# COMP5046
# Natural Language Processing

Lecture 3: Word Classification and Machine Learning

Semester 1, 2020
School of Computer Science
The University of Sydney, Australia

**Dr Caren Han**
Caren.Han@sydney.edu.au

THE UNIVERSITY OF
SYDNEY

# LECTURE PLAN

**Lecture 3: Word Classification and Machine Learning**

1. **Previous Lecture: Word Embedding Review**
2. Word Embedding Evaluation
3. Deep Neural Network for Natural Language Processing
    1. Perceptron and Neural Network (NN)
    2. Multilayer Perceptron
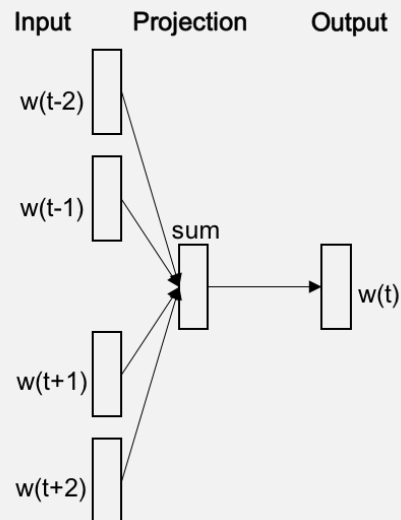    3. Applications
4. Next Week Preview
    See how the Deep Learning can be used for NLP
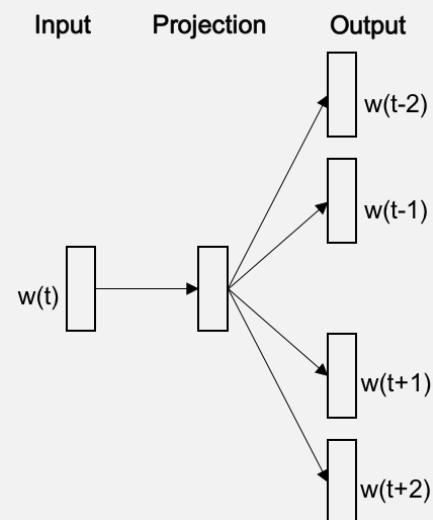    - Text Classification, etc.

# Previous Lecture Review

## Word2Vec Models

| CBOW | Skip-gram |
|---|---|
| *Predict center word from (bag of) context words* | *Predict context words given center word* |

