# Word2Vec and Doc2Vec

**Nagiza F. Samatova,** samatova@csc.ncsu.edu
**Professor, Department of Computer Science**
**North Carolina State University**

# Word2Vec is a tool to create vector space representations of words

- **"embeddings" is a term often used instead of "vector space representation":**
  - Embedding means representation of an object embedded (placed) in a vector space.
- **Word2Vec produces vector space representations of words that:**
  - Capture co-occurrence of words in text
  - Captures sematic meaning by accounting for word order
- **Word2Vec is a shallow neural network**
  - Input layer, hidden layer, and output layer
  - The number of neurons in the hidden layer is the number of dimensions in the vector space
  - The neural net in Word2Vec is only used to obtain the weights, it is not used in the normal way – e.g. to perform classification with forward propagation.
- **Word2Vec has two different versions:**
  - continuous bag of words and
  - skip-gram

# Ex: Semantic Meaning of Relationships in Vector Space

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?

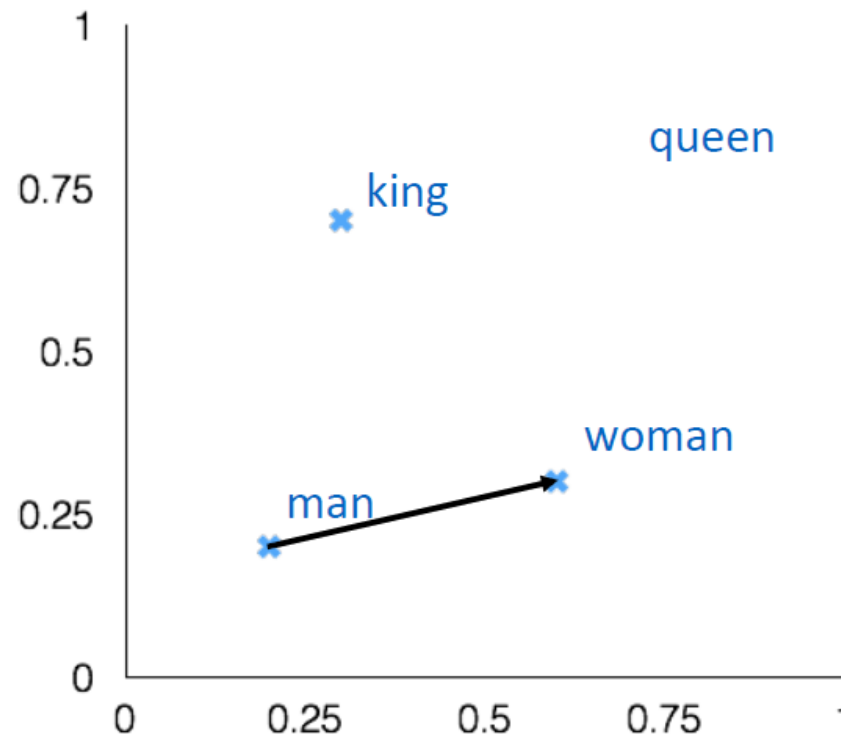$$d = \arg\max_{x} \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\|}$$

**Word Analogies**

man:woman :: king:?

| | | |
|---|---|---|
| + | king | [ 0.30 0.70 ] |
| - | man | [ 0.20 0.20 ] |
| + | woman | [ 0.60 0.30 ] |
| | queen | [ 0.70 0.80 ] |

http://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf (Author of GloVe)

# word2vec: Input, output, parameters, training...

- **Input**
  - Corpus – the text that you train word2vec model on (e.g., the entire Wikipedia, your Inbox, all Harry Potter books, etc.)
  - Vocabulary (it is extracted by word2vec but could be an input)
- **Parameters**
  - Number of dimensions of the target vector space – N:
    - Usually small: 64, 128, ..., under 1,000
    - Unlike bag of words: ~100,000 of words, with each word as a dimension
  - context window size:
    - Usually: 5-11 words
  - others
- **Output**
  - Vectors in the N-dimensional space (Euclidean space)
    - Each word is represented by an array of N numbers

# Continuous Bag of Words (CBOW): One-word context

We assume that there is only one word considered per context.

Model will predict one target word given one context word

# Neural network: One-word context

The input vector $x = \{x_1, \dots, x_V\}$ is a one-hot encoded vector
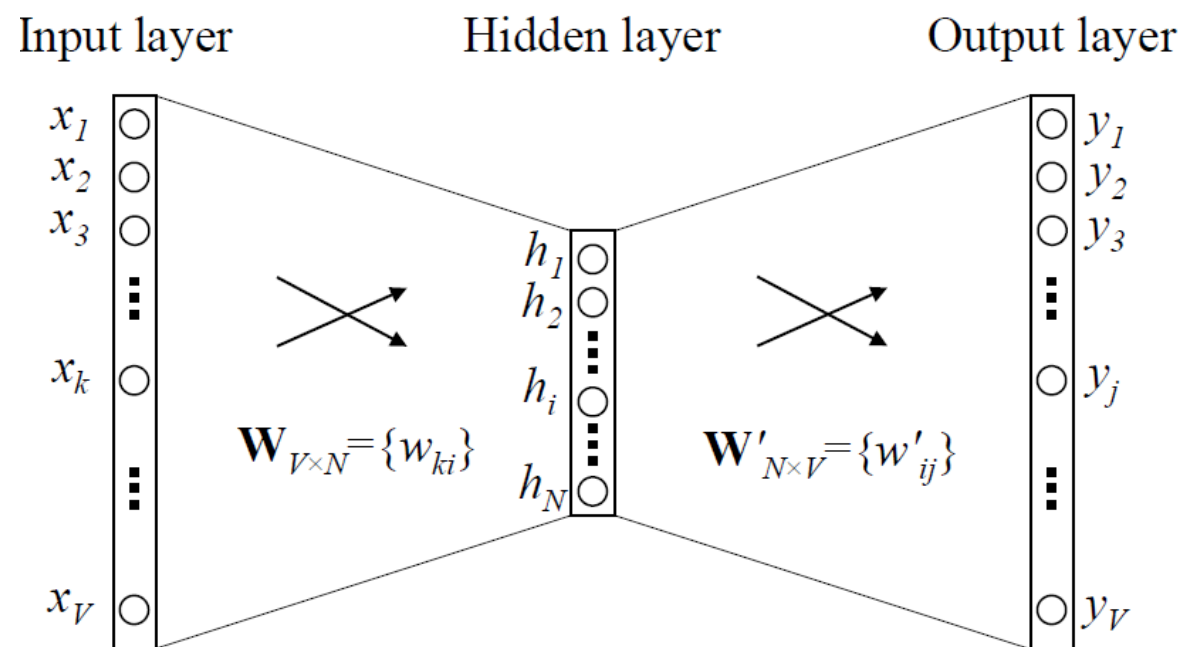
E.g. "eats" -> [0 0 0 1 0 0 0]

$V$ is the vocabulary size

$N$ is size of the hidden layer

The weights between the input and the hidden layer can be represented by a $V \times N$ matrix $W$

Each row of $W$ is the $N$-dimensional vector $v_w$

$v_w$ is the vector representation of associated word of the input layer.



