

2 Understanding word2vec

2.a

assume $x \in R^d$ and let $\vec{c} \in d$ be a vector with constant entries.

$$\text{softmax}(\vec{x} + \vec{c})_i = \frac{\exp(x_i + c)}{\sum_{j=1}^d \exp(x_j + c)} = \frac{\exp x_i \exp c}{\exp(c) \sum_{j=1}^d \exp(x_j)} = \frac{\exp x_i}{\sum_{j=1}^d \exp(x_j)} = \text{softmax}(\vec{x})$$

2.b

let y be one-hot vector.

$$-\sum_{w \in W} y_w \log(\hat{y}_w) = -\sum_{w \neq o} 0 * \log(\hat{y}_w) + (-\log(\hat{y}_o)) = \log(\hat{y}_o)$$

2.c

$$\begin{aligned} J(\vec{V}_c, o, U) &= -\log\left(\frac{\exp(\mu_o^T \vec{V}_c)}{\sum_{w \in W} \exp(\mu_w^T \vec{V}_c)}\right) = \log\left(\sum_{w \in W} \exp(\mu_w^T \vec{V}_c) - \mu_o^T \vec{V}_c\right) \\ \nabla_{\vec{v}_c}(J) &= -\mu_o + \frac{1}{\sum_{w \in W} \exp(\mu_w^T \vec{V}_c)} \sum_{r \in W} \mu_r^T \exp(\mu_r^T \vec{V}_c) = \sum_{r \in W} \mu_r^T \frac{\exp(\mu_r^T \vec{V}_c)}{\sum_{w \in W} \exp(\mu_w^T \vec{V}_c)} = \\ &\sum_{r \in W} P(o = r) | C = c \rangle \mu_r^T - \mu_o^T = \sum_{r \in W} (\hat{\vec{y}})^T \mu_r^T - \mu_o^T = \\ &E[\vec{\mu}]_{\vec{\mu} \sim P(O|C=c)} - \vec{\mu}_0 \end{aligned}$$

2.d

for $w = 0$:

$$\nabla_{\vec{\mu}_w} = \vec{V}_c \exp(\vec{\mu}_w^T \vec{V}_c) \frac{1}{\sum_{r \in W} \exp(\vec{\mu}_r^T \vec{V}_c)} - \vec{V}_c = \vec{V}_c [(\hat{\vec{Y}}) - 1]$$

for $w \neq 0$:

$$\nabla_{\vec{\mu}_w} = \vec{V}_c (\hat{\vec{Y}})_w$$

2.e

since 6 is applied element-wise, the entries of the derivative vector will be the derivative of σx_i w.r.t for each entry i. thus, will calculate the derivative for single element x_i :

$$\begin{aligned} \frac{\partial \sigma(x_i)}{\partial x_i} &= -\frac{1}{(1 + \exp(-x_i))^2} (-\exp^{-x_i}) = -\frac{1 + (-\exp^{-x_i})}{(1 + \exp(-x_i))^2} - \frac{-1}{(1 + \exp(-x_i))^2} = \\ \sigma(x_i) - \sigma^2(x_i) &= \sigma(x_i)(1 - \sigma(x_i)) = \sigma(x_i)\sigma(-x_i) \end{aligned}$$

2.f

$$\begin{aligned}\nabla_{\vec{\mu}_o}[J_{neg}] &= -[\frac{1}{\sigma(\vec{\mu}_o^T V_c)} \sigma(\vec{\mu}_o^T V_c) \sigma(-\vec{\mu}_o^T V_c) V_c] - [\sum_{k=1}^K \frac{1}{\sigma(\vec{\mu}_k^T V_c)} \sigma(\vec{\mu}_k^T V_c) \sigma(-\vec{\mu}_k^T V_c) \vec{O}_{o \notin (\mu_1, \dots, \mu)} = \\ &= -\sigma(-\vec{\mu}_o^T V_c) V_c\end{aligned}$$

for $j \in (1, 2 \dots K)$:

$$\begin{aligned}\nabla_{\vec{\mu}_j}[J_{neg}] &= -[\frac{1}{\sigma(\vec{\mu}_o^T V_c)} \sigma(\vec{\mu}_o^T V_c) \sigma(-\vec{\mu}_o^T V_c) \vec{O}] - [\sum_{k=1}^K \frac{1}{\sigma(\vec{\mu}_k^T V_c)} \sigma(\vec{\mu}_k^T V_c) \sigma(-\vec{\mu}_k^T V_c) \frac{\partial}{\partial \vec{\mu}_j} [-\mu_k^T \vec{V}_c] = \\ &= +\sigma(-\vec{\mu}_j^T V_c) V_c\end{aligned}$$