# Previous Lecture Review

## Word2Vec

**Sentence: "Sydney is the state capital of NSW"**

*Using window sliding, develop the training data*

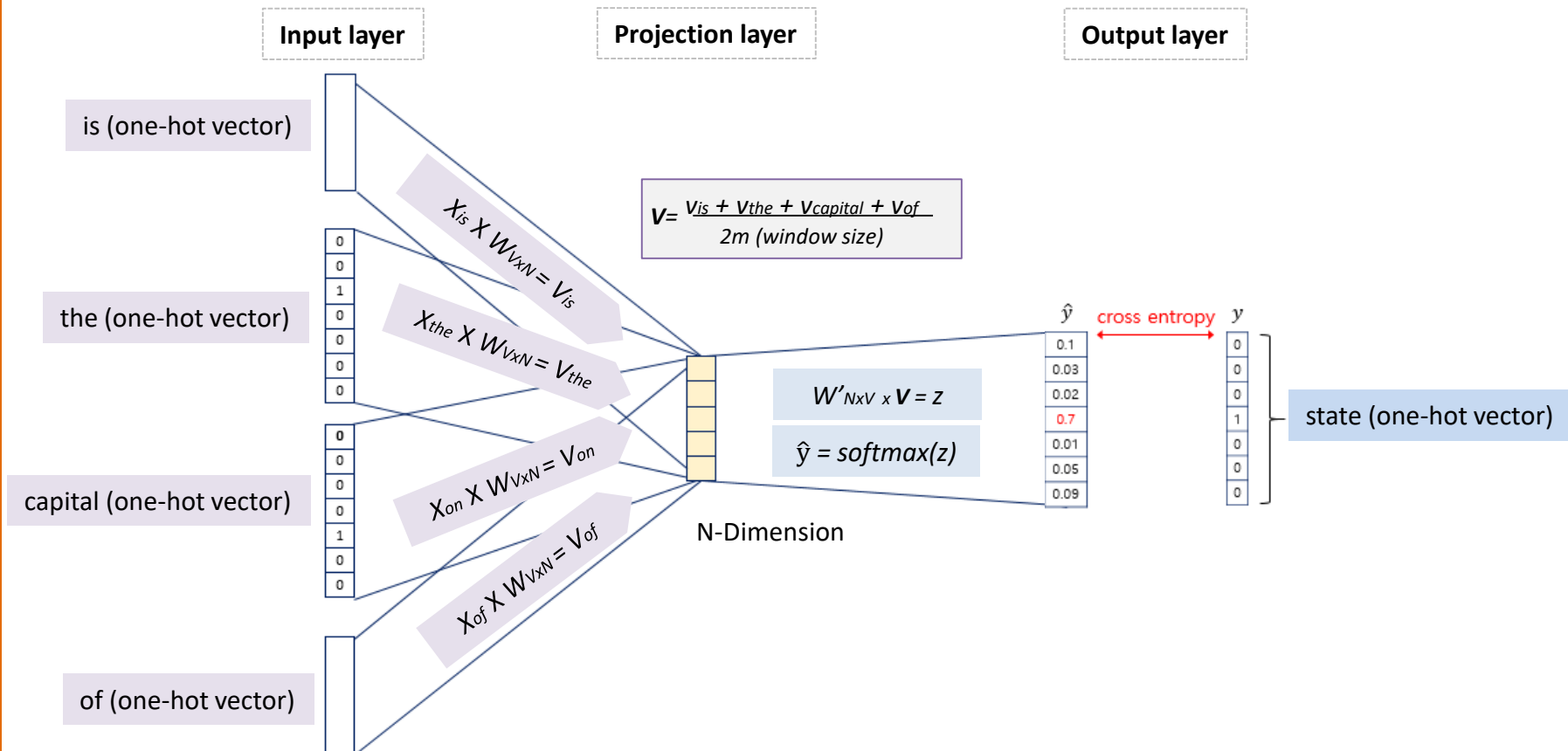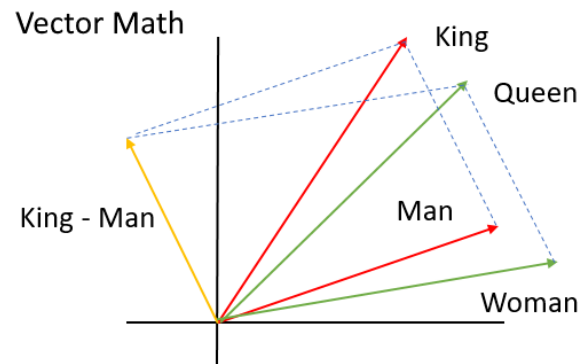| Center word | Context ("outside") word | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| [1,0,0,0,0,0,0] | [0,1,0,0,0,0,0], [0,0,1,0,0,0,0] | Sydney | is | the | state | capital | of | NSW |
| [0,1,0,0,0,0,0] | [1,0,0,0,0,0,0], [0,0,1,0,0,0,0], [0,0,0,1,0,0,0] | Sydney | is | the | state | capital | of | NSW |
| [0,0,1,0,0,0,0] | [1,0,0,0,0,0,0], [0,1,0,0,0,0,0] [0,0,0,1,0,0,0], [0,0,0,0,1,0,0] | Sydney | is | the | state | capital | of | NSW |
| [0,0,0,1,0,0,0] | [0,1,0,0,0,0,0], [0,0,1,0,0,0,0] [0,0,0,0,1,0,0], [0,0,0,0,0,1,0] | Sydney | is | the | state | capital | of | NSW |
| [0,0,0,0,1,0,0] | [0,0,1,0,0,0,0], [0,0,0,1,0,0,0] [0,0,0,0,0,1,0], [0,0,0,0,0,0,1] | Sydney | is | the | state | capital | of | NSW |
| [0,0,0,0,0,1,0] | [0,0,0,1,0,0,0], [0,0,0,0,1,0,0] [0,0,0,0,0,0,1] | Sydney | is | the | state | capital | of | NSW |
| [0,0,0,0,0,0,1] | [0,0,0,0,1,0,0], [0,0,0,0,0,1,0] | Sydney | is | the | state | capital | of | NSW |

■ Center word

■ Context ("outside") word

# Previous Lecture Review

## CBOW Process

Predict center word from (bag of) context words

**Sentence: "Sydney is the state capital of NSW"**



| Input layer | Projection layer | Output layer |

is (one-hot vector)

the (one-hot vector)

capital (one-hot vector)

of (one-hot vector)