Input layer     Hidden layer     Output layer

$x_1$, $x_2$, $x_3$, $x_k$, $x_V$

$h_1$, $h_2$, $h_i$, $h_N$

$y_1$, $y_2$, $y_3$, $y_j$, $y_V$

$\mathbf{W}_{V \times N} = \{w_{ki}\}$     $\mathbf{W}'_{N \times V} = \{w'_{ij}\}$

The weights between the hidden and the output layer can be represented by a $V \times N$ matrix $W'$

# CBOW: Bag of words context

We assume that there several words considered per context.

Model will predict one target word given many context word

context                target

**Cookie monster eats** ⟶ **cookies**

# Continuous bag-of-words neural network architecture
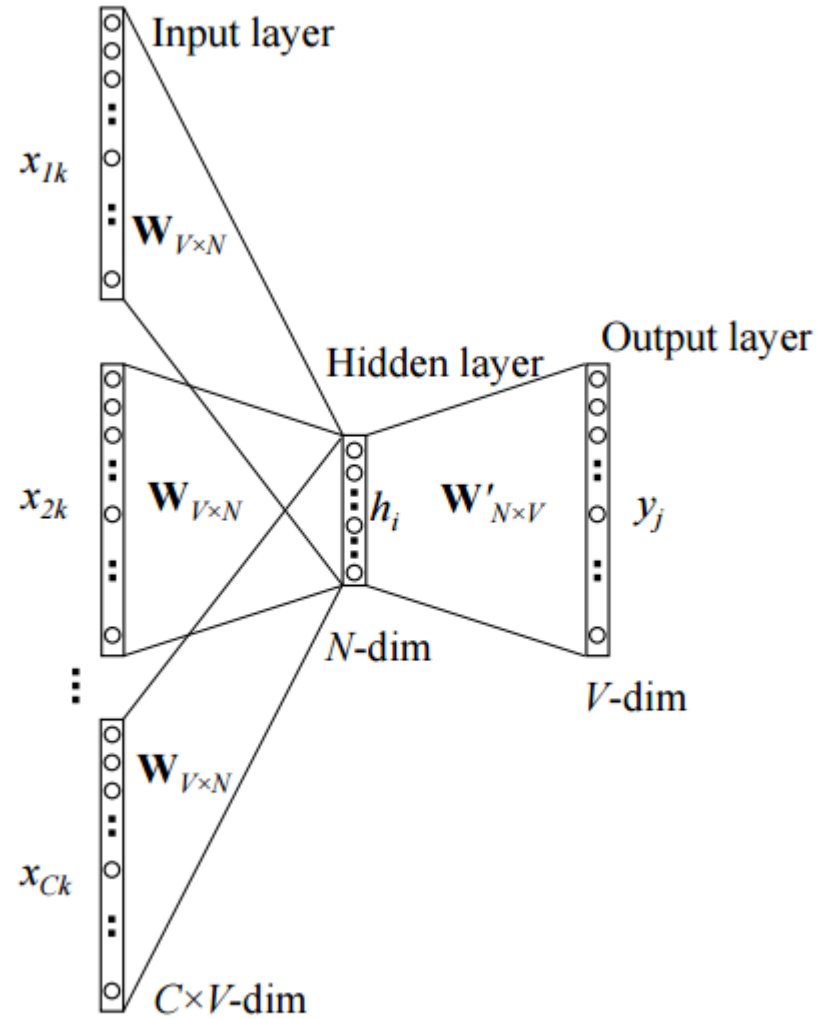


Figure 2: Continuous bag-of-word model

# Skip-Gram Model: One-word target -> context

We assume that there is only one word considered in target and many (window-size) in context.

Model will predict many context words given one target word

target                                  context

**eats**    ⟶    **Cookie monster ___ cookies**

- Skip-gram example:
  - 3-grams:
    - Cookie, monster, cookies
    - Cookie, eats, cookies
    - monster, eats, cookies

# Skip-Gram with Negative Sampling

Idea: don't just use 1's for the words you see in the context also use 0 for the words you **don't see** in the context

target              context

**eats** ⟶ **Frobenius norm eats cookies**

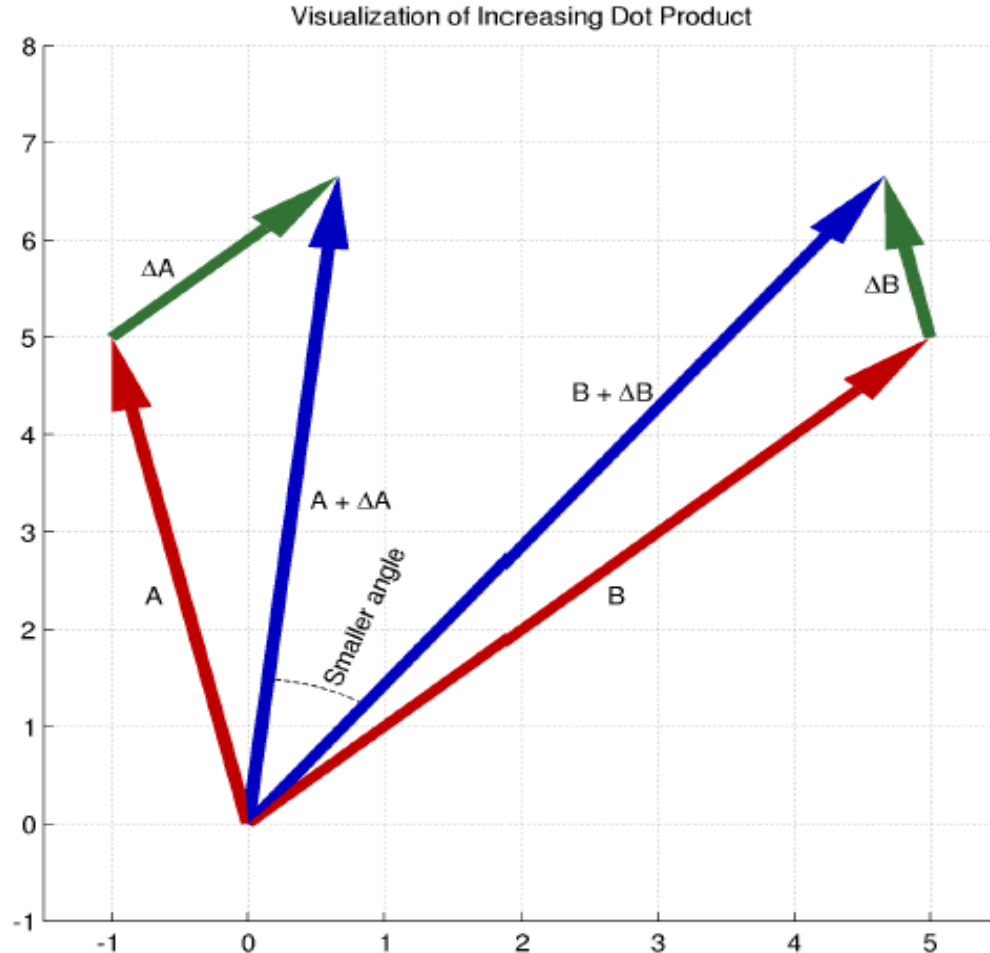$$\|A\|_F = \left( \sum_{i=1}^{m} \sum_{j=1}^{} a_{ij}^2 \right)^{1/2}$$

