

1.1 N-Grams

$$P(w_1, w_2, w_3, \dots, w_n) = P(w_1) P(w_2 | w_1) P(w_3 | w_{1:2}) \dots P(w_n | w_{1:n-1})$$
$$= \prod_{k=1}^n P(w_k | w_{1:k-1})$$

\Downarrow "Markov Assumption"

$$\approx \prod_{k=1}^n P(w_k | w_{k-1}) \quad (\text{e.g. for bigram})$$

* MLE (Maximum Likelihood Estimation)

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{\sum_w C(w_{n-1} w)} = \frac{C(w_{n-1} w_n)}{C(w_{n-1})} \dots (\text{bigram})$$

general
 \Rightarrow

$$\underbrace{P(w_n | w_{n-N+1:n-1})}_{\text{"relative frequency"}} = \frac{C(w_{n-1})}{C(w_{n-N+1:n-1})} \dots \text{N-Gram}$$

* Smoothing

\hookrightarrow to avoid assigning zero probabilities to unseen events

1) Laplace Smoothing

: add 1 to every count $\Rightarrow P(w_n | w_{n-N+1:n-1}) = \frac{C(w_{n-N+1:n-1} w_n) + 1}{\sum_{k=1}^V (C(w_{n-N+1:n-1} w_k) + 1)}$

$$= \frac{C(w_{n-N+1:n-1} w_n) + 1}{C(w_{n-N+1:n-1}) + V}$$

where V is,
the number of tokens

2) Add-k Smoothing

: add k to every count (Similar to Laplace Smoothing)

\Rightarrow Sharp change in probabilities may occur because too much probability mass may move to zeros

3) Backoff

1.2 Naive Bayes Classifiers

* Generative vs Discriminative Classifiers

- Generative: Build a model of how a class could generate some input data (e.g. Naive Bayes)
- Discriminative: Learn useful features from the input to discriminate between classes (e.g. Logistic Regression)

The class of a document (d): $\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c|d) = \underset{c \in C}{\operatorname{argmax}} \frac{P(d|c)P(c)}{P(d)}$ (Bayes' rule)

$$= \underset{c \in C}{\operatorname{argmax}} P(d|c)P(c)$$

W.L.O.G. for a document of features f_1, f_2, \dots, f_n , $\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(f_1, f_2, \dots, f_n | c) P(c)$

* Assumptions

- bag-of-words assumption: position doesn't matter
- naive Bayes assumption: feature independence $P(f_1, f_2, \dots, f_n | c) = P(f_1 | c) \dots P(f_n | c)$

$$\Rightarrow C_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{i \in \text{positions}} P(w_i | c)$$

} Linear Function of input features

$$C_{NB} = \underset{c \in C}{\operatorname{argmax}} \log P(c) + \sum_{i \in \text{positions}} \log P(w_i | c)$$

(from training data)

$$\hat{P}(c) = \frac{N_c}{N_{\text{doc}}} \quad \hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)}$$

* How to improve Naive Bayes Classifiers?

1) Delete Duplicates \leadsto binary Naive Bayes

: whether a word occurs or not may be more important than its frequency

2) Negation

e.g. $\left\{ \begin{array}{l} \text{didn't like this movie, but I} \\ \quad \downarrow \\ \text{didn't NOT-like NOT-this NOT-movie, but I} \end{array} \right.$

3) Use existing sentiment lexicons

1.3 Logistic Regression

* Generative vs Discriminative model

\Rightarrow generative는 $\hat{c} = \operatorname{argmax}_c P(d|c)P(c)$ 로, $P(d|c)$ 를 구하려 하지만,
discriminative는 $P(c|d)$ 를 바로 구하려 한다.

* When to use Logistic Regression : Correlation

\Rightarrow Naive Bayes는 independence assumption 때문에 correlation이 강한 정보에 적합
 \leadsto larger dataset에는 logistic regression이 더 나음.