$X_{is} \times W_{VxN} = V_{is}$

$X_{the} \times W_{VxN} = V_{the}$

$X_{on} \times W_{VxN} = V_{on}$

$X_{of} \times W_{VxN} = V_{of}$

$$V = \frac{V_{is} + V_{the} + V_{capital} + V_{of}}{2m \ (window \ size)}$$

$W'_{NxV} \times V = z$

$\hat{y} = softmax(z)$

N-Dimension

$\hat{y}$  cross entropy  $y$

| $\hat{y}$ |
|---|
| 0.1 |
| 0.03 |
| 0.02 |
| 0.7 |
| 0.01 |
| 0.05 |
| 0.09 |

| $y$ |
|---|
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |

state (one-hot vector)

**LECTURE PLAN**

**Lecture 3: Word Classification and Machine Learning**

1. Previous Lecture: Word Embedding Review
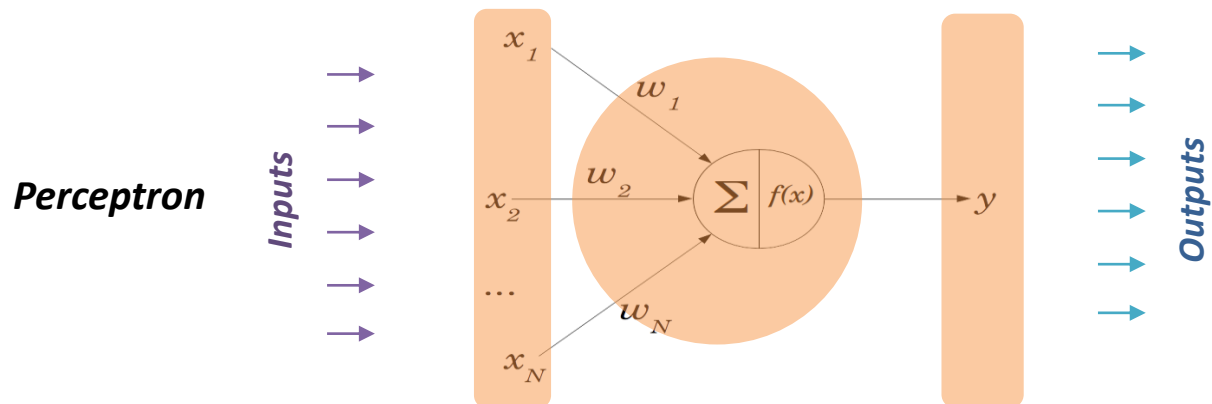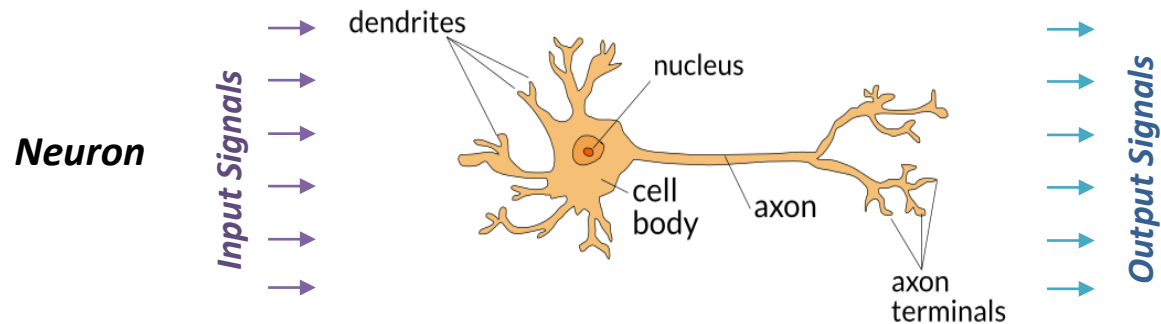2. **Word Embedding Evaluation**
3. Deep Neural Network for Natural Language Processing
   1. Perceptron and Neural Network (NN)
   2. Multilayer Perceptron
   3. Applications
4. Next Week Preview
   See how the Deep Learning can be used for NLP
   - Text Classification, etc.

