

Handout 7: Language Models

1. Simple generative model of text

- a. Choose a random token to start. (Why choose from tokens rather than types?)

```
1 >>> from nltk.corpus import gutenberg
2 >>> emma = gutenberg.words('austen-emma.txt')
3 >>> import random
4 >>> w = random.choice(emma)
5 >>> w
6 u'mind'
```

- b. Choose a random bigram with that word on the left

```
1 >>> w = random.choice([r for (l,r) in bigrams(emma) if l == w])
2 >>> w
3 u'like'
```

- c. Repeat

```
1 >>> w = random.choice([r for (l,r) in bigrams(emma) if l == w])
2 >>> w
3 u'to'
```

- d. Write a function called `brute_generate` to generate a random string of length n . It should return a generator.

- e. `brute_generate` is very slow! Can we do better?

2. Sampling from a frequency distribution

- a. Consider a simple case:

```
1 >>> tokens = ['a', 'a', 'a', 'a', 'b', 'b', 'c']
2 >>> ud = FreqDist(tokens)
3 >>> ud.N()
4 7
```

- b. Choose a target in `range(len(tokens))`.

```
1 >>> tgt = random.randrange(ud.N())
2 >>> tgt
3 4
4 >>> tokens[tgt]
5 'b'
```

- c. Compute that from `ud`. Is `tgt` in the “a” range?

```
1 >>> runningtotal = ud['a']
2 >>> runningtotal
3 4
4 >>> tgt < runningtotal
```

```

5      False
6      >>> runningtotal += ud['b']
7      >>> runningtotal
8      6
9      >>> tgt < runningtotal
10     True

```

d. So:

```

1      def fdchoose (fd):
2          tgt = random.randrange(fd.N())
3          rtot = 0
4          for key in fd:
5              rtot += fd[key]
6              if tgt < rtot:
7                  return key

```

3. A language model

a. Example:

```

1      >>> lm = (FreqDist(emma), ConditionalFreqDist(bigrams(emma)))

```

b. Define `generate` to use the language model. After the first word, we want to sample from the distribution of words following w . How do we do that?

c. Usage:

```

1      >>> ' '.join(generate(model, 10))
2      u'for beauty happened which was pained by a regard to her'

```

4. A language model defines a probability distribution over texts

a. Suppose we train (estimate) a language model from our running example text `abbdabdbbd`. Model M_1 :

	a	b	d
ud	0.2	0.5	0.3
a	0	1	0
b	0	0.4	0.6
d	0.667	0.333	0

b. What is the probability that we choose “b” as the first word?

c. What is the probability that we choose “b” as the next word?

d. So what is the probability of starting “b b”?

e. What is the probability of generating the text “b b d a”?

5. Comparing language models

- a. An alternative language model, M_2 :

	a	b	d
ud	0.2	0.5	0.3
a	0.2	0.5	0.3
b	0.2	0.5	0.3
d	0.2	0.5	0.3

- b. What is the probability of M_2 generating “b b d a”?
- c. Which model is more **likely**, given the observation “b b d a”?

6. Bayes’ Rule

- a. Def. of cond. prob.: $p(M|D) = p(M, D)/p(D)$
- b. Chain rule: $p(M, D) = p(M) p(D|M)$
- c. Hence:

$$\underbrace{p(M|D)}_{\text{posterior}} = \frac{\underbrace{p(M)}_{\text{prior}} \cdot \underbrace{p(D|M)}_{\text{likelihood}}}{p(D)}$$

- d. Likelihood is probability of data given model.
- e. If models have equal prior probability, highest-likelihood model is the same as highest-posterior model

7. A technical issue

- a. Computing language-model likelihoods involves multiplying small numbers together. We may run into problems of **underflow**.

```

1  >>> p = 1
2  >>> for n in [.1] * 1000: p *= n
3  >>> p
4  0.0
5  >>> q = 1
6  >>> for n in [.2] * 1000: q *= n
7  >>> q
8  0.0
9  >>> q > p
10 False

```

- b. The fix: use logarithms

$$p = 10^{\log p}$$

$$pq = 10^{\log p} 10^{\log q} = 10^{\log p + \log q}$$

$$\log pq = \log p + \log q$$

$$\log p_1 q_1 > \log p_2 q_2 \quad \text{if and only if} \quad p_1 q_1 > p_2 q_2$$

c. Example

```

1 >>> import math
2 >>> math.log10(.1 * .1 * .1)
3 -3.0
4 >>> math.log10(.1) + math.log10(.1) + math.log10(.1)
5 -3.0
6 >>> math.log10(.1) * 1000
7 -1000.0
8 >>> math.log10(.2) * 1000
9 -698.9700043360187

```

d. Log probability of “b b d a” according to M_1 and M_2

```

1 >>> lp1 = math.log10(.5) + math.log10(.4) + math.log10(.6) + math.log10(2/3)
2 >>> lp1
3 -1.0969100130080565
4 >>> lp2 = math.log10(.5) + math.log10(.5) + math.log10(.3) + math.log10(.2)
5 >>> lp2
6 -1.8239087409443187
7 >>> lp1 > lp2
8 True

```

8. Log probability as **cost**

- a. We can think of $(-\log p(y|x))$ as the cost that the model pays if the bigram (x, y) occurs in testing.

