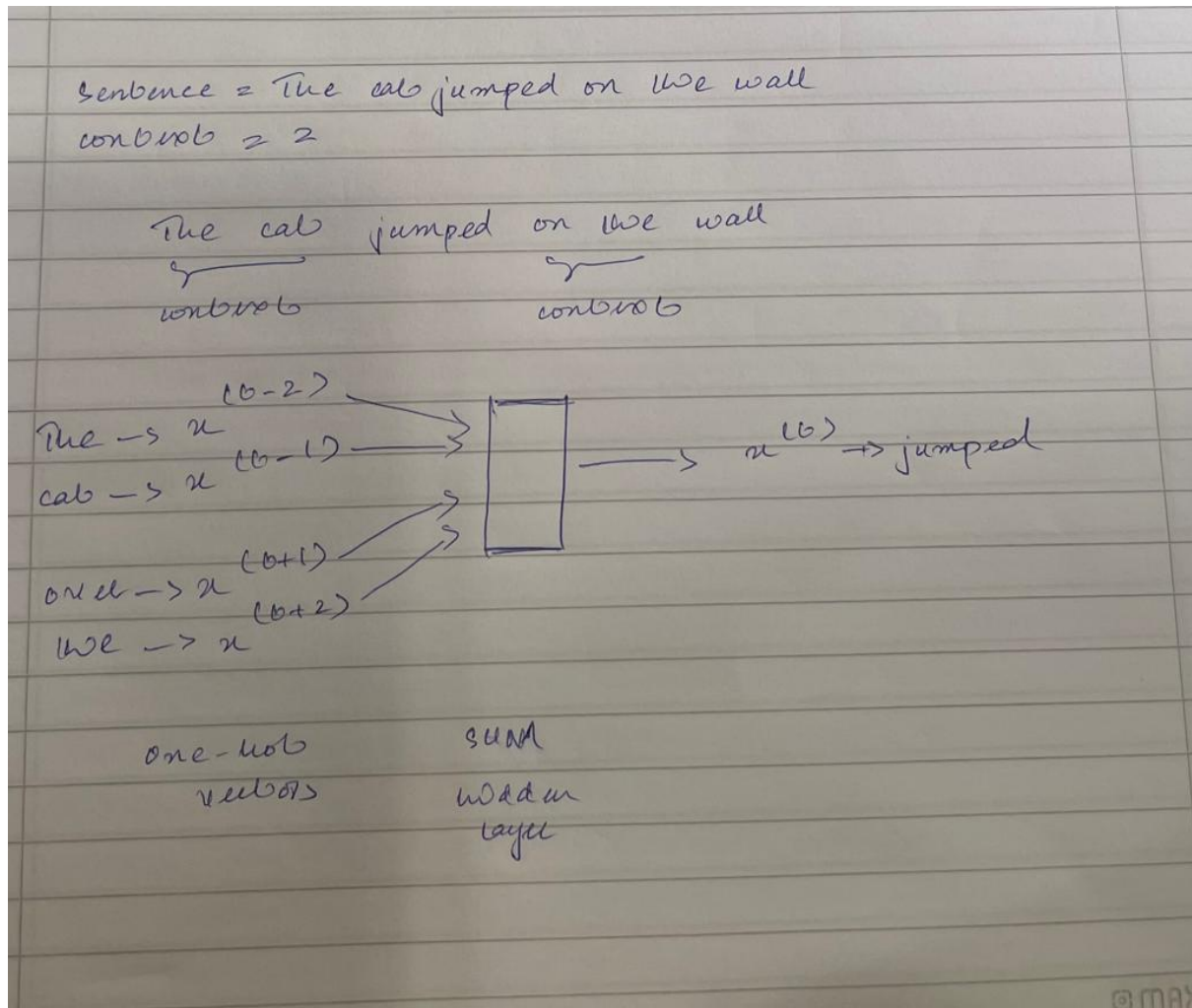**Question # 02**

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

**Annotated CBOW Architecture**

**CBOW Pseudocode**

```
CBOW Pseudocode
initialize training_set <- Word2Vec (brain)

for sentence in training_set:
    set sent = token
    set sentence = lower (sentence)
    set sentence = remove_punc (sent)
    set sentence = remove_line (sentence)


initialize Model CBOW (vocab_size, embed_dim)


# Training
initialize EPOCHS = N
for epoch in EPOCHS do
    initialize X & initialize Y
    for word in sentence do
        initialize conbvob = [word - CONTEXT_SIZE] +
                             [word + CONTEXT_SIZE]
        center_word = word
        X. add (conbvob)
        Y. add (center_word)
    end for
```

```
    output = model (X)
    # Quantitative evaluation
    loss = model.calc_loss ()
    model.update_weights ()
end for


# Qualitative evaluation
embedding ("Queen") = embedding ("king") - embedding ("man")
                                + embedding ("woman")
```

What is the computational bottleneck operation in word2vec? What can you adopt other than 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 the training will partially update the model weights.

The computational bottleneck operation in the word2vec model is the final softmax layer. The softmax operation sums over the entire vocabulary making it the most expensive. Other than negative sampling a hierarchical softmax operation can be adopted which modifies the architecture of the softmax layer to improve its efficiency. The approach can yield up to 5x speedups.