# Word Embedding Evaluation

## How to evaluate word vectors?

| Type | How to work / Benefit |
|---|---|
| Intrinsic | Evaluation on a specific/intermediate subtask |
| | • Fast to compute<br>• Helps to understand that system<br>• Not clear if really helpful unless correlation to real task is established |
| Extrinsic | Evaluation on a real task |
| | • Can take a long time to compute accuracy<br>• Unclear if the subsystem is the problem or its interaction or other subsystems |

# Word Embedding Evaluation

## Intrinsic word vector evaluation

**Word Vector Analogies**

$$a \longleftrightarrow b \quad :: \quad c \longleftrightarrow ???$$
$$man \longleftrightarrow women :: king \longleftrightarrow ???$$

- Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy questions

# Word Embedding Evaluation

## Intrinsic word vector evaluation

**Word Vector Analogies**

**King – Man + Woman = ?**

| No | Dataset | Type | Result |
|----|---------|------|--------|
| 1 | TED Script | word2vec CBOW | President |
| 2 | | word2vec Skip-gram | Luther |
| 3 | | fastText CBOW | Kidding |
| 4 | | fastText Skip-gram | Jarring |
| 5 | Google News | word2vec CBOW | queen |
| 6 | | word2vec Skip-gram | queen |

# Word Embedding Evaluation

## Intrinsic word vector evaluation

### Evaluation Result Comparison

The Semantic-Syntatic word relationship tests for understanding of a wide variety of relationships as shown below. Using 640-dimensional word vectors, a skip-gram trained model achieved 55% semantic accuracy and 59% syntatic accuracy.

Table 3: *Comparison of architectures using models trained on the same data, with 640-dimensional word vectors. The accuracies are reported on our Semantic-Syntactic Word Relationship test set, and on the syntactic relationship test set of [20]*

| Model Architecture | Semantic-Syntactic Word Relationship test set | | MSR Word Relatedness Test Set [20] |
|---|---|---|---|
| | Semantic Accuracy [%] | Syntactic Accuracy [%] | |
| RNNLM | 9 | 36 | 35 |
| NNLM | 23 | 53 | 47 |
| CBOW | 24 | 64 | 61 |
| Skip-gram | 55 | 59 | 56 |

# Previous Lecture Review

## CBOW Process

Predict center word from (bag of) context words

**Sentence: "Sydney is the state capital of NSW"**

| Input layer | Projection layer | Output layer |
| --- | --- | --- |

is (one-hot vector)

the (one-hot vector)

capital (one-hot vector)

of (one-hot vector)

$X_{is} \times W_{VxN} = V_{is}$

$X_{the} \times W_{VxN} = V_{the}$

$X_{on} \times W_{VxN} = V_{on}$

$X_{of} \times W_{VxN} = V_{of}$

$$V = \frac{V_{is} + V_{the} + V_{capital} + V_{of}}{2m \ (window\ size)}$$

$W'_{NxV} \times V = z$

$\hat{y} = softmax(z)$

N-Dimension

$\hat{y}$  ←cross entropy→  $y$

| $\hat{y}$ |
| --- |
| 0.1 |
| 0.03 |
| 0.02 |
| 0.7 |
| 0.01 |
| 0.05 |
| 0.09 |

| $y$ |
| --- |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |

state (one-hot vector)

- *How does this weight work?*
- *What is the Softmax?*
- *What is the Cross Entropy?*

THE UNIVERSITY OF SYDNEY

**Lecture 3: Word Classification and Machine Learning**

1. Previous Lecture: Word Embedding Review
2. Word Embedding Evaluation
3. **Deep Neural Network for Natural Language Processing**
   1. Perceptron and Neural Network (NN)
   2. Multilayer Perceptron
   3. Applications
4. Next Week Preview
   See how the Deep Learning can be used for NLP
   - Text Classification, etc.

# Deep Learning for NLP

## Deep Learning with Neural Network

### Neuron and Perceptron

# Deep Learning for NLP

## Deep Learning with Neural Network

**Inputs and Outputs (Labels) for Natural Language Processing**

| $x_i$ | Inputs | **Features**<br>words (indices or vectors!), context windows, sentences, documents, etc. |
|---|---|---|
| $y_i$ | Outputs (labels) | **What we try to predict/classify**<br>• E.g. word meaning, sentiment, name entity |

*Perceptron*

*Inputs*

$x_1$

$w_1$

$x_2$

$w_2$

$\Sigma$ $f(x)$

$y$

…

$w_N$

$x_N$

*Outputs*

## Deep Learning with Neural Network

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

Lisa, give me an apple. I will give you three bananas then!