# Skip-gram overview

- **Each word corresponds to two vectors, one is the target word and the other as a context word.**
- **The algorithm starts with random vectors for the context and zero vectors for the target words.**
- **Each time one word occurs in the context of the other, the context vector and target vector are modified slightly so that their product is slightly larger.**

Dissertation by Eric D. Moyer
"What machines understand about personality words after reading the news"

# Intuition: minimization of context-target dot-product



Visualization of Increasing Dot Product

- Each time one word occurs in the context of the other, the context vector and target vector are modified slightly so that their dot product is slightly larger.

- the dot product between a word and a target will be proportional to the probability of that word appearing in the context of that target when compared to the dot products of other words also appearing in that context.

- This means that words will cluster around their most common contexts.

# Training skip-gram model on a tiny example

Corpus: Mary had a little lamb, little lamb little, little lamb Mary had a little lamb, little lamb little, little lamb Mary had a little lamb, little lamb little, little lamb Mary had a little lamb, little lamb little, little lamb

Vocabulary (V=5): Mary, had, a, little, lamb

Dimensionality: N = 2

Context window size: w = 1

# Training skip-gram model

Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb

Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb

Mary had a little lamb, little lamb, little lamb, little lamb
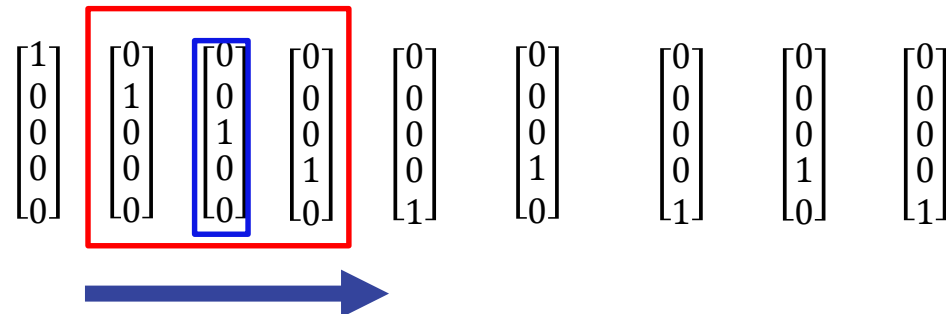Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
Mary had a little lamb, little lamb, little lamb, little lamb
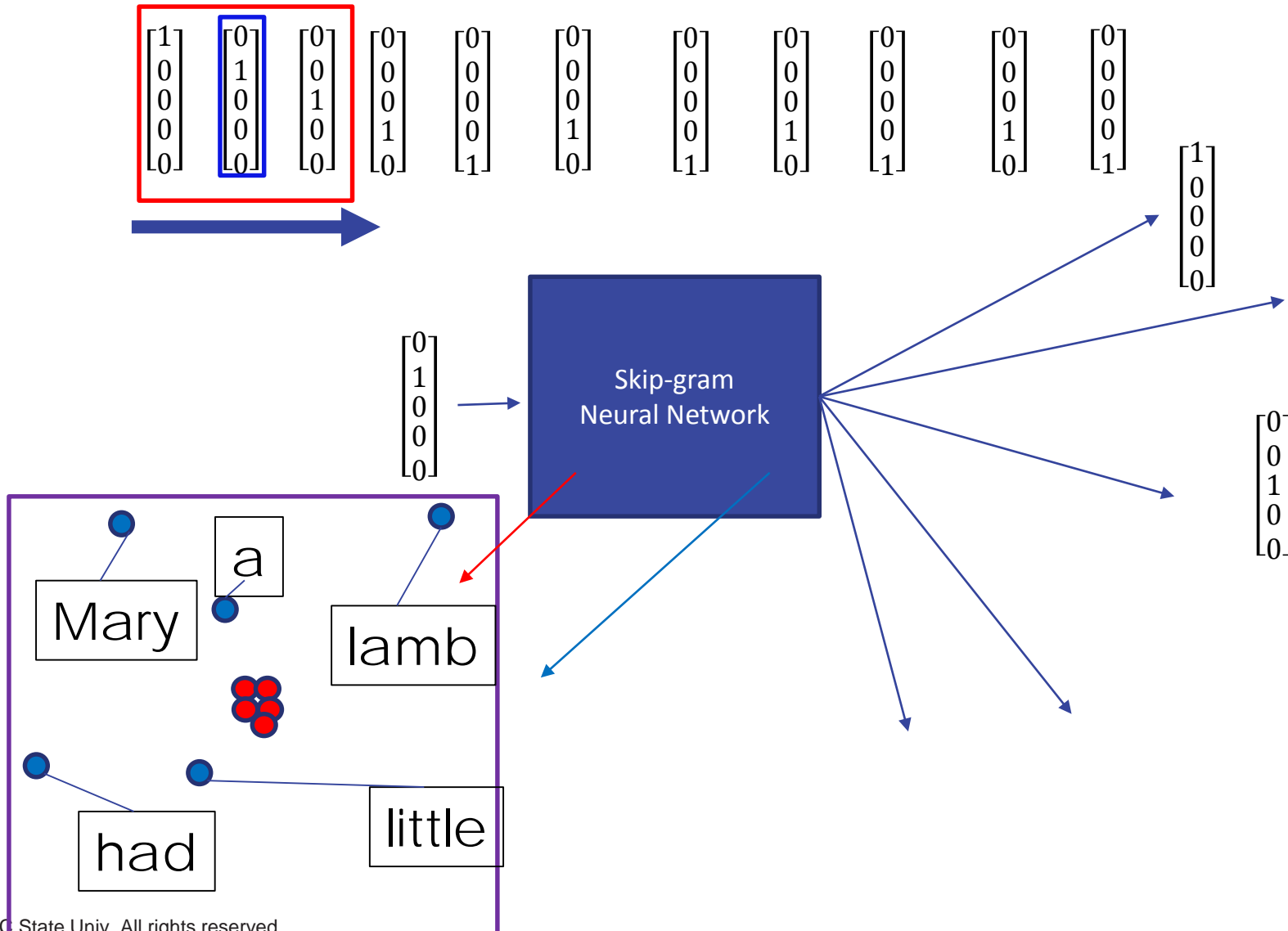
context

target

# One-hot encoding

| Index | Word | One-hot Encoding vector |
|-------|------|-------------------------|
| 0 | Mary | [1,0,0,0,0] |
| 1 | had | [0,1,0,0,0] |
| 2 | a | [0,0,1,0,0] |
| 3 | little | [0,0,0,1,0] |
| 4 | lamb | [0,0,0,0,1] |