Bet:	p	.1	.2	.3	.4	.5	.6	.7	.8	.9
Pay:	$-\log p$	1	.70	.52	.40	.30	.22	.15	.10	.05

- b. For test data “b b d a”:

	b	b	d	a	Total
M_1 pays:	.30	.40	.22	.18	1.10
M_2 pays:	.30	.30	.52	.70	1.82

- c. The total cost is called the **cross entropy** of the model with the testing data

9. Overfitting

- a. Model M_1 **fits** the training data better: it extracts more information, assigns higher probabilities to the training examples
- b. It never saw **ac** in training, so it assumes that **ac** is impossible. What if that assumption is wrong?
- c. What is the cost for a bet $p(x) = 0$, if x actually occurs in testing?

- d. What are the model likelihoods and cross entropies if the **test data** is “b b a d”?

		<i>b</i>	<i>b</i>	<i>a</i>	<i>d</i>	
M_1	p	.5	.4	0	0	0
	$-\log p$.30	.40	∞	∞	∞
M_2	p	.5	.5	.2	.3	.015
	$-\log p$.30	.30	.70	.52	1.82

10. Fit and simplicity

- Model M_1 **overfits** the training data: it bets too heavily that the future will look exactly like the past.
- Model M_2 **underfits** the training data: it throws away too much information.
- There is a trade-off between **fit** and **simplicity**. Recording less information makes for a smaller model, but it also fits the training data less well.

11. Smoothing: finding a better trade-off

- We could beat M_2 on both test sets if we model bigram probabilities (like M_1), but hedge our bets
- Pretend every bigram occurs ϵ more often than it actually occurs. Smoothed counts for training **abbdabdbbd**:

	<i>a</i>	<i>b</i>	<i>d</i>	Sum
<i>a</i>	ϵ	$2 + \epsilon$	ϵ	$2 + 3\epsilon$
<i>b</i>	ϵ	$2 + \epsilon$	$3 + \epsilon$	$5 + 3\epsilon$
<i>d</i>	$2 + \epsilon$	$1 + \epsilon$	ϵ	$3 + 3\epsilon$

- General formula. Let n be the number of different types.

$$\tilde{p}(y|x) = \frac{c(x,y) + \epsilon}{c(x) + n\epsilon} \quad \tilde{p}(x) = \frac{c(x) + n\epsilon}{N + n^2\epsilon}$$

12. For our example, with $\epsilon = .5$:

- Model M_3 , probabilities and costs:

	<i>a</i>	<i>b</i>	<i>d</i>		<i>a</i>	<i>b</i>	<i>d</i>
ud	3.5/14.5	6.5/14.5	4.5/14.5	ud	.62	.35	.51
<i>a</i>	.5/3.5	2.5/3.5	.5/3.5	<i>a</i>	.85	.15	.85
<i>b</i>	.5/6.5	2.5/6.5	3.5/6.5	<i>b</i>	1.11	.41	.27
<i>d</i>	2.5/4.5	1.5/4.5	.5/4.5	<i>d</i>	.26	.48	.95

- Costs on test sets:

		M_3	M_1	M_2
<i>bbda</i>	$.35 + .41 + .27 + .26 =$	1.29	1.10	1.82
<i>bbad</i>	$.35 + .41 + 1.11 + .85 =$	2.72	∞	1.82

Classes

13. Defining LangModel as a class

```
1 class LangModel (object):
2
3     def __init__ (self, text):
4         self.unigram = FreqDist(text)
5         self.bigram = ConditionalFreqDist(bigrams(text))
```

14. Creating an instance:

```
1 >>> lm = LangModel('abbdabdbbd')
2 >>> lm
3 <LangModel instance>
4 >>> lm.unigram
5 <FreqDist>
6 >>> lm.unigram.freq('b')
7 0.5
8 >>> lm.bigram['b'].freq('d')
9 0.6
```

15. Defining a method:

```
1 class LangModel (object):
2     ...
3
4     def cost (self, test):
5         c = math.log10(self.unigram.freq(test[0]))
6         for i in range(1, len(test)):
7             c += math.log10(self.bigram[test[i-1]].freq(test[i]))
8         return c
```

16. Calling it:

```
1 >>> lm = LangModel('abbdabdbbd')
2 >>> lm.cost('bbda')
3 -1.0969100130080565
```