**Negative Sampling:** As mentioned before, summation over |V| is computationally huge, practically |V| can be in millions. In each training step, instead of looping over the entire vocabulary, for each positive sample (word, context pair belonging to the corpus), the negative sampling algorithm generates several negative samples (word, context pair belonging to a false or negative corpus). The negative samples are sampled from a noise distribution $P_n(w)$. The algorithm maximizes the probability of a word and context being in the corpus data if it is, otherwise maximizes the probability of it not being in the corpus data if it is not.

Negative sampling Pseudocode

Input: training set $s \rightarrow (w, c)$ # word and context pair

loop
  $s_{batch} \leftarrow$ sample $(s, b)$ # sample minibatch of size $b$
  $T_{batch} \leftarrow$ null # batch that will contain negative and positive samples
  # for loop to generate negative samples
  for pos_sample in $s_{batch}$ do
    neg_sample $\leftarrow$ sample $(s'(w', c's))$ # sample corrupted pairs
    $T_{batch} \leftarrow T_{batch} + \{(pos\_sample, neg\_sample)\}$
  end for

  If model is skip_gram:
    # Update embeddings w.r.t following objective
    $$J(\theta) \rightarrow [-\log \sigma(u_{c-m+j}^T \cdot v_c) - \sum_{t=1} \log \sigma(-\tilde{u}_{t_c}^T \cdot v_c)]$$