## Mary had a little lamb, little lamb, little lamb, little lamb

$$\begin{bmatrix}1\\0\\0\\0\\0\end{bmatrix} \begin{bmatrix}0\\1\\0\\0\\0\end{bmatrix} \begin{bmatrix}0\\0\\1\\0\\0\end{bmatrix} \begin{bmatrix}0\\0\\0\\1\\0\end{bmatrix} \begin{bmatrix}0\\0\\0\\0\\1\end{bmatrix} \begin{bmatrix}0\\0\\0\\1\\0\end{bmatrix} \begin{bmatrix}0\\0\\0\\0\\1\end{bmatrix} \begin{bmatrix}0\\0\\0\\1\\0\end{bmatrix} \begin{bmatrix}0\\0\\0\\0\\1\end{bmatrix}$$
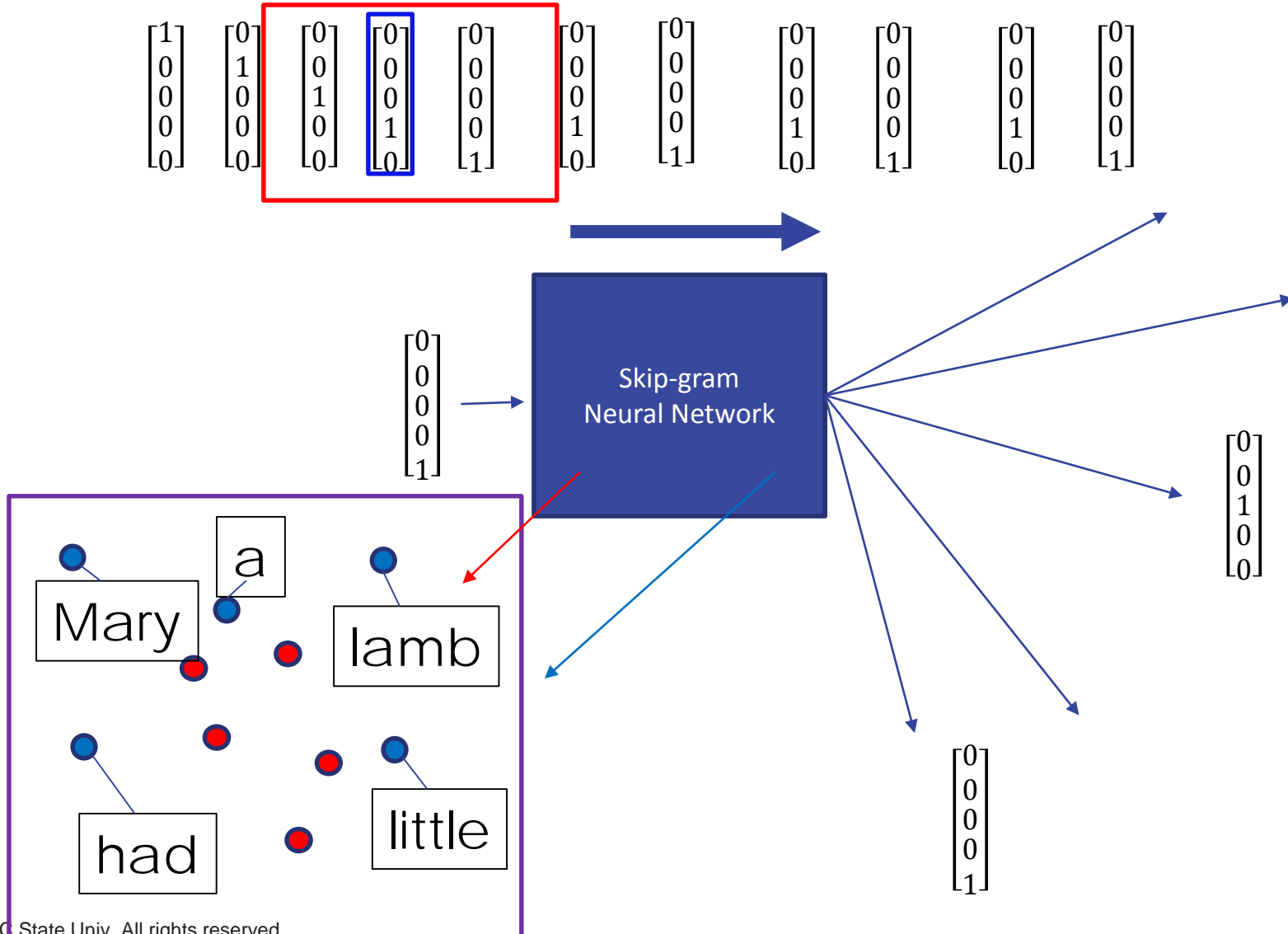
# Initialized weights in neural net and start

Mary had a little lamb, little lamb, little lamb, little lamb
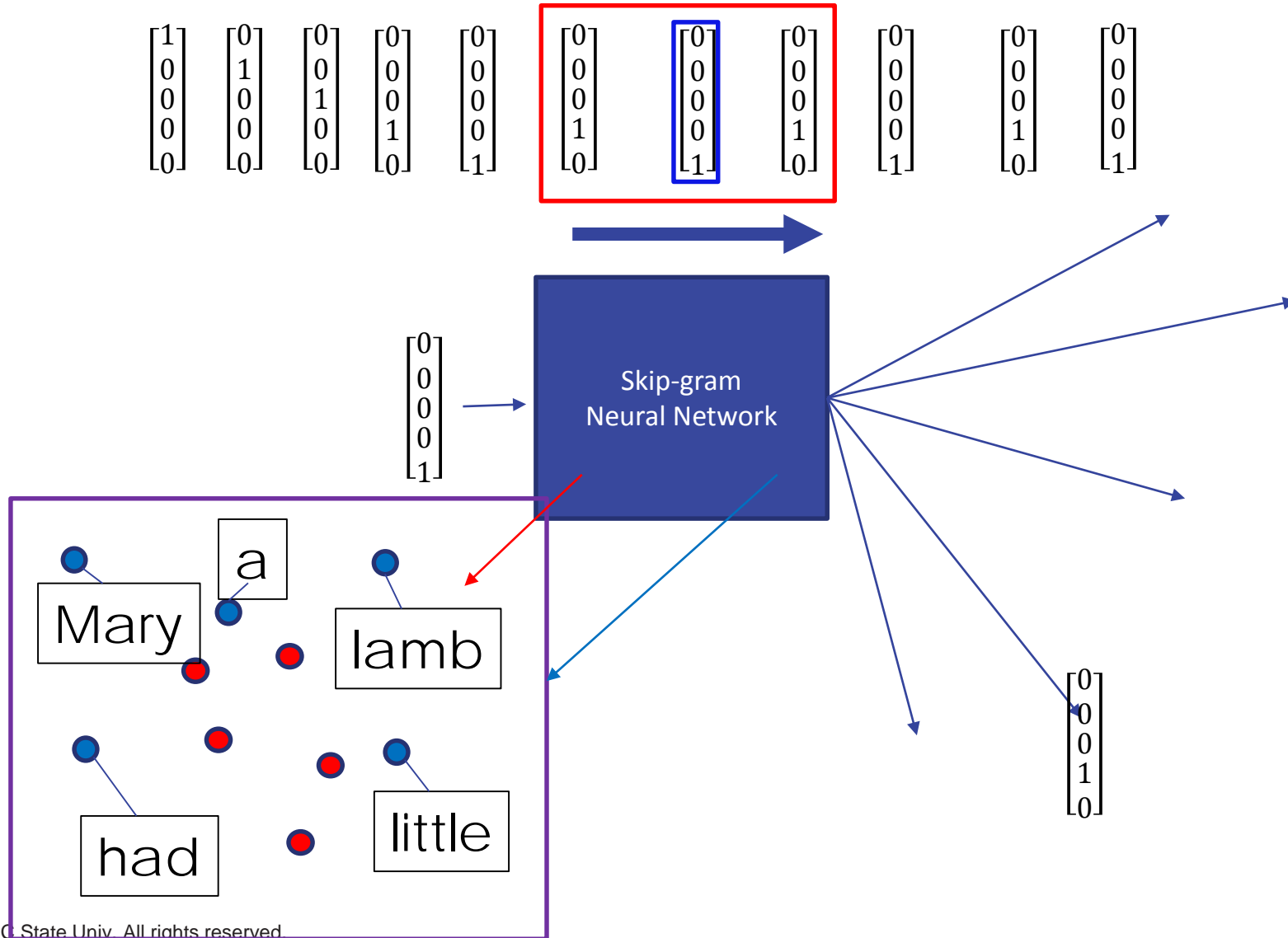
# Training neural network

