



Basic statistical language modeling

CS 780/880 Natural Language Processing Lecture 7

Samuel Carton, University of New Hampshire

Last lecture

Key idea: Dimension reduction

Concepts

- Dimensionality of data
- Variance of data
- Principle components
- Matrix factorization
- SVD and PCA
- Application to clustering

Toolkits

- Scikit-learn for SVD



Probability review



Random variables

A random variable X can take different values depending on chance

Notation:

- $p(X = x)$ is the probability that r.v. X takes value x
 - $p(x)$ is shorthand for the same
- $p(X)$ is the distribution over values X can take (a function)

Example: flipping a coin; $P(X = \text{heads}) = P(X = \text{tails}) = 0.5$



Discrete distributions

A discrete distribution enumerates the values a random variable can take and how likely each one is

Examples:

$p(\text{flipping a coin}) = [0.5, 0.5]$

$p(\text{rolling a die}) = [.167, .167, .167, .167, .167, .167]$

$p(\text{flipping a rigged coin}) = [0.25, 0.75]$

$p(\text{rolling a die}) = [0.1, 0.1, 0.1, 0.1, 0.1, 0.5]$

How does **entropy** relate to the values in the discrete distribution?



Joint probability and product rule

The **joint probability** of two random variables X and Y describes the total chance they take on a particular pair of values: $p(X = x, Y = y)$

If X and Y are **independent**, then $p(X = x, Y = y) = p(X = x) * p(Y = y)$

Example: two coin flips X and Y . $p(X = \text{heads}, Y = \text{heads}) = 0.5 * 0.5 = 0.25$



Conditional probability

If X and Y are **dependent**, then you have to think of the probability of X given Y:
 $p(X = x \mid Y = y)$

In this case, the joint probability of X and Y is $p(Y=y) * p(X = x \mid Y = y)$

Example: Weather is 50% sunny and 50% cloudy; I am 25% likely to run when sunny and 10% likely when rainy.

$$P(\text{run}|\text{sunny}) = \mathbf{.25}$$

$$P(\text{run, sunny}) = 0.5 * 0.25 = \mathbf{0.125}$$

$$P(\text{run}) = P(\text{run, sunny}) + P(\text{run, cloudy}) = 0.5 * 0.25 + 0.5 * 0.1 = \mathbf{0.175}$$

How does **mutual information** relate to dependence versus independence?



Chain rule

If we generalize to N joint random variables, we end up with the **chain rule**

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_2, X_1)\dots P(X_n|X_1, \dots, X_{n-1}) \\ &= P(X_1) \prod_{i=2}^n P(X_i|X_1 \dots X_{i-1}) \end{aligned}$$



Conditional probability table

With two variables X and Y , we can summarize their joint distribution with a **conditional probability table**

Each cell is $P(X=\text{column} \mid y = \text{row})$

Example:

		X	
		Run	Don't run
Y	Sunny	0.25	0.75
	Cloudy	0.1	0.9



Maximum-likelihood probabilistic modeling

Whenever we build a probabilistic model of some phenomenon, we are deciding to fit it within some probabilistic form, and then finding the **most likely** parameters of our model to fit the data.

Example:

- Model: the weather is sunny with probability p , and cloudy with probability $1-p$
- Data: 10 days; [sunny, cloudy, sunny, sunny, cloudy, sunny, sunny, cloudy, sunny, sunny]
- MLE estimate of p :



Maximum-likelihood probabilistic modeling

Whenever we build a probabilistic model of some phenomenon, we are deciding to fit it within some probabilistic form, and then finding the **most likely** parameters of our model to fit the data.

Example:

- Model: the weather is sunny with probability p , and cloudy with probability $1-p$
- Data: 10 days; [sunny, cloudy, sunny, sunny, cloudy, sunny, sunny, cloudy, sunny, sunny]
- MLE estimate of p : 0.7

How did we do that?



