

**CPSC 477**

**Spring 2019**

**HW3**

**Due Friday, March 29 by 4 PM**

<b>Problem</b>	<b>Grade</b>
QUESTION 1 (Language Modeling)	/10
QUESTION 2 (Recurrent Neural Networks)	/10
QUESTION 3 (Morphology)	/10
QUESTION 4 (Text Similarity)	/10
QUESTION 5 (Dimensionality Reduction)	/10
QUESTION 6 (Naïve Bayes Classifier)	/10
QUESTION 7 (Training Neural Networks)	/10
QUESTION 8 (Word Embeddings)	/10
<b>TOTAL</b>	<b>/80</b>

**Submission Instructions**

1. Submit each solution on a separate sheet.
2. Include your **name**, **netID** and **problem number** on each page.

$$p(a|bc) = \frac{p(cb|ac)p(a|c)}{p(cb|c)}$$

### QUESTION 1 (Language Modeling)

You are in a noisy coffee shop, diligently studying for your midterm, and your friend is trying to get your attention, using only a two words vocabulary. She has said a sentence but you couldn't hear one of the words:

$$(w_1 = 'hi'; w_2 = 'yo'; w_3 = ???; w_4 = 'yo')$$

Assume that your friend was generating words from this first-order Markov model:

$$\begin{aligned} p('hi'|'hi') &= 0.7 & p('yo'|'hi') &= 0.3 \\ p('hi'|'yo') &= 0.5 & p('yo'|'yo') &= 0.5 \end{aligned}$$

Given these parameters, what is the posterior probability of whether the missing word is 'hi' or 'yo'? Show your work.

want  $P(w_3|w_1, w_2, w_4)$ , by Markov assumption  $p(w_3|w_2, w_4)$

$\propto P(w_3|w_2) P(w_4|w_2, w_3)$

(by Markov assumption)  $\propto P(w_3|w_2) P(w_4|w_3)$

$$\begin{aligned} \frac{P(w_3=hi)}{P(w_3=yo)} &= \frac{P(hi|yo) \cdot P(yo|hi)}{P(yo|yo) \cdot P(yo|yo)} \\ &= \frac{0.5 \times 0.3}{0.5 \times 0.5} \\ &= \frac{3}{5} \end{aligned}$$

?  $P(w_3=hi) = 0.375$   
 $P(w_3=yo) = 0.625$

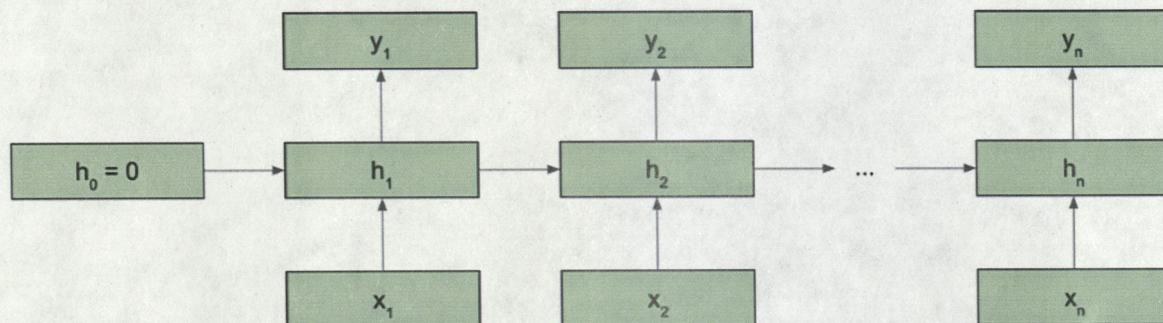
## QUESTION 2 (Recurrent Neural Networks)

Your goal will be to set up weights for a simple recurrent network (SRN and also known as an Elman Network) to solve two simple sequence tasks. The network reads in a sequence  $x_1, \dots, x_t$  and produces outputs  $y_1, \dots, y_t$ . We consider the final output  $y_t$  to be the answer given by the network.