Mary had a little lamb, little lamb, little lamb, little lamb
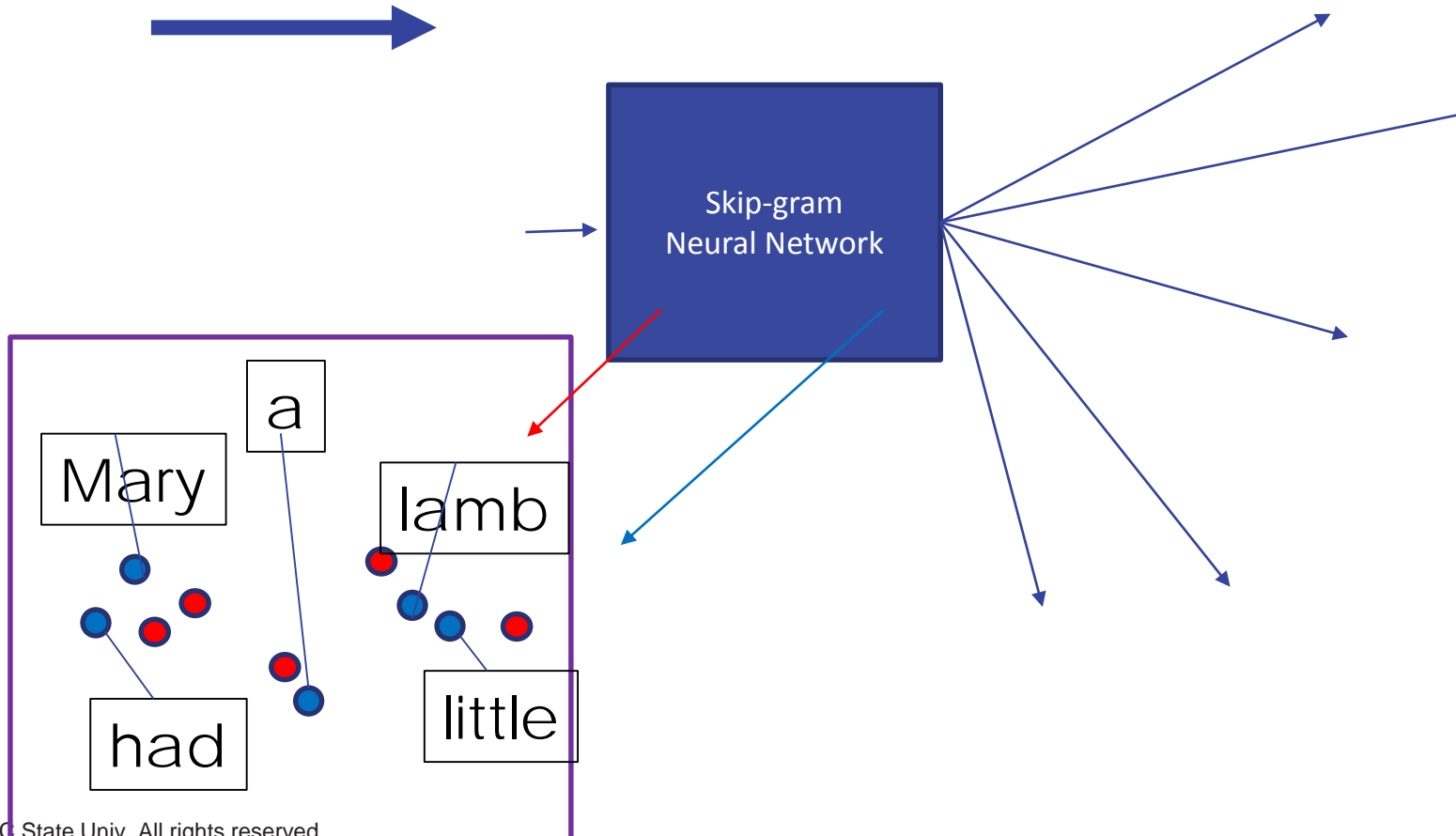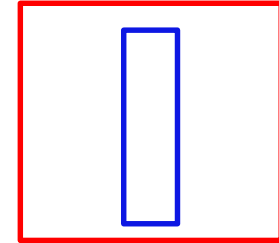
# Training neural network

Mary had a little lamb, little lamb, little lamb, little lamb

# Finished training neural network

Mary had a little lamb, little lamb, little lamb, little lamb

# Skip-gram model neural net

- The **target** word is in the **input** layer.
- The **context** words are in the **output** layer

$v_{w_I}$ is the input vector of the only word on the input layer

For a given input word, and assuming its index in the vocabulary $k$, one-hot input vector representation is $x$.

$$x_i = \begin{cases} 1 \; for \; i = k \\ 0 \; for \; i \neq k \end{cases}$$

$$v_{w_I} := \boldsymbol{h} = \boldsymbol{x}^T \boldsymbol{W} = \boldsymbol{W}_{(k,\cdot)}$$



Figure 3: The skip-gram model.

