

CS159 - Assignment 2a: Solution

1. Unigram probabilities:

	MLE prob
$p(a)$	6/12
$p(b)$	5/12
$p(c)$	1/12

2. The bigram probabilities:

	MLE prob
$p(a a)$	1/5
$p(b a)$	3/5
$p(c a)$	1/5
$p(a b)$	2/3
$p(b b)$	1/3
$p(c b)$	0
$p(a c)$	1/1
$p(b c)$	0
$p(c c)$	0

3. Interpolated bigram probabilities with $\lambda = 1$:

	MLE prob
$p(a a)$	2/8
$p(b a)$	4/8
$p(c a)$	2/8
$p(a b)$	3/6
$p(b b)$	2/6
$p(c b)$	1/6
$p(a c)$	2/4
$p(b c)$	1/4
$p(c c)$	1/4

4. Adding in <UNK>

<UNK> a <UNK> b
a <UNK> a b
b a b a

5. Updated unigram and bigram probabilities

	MLE prob
$p(a)$	5/12
$p(b)$	4/12
$p(\langle\text{UNK}\rangle)$	3/12

The bigram probabilities:

	MLE prob
$p(a a)$	0
$p(b a)$	2/4
$p(\langle\text{UNK}\rangle a)$	2/4
$p(a b)$	2/2
$p(b b)$	0
$p(\langle\text{UNK}\rangle b)$	0
$p(a \langle\text{UNK}\rangle)$	2/3
$p(b \langle\text{UNK}\rangle)$	1/3
$p(\langle\text{UNK}\rangle \langle\text{UNK}\rangle)$	0

6. Backoff model with $D = 0.5$:

- a

$$\text{reserved_mass}(a) = (2 * 0.5)/4 = 1/4$$

Since there is only one unseen bigram, it will get all of the reserved mass (i.e. 1/4).

If you wanted to do the entire calculation (e.g. for practice):

$$\text{denominator} = 1 - \sum_{x:\text{count}(ax)>0} p(x) = 1 - (p(b) + p(\langle\text{UNK}\rangle)) = 1 - (4/12 + 3/12) = 5/12$$

$$\alpha(a) = \frac{1/4}{5/12} = 12/20 = 3/5$$

- b

$$\text{reserved_mass}(b) = (1 * 0.5)/2 = 1/4$$

$$\text{denominator} = 1 - \sum_{x:\text{count}(bx)>0} p(x) = 1 - p(a) = 1 - 5/12 = 7/12$$

(Note this is the same as $p(b) + p(\langle\text{UNK}\rangle)$, but for efficiency, since the number of words that do occur following a particular word is usually much, much smaller than the number that don't, programmatically we calculate it using the 1 minus formulation).

$$\alpha(b) = \frac{1/4}{7/12} = 12/28 = 3/7$$

- $\langle\text{UNK}\rangle$

$$\text{reserved_mass}(\langle\text{UNK}\rangle) = (2 * 0.5)/3 = 1/3$$

Since there is only one unseen bigram, it will get all of the reserved mass.

If you wanted to do the entire calculation (e.g. for practice):

$$\text{denominator} = 1 - \sum_{x:\text{count}(\langle\text{UNK}\rangle x)>0} p(x) = 1 - (p(a) + p(b)) = 1 - (5/12 + 4/12) = 3/12 = 1/4$$

$$\alpha(\langle\text{UNK}\rangle) = \frac{1/3}{1/4} = 4/3$$

Finally, now that we have the α s, we can calculate the smoothed bigram probabilities. For those that occurred, we simply discount the count. For those that did not occur, we calculate the probability as *alpha* times the unigram probability of the word.

	eqn	prob
$p(a a)$	$3/5 * 5/12$	$2/8$
$p(b a)$	$(2-0.5)/4$	$3/8$
$p(\langle \text{UNK} \rangle a)$	$(2-0.5)/4$	$3/8$
$p(a b)$	$(2-0.5)/2$	$3/4$
$p(b b)$	$3/7 * 4/12$	$4/28$
$p(\langle \text{UNK} \rangle b)$	$3/7 * 3/12$	$3/28$
$p(a \langle \text{UNK} \rangle)$	$(2-0.5)/3$	$3/6$
$p(b \langle \text{UNK} \rangle)$	$(1-0.5)/3$	$1/6$
$p(\langle \text{UNK} \rangle \langle \text{UNK} \rangle)$	$4/3 * 3/12$	$2/6$