

What is Statistical Language Model?

- Statistical language models are probability $P(w_1, w_2, \dots, w_n)$ defined on sequences of words w_1, w_2, \dots, w_n
- Many NLP applications
 - speech recognition, machine translation and information retrieval
- Estimating the probability becomes expensive and infeasible for long sequences (not enough data/storage)
- Smoothed N-gram models based on unigrams, bigrams, and trigrams are often used to approximate the language model:

$$P(w_1, w_2, \dots, w_n) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_2) \dots \times P(w_n|w_{n-1})$$

How to derive a language model

1. Get a corpus of texts
2. Break texts into sentences, sentences into tokens
3. Count the frequency of tokens, pairs, triples
4. Apply probability theory and information theory to estimate model parameters (prob. of tokens, pairs, triples)
 - MLE (Maximal Likelihood Estimation) = based on frequency, maximizing the prob. of sample, but gives zero prob. to unseen events
 - Smoothing = Dealing with low-count and unseen events (giving them non-zero probability)

Example – Case Restoration

- Input:
 - THE U.S. MILITARY ISSUED A PUBLIC APOLOGY TO THE PEOPLE OF A SHIITE MUSLIM NEIGHBORHOOD IN BAGHDAD ON THURSDAY ...
- Output:
 - The U.S. military issued a public apology to the people of a Shiite Muslim neighborhood in Baghdad on Thursday ...
- Many Likely Candidates for Output
 - the u.s. military issued a public apology to the people of a shiite muslim neighborhood in baghdad on thursday ...
 - The u.s. military issued a public apology to the people of a shiite muslim neighborhood in baghdad on thursday ...
- Which is more likely = Max Probability
 - P(the u.s. military issued a public apology to the people of a shiite muslim neighborhood in maghdad on thursday ...)
 - P(The u.s. military issued a public apology to the people of a shiite muslim neighborhood in maghdad on thursday ...)

Example – Filling in Articles and Prepositions

- In Alan Meryers (2005) Gateways to Academic Writing pp. 277, learners are asked to fill articles and prepositions in the blanks.
- The model T Ford was a fragile-looking automobile, but it became the most popular car in history. Henry Ford sold 16 million Model Ts _____ the years 1908 and 1928. The Model T is _____ immediate best-seller, not only because of its low price, but because it was _____ powerful car.

Example – using articles and prepositions correctly

- In Gateways to Academic Writing (Meryers 2005, pp. 277), learners are asked to fill articles and prepositions in the blanks.
- The model T Ford was a fragile-looking automobile, but it became the most popular car in history. Henry Ford sold 16 million Model Ts in the years 1908 and 1928. The Model T is an immediate best-seller, not only because of its low price, but because it was the powerful car.

How to Estimate Probability

- Counting words – $(w, \text{count}(w))$ (done in Lab 1)
- MLE (maximal likelihood estimator)
 - For all words w with count $i \rightarrow P(w) = i / N$
 - For all unseen words $\rightarrow P(w) = 0$

Smoothing

- Good-Turing Estimation
 - V = Vocabulary Size
1. Counting word types with counts 1, 2, 3, ...
 - $N = N_1 + 2 N_2 + 3 N_3 + \dots$
 - $N_0 = V - N_1 - N_2 - N_3 - \dots$ (for 1gram)
 - $N_0 = V^2 - N_1 - N_2 - N_3$ (N_1, N_2, N_3 are different for 2gram)
 2. Compute $0^* = 1 \times N_1 / N_0$
$$r^* = (r+1) N_{r+1} / N_r, \quad \text{for } r = 1, k-1$$
$$r^* = r, \quad \text{for } r \geq k$$
 3. Adjusting counts and prob. with r^*
 - For unseen word w ($r = 0$)
$$r^* = 0^* = N_1 / N_0$$
$$P(w) = r^* / N = N_1 / (N_0 N)$$
 - For seen word w with count $r < k$ ($k = 10$)
$$r^* = (r+1) N_{r+1} / N_r$$
$$P(w) = r^* / N = (r+1) N_{r+1} / (N_r N)$$

Re-normalization

- All probabilistic values must sum to 1
- Adjusted total
 - $N' = N_1 + 2 N_2 + 3 N_3 + \dots + (k-1) N_{k-1} +$
 $k N_k + k N_k +$
 $(k+1) N_{k+1} + \dots$
- Prob sum to $N' / N = (N + k N_k) / N$ (not 1.0)
- Normalization
 - $P'(w) = N P(w) / (N + k N_k)$
- Do the same for bigram (w, w')
- Or directly
 - $P(w) = r^* / (N + k N_k) = N_{r+1} / N_r / (N + k N_k)$ for $r = 0, 1, \dots, k-1$
 - $P(w) = r / (N + k N_k) = r / (N + k N_k)$ for $r \geq k$

Estimating $P(w)$ and $P(w'|w)$

- $P(w) = r^*/N$
 - $r = \text{count}(w)$, go from r to r^* (for unigram)
- $P(w, w') = r^*/N$
 - $r = \text{count}(w, w')$, go from r to r^* (for bigram)
- $P(w'|w) = P(w', w) / P(w)$

Language Modeling Tools

- SRI Language Modeling Toolkit
 - <http://www.speech.sri.com/projects/srilm/>
- Tutorial
 - Introduction to SRILM by Berlin Chen
 - http://berlin.csie.ntnu.edu.tw/Courses/2005F-SpeechRecognition/Lectures2005F/SP2005F_Lecture08_SRILM%20Tutorial.pdf
- Python Wrapper for SRILM
 - <http://www.isi.edu/~chiang/software/psrilm.tgz>

Example of Running SRILM

- INPUT (科技想要什麼？)
 - 寫給 台灣 的 讀者 凱文 · 凱利
 - 這 本 書 ， 可 以 說 是 四 十 年 前 在 台 灣 扎 下 了 根 基
 - 一 九 七 二 年 ， 二 十 歲 的 我 離 開 美 國 紐 澤 西 ， 結 束 平 淡 的 生 活 ， 展 開 長 長 的 旅 程 ， 來 到 一 個 完 全 不 同 的 世 界 ， 也 就 是 當 時 正 開 始 面 臨 轉 變 的 台 灣
 - 我 第 一 次 出 國 ， 目 的 地 就 是 台 灣
 - 眼 前 所 見 ， 令 我 震 驚 不 能 自 己
 - 那 個 時 候 ， 最 主 要 的 交 通 工 具 是 腳 踏 車 ， 現 在 的 新 北 市 則 是 一 片 片 稻 田
 - 我 親 眼 見 證 ， 這 塊 活 躍 的 土 地 急 速 現 代 化 ， 幾 乎 每 天 都 在 改 變 ， 從 第 三 世 界 國 家 提 升 到 名 列 世 界 富 國
 - 台 灣 ， 達 成 了 不 可 能 的 任 務

Example of Running SRILM

- Language Model

- \data\

- ngram 1=11589

- ngram 2=71201

- ngram 3=8893

- \1-grams:

- -5.149576 一共 -0.1150313

- -4.149576 一再 -0.2703053

- -3.516108 一切 -0.3589188

- -3.631062 一千 -0.5536026

~ $\log_{10}(24/11589)$

(back off weight = 『除非 一千』 查不到時的權重)

Example of Running SRILM

- \2-grams:
- -0.6187539 一 共 有
- -0.4609967 一 再 出 現 0.00103879
- -1.618754 一 再 提 出
- ...