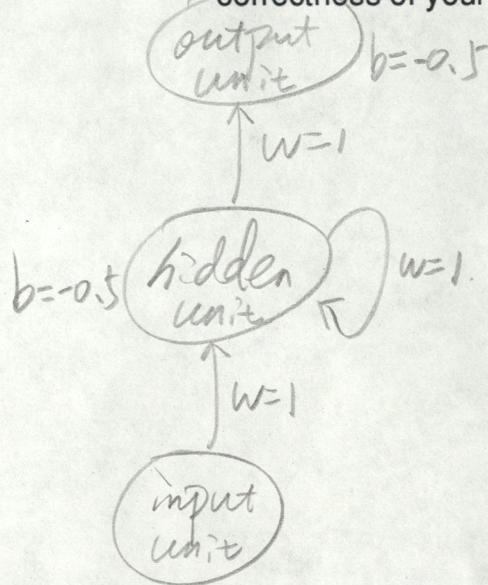
We will assume a threshold activation function:

$$f(x) = 1 \text{ if } x > 0; \text{ otherwise } f(x) = 0.$$

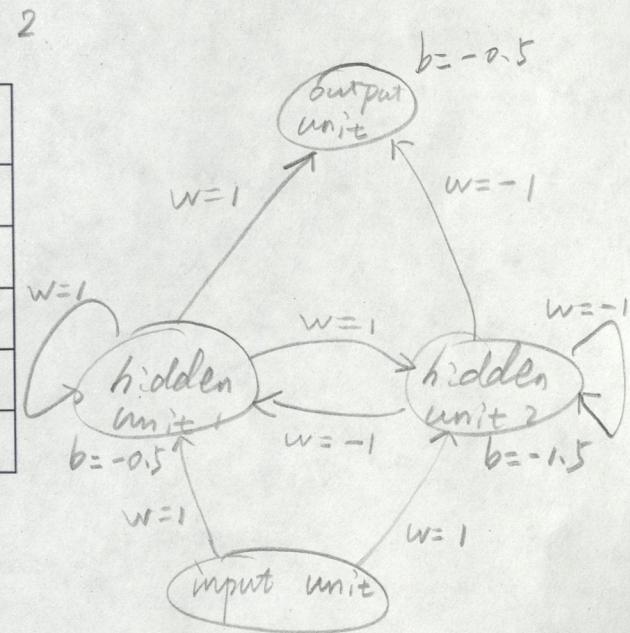
The network has one simple recurrent layer followed by a linear layer (shown in below diagram). The dimension of both the input to the network and its output should be 1.



1. Construct an SRN with 1 hidden unit that reads a sequence of 1s and 0s and outputs 1 if any of  $x_1, \dots, x_t$  had the value 1. Justify the correctness of your network.
2. Construct an SRN with 2 hidden units to solve the parity task (read in a sequence of 1s and 0s and return the sum mod 2, shown in the below table). Justify the correctness of your network.



$x$	$y$
0	0
1	1
010	1
0110	0
10101	1



Justification (1)

time

input 0 0 1 0 1 1 0

unit 0 0 1 1 1 1 1

output 

0	0	1	1	1	1	1
---	---	---	---	---	---	---

at each time  $i$

$$U^i = \text{input}^i + U^{i-1} - 0.5$$

$$\text{output}^i = U^i - 0.5$$

"OR" of 0, 00, 001, 0010,  
00101, 001011, 0010110, respectively.

∴ verified

Justification (2)

time

1 2 3 4 5

input

1 0 1 0 1

unit 1

1 1 1 0 1

unit 2

0 0 1 0 0

output 

1	1	0	0	1
---	---	---	---	---

at each time  $i$ ,

$$U_1^i = \text{input}^i + U_1^{i-1} - U_2^{i-1} - 0.5$$

$$U_2^i = \text{input}^i + U_1^{i-1} - U_2^{i-1} - 1.5$$

$$\text{output}^i = U_1^i - U_2^i - 0.5$$

Parity of 1, 10, 101, 1010, 10101, respectively.

∴ verified

### QUESTION 3 (Morphology)

Refer to the slides for lecture #143 for background on this question.

1. Consider the English verb *drink*. Give two inflected forms of *drink*, and two derivational forms of *drink*.

inflectional form: *drinks, drank*  
derivational form: *drinkable, undrinkable*.