It is the $k$-th row of $\boldsymbol{W}$, which is the input->hidden weight matrix
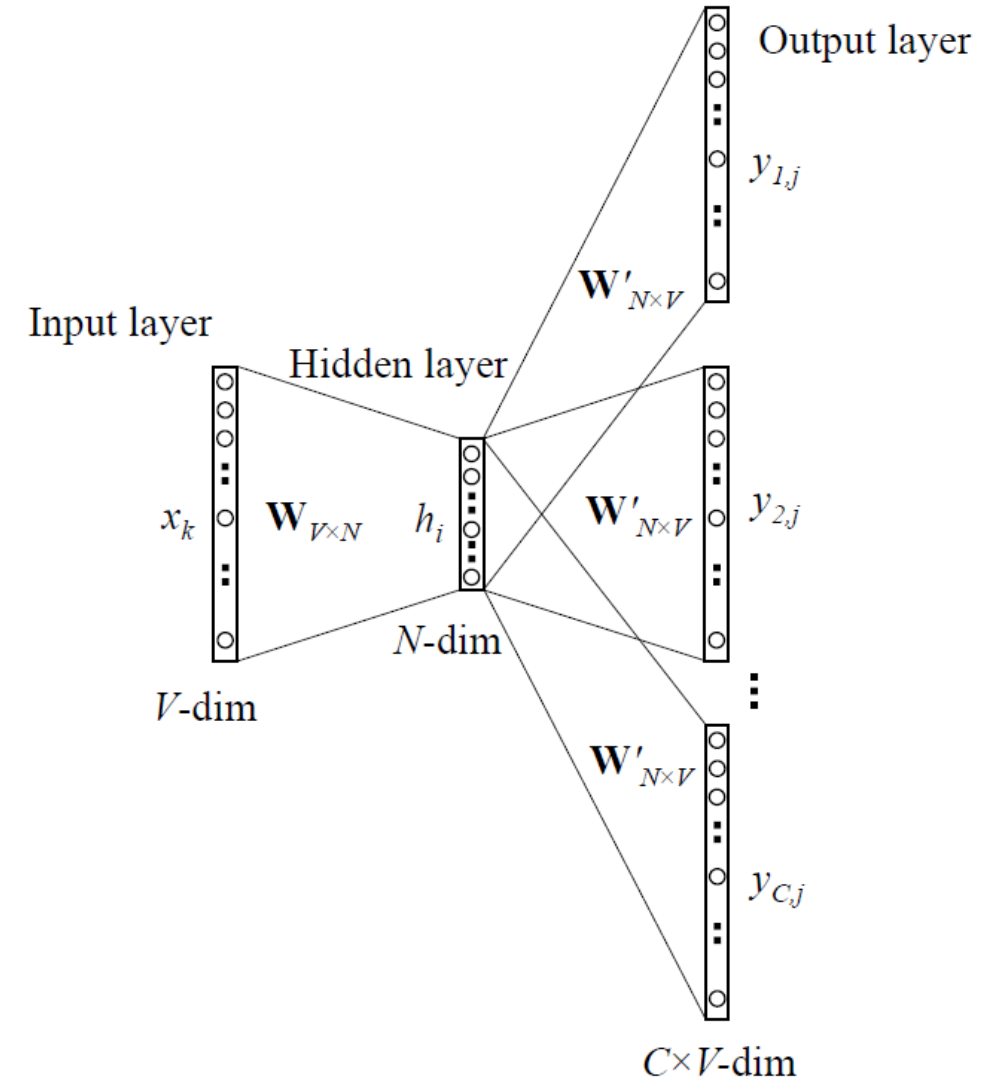
# Skip-gram model:  details

One the output layer we output $C$ multinomial distributions

Each output computed using the same hidden->output weight matrix $W'$ and softmax link function:

$$y_{c,j} = \frac{\exp(u_{c,j})}{\sum_{j'=1}^{V} \exp(u'_j)}$$

- Also note that we hope that:
- $y_{c,j} = p(w_{c,j} = w_{O,c}|w_I)$
- $w_{c,j}$ is the $j$-th word on the $c$-th panel of the output layer
- $w_{O,c}$ is the actual $c$-th word in the output context words
- conditional probability of observing context word given a target word

Input layer

Hidden layer

Output layer

$x_k$  $\mathbf{W}_{V \times N}$  $h_i$  $\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$y_{1,j}$

$y_{2,j}$

$y_{C,j}$

$V$-dim

$N$-dim

$C \times V$-dim

# Skip-gram model: output

$$y_{c,j} = \frac{\exp(u_{c,j})}{\sum_{j'=1}^{V} \exp(u'_j)}$$

$$u_j = v'^T_{w_j} \cdot h \text{ for } c = 1,2,\dots,C$$

- $v'^T_{w_j}$ is the output vector of the $j$-th word in the vocabulary, $w_j$
- $v'^T_{w_j}$ is a column from the hidden->output matrix, $W'$



Input layer

Hidden layer

Output layer

$x_k$   $\mathbf{W}_{V \times N}$   $h_i$

$V$-dim

$N$-dim

$\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$\mathbf{W}'_{N \times V}$

$y_{1,j}$

$y_{2,j}$

$y_{C,j}$

$C \times V$-dim

# Python implementation of word2vec

- https://radimrehurek.com/**gensim**/models/word2vec.html
  https://rare-technologies.com/word2vec-tutorial/

```
>>> model = Word2Vec(sentences, size=100, window=5, min_count=5, workers=4)
```

```
>>> model.most_similar(positive=['woman', 'king'], negative=['man'])
[('queen', 0.50882536), ...]

>>> model.doesnt_match("breakfast cereal dinner lunch".split())
'cereal'

>>> model.similarity('woman', 'man')
0.73723527

>>> model['computer']  # raw numpy vector of a word
array([-0.00449447, -0.00310097,  0.02421786, ...], dtype=float32)
```

# Doc2Vec: Distributed Representations of Sentences and Documents

- Doc2Vec has 2 models:
  - Paragraph Vector Distributed Bag of Words (PV-DBOW)
  - Paragraph Vector Distributed Memory (PV-DM)

