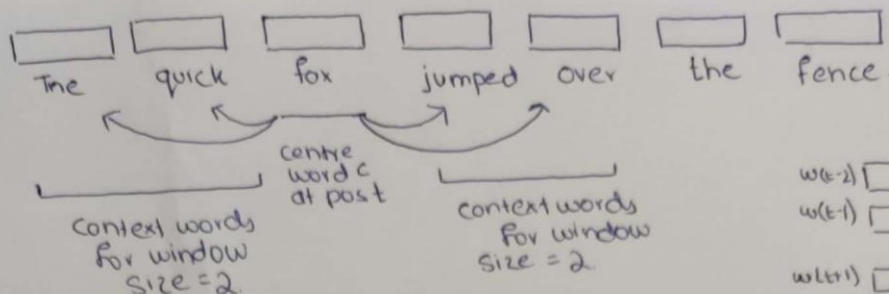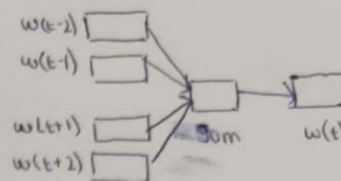# MID TERM NLP 702
# PSUEDO CODES
## Muhammad Umar Salman - 21010241

## Question 2

Assume you are asked to implement CBOW. Please choose a sentence and some sample training set, draw an annotated CBOW model architecture and write a pseudo code. Make sure you show how you prepare the training samples, translate the training window and qualitatively and quantitatively evaluate your embedding implementation. (4 points)

The [ ]  quick [ ]  fox [ ]  jumped [ ]  over [ ]  the [ ]  fence. [ ]

CBOW

centre word c at post

Context words for window Size = 2

context words for window Size = 2

$w(t-2)$ [ ]
$w(t-1)$ [ ]
$w(t+1)$ [ ]
$w(t+2)$ [ ]

Som [ ] $w(t)$

Input    Projection  Output

# Preparing Training Samples

```
window_size = 2
text ← load_data()
text ← text.lower()
text ← regex.replace (punctuation,' ', text)
text ← text.split()
train, val, test ← train.test_split(text).
Iterate: for token in train:
            if length(token)) > 2*window-Size + 1
            then
                for index, word in enumerate (window-size, length(token) - window size)
                then
                    pre-window ← [words for twords in token before centre
                                  word of size 2]
                    post_window ← [words for words in token after centre
                                   word of size window-size]
                    context-words ← pre_window + post_window
                    context-words ← [get_encoding (w) for w in context_words]
                    label_word ← [get_encoding (centre_word)
        end
            return context_words, label_word.
```

In Translate Training window with more detail.

# Translate training window

```
token ←[ the, quick, fox, jumped, over, the, fence]     # take w as token.
                                                          in psuedo_code.
for w in Train:
    if len(w) > 2*window-size +1
    then.
        for i in range (window to length(w) -window-size)
                       size
            pre-window ← [w(i-j-1) for j in window-size]
            post_window ← [w(i+j+1) for j in window size]
            context_words ← pre-window + post-window
            centre-word ← w(i)
```

e.g
context [the, quick, jumped, over]     centre → fox
context [quick, fox, over, the]        centre → jumped
context [ fox, jumped, the, fence]     centre → over.

What is the computational bottleneck operation in Word2Vec? What can you adopt other than negative sampling? Assume you decide to adopt negative sampling. Briefly describe what is negative sampling. Please show how you would prepare the training data for such a task and write the pseudo code to show how training will partially updates the model weights. (4 points)

- The computational bottleneck for Word2Vec is the Softmax operation. Although the computation is that of a dot product. But because for all pairs with the centre word we each time normalize over the entire vocabulary, it becomes very expensive and hence is the bottle neck.

- The 2 improvements other than Negative sampling for word2vec is 1) Word pairs and phrases, where we prefer phrases of infrequent words. Although this increases the vocabulary size it decreases the training expense. 2) To Subsample frequent words to decrease no. of training examples.

- Negative Sampling defines a new objection function on top of the existing skip gram model which maximizes the similarity of in context words compared to out of context words which it aims to minimize. Instead of minimizing all words in the dictionary, the model causes each training sample to update only a small percentage of the models weight. It randomly selects k negative samples and our one true sample and trains a logistic regression instead of a softmax across the entire vocabulary.

```
k = 5
text = load_text_data()
new_text = [line.split() for line in text]

function get_samples (k)
        negword = []
        while negwords < k
              negwords.append (Vocab(random.randint))
              neglabels = [0] * length(negwords).
        end
        return negword, neglabels.

for N epochs
      for sentence in new_text
            for j, word in sentence.
                  u, u_labels = get_samples (k)

                  U = [getOHV (i) for i in u]

                  c = word
                  c_label = 1
                  v = sentence[i+1]
```