$$\begin{aligned}\nabla_{\vec{V}_c}[J_{neg}] &= -[\frac{1}{\sigma(\vec{\mu}_o^T V_c)} \sigma(\vec{\mu}_o^T V_c) \sigma(-\vec{\mu}_o^T V_c) V_c] - [\sum_{k=1}^K \frac{1}{\sigma(\vec{\mu}_k^T V_c)} \sigma(\vec{\mu}_k^T V_c) \sigma(-\vec{\mu}_k^T V_c) \vec{\mu}_k = \\ &\quad \sum_{k=1}^K \sigma(\vec{\mu}_k^T \vec{V}_c) \vec{\mu}_k - \sigma(\vec{\mu}_o^T \vec{V}_c) \vec{\mu}_o\end{aligned}$$

In naïve soft-max one has to update the parameters - going through the entire corpus in each iteration, whereas here, we only consider k "additional" vecotrs.

2.g.i

$$\frac{\partial J_{skip-gram}}{\partial U} = \sum_{-M \leq j \leq M; j \neq o} \frac{\partial (\vec{V}_c, W_{t+j}, U)}{\partial U}$$

2.g.ii

$$\frac{\partial J_{skip-gram}}{\partial \vec{V}_c} = \sum_{-M \leq j \leq M; j \neq o} \frac{\partial (\vec{V}_c, W_{t+j}, U)}{\partial \vec{V}_c}$$

2.g.iii

$$\frac{\partial J_{skip-gram}}{\partial \vec{V}_w} = 0$$

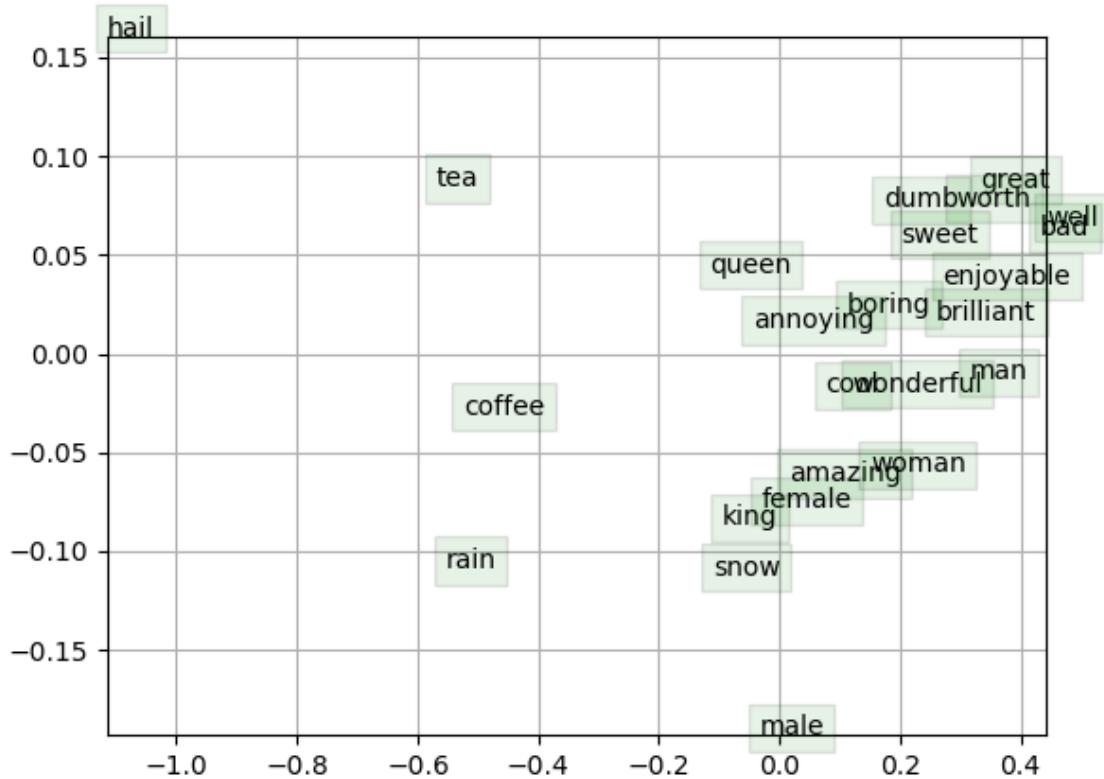


Figure 1: Word Vectors SVD plot

3 Implementing word2vec

3.e

In figure 1 we can see several interesting phenomena, regarding the clusters, trends and algebraic operations.