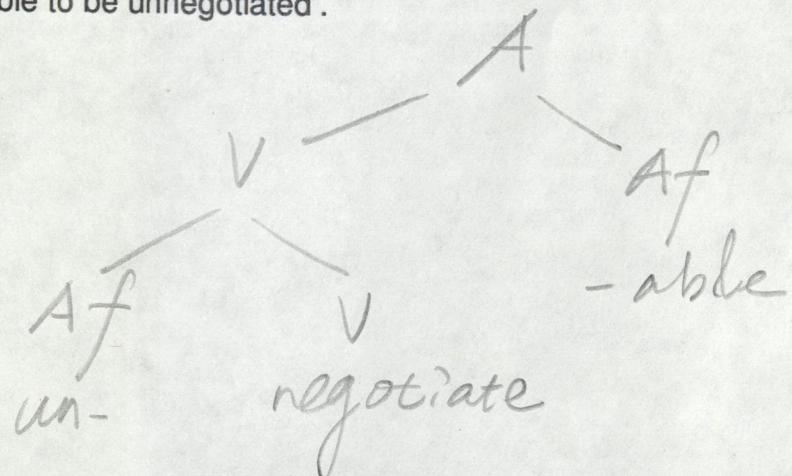
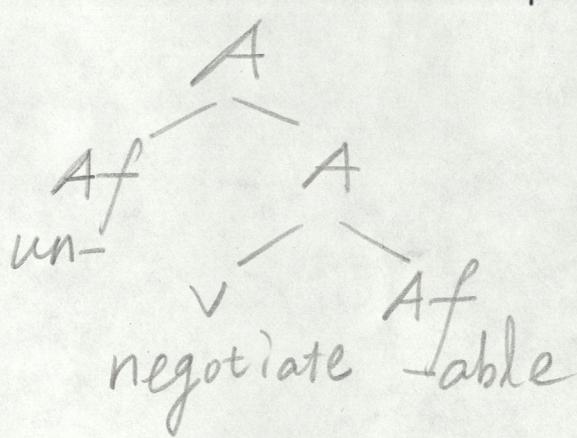
2. Identify all the morphemes in *infallibilities*. For each morpheme, say whether it is a prefix, suffix, or root.

prefix: *in-*

root: *fall*

suffix: *-ible, -ity, -ies*

3. Consider the word *unnegotiable*. Draw and label two morphological trees for this word, one which corresponds to the meaning 'not able to be negotiated', and another which corresponds to 'able to be unnegotiated'.



#### QUESTION 4 (Text Similarity)

We first need to introduce a function that compares the similarity of two sentences. Recall the definition for the cosine distance between two vectors  $x$  and  $y$ :

$$\text{cosine}(x, y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

*Part A*  
**Part A:** How would you define a mapping from sentences  $s$  to vectors  $f(s)$  such that the cosine distance  $\text{cosine}(f(s_1), f(s_2))$  is a measure of the similarity between two sentences  $s_1$  and  $s_2$ ? We will give full marks for any reasonable mapping.

*Let  $f(s)$  be a vector of count for each word in the vocabulary.*  
**Part B:** Now, we will define a dynamic programming algorithm that computes sentence alignment of two translated documents. The alignment score is the sum of the similarity scores of the aligned sentences. Our goal is to find an alignment with the highest score.

We will consider alignments of the following form:

- a sentence can be aligned to an empty sentence. This happens when one of the translators omits a sentence.
- a sentence can be aligned to exactly one sentence.
- a sentence can be aligned to two sentences. This happens when one of the translators either breaks or joins sentences.

Our sentence alignment algorithm recursively computes the alignment matrix  $F$  indexed by  $i$  and  $j$ , one index for each sentence. The value stored in  $F(i, j)$  is the score of the best alignment between the first  $i$  sentences of  $x$  and the first  $j$  sentences of  $y$ .  $s(x_i, y_j)$  is the similarity between sentence  $x_i$  and sentence  $y_j$ .

Based on the above information, please

Define  $F(0, 0)$ ,  $F(i, 0)$  and  $F(0, j)$ .

Define  $F(i, j)$  for  $i > 0$  and  $j > 0$ .