MLE with conditional probability

Model:

- the weather is sunny with probability p and cloudy with probability $1-p$
- I run with probability r_s when it is sunny and probability r_c when it is cloudy

Data: 10 runs

sunny	cloudy	cloudy	sunny	cloudy	sunny	cloudy	sunny	sunny	cloudy
run	no run	no run	no run	no run	no run	no run	run	no run	run

Counts:

	Run	Don't run
Sunny	2	4
Cloudy	1	3

$p =$
 $r_s =$
 $r_c =$



MLE with conditional probability

Model:

- the weather is sunny with probability p and cloudy with probability $1-p$
- I run with probability r_s when it is sunny and probability r_c when it is cloudy

Data: 10 runs

sunny	cloudy	cloudy	sunny	cloudy	sunny	cloudy	sunny	sunny	cloudy
run	no run	no run	no run	no run	no run	no run	run	no run	run

Counts:

	Run	Don't run
Sunny	2	4
Cloudy	1	3

$$p = 2+4 / (2+4+1+3) = 6/10 = .6$$

$$r_s = 2 / (2+4) = 2/6 = 0.33$$

$$r_c = 1 / (1+3) = 1/4 = 0.25$$



Probabilistic language modeling



Unigram model

Basic idea: probability of a given word w depends **only** on its overall frequency within the corpus

- So the probability of a given text is the product of the individual word probabilities

$$P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)})$$

Similar to bag-of-words in that it doesn't respect word order, but bag-of-words isn't explicitly probabilistic



Bigram model

Basic idea: the probability of word i depends **only** on word $i-1$

Example: “I am” is more likely than “I is”

$$P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$$



What can we do with a language model?

Two main things:

1. Generate new text
2. Assess the likelihood of existing text



The corpus

```
1 review_0 = "The film was a delight--I was riveted."  
2 review_1 = "It's the most delightful and riveting movie."  
3 review_2 = "It was a terrible flick, the worst I have ever seen."  
4 review_3 = "I have a feeling the film was recut poorly."  
5  
6 reviews = [review_0, review_1, review_2, review_3]
```



Preprocessing

```
1 raw_token_seqs = [word_tokenize(review.lower()) for review in reviews]
2 raw_token_seqs

[['the', 'film', 'was', 'a', 'delight', '--', 'i', 'was', 'riveted', '.'],
 ['it', "'s", 'the', 'most', 'delightful', 'and', 'riveting', 'movie', '.'],
 ['it',
  'was',
  'a',
  'terrible',
  'flick',
  ',',
  'the',
  'worst',
  'i',
  'have',
  'ever',
  'seen',
  '.'],
 ['i', 'have', 'a', 'feeling', 'the', 'film', 'was', 'recut', 'poorly', '.']]
```



Preprocessing

```
1 # For convenience, we are going to add an imaginary [START] and [END] token at the
2 # beginning and end of each sequence. You'll see why in a little bit.
3
4 token_seqs = [['[START]'] + seq + ['[END]'] for seq in raw_token_seqs]
5 token_seqs
```

```
[['[START]',
  'the',
  'film',
  'was',
  'a',
  'delight',
  '--',
  'i',
  'was',
  'riveted',
  '.',
  '[END]'],
 ...]
```



Counting

```
1 token_counts = {}
2 for token_seq in token_seqs:
3     for i in range(len(token_seq)):
4         prev_token = token_seq[i]
5         if prev_token not in token_counts:
6             token_counts[prev_token] = {}
7         if i < len(token_seq)-1: #if there is a next token, add it to the counts dict
8             next_token = token_seq[i+1]
9
10            if next_token not in token_counts[prev_token]:
11                token_counts[prev_token][next_token] = 1
12            else:
13                token_counts[prev_token][next_token] += 1
14 token_counts
```

```
{'[START]': {'the': 1, 'it': 2, 'i': 1},
'the': {'film': 2, 'most': 1, 'worst': 1},
'film': {'was': 2},
'was': {'a': 2, 'riveted': 1, 'recut': 1},
'a': {'delight': 1, 'terrible': 1, 'feeling': 1},
'delight': {'--': 1},
'--': {'i': 1},
'i': {'was': 1, 'have': 2},
'riveted': {'.': 1},
'.': {'[END]': 4},
'[END]': {},
...}
```