* Sigmoid Classifier

$z = \left(\sum_{i=1}^n w_i x_i \right) + b$ 를 0~1 range로 mapping

$$\Rightarrow \sigma(z) = \frac{1}{1 + \exp(-z)} \approx \begin{cases} p(y=1) = \sigma(w \cdot x + b) \\ p(y=0) = 1 - \sigma(w \cdot x + b) \end{cases}$$

* Softmax Regression (Multinomial Logistic Regression)

\Rightarrow multiclass 일 때는 output y 가 one-hot vector

$$\operatorname{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)} \approx p(y_k=1|x) = \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \quad (K \text{ is the number of classes})$$

$$\Rightarrow \hat{y} = \operatorname{softmax}(Wx + b)$$

* Learning in Logistic Regression

1) Loss Function (Cost Function)

: a metric for how close the current label (\hat{y}) is to the true gold label (y). (e.g. cross-entropy)

2) Gradient Descent

: optimization algorithm for updating the weights

1) Cross-Entropy Loss (binary classifier)

$L(\hat{y}, y)$ = How much \hat{y} differs from y , where $\hat{y} = \sigma(w \cdot x + b)$

$$\approx p(y|x) = \hat{y}^y (1-\hat{y})^{1-y} \quad (\because \text{Bernoulli Distribution})$$

따라서, $\log p(y|x) = y \log \hat{y} + (1-y) \log (1-\hat{y})$ should be maximized

$$\Leftrightarrow \underline{L_{CE} = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]} \text{ should be Minimized}$$

$$\Rightarrow L_{CE}(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1-y) \log (1 - \sigma(w \cdot x + b))]$$

2) Gradient Descent

↳ Goal: Minimize the loss (averaged over all examples)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}; \theta), y^{(i)})$$

* (Logistic Regression에서는 loss function이 convex라서 무조건 minimum이 보장되는 반면,
multi-layer neural network에서는 local minima에 stuck될 수 있다.)

$$\Rightarrow W^{t+1} = W^t - \eta \frac{d}{dW} L(f(x; W), y) \quad (\eta \text{ is learning rate})$$

$$\Downarrow$$

$$\nabla L(f(x; \theta), y) = \begin{bmatrix} \frac{\partial}{\partial w_1} (f(x; \theta), y) \\ \frac{\partial}{\partial w_2} (f(x; \theta), y) \\ \vdots \\ \frac{\partial}{\partial b} (f(x; \theta), y) \end{bmatrix}$$

$$\frac{\partial L_{CE}(\hat{y}, y)}{\partial w_j} = (\hat{y} - y) x_j$$

derivation

$$\left\{ \begin{aligned} \frac{\partial L_{CE}}{\partial w_j} &= \frac{\partial}{\partial w_j} (-y \log \sigma(wx+b) + (1-y) \log (1-\sigma(wx+b))) \\ &= -\frac{y}{\sigma(wx+b)} \cdot \frac{\partial}{\partial w_j} \sigma(wx+b) + \frac{1-y}{1-\sigma(wx+b)} \cdot \frac{\partial}{\partial w_j} \sigma(wx+b) \\ &= \frac{\partial}{\partial w_j} \sigma(wx+b) \left[\frac{\sigma(wx+b) - y}{\sigma(wx+b)(1-\sigma(wx+b))} \right] \\ &= \cancel{\sigma(wx+b)} (1-\cancel{\sigma(wx+b)}) \cdot \frac{\partial}{\partial w_j} (wx+b) \cdot \left[\frac{\sigma(wx+b) - y}{\cancel{\sigma(wx+b)}(1-\cancel{\sigma(wx+b)})} \right] \\ &= (\sigma(wx+b) - y) x_j \\ &= (\hat{y} - y) x_j \end{aligned} \right.$$