    # where $\theta$ are the parameters of the model
    $$\theta_{new} = \theta_{old} - \alpha \frac{\partial J(\theta)}{\partial \theta}$$
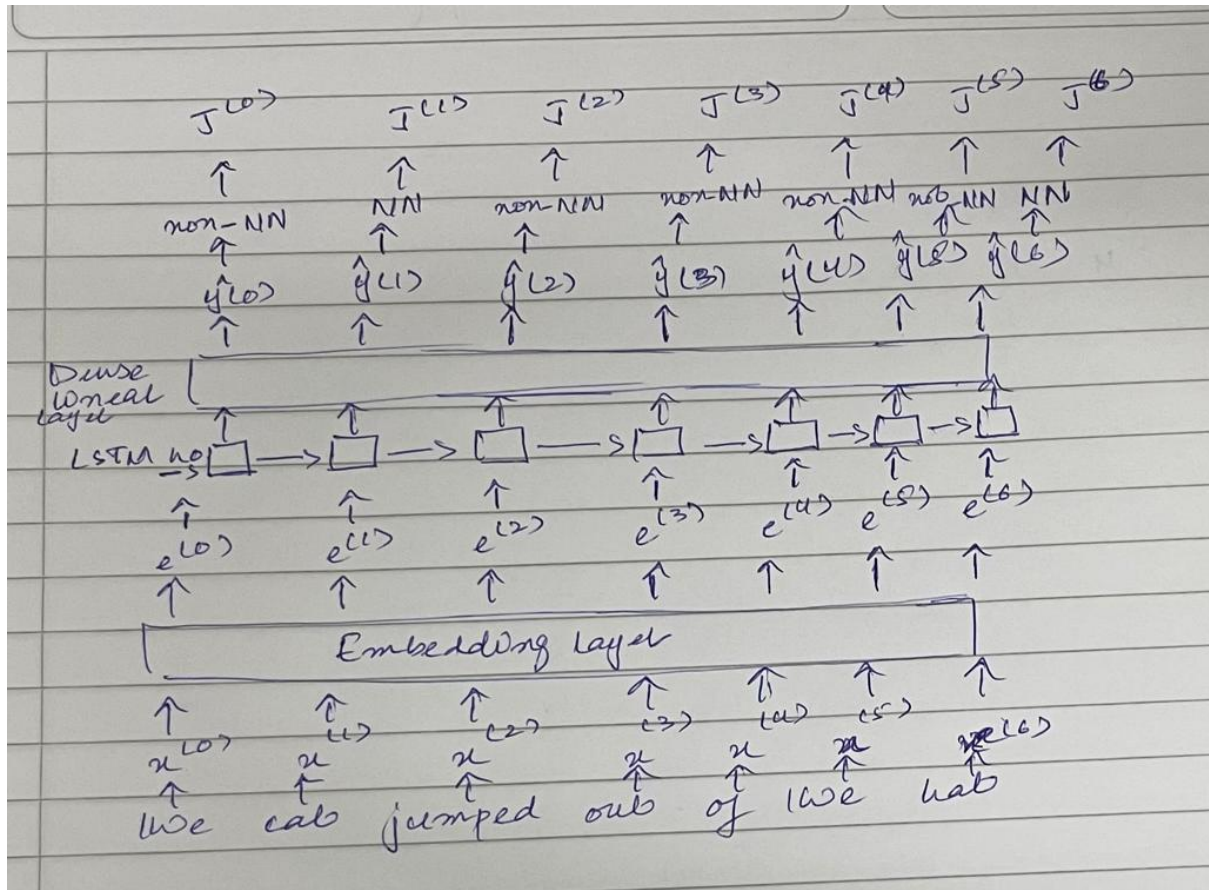
  else if model is cbow:
    # Update embeddings w.r.t following objective
    $$J(\theta) = [-\log \sigma(u_c^T \cdot \hat{v}) - \sum_{t=1} \log \sigma(-\tilde{u}_t^T \cdot \hat{v})]$$

end loop

## Question # 05

Please annotate an LSTM architecture to implement a 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 regularization in
the NER cost function? Write a pseudocode showing how you handle different vector sizes and how
you read LSTM output in NER.



Neural network algorithms in general tend to be complex and prone to overfitting. Including a
regularization term in the NER cost function, in-turn helps to reduce the complexity of the learning
algorithm and reduce the variance in the model. The result is that the model becomes more general in nature.

```
for label in output:
    if label[0] > label[1]:
        print "Not an NN"
    else
        print "NN"
```

## Question # 06

Sketch a block based diagram and write a pseudocode 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 embeddings, positional embeddings and (bonus) a retrofitting of embeddings to existing lexicon as input.

**Pseudocode**

```
Pseudocode Hybrid Event Detection Approach


# sentence encoding
func sentence_encoding ( s )
    Initialize sentence = S # Array of words
    for word in sentence do
            embedding_array.add(get_embeddings ( word ))
    end for


    forward, backward = biLSTM ( embedding_array )
    hidden   =   concat ( forward, backward )
    output   =   linear ( hidden )
    output =   sigmoid ( output )


end func


    Initialize train_set = T
    for sentence in train_set do
        for word in sentence do
                embedding_array . add ( get_embeddings ( word ))
        end for
        retro_embed = retrofit_embeddings ( embedding_array )
        pos_embed = get_positional_embeddings ( sentence )
        embeddings = concat ( retro_embed, pos_embed )
        conv_output =      apply_conv_filters ( embeddings )
        max_output =     apply_max_pooling ( conv_output )
        sent_encod =     sentence_encoding ( sentence )
        output =     linear ( sent_encod, max_output )
        output  =   softmax ( output )
```