**https://cs.stanford.edu/~quocle/paragraph_vector.pdf**

**Distributed Representations of Sentences and Documents**

**Quoc Le**                                                QVL@GOOGLE.COM
**Tomas Mikolov**                                          TMIKOLOV@GOOGLE.COM
Google Inc, 1600 Amphitheatre Parkway, Mountain View, CA 94043
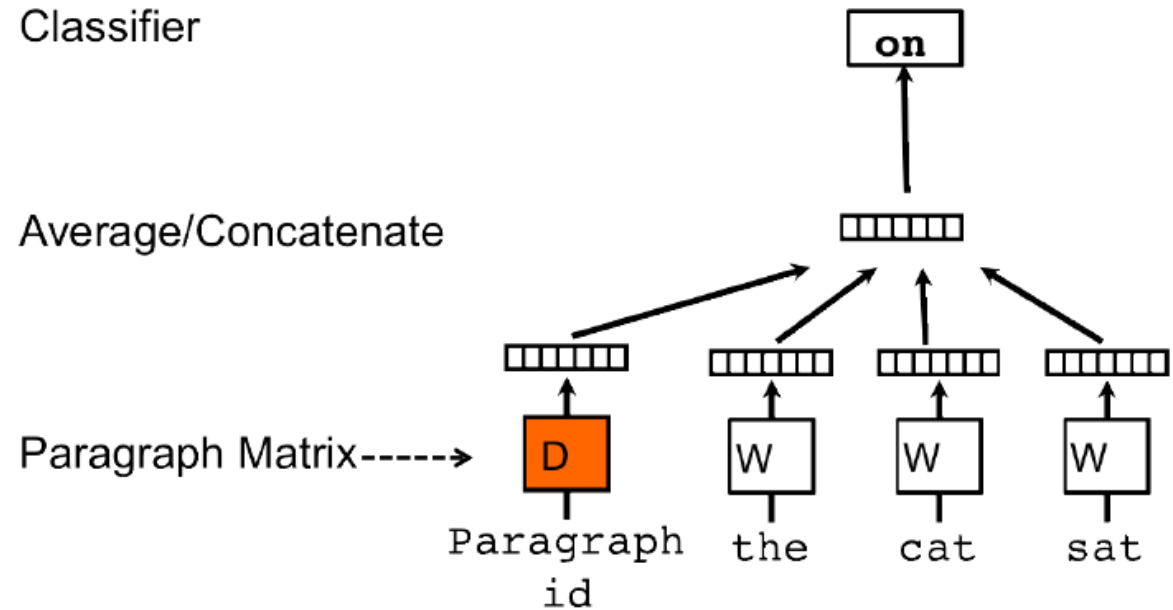
# Doc2Vec

- "In our model, the vector representation is trained to be useful for predicting words in a paragraph.
- ...we concatenate the paragraph vector with several word vectors from a paragraph and predict the following word in the given context.
- Both word vectors and paragraph vectors are trained by the stochastic gradient descent and backpropagation.
- While paragraph vectors are unique among paragraphs, the word vectors are shared.
- At prediction time, the paragraph vectors are inferred by fixing the word vectors and training the new paragraph vector until convergence."

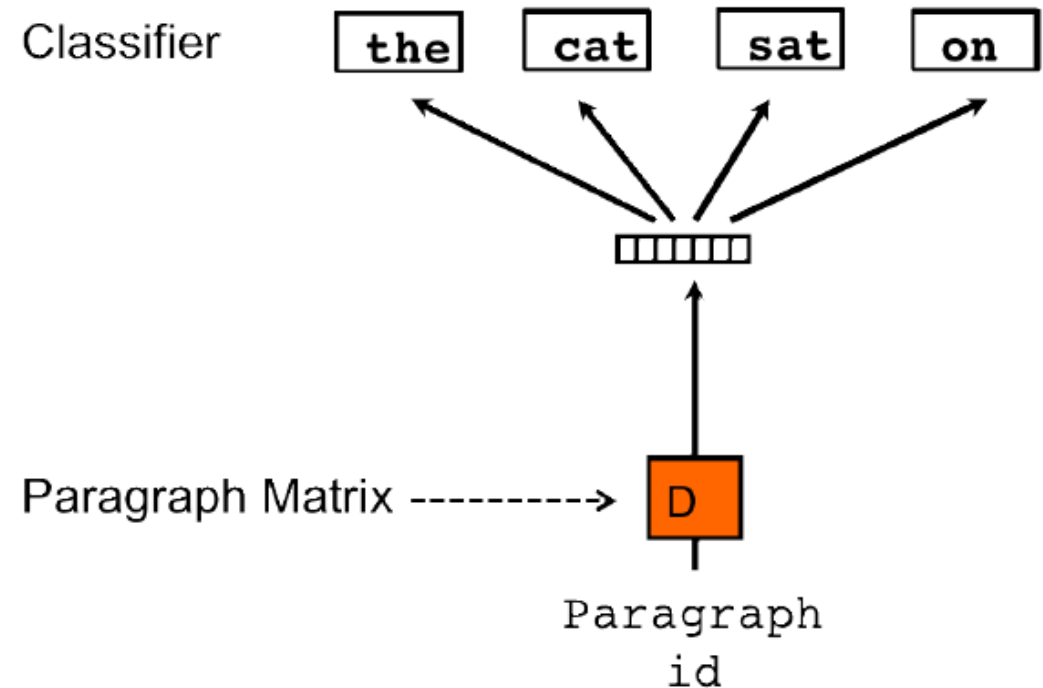# Paragraph vector: distributed memory

- every **paragraph** is mapped to a unique **vector**, represented by a column in matrix D
- and every **word** is also mapped to a unique **vector**, represented by a column in matrix W.
- The paragraph vector and word vectors are averaged or **concatenated** to predict the next word in a context.

# Paragraph vector: distributed bag of words

- …ignore the context words in the input, but force the model to predict words randomly sampled from the paragraph in the output.
- In reality, what this means is that at each iteration of stochastic gradient descent, we sample a text window, then sample a random word from the text window and form a classification task given the Paragraph Vector.

Classifier: the | cat | sat | on

Paragraph Matrix --------→ D

Paragraph id

# Doc2Vec usage with gensim package in Python

https://radimrehurek.com/gensim/models/doc2v

```
>>> model = Doc2Vec(documents, size=100, window=8, min_count=5, workers=4)
```

```
>>> trained_model.n_similarity(['sushi', 'shop'], ['japanese', 'restaurant'])
0.615404665561049689

>>> trained_model.n_similarity(['restaurant', 'japanese'], ['japanese', 'restaurant'])
1.0000000000000004
```

- class gensim.models.doc2vec.Doc2Vec(documents=None, size=300, alpha=0.025, window=8, min_count=5, max_vocab_size=None, sample=0, seed=1, workers=1, min_alpha=0.0001, dm=1, hs=1, negative=0, dbow_words=0, dm_mean=0, dm_concat=0, dm_tag_count=1, docvecs=None, docvecs_mapfile=None, comment=None, trim_rule=None, **kwargs)

# Project: Sentiment analysis with Doc2Vec

- **perform sentiment analysis over IMDB movie reviews and Twitter data.**
- **build BOW, Doc2Vec, Word2Vec models with labeled training and testing data to evaluate the model.**
- **classify tweets or movie reviews as either positive or negative given labeled training For classification, you will experiment with logistic regression as well as a Naive Bayes classifier from python's well-regarded machine learning package scikit-learn.**

# Project steps

## 1. Load datasets

```
34  def main():
35      train_pos, train_neg, test_pos, test_neg = load_data(path_to_data)
36
```

## 2. Train Doc2Vec model

```
62      if method == "d2v":
63          train_pos_vec, train_neg_vec, test_pos_vec, test_neg_vec =
            feature_vecs_DOC(train_pos, train_neg, test_pos, test_neg)
64          #filename = './'+path_to_data+'train_pos_vec_d2v.txt'
```

