

Aa

:-

Suggested Title: New Note [Edit](#)

September 19, 2025 at 20:52



...

Problem Set 1

1. Word2Vec

word o and c \rightarrow predict $P(O=o | C=c)$

context window of c : W_c

2 2D Vectors : \vec{u}_w outside vector
w is outside word

\vec{v}_w center vector
w is center word

\vec{u}_k outside vectors word indexed by k

\vec{v}_k center ...

U k^{th} col of U and V are \vec{u}_k and
V \vec{v}_k

Soft max

$$P(O=o | C=c) = \frac{\exp(\vec{u}^T \vec{c})}{\sum_{w \in \text{Vocab}} \exp(\vec{u}^T \vec{v}_w)}$$

To minimize $-\log$ / maximize log likelihood

To minimize $-\log$ / maximize log likelihood

$$-\sum_{c \in C} \sum_{o \in \text{vocab}} \log P(O=o | C=c)$$

Single pair of words c and o loss:

$$\mathcal{L}(\vec{v}_c, o, u) = -\log P(O=o | C=c)$$

\mathcal{L} : $y=1$ at position of word o , 0 everywhere else

$$\Rightarrow y_o = 1, y_w = 0 (w \neq o)$$

$$-\sum_{w \in \text{vocab}} y_w \log \hat{y}_w$$

$$= -[y_o \log \hat{y}_o + \sum_{w \neq o} y_w \log \hat{y}_w]$$

$$= -[\cancel{y_o}^1 \log \hat{y}_o + \sum_{w \neq o} 0 \cdot \log \cancel{\hat{y}_w}^0]$$

$$= -\log \hat{y}_o$$

Since: $\hat{y}_o = P(O=o | C=c)$

Since: $\hat{y}_o = P(O=o | C=c)$

$$\Rightarrow -\log \hat{y}_o = -\log P(O=o | C=c)$$

↑
this part is $L(\vec{v}_c, o, u)$

So all 3 parts are equal

1.2 Gradient of loss function

$$L(\vec{v}_c, o, u) = -\log P(O=o | C=c)$$

$$= -\log \left(\frac{\exp(\vec{u}_o^\top \vec{v}_c)}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \right)$$

$$= -\vec{u}_o^\top \vec{v}_c + \log \left(\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) \right)$$

$$\nabla_{\vec{v}_c} (-\vec{u}_o^\top \vec{v}_c) = -\vec{u}_o$$

⚠

↙

$$\nabla_{\vec{v}_c} \log \left(\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) \right)$$

⚠

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \nabla_{\vec{v}_c} \left(\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) \right)$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) \cdot \vec{u}_w$$

$$= \sum_{w \in \text{vocab}} \frac{\exp(\vec{u}_w^\top \vec{v}_c)}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \vec{u}_w$$

$$= \sum_{w \in \text{vocab}} \hat{y}_w \vec{u}_w$$

$$\Rightarrow \nabla_{\vec{v}_c} L = -\vec{u}_0 + \sum_{w \in \text{vocab}} \hat{y}_w \vec{u}_w$$

U : matrix, col k is \vec{u}_k

\vec{y} : 1 at 0, 0 everywhere else

$$U\vec{y} = \vec{u}_0 \cdot 1 + \sum_{w \neq 0} \vec{u}_w \cdot 0 = \vec{u}_0$$



Aa

o-



...



Search



$\hat{\vec{y}}$ is p vector, w is \hat{y}_w

$$U\vec{y} = \sum_{w \in \text{vocab}} \vec{u}_w \cdot \vec{y}_w$$

w \in vocab

$$\Rightarrow \nabla_{\vec{v}_c} L = U\vec{y} - U\hat{\vec{y}}$$

1.3 Improvement

$$\vec{v}_c - \nabla_{\vec{v}_c} L$$

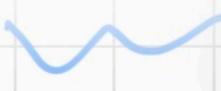
$$= \vec{v}_c - \alpha \nabla_{\vec{v}_c} L$$

$$= \vec{v}_c - \alpha (U\vec{y} - U\hat{\vec{y}})$$

$$= \vec{v}_c - \alpha U\vec{y} + \alpha U\hat{\vec{y}}$$

$\alpha U\vec{y}$

- $\alpha U\vec{y}$: subtract weighted average of all vectors to move \vec{v}_c away from average of all other outside vectors to reduce



Aa

o-



...