3.e.i clusters

The most distinct cluster appears on the top right corner of the plot. It includes words such as ("great", "sweet", "dumb") and others, all can be seen in figure 1. All words in this cluster are adjectives. When looking at the left cluster, which is less distinctive, since words are further away from each other, we can see nouns (tea, coffee, rain, hail). The reasons why these words are not clustered as well could be several: their appearances on the text might be too little, or biased (men drink tea, women drink coffee), or the sub dimension we are looking at omits their relationship. In the middle, we can see a sub-cluster when looking at x values - king, queen, male, female. These words are considered as nouns, but they somewhat describe an object (king as a title for a man). We might be enforcing here what we wish to see, but it makes sense in the (x,y) coordinate system and according to the word2vec article.

3.e.ii trends

When looking at x-axis trends, we see the change from nouns (hail, tea, coffee) to adjectives (great, well, bad). In the middle we see titles (king, queen, etc.). For some reason, the word snow appears closer to king and female, than to rain. This of course doesn't support our conclusions, but we suspect it to be an outlier, brought-forth by biased text or the specific dim-reduction applied in SVD. On the y-axis, we couldn't see any clear trends.

3.e.iii algebraic operations

Like word2vec suggested, we tried to perform some algebraic operations. Below is an evaluation of word-coordinates using the grid we added to the plot. lets take a look at the words:

$$\text{king} : (x = -0.10, y = -0.09)$$

$$\text{male} : (x = -0.03, y = -0.19)$$

$$\text{female} : (x = -0.03, y = -0.08)$$

$$\text{queen} : (x = -0.12, y = 0.04)$$

$$\text{result} = \text{king} - \text{male} + \text{female} = (-0.1, 0.02)$$

when using L_2 as distance function, annotated d,

$$d(\text{result}, \text{queen}) = 0.02$$

thus makes it the closest vector to results.

(4,4)

Berücks-Solution:

- Formulate as an Optimization Problem:

$$\underset{\theta}{\text{Max}} \left[\sum_{(c,o) \in b} \#(c,o) \cdot \log [P_{\theta}(o|c)] \right]$$

Subject to:

~~Normal~~ $\forall c: \sum_{o' \in b} P_{\theta}(o'|c) = 1 \Leftrightarrow$

$$\Leftrightarrow \left[1 - \sum_{o' \in b} P_{\theta}(o'|c) \right] = 0$$

\Rightarrow • $L(\theta, \lambda) = \sum_{(c,o) \in b} \#(c,o) \cdot \log [P_{\theta}(o|c)] +$

$$+ \sum_{c \in b} \lambda_c \left[1 - \sum_{o' \in b} P_{\theta}(o'|c) \right]$$

$f(c,o):$

$$\frac{\partial L}{\partial P_{\theta}(o|c)} = \frac{\#(c,o)}{P_{\theta}(o|c)} - \lambda_c = 0 \text{ iff}$$

$$P_{\theta}(o|c) = \frac{\#(c,o)}{\lambda_c}$$

1.



\Rightarrow From our constraint:

$$\forall C : \sum_{O \in \Omega} p_{\theta}(O|C) = 1 \text{ iff } \sum_{O \in \Omega} \left[\frac{\#(C, O)}{\lambda_C} \right] = 1$$

iff $\lambda_C = \sum_{O \in \Omega} \#(C, O)$

~~From the previous slide~~ $\Rightarrow \forall (C, O) \in \Omega :$

$$p_{\theta_{opt}}(O|C) = \frac{\#(C, O)}{\sum_{O' \in \Omega} \#(C, O')}$$

Since $p(O|C) \geq 0$ if $O \in \Omega$ and $\sum_{O \in \Omega} p(O|C) = 1$, we have $\#(C, O) \leq \lambda_C$.

1/2

H, \mathcal{A}

Fish-Solution:

- Define 2 probability distributions over \mathcal{D} :

$$\underline{H(c, o) \in \mathcal{D}}: P[(c, o)] = \frac{\#(c, o)}{|\mathcal{D}|} ; \underline{Q[(c, o)]} = \frac{P_\theta(o|c) \cdot \#(c)}{|\mathcal{D}|}$$

$$\#(c) \stackrel{!}{=} \sum_{o'} \#(c, o')$$

- Re-formulate opt problem:

$$\underset{\theta}{\text{argmax}} [J(Q)] = \underset{\theta}{\text{argmin}} \left\{ - \sum_{(c, o) \in \mathcal{D}} \#(c, o) \cdot \log [P_\theta(o|c)] \right\} =$$

$$= \underset{\theta}{\text{argmin}} \left\{ - \sum_{(c, o) \in \mathcal{D}} \left[\frac{\#(c, o)}{|\mathcal{D}|} \right] \cdot \left[\log [P_\theta(o|c)] + \log [\#(c)] - \log [|\mathcal{D}|] \right] \right\} =$$

$$= \underset{\theta}{\text{argmin}} \left\{ - \sum_{(c, o) \in \mathcal{D}} \left[\frac{\#(c, o)}{|\mathcal{D}|} \right] \cdot \log \left[\frac{P_\theta(o|c) \cdot \#(c)}{|\mathcal{D}|} \right] \right\} =$$

$$= \underset{\theta}{\text{argmin}} \left[H(P, Q) \right] = \underset{\theta}{\text{argmin}} [H(P) + KL(P||Q)] =$$

Cross-entropy!