Input (x)

1 🍎

3 🍌
Output (y)

# Deep Learning for NLP

## Deep Learning with Neural Network - Model

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

$y = 3x$

Input (x)

1 🍎

3 🍌

Output (y)

Weight

**3**

Input (x)

Output (y)

$$y = Wx$$

**Deep Learning for NLP**

## Deep Learning with Neural Network - Model

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

Guess how much I will give you back!

# Deep Learning for NLP

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

Guess how much I will give you back!

1 🍎

0 🍌

5 🍎

16 🍌

6 🍎

20 🍌

3 🍎

? 🍌

$$y = W x$$

*What is **W** then?*

# Deep Learning for NLP

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

Guess how much I will give you back!

Data
1 🍎
0 🍌
5 🍎
16 🍌
6 🍎
20 🍌

$$y = Wx$$

*What is **w** then?*

# Deep Learning for NLP

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

**Data**

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

$$y = Wx$$

*What is **W** then?*

# Deep Learning for NLP

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*



| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

$$y = W x$$

*What is **W** then?*

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

$$y = W x$$

*What if W is **3**?*

  3   =  3 x 1
  15  =  3 x 5
  20  =  3 x 6

**Deep Learning for NLP**

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

*weight*   *bias*

$$y = Wx + b$$

*Weight is not enough…*

**Deep Learning with Neural Network - Parameter**

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

weight      bias

$$y = W x + b$$

*How can we find the parameters, **w** and **b**?*

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

$$y = Wx + b$$

weight — $W$ ; bias — $b$

*How can we find the parameters, w and b?*

# Deep Learning for NLP

## Deep Learning with Neural Network - Parameter

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

| x 🍎 | y 🍌 |
|------|------|
| 1    | 0    |
| 5    | 16   |
| 6    | 20   |

*Model* $y = Wx + b$

weight — $W$
bias — $b$

*How can we find the parameters, $w$ and $b$?*

**Deep Learning with Neural Network - Cost**

*Actual Data*

weight        bias

$$y = ?\ x + ?$$

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

Weight    Bias

Input (x)    ?    ?    Output (y)

*Model Ex#1*

weight        bias

$$\hat{y} = 1\ x + 0$$

|  | predicted | actual |
|---|---|---|
| x 🍎 | ŷ 🍌 | y 🍌 |
| 1 | 1 | 0 |
| 5 | 5 | 16 |
| 6 | 6 | 20 |

*Model Ex#2*

weight        bias

$$\hat{y} = 2\ x + 2$$

|  | predicted | actual |
|---|---|---|
| x 🍎 | ŷ 🍌 | y 🍌 |
| 1 | 4 | 0 |
| 5 | 12 | 16 |
| 6 | 14 | 20 |

🤔 *Which one is closer?*

## Deep Learning with Neural Network - Cost

### Actual Data

weight     bias

$$y = \boxed{?}\ x + \boxed{?}$$

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

**Weight**    **Bias**

[?] — [?] →

Input (x)      Output (y)

### Model Ex#1

weight     bias

$$\hat{y} = \boldsymbol{1}\ x + \boldsymbol{0}$$

| | predicted | actual | cost |
|---|---|---|---|
| x 🍎 | $\hat{y}$ 🍌 | y 🍌 | $(y-\hat{y})^2$ |
| 1 | 1 | 0 | |
| 5 | 5 | 16 | |
| 6 | 6 | 20 | |

### Model Ex#2

weight     bias

$$\hat{y} = \boldsymbol{2}\ x + \boldsymbol{2}$$

| | predicted | actual | cost |
|---|---|---|---|
| x 🍎 | $\hat{y}$ 🍌 | y 🍌 | $(y-\hat{y})^2$ |
| 1 | 4 | 0 | |
| 5 | 12 | 16 | |
| 6 | 14 | 20 | |

😎 *Let's calculate the cost(loss)!*

**Mean Squared Error (MSE)**   $C(w,b) = \sum(y_n - \hat{y}_n)$

$$n \in \{0,1,2\}$$

# Deep Learning for NLP

## WAIT! Loss Function? Cost Calculation?

**WAIT! Loss Function? Cost Calculation?**



*1) Mean Squared Error (MSE): measures the average of the squares of the errors*

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \tilde{y}_i)^2$$

*2) Cross Entropy: calculating the difference between two probability distributions*

$$L_{\text{cross-entropy}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_i y_i \log(\hat{y}_i)$$

