

# Natural Language Processing

## Assignment - 2

a)  $y$  is the true distribution and the predicted distribution is  $\hat{y}$ .

$$\therefore (\hat{y})_k = \frac{\sup(\mu_k^T v_k)}{\sum_{i=1}^V \sup(\mu_i^T v_k)}$$

Where  $\mu_k$  is the  $k$ th row of  $U$ .  
~~or also called the content vector for word with index  $k$ .~~

Suppose the true outside word is 0

$$\hat{y} \therefore \hat{y}(y)_0 = 1$$

$$(y)_w = 0 \quad \text{if } w \neq 0$$

$$\therefore \text{cross entropy} = - \sum_{w \in V_{\text{vocab}}} y_w \log(\hat{y}_w)$$

$$= - \sum_{w \neq 0} y_w \log(\hat{y}_w) - \sum_0 y_0 \log(\hat{y}_0)$$

$$y_w = 0 \quad \forall w \neq 0$$

$$\Rightarrow = -y_0 \log(\hat{y}_0)$$

$$y_0 = 1$$

$$\Rightarrow = -\log(\hat{y}_0)$$

Hence proved.

$$h) \quad J = -\log P(C=0 | C=2)$$

$$= - \log \frac{\sup(\mu_0^T v_c)}{\sum_{w \in V_{\text{node}}} \sup(\mu_w^T v_c)}$$

$$= -\mu_0^T V_c + \log \left( \sum_w \exp(\mu_w^T V_c) \right)$$

$$\frac{\partial J}{\partial v_c} = -\mu_0 + \frac{1}{\sum_w \sigma_{\text{out}}(\mu_w^T v_c)} \times \sum_w \mu_w \sigma_{\text{out}}(\mu_w^T v_c)$$

$$\frac{\partial J}{\partial v_c} = -\mu_0 + \sum_w \mu_w \frac{\sigma_{\text{out}}(\mu_w^T v_c)}{\sum_w \sigma_{\text{out}}(\mu_w^T v_c)}$$

$$\frac{\partial J}{\partial v_c} = \sum_w \mu_w \hat{y}_w - \mu_0$$

$$\boxed{\frac{\partial J}{\partial v_c} = E[\mu_w] - \mu_0}$$

c) (case 1)  $w \neq 0$

$$J = -\mu_0^T v_c + \log \left( \sum_k \mu_k^T v_c \right)$$

$$\frac{\partial J}{\partial \mu_w} = 0 + \frac{1}{\sum_k \mu_k^T v_c}$$

$$J = -\mu_0^T v_c + \log \left( \sum_k \mu_k^T v_c \right)$$

$$\frac{\partial J}{\partial \mu_w} = 0 + \frac{1}{\sum_k \mu_k^T v_c} \times v_c \mu_w^T v_c$$

$$\frac{\partial J}{\partial \mu_w} = v_c \hat{y}_w$$

Case 2:  $w = 0$

$$J = -\mu_0^T v_c + \log \left( \sum_w \exp(\mu_w^T v_c) \right)$$

$$\frac{\partial J}{\partial \mu_0} = -v_c + \frac{1}{\sum_w \exp(\mu_w^T v_c)} \times \sum_w \frac{\partial}{\partial \mu_0} \exp(\mu_w^T v_c)$$

$$= -v_c + \frac{1}{\sum_w \exp(\mu_w^T v_c)} \times \cancel{\sum_w} v_c \exp(\mu_0^T v_c)$$

$$\boxed{\frac{\partial J}{\partial \mu_0} = -v_c (1 - \hat{y}_0)}$$



## Coding Part

c) In the plot I can see several clusters of words that make sense, such as tea & coffee, enjoyable & annoying, woman & female, etc.

However many things don't make sense, like how male is away from man, worth is close to man.

A boy & woman are close, which does not make sense. We can also see the bias of the model, by observing that the word dumb is closer to queen than king.