Entropy!

KL-Divergence!

$$= \underset{\theta}{\text{argmin}} [KL[P||Q]] = ?$$

JK (3)

• As $kL \geq 0 \Rightarrow$ setting: $q = p$

minimizes the objective. \Rightarrow

$$p = q \text{ iff } \forall (c, o) \frac{\#(c, o)}{|D|} = p_o(o|c) \cdot \frac{\#(o|c)}{|D|}$$

$$\text{iff } \frac{\#(c, o)}{\#(c)} = p_o(o|c) \text{ iff } \boxed{p_o(o|c) = \frac{\#(c, o)}{\#(c, o)}}.$$

F.C.J.

(b)

Let $V = \{a, b, c, d\}$ be the vocabulary, and $G = \{"aa", "bb", "cc", "dd"\}$ be the corpus.

→ If we use a window size of ($M=1$), it is clear that each word will appear in D - both as ^ainside words and as an outside-word. Also, no word from V occurs next to a different word from V in the set of tuples D .
[EXCEPT maybe for the "blank" word],

$$\Rightarrow \text{Therefore, } P_1(o|c) = \frac{\#\{o\}}{\sum_{o'} \#\{c,o'\}}$$

a well defined probability for all words in the vocabulary, yet, it ~~assumes~~, for example

$$P_1(a=o | c=b) = \frac{\#\{a,b\}}{1} = \frac{0}{1} = 0$$

$$\text{Whereas, } P_2(o|c) = \frac{e^{\sum_{w \in o} M_w \cdot V_w}}{\sum_{w \in G} e^{\sum_{w \in w} M_w \cdot V_w}} > 0$$

$$\Rightarrow \text{SPECIFICALLY: } P_1(a=o | c=b) = 0 \neq P_2(o=a | c=b)$$

(1.)

\Rightarrow In this instance - the two probabilities differ $\Rightarrow P_2 \neq P_1 \Rightarrow P_2$ cannot approximate solve the optimization problem [AS shown in (a)].

(g.c.d.)

(f.)

⑤ Paraphrase-detection:

- (a) $\begin{cases} \cdot \underline{x}_1, \underline{x}_2 = \text{the input of sentence} \\ \cdot \underline{x}_1, \underline{x}_2 \stackrel{\text{pair}}{\neq} \text{a vector representation for the pair} \end{cases}$

of sentences obtained with so named ~~algorithm~~,
while $\underline{x}_i \in \mathbb{R}^d$.

$$P(\text{Pair is a paraphrase} | \underline{x}_1, \underline{x}_2) = \sigma(\text{relu}^T(\underline{x}_1) \cdot \text{relu}^T(\underline{x}_2))$$

\Rightarrow note that $\forall \underline{x}, \underline{y} \in \mathbb{R}^d : \text{relu}^T(\underline{x}) \cdot \text{relu}^T(\underline{y}) \geq 0$

$$\Rightarrow \sigma(\text{relu}^T(\underline{x}_1) \cdot \text{relu}^T(\underline{x}_2)) = 1$$

due to the behavior of the
~~sigmoid~~ sigmoid function!

\Rightarrow therefore, prediction will always be positive

~~because it's a linear function~~

\Rightarrow accuracy of $\boxed{\frac{1}{4}}$.

(b)

1

(b) An easy fix would be - to use

a different activation for example

the composition: $\text{Sigmoid} \circ \text{tanh}(\text{linear-product})$

It is also possible to consider introducing a bias b - parameter $\epsilon > 0$ s.t.

$$\text{Prediction} = \text{Sigmoid}[\text{relu}^T(\underline{x}_1) \cdot \text{relu}(\underline{x}_2) + b] =$$

\Rightarrow This would allow us to keep the benefits

of using sparse vectors while changing

us classifying negative examples.

[The parameter ϵ should be learned as well].

• The simplest solution of all is to

remove the linear-activation altogether:

$$\Rightarrow \boxed{(\text{Sigmoid}) \circ [\text{linear-product}]}$$

Using tanh might result in squeezing our values toward the edges, that in turn will cause the PDF to be very decisive – values will be very close to either 0 or 1. And alternative solution, will be replacing ReLU with leaky ReLU, which holds representation for negative values as well, thus will solve the problem presented, while keeping the values spread as possible.