**Part C:** Next, we will modify the alignment score to be the same as before, but to include a fixed penalty  $p$  each time a sentence is aligned to an empty sentence.  $p$  is a parameter,

which is  $\geq 0$ , chosen by the user of the algorithm. Describe a modified dynamic programming method that takes this new penalty into account.

Part B.

$$F(0,0)=0, \quad F(i,0)=0, \quad F(0,j)=0.$$
$$F(i,j) = \max \left\{ \begin{array}{l} F(i-1, j-1) + S(i, j), \\ F(i-1, j), \\ F(i, j-1), \\ F(i-2, j-1) + S(i-1, j) + S(i, j), \\ F(i-1, j-2) + S(i, j-1) + S(i, j) \end{array} \right\}$$

Part C

$$F(0,0)=0, \quad F(i,0)=-i \cdot p, \quad F(0,j)=-j \cdot p$$
$$F(i,j) = \max \left\{ \begin{array}{l} F(i-1, j-1) + S(i, j), \\ F(i-1, j) - p, \\ F(i, j-1) - p, \\ F(i-2, j-1) + S(i-1, j) + S(i, j), \\ F(i-1, j-2) + S(i, j-1) + S(i, j) \end{array} \right\}$$

**QUESTION 5 (Dimensionality reduction and SVD)**

(a) Prove that the left singular vectors of a matrix  $A$  are the right singular vectors of  $A^T$

let  $A = U \Sigma V^T \rightarrow U$  is the left singular vectors of  $A$   
 then  $A^T = V \Sigma^T U^T \rightarrow U$  is the right singular vectors of  $A^T$

(b) Calculate the covariance matrix for the matrix  $M$  shown below.

covariance matrix  $M = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 \end{bmatrix}$

for each column, minus  
each value by the column  
mean

$A = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ -\frac{2}{3} & -\frac{2}{3} & 0 & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & 0 & -\frac{2}{3} \end{pmatrix}$

$\begin{pmatrix} x_1 & x_2 & x_3 & x_4 \end{pmatrix}$

mean  $\frac{2}{3} \quad \frac{2}{3} \quad 0 \quad \frac{2}{3}$  the covariance matrix

$= \frac{1}{n} A^T A$ , where  $n=3$

(c) Find the singular values of the matrix  $M$

$$M^T M = \begin{pmatrix} 2 & 2 & 0 & 1 \\ 2 & 2 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 2 \end{pmatrix}$$

let  $\det(M^T M - \lambda I) = 0$ ,

$$\lambda = 0, 3 + \sqrt{3}, 3 - \sqrt{3}$$

$\therefore$  singular values of  $M$  is  $0, \sqrt{3+\sqrt{3}}, \sqrt{3-\sqrt{3}}$ ,  
 or equivalently,  $0, 2.175, 1.126$ .

### QUESTION 6 (Naïve Bayes Classifier)

Suppose you are given the following set of data with three Boolean input variables  $a, b$ , and  $c$ , and a single Boolean output variable  $K$ .

$$P(K=1 | a=1, b=1, c=0)$$

$$= P(a=1 | K=1) \cdot P(b=1 | K=1) \cdot P(c=0 | K=1)$$

$$P(a=1, b=1, c=0)$$

$$= \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{8}$$

$$= \frac{\frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{8} + \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{5}{8}}{52}$$

$$= \frac{25}{52}$$

$$= 0.48$$

<b>a</b>	<b>b</b>	<b>c</b>	<b>K</b>
1	0	1	0
1	1	1	1
0	1	1	0
1	1	0	0
1	0	1	0
0	0	0	1
0	0	0	1
0	0	1	0

	<b>a</b>	<b>b</b>	<b>c</b>
$K=0$	0	1	0
$K=0$	$\frac{3}{5}$	$\frac{3}{5}$	$\frac{2}{5}$
$K=1$	$\frac{2}{3}$	$\frac{1}{3}$	$\frac{2}{3}$

$$(b, P(K=0 | a=1, b=1))$$

$$= P(a=1 | K=0) \cdot P(b=1 | K=0) \cdot P(K=0)$$

$$= \frac{\frac{3}{5} \times \frac{2}{5} \times \frac{5}{8}}{P(a=1, b=1)}$$

$$= \frac{\frac{3}{5} \times \frac{2}{5} \times \frac{5}{8}}{\frac{3}{5} \times \frac{2}{5} \times \frac{5}{8} + \frac{1}{3} \times \frac{1}{3} \times \frac{3}{8}}$$

$$= \frac{18}{23}$$

$$= 0.78$$

(For parts (a) and (b), assume we are using a Naïve Bayes classifier to predict the value of  $K$  from the values of the other variables.)

