

Language Models: Auto-Complete

Ex. Typing I like a cat (where the next word cat is recommended because it has the highest probability (in our model) given by a sequence "I like a")

I. Load and Preprocess Data

- (1) Doc (e.g. en-US.twitter.txt) → split ("\\n") → sentences
 - (2) Each sentence → tokenize → list of tokens (= words)
 - (3) Split to train / test data
 - (4) GetTokens in train data that has \geq threshold times
otherwise set them to <UNK> (unknown)
- Output: train data → Ex: [['sky', 'is', 'blue', '.'],
test data ['love', 'are', 'green', '.'],
['unk', 'are', 'unk'], ...]

II. N-gram based Language Models

- Assumption: The probability of next word depends on only the previous N-grams

Ex. I like a cat with N-grams = 2

↳ (I like), (like a), (a, cat)

- With a list of words: $w_1 w_2 \dots w_t, \dots, w_M$

With N-grams $P(w_t | w_{t-1} \dots w_{t-N}) = \frac{C(w_{t-1} \dots w_{t-N}, w_t)}{C(w_{t-1} \dots w_{t-N})}$ (Eq. 1)

(...): number of occurrence of the given sequence.

- Add start and end token

Ex. with N-grams = 2

$\langle \rangle \langle \rangle$ I like a cat $\langle \rangle \langle \rangle$

equal to N-grams

only one token with any N-grams

- At (Eq. 1) If count of N-grams = 0 (mean this N-grams does not exist in training data)

$$\Rightarrow (Eq. 1) = \frac{0}{0} ?$$

↳ K-smoothing: $\hat{P}(w_t | w_{t-1} \dots w_{t-N}) = \frac{C(w_{t-1} \dots w_{t-N}, w_t) + k}{C(w_{t-1} \dots w_{t-N}) + k \cdot V}$ (Eq. 2)

where V is size of training vocab
(a set of unique words in training vocab)

$$k = 1 \text{ or } 2, M \text{ or } 3, \dots$$

III. Perplexity

$$PP(W) = \sqrt[M]{\prod_{t=N+1}^M P(w_t | w_{t-1} \dots w_{t-N})} \quad (Eq. 3)$$

where N is N-grams

M: length of the sentence.

t: is number starting from 1 (NOT 0)

Interpret: The higher Probability → the lower PP(W) (it is better because it is more natural)

IV. Build an auto-complete system based on N-grams

(1) Build n-grams and n-plus1-grams dictionary

Ex. n-grams = { ('i', 'like'): 2000, n-plus1-grams = { ('i', 'like', 'a'): 1800,
('like', 'a'): 99, ('a', 'cat'): 888,
<'>', '<>': 10, }
<'>', 'i': 100, }
<'cat', '</>': 666 }

(2) Calculate probability of a word given by prior N-words (N-grams)

Using (Eq. 2)

def estimate_probability (word, previous_n-grams, n-grams, n-plus1-grams, vocabulary_size, k=1)
↓ ↓ ↓ ↓ ↓ ↓
Ex: ('cat', 'I like' bi-gram, tri-gram) V, K = 1
↓ ↓ ↓ ↓ ↓ ↓
return probability (a value. Ex. 0.227)

(3) Calculate probability of all words in vocab given by prior N-words (N-grams)

def estimate_probabilities (previous_n-grams, n-grams, n-plus1-grams, vocabulary_size, k)
↓ ↓ ↓ ↓ ↓ ↓
Ex: 'I like' bi-gram, tri-gram, vocabulary, k
just using (2) for every word in vocabulary
↓ ↓ ↓ ↓ ↓ ↓
return probabilities (a dict). Ex: probabilities = { 'cat': 0.272, A set of unique words in training data.
'a': 0.09,
'like': 0.09,
'dog': 0.09, }
↓ ↓ ↓ ↓ ↓ ↓

(4) Calculate PP(W) using (Eq. 3)

(5) Suggest a word

Similar to (3) but may add some conditions to choose word from probabilities dict.
- Ex: start with character(s)

(6) Multiple suggestions

Similar to (5) but run for a series of n-grams and for each pair of (n-grams and n-plus1-grams) suggest a word.