*softmax*

**WAIT! Softmax!!**

**Neural Network**

**Input**

*1. Define N Class*

*5. Score given classes*

*2. Construct Network*

*4. Classification*

*3. Input*

# Deep Learning for NLP

## WAIT! Softmax!!



nice   food   state   big

1. Define N Class

5. Score given classes

2. Construct Network

4. Classification

Input

3. Input

is        the        capital        of

# Deep Learning for NLP

**WAIT! Softmax!!**

# Deep Learning for NLP

## WAIT! Softmax!!

*Depends on the dataset*

| Problem type | Last-layer activation | Loss (Cost) function | Example |
|---|---|---|---|
| Binary classification | sigmoid | Binary Cross Entropy | Sentiment analysis (Positive/Negative) |
| Multi-class, single-label classification | softmax | Categorical Cross Entropy | Part-of-Speech tagging Named Entity Recognition |
| Multi-class, multi-label classification | sigmoid | Binary Cross Entropy | Multi-topic classification, one can have multiple topics |
| Regression to arbitrary values | None | MSE (Mean Squared Error) | Predict house price |
| Regression to values between 0 and 1 | sigmoid | MSE or Binary Cross Entropy | Engine health assessment where 0 is broken, 1 is new |

## Deep Learning with Neural Network - Cost

**Actual Data**

weight        bias

$$y = ?\ x + ?$$

| x 🍎 | y 🍌 |
|------|------|
| 1    | 0    |
| 5    | 16   |
| 6    | 20   |



Weight    Bias

? ?

Input (x)        Output (y)

**Model Ex#1**

weight        bias

$$\hat{y} = 1\ x + 0$$

|      | predicted | actual | cost |
|------|-----------|--------|------|
| x 🍎 | $\hat{y}$ 🍌 | y 🍌 | $(y-\hat{y})^2$ |
| 1    | 1         | 0      |      |
| 5    | 5         | 16     |      |
| 6    | 6         | 20     |      |

**Model Ex#2**

weight        bias

$$\hat{y} = 2\ x + 2$$

|      | predicted | actual | cost |
|------|-----------|--------|------|
| x 🍎 | $\hat{y}$ 🍌 | y 🍌 | $(y-\hat{y})^2$ |
| 1    | 4         | 0      |      |
| 5    | 12        | 16     |      |
| 6    | 14        | 20     |      |

😎 *Let's calculate the cost(loss)!*

**Mean Squared Error (MSE)**  $$C(w,b) = \sum (y_n - \hat{y}_n)$$
$$n \in \{0,1,2\}$$

## Deep Learning with Neural Network - Cost

**Actual Data**

$$y = ?\ x + ?$$

weight · bias

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

Input (x) → Weight **?** → Bias **?** → Output (y)

**Model Ex#1**

$$\hat{y} = 1\ x + 0$$

weight · bias

| x 🍎 | predicted $\hat{y}$ | actual y 🍌 | cost $(y-\hat{y})^2$ |
|---|---|---|---|
| 1 | 1 | 0 | 1 |
| 5 | 5 | 16 | 121 |
| 6 | 6 | 20 | 196 |

$C(1,0) =$ **318**

**Model Ex#2**

$$\hat{y} = 2\ x + 2$$

weight · bias

| x 🍎 | predicted $\hat{y}$ | actual y 🍌 | cost $(y-\hat{y})^2$ |
|---|---|---|---|
| 1 | 4 | 0 | 16 |
| 5 | 12 | 16 | 16 |
| 6 | 14 | 20 | 36 |

$C(2,2) =$ **68** 🏅 Winner

😎 *Let's calculate the cost(loss)!*

**Mean Squared Error (MSE)**

$$C(w,b) = \sum (y_n - \hat{y}_n)$$

$$n \in \{0,1,2\}$$

## Deep Learning with Neural Network - Cost

### Actual Data

weight — bias

$$y = ? \, x + ?$$

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

Weight — Bias

Input (x) — ? — ? — Output (y)

### Model Ex#1

weight — bias

$$\hat{y} = 1 \, x + 0$$

| x 🍎 | predicted $\hat{y}$ | actual y | $(y-\hat{y})^2$ |
|---|---|---|---|
| 1 | 1 | 0 | 1 |
| 5 | 5 | 16 | 121 |
| 6 | 6 | 20 | 196 |

$C(1,0) = $ 318

### Model Ex#2

weight — bias

$$\hat{y} = 2 \, x + 2$$

| x 🍎 | predicted $\hat{y}$ | actual y | $(y-\hat{y})^2$ |
|---|---|---|---|
| 1 | 4 | 0 | 16 |
| 5 | 12 | 16 | 16 |
| 6 | 14 | 20 | 36 |