## 3. Train logistic regression model with Doc2Vec vector-features

```
77          nb_model, lr_model = build_models_DOC(train_pos_vec,
            train_neg_vec)
```

## 4. Evaluate the model on test data

```
113     print(                    )
114     evaluate_model(lr_model, test_pos_vec, test_neg_vec, True)
115
```

# Training Doc2Vec steps:

## Step 1. Create TaggedDocument-s from lists of words

```
197    labeled_train_pos = [TaggedDocument(words, ["TRAIN_POS_" + str(i)])
       for i, words in enumerate(train_pos)]
198    labeled_train_neg = [TaggedDocument(words, ["TRAIN_NEG_" + str(i)])
       for i, words in enumerate(train_neg)]
199    labeled_test_pos = [TaggedDocument(words, ["TEST_POS_" + str(i)])
       for i, words in enumerate(test_pos)]
200    labeled_test_neg = [TaggedDocument(words, ["TEST_NEG_" + str(i)])
       for i, words in enumerate(test_neg)]
```

## Step 2. Initialize the model

```
203    model = Doc2Vec(min_count=1, window=10, size=100, sample=1e-4,
       negative=5, workers=4)
```

```
205    sentences = labeled_train_pos + labeled_train_neg + labeled_test_pos
        + labeled_test_neg
206    model.build_vocab(sentences)
```

## Step 3. Train the model

```
210    for i in range(5):
211        print("Training iteration %d" % (i))
212        random.shuffle(sentences)
213        model.train(sentences,total_examples=model.corpus_count, epochs=
           model.iter)
```

# Training Doc2Vec steps:

## Step 4. Extract vectors

```python
216    # Use the docvecs function to extract the feature vectors for the
       training and test data
217    train_pos_vec = [model.docvecs["TRAIN_POS_" + str(i)] for i in range
       (len(labeled_train_pos))]
218    train_neg_vec = [model.docvecs["TRAIN_NEG_" + str(i)] for i in range
       (len(labeled_train_neg))]
219    test_pos_vec = [model.docvecs["TEST_POS_" + str(i)] for i in range(
       len(labeled_test_pos))]
220    test_neg_vec = [model.docvecs["TEST_NEG_" + str(i)] for i in range(
       len(labeled_test_neg))]
221
222    # Return the four feature vectors
223    return train_pos_vec, train_neg_vec, test_pos_vec, test_neg_vec
224
```

# Build logistic regression model with Doc2Vec vectors

```python
342  def build_models_DOC(train_pos_vec, train_neg_vec):
343      """
344      Returns a GaussianNB and LosticRegression Model that are fit to the
         training data.
345      """
346      Y = ["pos"]*len(train_pos_vec) + ["neg"]*len(train_neg_vec)
347
348      # Use sklearn's GaussianNB and LogisticRegression functions to fit
         two models to the training data.
349      # For LogisticRegression, pass no parameters
350      X = train_pos_vec + train_neg_vec
351      nb_model = sklearn.naive_bayes.GaussianNB()
352      nb_model.fit(X, Y)
353      lr_model = sklearn.linear_model.LogisticRegression()
354      lr_model.fit(X, Y)
355      return nb_model, lr_model
```

# Evaluating the model

```
Doc2Vec
Training iteration 0
Training iteration 1
Training iteration 2
Training iteration 3
Training iteration 4
end of training
Naive Bayes
----------
predicted:      pos       neg
actual:
pos             5482      7018
neg             2262      10238
accuracy: 0.628800

Logistic Regression
-------------------
predicted:      pos       neg
actual:
pos             10745     1755
neg             1700      10800
accuracy: 0.861800
```

```python
def evaluate_model(model, test_pos_vec, test_neg_vec, print_confusion=False):
    """
    Prints the confusion matrix and accuracy of the model.
    """
    # Use the predict function and calculate the true/false positives
    and true/false negative.
    pos_results = list(model.predict(test_pos_vec))
    neg_results = list(model.predict(test_neg_vec))
    tp = pos_results.count("pos")
    tn = neg_results.count("neg")
    fn = pos_results.count("neg")
    fp = neg_results.count("pos")
    accuracy = float(tp + tn) / (tp + tn + fp + fn)

    if print_confusion:
        print("predicted:\tpos\tneg")
        print("actual:")
        print("pos\t\t%d\t%d" % (tp, fn))
        print("neg\t\t%d\t%d" % (fp, tn))
    print("accuracy: %f" % (accuracy))
```

**If you would like to write your own implementation of word2vec's neural net …**

# Error function for Skip-gram model

$$E = -\log p\left(w_{O,1}, w_{O,2}, \ldots, w_{O,C} \middle| w_I\right) =$$

$$-\log \prod_{c=1}^{C} \frac{\exp(u_{c,j_c^*})}{\sum_{j'=V}^{V} \exp(u_j')} =$$

$$-\sum_{c=1}^{C} u_j^* + C \cdot \log \sum_{j'=1} \exp(u_{j'})$$

$j_c^*$ **is the index of the actual** $c$**-th output context word in the vocabulary**

**Now we need the gradient for the stochastic gradient descent**

# Prediction errors and gradients

Taking the derivative of $E$ with regard to the net input of every node on every panel of the output layer, $u_{c,j}$ and setting it as prediction error:

Defining $EI_i$ as a sum of prediction errors over all context words:

Taking the derivative of $E$ with regard to the hidden->output matrix $W'$ with elements $w'_{ij}$

Thus we obtain the update equation for the hidden->output matrix $W'$

$$e_{c,j} = \frac{\partial E}{\partial u_{c,j}} = y_{c,j} - t_{c,j}$$

$$EI_j = \sum_{c=1}^{C} e_{c,j}$$

$$\frac{\partial E}{\partial w'_{ij}}$$

$$= \sum_{c=1}^{C} \frac{\partial E}{\partial u_{c,j}} \cdot \frac{\partial u_{c,j}}{\partial w'_{ij}} = EI_j \cdot h_i$$

$${w'_{ij}}^{new} = {w'_{ij}}^{old} - \eta \cdot EI_j \cdot h_i$$

# Update Equations for Skip-gram Model

$$\boldsymbol{v'_{w_j}}^{new} = \boldsymbol{v'_{w_j}}^{old} - \eta \cdot EI_j \cdot \boldsymbol{h} \qquad \text{update of "context" vectors}$$

$$\boldsymbol{v}_{w_I}^{new} = \boldsymbol{v}_{w_I}^{old} - \eta \cdot \boldsymbol{EH} \qquad \text{update of "target" vectors}$$

$$EH_i = \sum_{j=1}^{V} EI_j \cdot w'_{ij}$$

$e_{c,j} := y_{c,j} - t_{c,j}$ is the prediction error on the node

$\boldsymbol{EI} = \{EI_1, \dots, EI_V\}$ is a $V-$dimensional vector, the sum of all prediction errors over all context words

$$EI_j = \sum_{c=1}^{C} e_{c,j}$$