(a) According to the naive Bayes Classifier, what is  $p(K = 1 | a = 1, b = 1, c = 0)$  ?

(b) According to the naive Bayes Classifier, what is  $p(K = 0 | a = 1, b = 1)$  ?

(c) Given  $P(Z|X) = 0.7$  and  $P(Z|Y) = 0.4$ .

Do you have enough information to compute  $P(Z|X, Y)$ ? If not, write "not enough info".

If so, compute the value of  $P(Z|X, Y)$  from the above information

Not enough info.

(e) Given  $P(Z, X) = 0.2$ ,  $P(X) = 0.3$ ,  $P(Y) = 1$

Do you now have enough information to compute  $P(Z|X, Y)$ ? If not, write "not enough info". If so, compute the value of  $P(Z|X, Y)$  from the above information.

$$P(Z|X, Y) = P(Z|X) \text{ since } P(Y) = 1$$

$$= \frac{P(Z, X)}{P(X)} = \frac{2}{3}$$

Enough info

### QUESTION 7 (Training Neural Networks)

This question is to train a three-layer neural network for classification task. Let  $H_1$  denote the number of hidden units, let  $D$  be the dimension of the input  $X$ , and let  $K$  be the number of classes. The three-layer network has the following form:

$$\begin{aligned} h_1 &= \text{ReLU}(W_1 X + b_1) \\ h_2 &= \text{ReLU}(W_2 h_1 + b_2) \\ p &= \text{Softmax}(W_3 h_2 + b_3) \end{aligned}$$

where the parameters of the network are  $W_1 \in R^{H_1*D}$ ,  $b_1 \in R^{H_1}$ ,  $W_2 \in R^{H_2*H_1}$ ,  $b_2 \in R^{H_2}$ ,  $W_3 \in R^{K*H_2}$ ,  $b_3 \in R^K$ .

$\text{ReLU}$  is the “rectified linear unit”  $\text{ReLU}(x) = \max\{0, x\}$  applied componentwise, and  $\text{SoftMax}$  maps a d-vector to the probability simplex according to

$$\text{Softmax}(v)_k = \frac{\exp(v_k)}{\sum_{j=1}^K \exp(v_j)}$$

For a given input  $X$ , this specifies how to calculate the probabilities  $P(Y = k|X)$  for  $k = 1, 2, \dots, K$ .

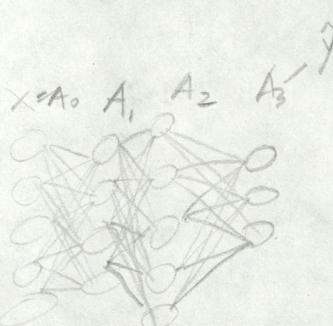
For a given training set  $\{(X_i, Y_i)\}$ , consider the loss function:

$$L = \frac{1}{n} \sum_{i=1}^n -\log p(Y = Y_i | X_i) + \frac{\lambda}{2} (\|W_1\|^2 + \|W_2\|^2 + \|W_3\|^2)$$

where  $\|A\|^2$  means the sum of the squares of the entries of the matrix  $A$ .

Give formulas for the derivatives for each layer of the network  $j = 1, 2, 3$ .

$$\frac{\partial L}{\partial W_j}, \frac{\partial L}{\partial b_j}$$



Hints:

i. You should ignore the fact that  $\text{ReLU}(x)$  is not differentiable at  $x_i = 0$ .

ii. You might try to first do the calculations with  $\text{ReLU}()$  replaced by the identity function.