$C(2,2) = $ 68  🏅 *Winner*

😎 *Let's calculate the costs and get the lowest one!*

*arg min* **C(w,b)**

w,b∈[-∞,∞]

## Deep Learning with Neural Network - Optimizer

*Actual Data*

*weight*      *bias*

$$y = ?\, x + ?$$

| x 🍎 | y 🍌 |
|------|------|
| 1    | 0    |
| 5    | 16   |
| 6    | 20   |

*Weight*    *Bias*

◯ — ? — ? → ◯

*Input (x)*       *Output (y)*

$y = Wx + b$

**weight**

**The lowest!**

$w, b = 4, -4$
$C(w_0, b_0) = 0$

2

2      **bias**

😎 *Let's calculate the costs and get the lowest one!*

$$arg\ min\ C(w, b)$$

$w, b \in [-\infty, \infty]$

## Deep Learning with Neural Network - Optimizer

*weight*    *bias*

$$y = 4\ x - 4$$

| x 🍎 | y 🍌 |
|------|------|
| 1    | 0    |
| 5    | 16   |
| 6    | 20   |

**weight**

$y = Wx + b$

**The lowest!**

$w,b = 4,-4$
$C(w_0, b_0) = 0$

2

2    **bias**

**Backpropagation (weight update)**

Weight    Bias    **Loss Function**

?    ?

Input (x)    Predicted (ŷ)  Output (y)

*arg min* **C(w,b)**

w,b∈[-∞,∞]

## Deep Learning with Neural Network - Optimizer

weight     bias

$$y = 4 x - 4$$

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |



**Oh, Wait….**
*Do we need to calculate all cost for all options of* **W and B**?

**The lowest!**
$w, b = 4, -4$
$C(w_0, b_0) = 0$

*Expensive to compute (hours or days)*

$$arg\ min\ C(w, b)$$

$$w, b \in [-\infty, \infty]$$

## Finding the Optimal weight and bias – Gradient Descent



***There are different types of Gradient descent optimization algorithms:***
*Batch Gradient Descent, Stochastic Gradient Descent, Momentum, Adam, etc.*

# Deep Learning for NLP

## Finding the Optimal weight and bias – Gradient Descent



***There are different types of Gradient descent optimization algorithms:***
*Batch Gradient Descent, Stochastic Gradient Descent, Momentum, Adam, etc.*

# Deep Learning for NLP

## Finding the Optimal weight and bias – Gradient Descent



***There are different types of Gradient descent optimization algorithms:***
*Batch Gradient Descent, Stochastic Gradient Descent, Momentum, Adam, etc.*

**Deep Learning for NLP**

## Choose the optimal Learning Rate!

*Learning Rate: a hyperparameter that controls how much to change the model in response to the estimated error each time the model weights are updated.*



*new_weight = existing_weight — learning_rate * gradient*

*new_weight = existing_weight — learning_rate * (current_output – desired output) *gradient(current output) * existing_input*

# Deep Learning for NLP

## Deep Learning with Neural Network

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*



$y = 4x - 4$

1 🍎
0 🍌
5 🍎
16 🍌
6 🍎
20 🍌
3 🍎
8 🍌

# Deep Learning for NLP

## Deep Learning with Neural Network

$$y = w_1 \underset{\text{pixel(1,1)}}{x_1} + w_2 \underset{\text{pixel(1,2)}}{x_2} + w_3 x_3 + w_4 x_4 + \ldots + w_n x_n + b$$

Data



*Millions of **Parameters***
*Millions of **Samples***

# Deep Learning for NLP

## Deep Learning with Neural Network

*Vector1*     *Vector2*

$$y = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + \ldots + w_nx_n + b$$



*Data*

*Millions of **Parameters***

*Millions of **Samples***

# Deep Learning for NLP

## Deep Learning with Neural Network

*Input: x=number of apple given by Lisa*
*Output: y=number of banana received by Lisa*
*Parameters: Need to be estimated*

There is a limit of bananas I can give you

1 🍎

0 🍌

5 🍎

16 🍌

6 🍎

20 🍌

3 🍎

8 🍌

## Deep Learning with Neural Network

Nonlinear Neural Network

$$y = 4\,x - 4$$

weight      bias

*Data*

| x 🍎 | y 🍌 |
|---|---|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |

## Deep Learning with Neural Network

Nonlinear Neural Network

weight    bias

$$y = 2\,x + 3$$

*Data*

| x 🍎 | y 🍌 |
|:---:|:---:|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |
| 9 | 20 |
| 11 | 20 |