$$J(\theta) = \log \sigma \left( getOHV(v)^T \cdot getOHV(c) \right) + \sum_{L=1}^{k} E_{j \sim P(w)} \left[ \log \sigma (-u_L^T \cdot getOHV(c)) \right]$$

$$\theta_{new} = \theta_{old} - \eta \frac{\partial J(\theta)}{\partial \theta}$$
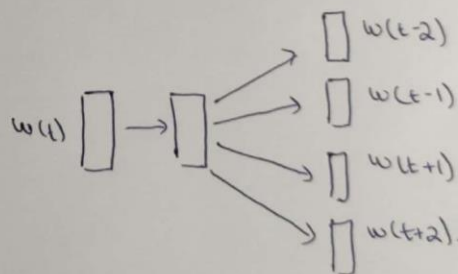
* ideally $\frac{U(w)^{3/4}}{z}$ but we did random.

```
end for  end for  end for
end for
```

The quick brown fox ....

| context | centre | output |
|---------|--------|--------|
| quick   | brown  | 1 |
| quick   | king   | 0 |
| quick   | brother| 0 |
| quick   | hire   | 0 |
| quick   | food   | 0 |
| quick   | orange | 0 |

K:5

input    projection    output
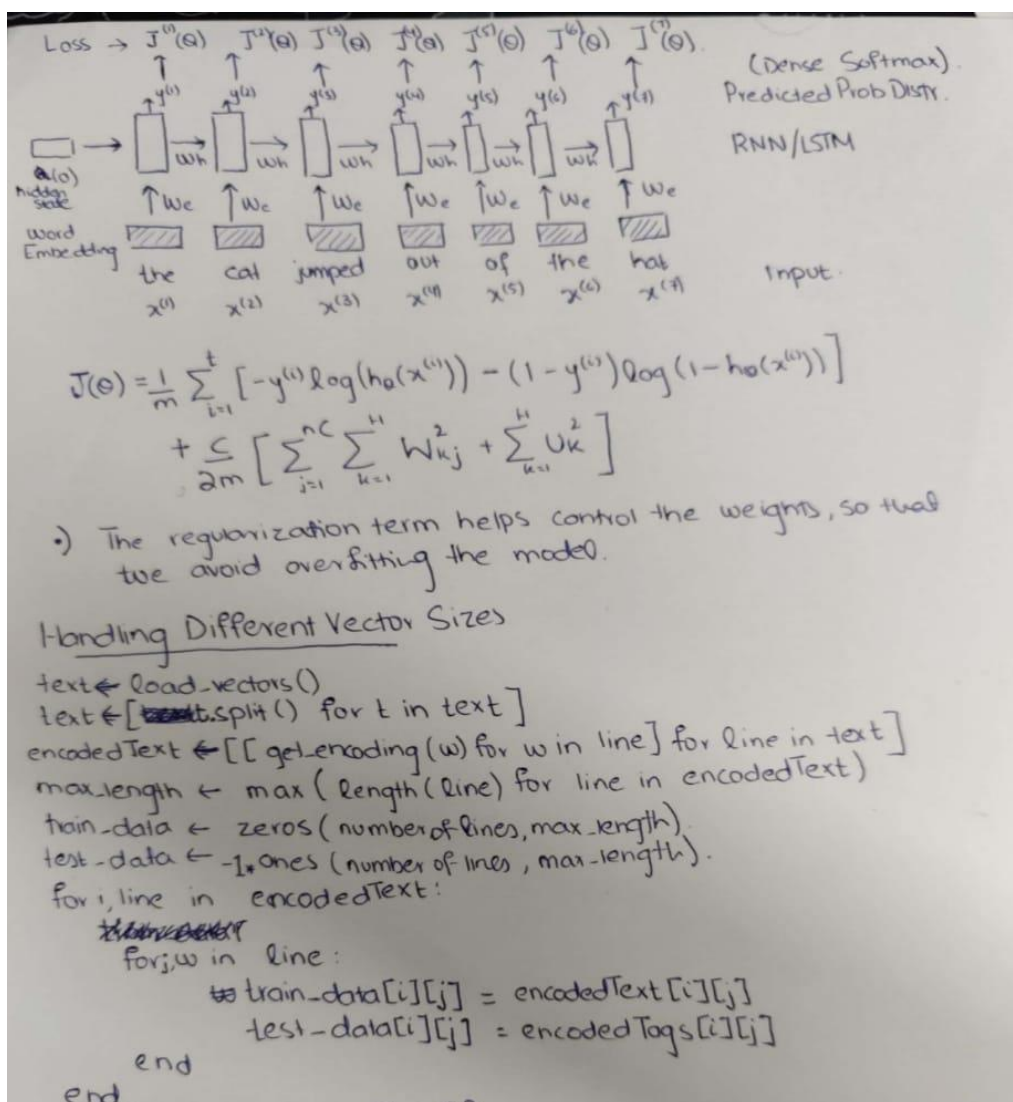
__Skip gram__

* Changed multinomial to binary classification problem
* change from prediction of word to get core words neighbours?

$$w(t) \Box \rightarrow \Box$$

$$\Box \ w(t-2)$$
$$\Box \ w(t-1)$$
$$\Box \ w(t+1)$$
$$\Box \ w(t+2).$$

# Question 5

Please annotate an LSTM architecture to implement an NER task (is a NN or not) for the following sentence: the cat jumped out of the hat. What do you achieve when you include a regularization term in the NER cost function? Write a pseudo code showing how you handle different vector sizes and how you read LSTM output in NER.

