data.



Overfitting versus generalisation

- For good generalization, we want two things:
 1. The train and test data to be representative of one another;
 2. The model to be optimally trained and not overfit to the data.



[15 minute break]

Working with N-gram models!

Team up!

Open exercises/week 3 in your course folder and start writing/running code!