Search



- $\alpha \nabla_{\vec{U}} \vec{y}$: subtract weighted average of all vectors to move \vec{v}_c away from average of all other outside vectors to reduce potential false positive words.

$\alpha \nabla_{\vec{U}} \vec{y}$: add true outside vector so that \vec{v}_c move towards the correct word.

1.4 Single pair C cancel O

$$\nabla_{\vec{u}_w} \mathcal{L} = \begin{cases} \vec{y}_w \vec{v}_c & \text{if } w \neq 0 \\ (\hat{y}_0 - 1) \vec{v}_c & \text{if } w = 0 \end{cases}$$

$$\mathcal{L} = -\log \frac{\exp(\vec{u}_0^\top \vec{v}_c)}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)}$$

$$= -\vec{u}_0^\top \vec{v}_c + \log \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)$$

when $w \neq 0$: $\nabla_{\vec{u}_w} (\vec{y}_w \vec{v}_c)$



Aa

o-



...

Search



1.4 Single pair C and O

$$\nabla_{\vec{u}_w} \mathcal{L} = \begin{cases} \vec{y}_w \vec{v}_c & \text{if } w \neq 0 \\ (\hat{y}_0 - 1) \vec{v}_c & \text{if } w = 0 \end{cases}$$

$$\mathcal{L} = -\log \frac{\exp(\vec{u}_0^\top \vec{v}_c)}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)}$$

$$= -\vec{u}_0^\top \vec{v}_c + \log \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)$$

when $w \neq 0$: $\nabla_{\vec{u}_w} (\vec{y}_w \vec{v}_c)$

$$= \nabla_{\vec{u}_w} \left(-\vec{u}_0^\top \vec{v}_c + \log \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) \right)$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \nabla_{\vec{u}_w} \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) (\vec{v}_c)$$

WEVOCAB

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c)} \cdot \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c) (\vec{v}_c)$$

$$= \hat{y}_w \vec{v}_c$$

when $w = 0$:

$$\nabla_{\vec{u}_0} (\hat{y}_0 - 1) \vec{v}_c$$

$$\nabla_{\vec{u}_0} (-\vec{u}_0^\top \vec{v}_c + \log \sum_{w \in \text{vocab}} \exp(\vec{u}_w^\top \vec{v}_c))$$

$$= -\vec{v}_c + \hat{y}_0 \vec{v}_c$$

$$= (\hat{y}_0 - 1) \vec{v}_c$$

Same as above

 $\nabla_{\vec{u}} L$

$$\vec{u} = \begin{bmatrix} \vec{u}_1 \\ \vec{u}_2 \\ \vdots \\ \vec{u}_0 \\ \vdots \\ \vec{u}_{|\mathcal{U}|} \end{bmatrix}$$

$$\nabla_{\vec{u}} L = \begin{bmatrix} \nabla_{\vec{u}_1} L \\ \nabla_{\vec{u}_2} L \\ \vdots \\ \nabla_{\vec{u}_0} L \\ \vdots \\ \nabla_{\vec{u}_{|\mathcal{U}|}} L \end{bmatrix} \rightarrow \begin{array}{l} \hat{y}_w \vec{v}_c \\ (\hat{y}_0 - 1) \vec{v}_c \end{array}$$

∇_{UL}

$$U = \begin{bmatrix} \vec{u}_1 \\ \vec{u}_2 \\ \vdots \\ \vec{u}_0 \\ \vdots \\ \vec{u}_{|U|} \end{bmatrix} \quad \nabla_{UL} = \begin{bmatrix} \nabla_{\vec{u}_1} L \\ \nabla_{\vec{u}_2} L \\ \vdots \\ \nabla_{\vec{u}_0} L \\ \vdots \\ \nabla_{\vec{u}_{|U|}} L \end{bmatrix} \rightarrow \hat{y}_w \vec{v}_c$$

$(\hat{y}_0 - 1) \vec{v}_c$

$$\nabla_{UL} = \vec{v}_c (\hat{\vec{y}} - \vec{y})^T$$

1.5 L^2 normalization

L^2 makes all words with same magnitude

but keeps the sign of the words.