Counting

```
1 count_df = pd.DataFrame.from_records(data=counts, columns=vocabulary).fillna(0.0)
2 count_df.index = vocabulary
3 count_df
```

	[START]	the	film	was	a	delight	--	i	riveted	terrible	flick	,	worst	have	ever	seen	feeling	recut	poorly
[START]	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
the	0.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
film	0.0	0.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
was	0.0	0.0	0.0	0.0	2.0	0.0	0.0	0.0	1.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
a	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	...	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
delight	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
--	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
i	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0
riveted	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0



Counts to probabilities

```
1 cpt_df = count_df.div(count_df.sum(axis=1), axis=0).fillna(0)
2 cpt_df
```

	[START]	the	film	was	a	delight	--	i	riveted	terrible	flick	,	worst	have	ever	seen	feeling	recut	poorly
[START]	0.0	0.25	0.0	0.000000	0.0	0.000000	0.0	0.25	0.00	0.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.00	0.0
the	0.0	0.00	0.5	0.000000	0.0	0.000000	0.0	0.00	0.00	0.0	...	0.000000	0.0	0.0	0.25	0.000000	0.0	0.0	0.000000	0.00	0.0
film	0.0	0.00	0.0	1.000000	0.0	0.000000	0.0	0.00	0.00	0.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.00	0.0
was	0.0	0.00	0.0	0.000000	0.5	0.000000	0.0	0.00	0.25	0.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.25	0.0
a	0.0	0.00	0.0	0.000000	0.0	0.333333	0.0	0.00	0.00	0.0	...	0.333333	0.0	0.0	0.00	0.000000	0.0	0.0	0.333333	0.00	0.0
delight	0.0	0.00	0.0	0.000000	0.0	0.000000	1.0	0.00	0.00	0.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.00	0.0
--	0.0	0.00	0.0	0.000000	0.0	0.000000	0.0	1.00	0.00	0.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.00	0.0
i	0.0	0.00	0.0	0.333333	0.0	0.000000	0.0	0.00	0.00	0.0	...	0.000000	0.0	0.0	0.00	0.666667	0.0	0.0	0.000000	0.00	0.0
riveted	0.0	0.00	0.0	0.000000	0.0	0.000000	0.0	0.00	0.00	1.0	...	0.000000	0.0	0.0	0.00	0.000000	0.0	0.0	0.000000	0.00	0.0



Generating text

Once we have a MLE estimate model, we can **generate** text by just sampling from our model one word at a time

We can randomly sample, or take the most probable word at each step

We can stop after N tokens, or when we hit some stopping condition (like a [STOP] token, or a “.”)



Generating text

```
1 random.seed(123)
2 prev_token = '[START]' # By defining this special token, we gave ourselves a
3 # place to start when we want to generate text
4 count = 0
5 while count < 20 and prev_token != '[END]':
6     probs = cpt_df.loc[prev_token] #remember we want to get the row, not the column
7     next_token = choice(vocabulary, p=probs)
8     print(next_token)
9     prev_token = next_token
10    count += 1
```

```
i
have
a
delight
--
i
was
a
feeling
the
most
delightful
and
riveting
movie
.
[END]
```



Generating text

```
1 # We can also just pick the maximum-likelihood next token at every step
2 prev_token = '[START]'
3 count = 0
4 while count < 20 and prev_token != '[END]':
5     probs = cpt_df.loc[prev_token]
6     next_token = vocabulary[np.argmax(probs)]
7     print(next_token)
8     prev_token = next_token
9     count += 1
```

```
it
was
a
delight
--
i
have
a
delight
--
i
have
a
delight
--
i
have
a
delight
--
```



Assessing text probability

Given our model, we can calculate the likelihood that a given text was produced by the model.