Loss → $J^{(1)}(\theta)$ $J^{(2)}(\theta)$ $J^{(3)}(\theta)$ $J^{(4)}(\theta)$ $J^{(5)}(\theta)$ $J^{(6)}(\theta)$ $J^{(7)}(\theta)$

(Dense Softmax). Predicted Prob Distr.

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$ $y^{(4)}$ $y^{(5)}$ $y^{(6)}$ $y^{(7)}$

RNN/LSTM

$a^{(0)}$ hidden state

$\uparrow W_e$ $\uparrow W_e$ $\uparrow W_e$ $\uparrow W_e$ $\uparrow W_e$ $\uparrow W_e$ $\uparrow W_e$

word Embedding

| the | cat | jumped | out | of | the | hat | Input. |

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$ $x^{(5)}$ $x^{(6)}$ $x^{(7)}$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{t} \left[ -y^{(i)} \log(h_\theta(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right]$$

$$+ \frac{c}{2m} \left[ \sum_{j=1}^{nc} \sum_{k=1}^{H} W_{kj}^2 + \sum_{k=1}^{H} U_k^2 \right]$$

·) The regularization term helps control the weights, so that we avoid overfitting the model.

## Handling Different Vector Sizes

text ← load_vectors()
text ← [text.split() for t in text]
encoded Text ← [[get_encoding(w) for w in line] for line in text]
max_length ← max(length(line) for line in encoded Text)
train_data ← zeros(number of lines, max_length)
test_data ← -1*ones(number of lines, max_length)
for i, line in encoded Text:
~~text.nrnnrnnrr~~
    for j, w in line:
      ~~to~~ train_data[i][j] = encoded Text[i][j]
      test_data[i][j] = encoded Tags[i][j]
  end
end

# Redd LSTM output for NER

X_train, y_train ← prepare_data()

Length_arr ← [len(rows) for row in x_train]    # list of all lengths.

masked_arr ← length_arr ≥ 0    # masking excludes padding which has -1 as labels.

Packed_seq ← pack_padded_seq (masked_arr, length_arr)

Lstm_out ← LSTM (Packed seq)

unpacked_seq ← pad_packed_seq (lstm_out)    # unpack output

linear_out ← Fully Connected Layer (unpacked_seq)

output_prob ← log_softmax (linear_out)

NER_Tag ← IDtoTag (max (output_prob))

# Question 6

Sketch a block-based diagram and write a pseudo code for a hybrid event extraction approach that uses 1) a bidirectional LSTM for sentence encoding and 2) a CNN for event detection that uses word embedding, a positional embedding and (bonus) a retrofitting of embeddings to existing lexicon as input

## CNN

word_embedding = load_pretrained_embeddings ()
pos_emb 1 = get_pos_embedding (1)
pos_emb 2 = get_pos_embedding (2)
embedding = concat (word_embedding, pos_emb 1, pos_emb 2)
convolution = conv1D (embedding, num_kernels, kernel size, padding, stride)
masked_conv = conv.masked_fill (mask, convolution)   # for 'PAD' tokens.
max_pool = maxpool (masked_conv)
all_masked_pool = concat ([max_pool 1; maxpool 2; ..... max_pool n;])

## LSTM

word_embedding = load_pretrained_embedding ()
packed_sequence = pack_padded_sequence (word_embedding, lengths)
               # pack sequence to ignore 'PAD' token.
lstm 1 = LSTM (packed_sequence, hidden_units, embed_dim)
LSTM 2 = LSTM ( " reversed ")
unpacked_seq1 = unpack_padded_sequence (lstm1).
output_lstm1 = linear Fully connected (hidden_units, num.classes, unpacked_seq1)
unpacked_seq2 = unpack_padded_sequence (LSTM2)
output_lstm2 = Fullyconnected (hidden_units, num.classes, unpacked_seq2)

output_concat = concat ([output_lstm1, output_lstm2])
sentence_embedding = [concat (output_concat) for all output_concats]

## JOINT

feature_vector = all_max_pool + sentence_embedding
output_FC = Fully connected (output_FC, hidden_units, num.classes)
final output = log_softmax (output_FC).