So if we only care about positivity or negativity of the words, L^2 will not impact the result. But if we need

word's magnitude, L^2 would make all words same mag, so very harmful.



Aa

o-



...

Search



2 Adam Optimizer and Dropout

$$\text{SGD} \quad \vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha \nabla_{\vec{\theta}} L(\vec{\theta}^{(t)})$$

limitation: too slow, might diverge.

same step for all, ignore freq & steepness
may jitter

2.1 Momentum

rolling average:

$$m^{(t+1)} = \beta_1 m^{(t)} + (1 - \beta_1) \nabla_{\vec{\theta}} L(\vec{\theta}^{(t)})$$

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha m^{(t+1)}$$

β_1 : hyperparameter 0 - 1 (usually 0.9)

$$f(\vec{\theta}) = \frac{1}{2} (1000 \theta_1^2 + \theta_2^2), \quad \vec{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

convergence tolerance $\epsilon = 0.001$

initial $\vec{\theta}^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, at most $i_{\max} = 1000$

CONVERGENCE CRITERIA $\epsilon = 0.001$

initial $\vec{\theta}^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, at most $i_{\max} = 1000$

SGD with $\alpha = 0.1$

SGD with $\alpha = 0.005$

SGD with $\alpha = 0.1 \quad \beta_1 = 0.9$

optimal ?

iterations ?

2.4 Dropout

\vec{h} : hidden layer with k elements.

randomly set to 0, with P_{drop} , then γ

$$\vec{h}_{\text{drop}} = \gamma \vec{d} * \vec{h}$$

\vec{d} same size \vec{h} , random

$$- 0 = P_{\text{drop}}$$

$$- 1 : 1 - P_{\text{drop}}$$

Aa

o-



...

Search



$$- | : 1 - P_{\text{drop}}$$

(a) γ to make \vec{h}_{drop} is \vec{h}

$$E_{P_{\text{drop}}} [(\vec{h}_{\text{drop}})_i] = \vec{h}_i$$

$$i = 1, 2, \dots, K$$

$$\text{prove } \gamma = \frac{1}{1 - P_{\text{drop}}}$$

since $d_i = 0$ or $d_i = 1$

$$E = \gamma * P(d_i=0) * h_i + \gamma * P(d_i=1) * h_i$$

$$= \gamma * 0 * h_i + \gamma * (1 - P_{\text{drop}}) * h_i$$

$$E = \gamma * h_i (1 - P_{\text{drop}}) = \vec{h}_i^T \vec{1}$$

$$\Rightarrow \gamma (1 - P_{\text{drop}}) = 1$$

$$\Rightarrow \gamma = \frac{1}{1 - P_{\text{drop}}}$$

If $P_{\text{drop}} = 0.5$ meaning half

Aa

Search



$$E = \gamma * h_i (1 - P_{drop}) = \cancel{h_i}^{\rightarrow} \gamma^2$$

$$\Rightarrow \gamma (1 - P_{drop}) = 1$$

$$\Rightarrow \gamma = \frac{1}{1 - P_{drop}}$$

b. If $P_{drop} = 0.5$, meaning half will be dropped, $\gamma = \frac{1}{1-0.5} = 2$, meaning doubling back. So $\gamma = \frac{1}{1-P_{drop}}$ makes sense.

C. Having dropout during training can help lower the overfitting chance.

We don't do dropout during evaluation because we need consistent input in order to compare the results to see if they remain consistent.