* Mini-batch Training

• Stochastic Gradient Descent: 1 example at a time

↑

• Mini-batch Training: train on a group of m examples
↳ Cost function is the average loss over the batch

$$\text{Cost}(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m L_{CE}(\hat{y}^{(i)}, y^{(i)})$$

$$\Rightarrow \frac{\partial \text{Cost}(\hat{y}, y)}{\partial W} = \frac{1}{m} (\hat{y} - y)^T X = \frac{1}{m} (\sigma(WX+b) - y)^T X$$

* Regularization

→ to avoid overfitting, add a regularization term, $R(\theta)$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) - \alpha R(\theta) \quad \rightarrow \text{to penalize large weights}$$

L2 Regularization

$$R(\theta) = \|\theta\|_2^2 = \sum_{i=1}^n \theta_i^2 \quad (\text{Euclidian distance})$$

Ridge Regression

- prefer weight vectors with many small weights
- corresponds to assuming that weights are distributed in a Gaussian distribution with $\mu=0$

L1 Regularization

$$R(\theta) = \|\theta\|_1 = \sum_{i=1}^n |\theta_i| \quad (\text{Manhattan Distance})$$

Lasso Regression

- prefer sparse solutions with larger weights, but many more zero weights
- viewed as a Laplace prior on the weights

3) Multinomial Logistic Regression

→ output: vector \hat{y} with K elements, where each element $\hat{y}_k = p(y_k=1|x)$

$$\begin{aligned}\Rightarrow L_{CE}(\hat{y}, y) &= - \sum_{k=1}^K y_k \log \hat{y}_k \\ &= - \log \hat{y}_c \quad (\text{where } c \text{ is the correct class}) \\ &= - \log \hat{p}(y_c=1|x) \\ &= - \log \frac{\exp(w_c \cdot x + b_c)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \quad (\text{where } K \text{ is the number of classes})\end{aligned}$$

$$\frac{\partial L_{CE}}{\partial w_{k,i}} = (\hat{y}_k - y_k) x_i = \left(\frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} - y_k \right) x_i$$

2.1. Vector Semantics (Idea)

Vector Semantics : learning representations of the meaning of words (embeddings)

→ Finding self-supervised ways to learn word representations is an important focus in NLP.

- the idea is to represent a word as a point in the multidimensional semantic space

⇒ by representing words as embeddings, classifiers can assign classes
as long as it has "similar meanings"

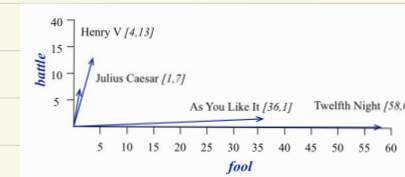
(naive Bayes / Logistic Regression에서는 특징 단어가 나온 횟수에 의존한 반면,
embedding 사용시 직접 나타나지 않아도, 비슷한 의미의 단어들로 고려)

2.2. Cosine Similarity

* Document Vectors

< Term - Document Matrix >

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

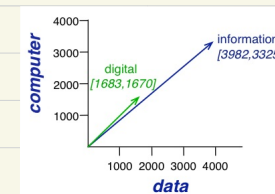


* Word Vectors

< Word - Word Matrix >

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...
digital	0	...	1670	1683	85	5	4	...
information	0	...	3325	3982	378	5	13	...

→ nearby words



$$\Rightarrow \text{Cosine}(v, w) = \frac{v \cdot w}{|v||w|} : \text{두 vector 사이 각도로 similarity 판별}$$

2.3 TF-IDF Weighting

* Problem

: word-word matrix에서 nearby에 많을수록 important,
반면에 too frequent words may be unimportant

1) Term Frequency : frequency of a word(t) in the document(d)

$$tf_{t,d} = \text{count}(t, d) \leadsto tf_{t,d} = \log_{10}(\text{count}(t, d) + 1)$$

"High term frequency \rightarrow more important"

2) Inverse Document Frequency : number of documents the word occurs in

$$idf_t = \log_{10}\left(\frac{N}{df_t}\right) \quad (N \text{ is the number of documents})$$