$$P_k = \text{softmax}(v)_k = \frac{e^{v_k}}{\sum_{j=1}^n e^{v_j}}$$

$$\begin{aligned} h_1 &= \text{ReLU}(w_1 x + b_1) \\ h_2 &= \text{ReLU}(w_2 h_1 + b_2) \\ v &= \text{softmax}(w_3^T h_2 + b_3) \end{aligned}$$

$$L = \frac{1}{n} \sum_{i=1}^n \left\{ -\log P(Y=Y_i(k)/X_i) \right\} + \frac{\lambda}{2} (||w_1||^2 + ||w_2||^2 + ||w_3||^2)$$

$$\frac{\partial L}{\partial v_k} = \frac{\partial L}{\partial P_k} \cdot \frac{\partial P_k}{\partial v_k} = \frac{1}{n} \sum_{k=1}^n \left( -\frac{1}{P_k} \right) \cdot P_k (1-P_k) = \frac{1}{n} \sum_{k=1}^n (P_k - 1)$$

$$\therefore \frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial v_k} \frac{\partial v_k}{\partial w_3} + \lambda ||w_3|| = \frac{1}{n} \sum_{k=1}^n (P_k - 1) h_2 + \lambda ||w_3|| \quad (1)$$

$$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial v_k} \frac{\partial v_k}{\partial b_3} = \frac{1}{n} \sum_{k=1}^n (P_k - 1) \quad (2)$$

$$\frac{\partial L}{\partial h_2} = \frac{\partial L}{\partial v_k} \frac{\partial v_k}{\partial h_2} = \begin{cases} \frac{1}{n} \sum_{k=1}^n w_3^T \cdot (P_k - 1) & \text{otherwise} \\ 0 & h_2 < 0 \end{cases}$$

$$\therefore \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial w_2} + \lambda ||w_2|| = \frac{1}{n} \sum_{k=1}^n w_3^T h_1^T \cdot (P_k - 1) + \lambda ||w_2|| \quad (3)$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial b_2} = \frac{1}{n} \sum_{k=1}^n w_3^T \cdot (P_k - 1) \quad (4)$$

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial h_1} = \begin{cases} \frac{1}{n} \sum_{k=1}^n w_2^T w_3^T \cdot (P_k - 1) & \text{otherwise} \\ 0 & h_1 < 0 \end{cases}$$

$$\therefore \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial w_1} + \lambda ||w_1|| = \frac{1}{n} \sum_{k=1}^n w_2^T w_3^T \cdot x^T \cdot (P_k - 1) + \lambda ||w_1|| \quad (5)$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial b_1} = \frac{1}{n} \sum_{k=1}^n w_2^T w_3^T \cdot (P_k - 1) \quad (6)$$

## QUESTION 8 (Word Embeddings<sup>1</sup>)

Word2Vec represents a family of word-embedding algorithms that are commonly used in a variety of contexts.

- (a) Let's consider a recommender system for an online music-streaming service. It includes information about a set of songs on how frequently users play them together (e.g., song X is commonly played together with song Y). Explain how you would use ideas similar to Word2Vec to recommend the next song to a user who has played some songs from this list.
- Treat user's listening queue as a sentence, with each word being a song. Then, train word2vec model to get high-dimensional vectors so that similar songs have weights that are closer together.*
- (b) Although pre-trained word vectors work quite well in many applications, it is sometimes better to "re-train" (i.e. continue to learn) the word vectors as parameters of our neural network. Explain why retraining the word vectors may hurt the model performance if the dataset for the specific task is too small.
- Pre-trained word vectors produce semantically related words to be located in the same part of the word space. When retrained with a small dataset, these words are shifted in the word space.*
- (c) Give two examples of how we can evaluate word-embedding vectors.
- ① word vector analogies:  $(a:b) :: (c:?)$
  - ② correlation evaluation: compare cosine similarity of word vectors with human's assessment on the
- (d) A simple way to produce word embeddings for a vocabulary of words is to create one-hot vectors for each word - a vector with 0s everywhere and only a single 1 in the position corresponding to the word's index in the vocabulary. Give at least two reasons why word2vec-produced word embeddings are better than these one-hot vectors.
- ① More computationally efficient, since word2vec word embeddings are lower dimensional, dense and doesn't grow with the vocabulary size. One-hot vectors are higher dimensional, sparse and grows with the vocabulary size.
  - ② word2vec word embeddings preserve the semantic meaning of words, i.e., semantically similar words have similar vectors. This is not the case with one-hot vectors.