Example: Bigram model

$$P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$$

$$p(\text{"this is a sentence ."}) = p(\text{this} | [\text{START}])p(\text{is} | \text{this})p(\text{a} | \text{is})p(\text{sentence} | \text{a})p(\text{.} | \text{sentence})$$



One token at a time

```
1 review_0 = "The film was a delight--I was riveted."  
2 review_1 = "It's the most delightful and riveting movie."  
3 review_2 = "It was a terrible flick, the worst I have ever seen."  
4 review_3 = "I have a feeling the film was recut poorly."  
5  
6 reviews = [review_0, review_1, review_2, review_3]
```

```
1 # If we know the previous token, we can calculate the probability of any given next token  
2 previous_token = previous_tokens[-1]  
3 print(f'Previous tokens: {previous_tokens}')  
4 for next_token in possible_next_tokens:  
5 | print(f'\tPossible next token: "{next_token}";\tProbability: {cpt_df.loc[previous_token, next_token]}')
```

```
Previous tokens: ['it', 'was', 'not', 'the', 'most', 'riveting']  
Possible next token: "movie"; Probability: 1.0  
Possible next token: "film"; Probability: 0.0  
Possible next token: "delight"; Probability: 0.0  
Possible next token: "."; Probability: 0.0  
Possible next token: "[END]"; Probability: 0.0
```



Over a whole sequence

```
1 # By applying the chain rule, we can calculate the probability of a whole sequence
2
3 # This is one of the sequences from the corpus
4 sequence = ['[START]', 'it', "'s", 'the', 'most', 'delightful', 'and', 'riveting', 'movie', '.', '[END]']
5
6 def sequence_likelihood(seq:list, cpt:pd.DataFrame):
7     cumulative_likelihood = 1.0
8     for i in range(len(seq)-1):
9         prev_token = seq[i]
10        next_token = seq[i+1]
11        prob = cpt.loc[prev_token, next_token]
12        cumulative_likelihood = cumulative_likelihood*prob
13    return cumulative_likelihood
14
15 print(f'Cumulative likelihood of entire sequence: {sequence_likelihood(sequence, cpt_df):.3f}')
```

Cumulative probability of entire sequence: 0.062



Over a whole sequence

Calculating cumulative log-likelihood is generally preferred to raw likelihood

- Why?

```
1 # This is one of the sequences from the corpus
2 sequence = ['[START]', 'it', "'s", 'the', 'most', 'delightful', 'and', 'riveting', 'movie', '.', '[END]']
```

```
1 def sequence_log_likelihood(seq:list, cpt:pd.DataFrame):
2     cumulative_log_likelihood = 0.0
3     for i in range(len(seq)-1):
4         prev_token = seq[i]
5         next_token = seq[i+1]
6         prob = cpt.loc[prev_token, next_token]
7         cumulative_log_likelihood = cumulative_log_likelihood + np.log(prob)
8     return cumulative_log_likelihood
9
10 print(f'Cumulative log-likelihood of entire sequence: {sequence_log_likelihood(sequence, cpt_df):.3f}')
```

```
Cumulative log-probability of entire sequence: -2.773
```



Perplexity

One way to assess the quality of a probabilistic language model is to calculate the average likelihood of the training data under that model.

This measure, called "perplexity", captures the intuition that a good language model of a corpus is one that would have been very likely to generate that corpus

```
1 normalized_sequence_likelihoods = [-sequence_log_likelihood(seq, cpt_df)/len(seq) for seq in token_seqs]
2 print(f'Length-normalized sequence negative log-likelihoods: {normalized_sequence_likelihoods}')
3 perplexity = np.mean(normalized_sequence_likelihoods)
4 print(f'Perplexity: {perplexity:.3f}')
```

Perplexity: 0.408

Can only really be compared to itself for a given corpus



Smoothing

Any model bigger than a unigram model suffers from **sparsity** issues that make certain sequences impossible, and screws with all the math

Example: “was delightful” is an impossible bigram in our toy corpus because those two words never happen to occur together, even though they very much could.

```
1 review_0 = "The film was a delight--I was riveted."  
2 review_1 = "It's the most delightful and riveting movie."  
3 review_2 = "It was a terrible flick, the worst I have ever seen."  
4 review_3 = "I have a feeling the film was recut poorly."  
5  
6 reviews = [review_0, review_1, review_2, review_3]
```

