

# NLP - Fall 2018 - Sample Midterm Exam

This material is ©Anoop Sarkar 2018.

Only students registered for this course are allowed to download this material.

Use of this material for “tutoring” is prohibited.

- (1) TrueCasing is the process of taking text with missing or unreliable case information and producing the proper case for each word, e.g. if the input looks like:

as previously reported , target letters were issued last month to  
michael milken , drexel 's chief of junk-bond operations ; mr. milken 's  
brother lowell ; cary maultasch , a drexel trader ; james dahl , a  
drexel bond salesman ; and bruce newberg , a former drexel trader .

Then the output of the TrueCasing program should be:

As previously reported , target letters were issued last month to  
Michael Milken , Drexel 's chief of junk-bond operations ; Mr. Milken 's  
brother Lowell ; Cary Maultasch , a Drexel trader ; James Dahl , a  
Drexel bond salesman ; and Bruce Newberg , a former Drexel trader .

Assume **we can only use** the following two probability distributions:

- A *translation probability*  $P(w | W)$  where  $w$  is the lowercase variant of the TrueCase word  $W$  (note that the TrueCase word might still be lowercase). The function `lower` can be used to lowercase a word, e.g. `“HAL900I”.lower() = “hal900I”`
- A *bigram probability*  $P(W | W')$ . A language model  $P(W_1, \dots, W_n)$  is used to provide the probability of a sentence. A bigram language model approximates the probability of a sentence as follows:

$$\Pr(W_1, \dots, W_n) \approx \prod_{i=1}^n P(W_i | W_{i-1})$$

Assume that  $W_{-1} = w_{-1} = \text{none}$  is a dummy word that begins each sentence.

- Assume that  $c(\cdot)$  gives the frequency of unigrams, bigrams, etc.
- a. Complete the following formula to provide a model of the TrueCasing task by using only the translation probability  $P(w | W)$  and the bigram probability  $P(W | W')$ :
- $$\begin{aligned} W_1^*, \dots, W_n^* &= \arg \max_{W_1, \dots, W_n} \Pr(W_1, \dots, W_n | w_1, \dots, w_n) \\ &= \text{provide this formula} \end{aligned}$$
- b. Using maximum likelihood, provide a formula to estimate the the translation probability parameters  $P(w | W)$  for lowercase words  $w$  and TrueCase words  $W$ . Assume you **only** have access to a sufficient amount of TrueCase text.
- c. Provide the equation that correctly computes add one smoothing for  $P(w | W)$ .
- d. Backoff smoothing for  $P(W_i | W_{i-1})$  is defined as follows:

$$P_{bo}(W_i | W_{i-1}) = \begin{cases} \frac{c^*(W_{i-1}, W_i)}{c(W_{i-1})} & \text{if } c(W_{i-1}, W_i) > 0 \\ \alpha(W_{i-1})P_{bo}(W_i) & \text{otherwise} \end{cases}$$

where  $c^*(W_{i-1}, W_i) = c(W_{i-1}, W_i) - D$  for some  $0 < D < 1$  and  $\alpha(W_{i-1})$  is chosen to make sure that  $P_{bo}(W_i | W_{i-1})$  is a proper probability. Provide the equation to compute  $\alpha(W_{i-1})$ . Assume that  $\sum_{W_i} P_{bo}(W_i) = 1$ .

(2) **Language Models**

For the CFG  $G$  given below:

$$S \rightarrow A \mid c$$

$$A \rightarrow B a$$

$$B \rightarrow b S$$

- a. What is the language  $L(G)$ ?
- b. Assign probabilities to each rule in the CFG above so that for each string  $w \in L(G)$ :

$$P(w) = \exp\left(\frac{|w| - 1}{2} \times \ln(0.3) + \ln(0.7)\right)$$

where,  $|w|$  is the length of string  $w$ ,  $\exp$  is exponentiation, and  $\ln$  is  $\log$  base  $e$ . Using an example, briefly explain *why* your PCFG provides the desired  $P(w)$  for any  $w$ .