"Low document frequency \rightarrow more important"

$$\Rightarrow \underline{w_{t,d} = tf_{t,d} \times idf_t}$$

\Rightarrow Represent a document with a Centroid document vector

$$d = \frac{w_1 + w_2 + \dots + w_k}{k} \quad (\text{given } k \text{ word vectors } w_i \ 1 \leq i \leq k)$$

2.4. PPMI Weighting

(PPMI: Positive Pairwise Mutual Information)

$$PMI(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)} \quad \begin{array}{l} \text{-----} \rightarrow \text{co-occur in the corpus} \\ \text{-----} \rightarrow \text{prior-expected to appear by chance} \end{array}$$

\hookrightarrow 각각 단어의 출현 비율로 예상한 것보다 실제로 두 단어가 함께 나타나는 횟수가 많으면 두 단어는 관련이 있다고 생각할 수 있다

\Rightarrow useful to find words that are strongly associated

$$PPMI(w, c) = \max\left(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0\right)$$

* Problem : Bias toward infrequent events (rare words \rightarrow high PMI)

\downarrow

Solution $\left\{ \begin{array}{l} \text{Power Count : } P(c) \text{ 대신 } P_\alpha(c) \text{ 사용} \leadsto P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_c \text{count}(c)^\alpha} \end{array} \right.$

$$\leadsto \underline{PPMI_\alpha(w, c) = \max\left(\log \frac{P(w, c)}{P(w)P_\alpha(c)}, 0\right)}$$

Laplace Smoothing α 는 주로 0.75

2.5. Word2vec

(Word vector는 sparse 하기엔, short dense vector (= embedding)로 변환하자
 \Rightarrow avoid overfitting + better capture synonymy \sim dimensions 50~1000)

* Static Embedding : the method learns 1 fixed embedding for each vocabulary

* Intuition

"instead of counting how often each word occurs near,
 train a classifier on a binary prediction task"

\sim take the weights as the word embeddings

\Rightarrow Self-Supervision : We can just use running text as training data (x labeling)

< SGNS : Skip Gram Negative Sampling >

\hookrightarrow method for computing embeddings

1) Intuition

- 1) If a word(c) is a context word to w, treat as positive(+)
- 2) Randomly sample other words in the lexicon to create negative(-) examples
- 3) Use logistic regression to train a classifier to distinguish (+/-)
- 4) Use the learned weights as embeddings

Similarity(c, w) \approx c · w (두 단어의 word vector의 dot product 높을수록 high similarity)

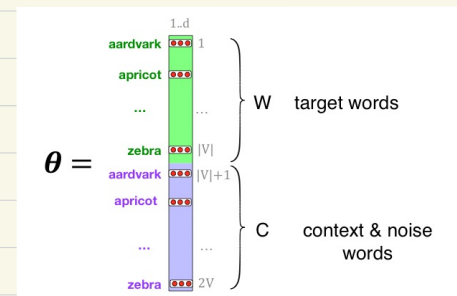
$$\rightarrow P(+|w, c) = \sigma(c \cdot w) \quad / \quad P(-|w, c) = \sigma(-c \cdot w)$$

2) Windowing

\rightarrow go over windows (multiple words)

$$\Rightarrow P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w) \quad (\leftarrow \text{assume that all context words are independent})$$

$$\Leftrightarrow \sum_{i=1}^L \log \sigma(c_i, w) \quad \text{how similar the context window is to the word}$$



|V| x |V| 에서 2|V| x d로 줄임.

2) Learning Skip-Gram Embeddings

... lemon, a (tablespoon of apricot jam, a pinch ...
 c1 c2 w c3 c4

positive examples +

w	c _{pos}
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

negative examples -

w	c _{neg}	w	c _{neg}
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

randomly sampled from the lexicon

based on $p_u(w)$

weighted unigram frequency