Solution: add some **smoothing** to the model which makes any bigram possible (if not likely)



Smoothing

```
1 # Big problem:
2 # We know this was an actual sequence:
3 real_sequence = ['[START]', 'it', "'s", 'the', 'most', 'delightful', 'and', 'riveting', 'movie', '.', '[END]']
4 # Our model rates this sequence as possible:
5 print(f'Cumulative likelihood of real sequence: {sequence_likelihood(real_sequence, cpt_df):.3f}')
```

Cumulative likelihood of real sequence: 0.062

```
1 # But what about this (very similar) sentence where "riveting" and "delightful" have been switched:
2 possible_sequence = ['[START]', 'it', "'s", 'the', 'most', 'riveting', 'and', 'delightful', 'movie', '.', '[END]']
3 # Our model rates this sequence as impossible!
4 print(f'Cumulative likelihood of possible sequence: {sequence_likelihood(possible_sequence, cpt_df):.3f}')
```

Cumulative likelihood of possible sequence: 0.000



Smoothing

```
1 smoothed_count_df = count_df+1  
2 smoothed_count_df
```

	[START]	the	film	was	a	delight	--	i	riveted	terrible	flick	,	worst	have	ever	seen	feeling	recut	poorly
[START]	1.0	2.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
the	1.0	1.0	3.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	1.0	1.0	1.0	2.0	1.0	1.0	1.0	1.0	1.0	1.0
film	1.0	1.0	1.0	3.0	1.0	1.0	1.0	1.0	1.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
was	1.0	1.0	1.0	1.0	3.0	1.0	1.0	1.0	2.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0
a	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	1.0	...	2.0	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0
delight	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
--	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	1.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
i	1.0	1.0	1.0	2.0	1.0	1.0	1.0	1.0	1.0	1.0	...	1.0	1.0	1.0	1.0	3.0	1.0	1.0	1.0	1.0	1.0
riveted	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.0	...	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0



Smoothing

```
1 # Now any sequence has at least a chance of being possible:  
2 print(f'Smoothed likelihood of possible sequence: {sequence_likelihood(possible_sequence, smoothed_cpt_df):.20f}')
```

Smoothed likelihood of possible sequence: 0.000000000000041044225

```
1 # And calculating the log-likelihoods isn't problematic  
2  
3 print(f'Smoothed log-likelihood of possible sequence: {sequence_log_likelihood(possible_sequence, smoothed_cpt_df):.3f}')
```

Smoothed log-likelihood of possible sequence: -28.522

```
1 # But the real sequences are still much more likely  
2 # (though much less so than in the unsmoothed model)  
3 print(f'Smoothed likelihood of real sequence: {sequence_likelihood(real_sequence, smoothed_cpt_df):.20f}')
```

Smoothed likelihood of real sequence: 0.000000000000656707607

```
1 # So the relative value is easier to compare when we look at the log-likelihoods:  
2  
3 print(f'Smoothed log-likelihood of real sequence: {sequence_log_likelihood(real_sequence, smoothed_cpt_df):.3f}')
```

Smoothed log-likelihood of real sequence: -25.749



N-gram models

$$\text{Unigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)})$$

$$\text{Bigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$$

$$\text{Trigram model: } P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)}, w^{(i-2)})$$

...and so on. But what's the problem? What stops us from conditioning on every previous token?



N-gram models

Unigram model: $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)})$

Bigram model: $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)})$

Trigram model: $P(w^{(1)} \dots w^{(N)}) = \prod_{i=1}^N P(w^{(i)} | w^{(i-1)}, w^{(i-2)})$

...and so on. But what's the problem? What stops us from conditioning on every previous token?

- data sparsity
- # of parameters



Concluding thoughts

Probabilistic modeling of text: count word occurrences and normalize to conditional probability distributions

Differing context sizes

- Unigrams: words occur independently
- Bigrams: words depend **only** on previous word
- Trigrams: words depend on previous two words
- etc.

Two important tasks

- Generate new text
- Assess text likelihood under model

Discussion question: how to do classification with these abilities?

