# MIDTERM EXAM

NAME: Shihab Muhtasim

ID: 21301610
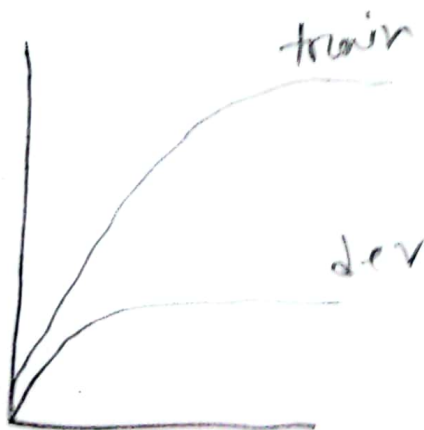
FALL 23

## Ans to qu 1(A)

Overfitting of a model is when a model is too complex for given data that it memorizes its data and fails to detect patterns. Again, using higher models for simple data can cause overfitting.

we can solve it by :

① Using simpler models or making the complex model simpler by using less polynomials.

② we can feed the model less data so that it tries to learn

train

dev

Overfitting performs good in training and bad during valid testing.

# Ans to or 1 (B)

TF : shallow appears 2 times.

But appears 7 times.

$$TF(shallow; d) = \log_{10}(C(shallow, d) + 1)$$

$$= \log_{10}(2 + 1) = 0.47$$

$$TF(but, d) = \log_{10}(7 + 1) = 0.9$$

$$IDF(shallow) = \log_{10}\left(\frac{N}{df_{shallow}}\right)$$

$$= \log_{10}\left(\frac{10000}{1000}\right)$$

$$= 1$$

$$IDF(but) = \log_{10}\left(\frac{10000}{9000}\right)$$

$$= 0.045$$

Weight of shallow = $0.47 \times 1$ = $0.47$

weight of but = $0.045 \times 0.9$ = $0.0405$

∴ shallow is more important than but in this doc.

## Ans to 1 (c)

Accuracy is bad performance metrix because it often works bad for imbalanced datasets. for example if a huge data is of positive class and bad classifier predicts even negative classes to positives classes still it can give good accuracy. Hence, its missleading

## Ans to or 2.

### (A)

Issues in english language tokenizations:

Tokenization is a process of splitting

a) into words.

① In case of punctuations it can be hard to tokenize a word

② If certain words are connected and have space between them. Ex: New york

③ Hyphen between words can cause problems.

## Ans to or 2B

A, Paris = $[3, 4, 0, 1]$    Oslo = $(1, 2, 1, 1)$
B, France = $[3, 2, 1, 0]$

A − B = $[0, 2, -1, 1]$

A − B + sweden = $[3, 3, 0, 2]$

$$\cos(\text{Oslo}, A-B+\text{sweden}) = \frac{3+6+0+2}{\sqrt{1+4+1+1} \ \sqrt{9+9+0+4}}$$

$$= 0.886$$

A − B + Norway = $[3, 4, 2, 3]$

$$\cos(\text{Oslo}, A-B+\text{Norway}) = \frac{3+8+2+3}{\sqrt{7} \ \times \ \sqrt{9+16+4+9}}$$

$$= 0.981$$

$$\text{oslo} = [1, 2, 1, 1)$$

$$A - B + \text{Finland} = [1, 3, 0, 1)$$

$$\cos\left(\overset{\text{oslo}}{A-B} + \text{Finland}\right) = \frac{1+6+0+1}{\sqrt{7}\ \sqrt{1+9+1}}$$

$$= 0.911$$

$$A - B + \text{Denmark} = [0, 2, 0, 3)$$

$$\cos(\text{oslo}, A-B+\text{Denmark}) = \frac{0+4+0+3}{\sqrt{7}\ \sqrt{4+9}}$$

$$= 0.7337$$

∴ Ans is oslo is to norway.

$$P = \cdot$$

## (A)

classifier A :

$W_A = [11 \ 3]^T = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$

$X = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$W_A^T = [11 \ 3]$

$W_A^T X + b_A = [11 \ 3] \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0$

$= 2 + 0 = 2$

$6(2) = \dfrac{1}{1 + e^{-2}} = 0.88 = y_A$

classifier B :

$W_B^T X + b_B = [1 \ 20] \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + (-1)$

$= 3 - 1 = 2$

$6(2) = 0.88 = y_B$

x's original table = 1

$$LCE = -\log\left(\hat{y}^{\,Y}(1-\hat{y})^{1-Y}\right)$$

$$= -\log\left(0.88^{1}\cdot(1-0.88)^{1-1}\right)$$

$$= 0.055$$

$\therefore$ Both classifier will have 0.055 cross entropy loss.

## Ans to or 3(B)

Although it may be hard but it is possible. We can use BOW, term document frequency and inverse doc freq to analyze the sentiments of the reviews from dataset and then collect symantic information for each review. The words that collec have most important in words in each review we can set them to positive or neg class either by human annotation or by using lexicon list of word and then run a classifier on that dataset.