*Underfitting Issue*

## Deep Learning with Neural Network

Nonlinear Neural Network

*y = ???*

*Data*

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |
| 9 | 20 |
| 11 | 20 |

*How to make this possible?*

*A: Use different linear functions depending on the value of x*

# Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

**Data**

| x 🍎 | y 🍌 |
|-----|-----|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |
| 9 | 20 |
| 11 | 20 |

weight   bias   weight   bias

$$y = (w_1 x + b_1)s_1 + (w_2 x + b_2)s_2$$

*If x < 6 and 0*          *If x >= 6 and 0*

*How to make this possible?*

*A: Use different linear functions depending on the value of x*

# Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

**Data**

| x 🍎 | y 🍌 |
|------|------|
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |
| 9 | 20 |
| 11 | 20 |

*weight*   *bias*   *weight*   *bias*

$$y = (w_1x + b_1)s_1 + (w_2x+b_2)s_2$$

*If x < 6 and 0*          *If x >= 6 and 0*



*How to make this possible?*

*A: Use different linear functions depending on the value of x*

$W_1 = 4, b_1 = -4, W_2 = 0, b_2 = 20$ 🍎

# Deep Learning for NLP

## Deep Learning with Neural Network



High variance — overfitting

High bias — underfitting

Low bias, low variance — Good balance

# Deep Learning for NLP

# Deep Learning for NLP

## Single Neuron VS Multilayer



$$\sum_{i=1}^{m}(w_i x_i) + bias$$

$$f(x) = \begin{cases} 1 & \text{if } \sum wx + b \geq 0 \\ 0 & \text{if } \sum wx + b < 0 \end{cases}$$

$\hat{y}$

**Inputs**  Weights  **Summation and Bias**  **Activation**  **Output**

Input Layer

Output Layer

Input Layer

Hidden Layer  Output Layer

## Application

Previous Word Embedding Models



NxP  NxPxH  HxV

*Neural Net Language Model (NNLM)*

- The first word embedding model!
- N – how many previous words will be checked
- Only care about the previous words
- Long Computation time! Really slow

H    HxH    HxV

*Recurrent Neural Net Language Model (RNNLM)*

- No need to setup the No of previous words
- No projection layer
- Only care about the previous word
- A bit faster than NNLM

# Deep Learning for NLP

## Application

Word2Vec with CBOW



```
Input layer          Projection layer          Output layer

is (one-hot vector)

                     V = (v_is + v_the + v_capital + v_of) / 2m (window size)

the (one-hot vector)

capital (one-hot vector)        W'_{NxV} x V = z

                                ŷ = softmax(z)

of (one-hot vector)

N-Dimension

ŷ    cross entropy   y

state (one-hot vector)
```

$X_{is} \times W_{VxN} = V_{is}$

$X_{the} \times W_{VxN} = V_{the}$

$X_{on} \times W_{VxN} = V_{on}$

$X_{of} \times W_{VxN} = V_{of}$

$V = \dfrac{v_{is} + v_{the} + v_{capital} + v_{of}}{2m \ (window\ size)}$

$W'_{NxV} \times V = z$

$\hat{y} = softmax(z)$

# Deep Learning for NLP

## Application

Word2Vec with Skip Gram

## Negative Sampling

$$new\_weight = existing\_weight - learning\_rate * gradient$$

# Deep Learning for NLP

## Application

Application #1: Embedding Pretraining

# Deep Learning for NLP

## Application

Application #2: Translation Rescoring

| Lisa | gosta | de | comer | Macas | e | bananas | **source** |

| Bart | does | to | eat | coconuts | and | bananas | *translation#1* |

| Lisa | likes | to | eat | apples | and | bananas | *translation#2* |

| Lisa | dislikes | to | drink | apples | and | bananas | *translation#3* |

# Deep Learning for NLP

## Application

Application #2: Translation Rescoring

| Lisa | gosta | de | comer | Macas | e | bananas | **source** |

**Lisa** likes to eat apples and bananas *translation#2*

0.2     0.1

## Application

Application #2: Translation Rescoring

| Lisa | gosta | de | comer | Macas | e | bananas | **source** |



| **Lisa** | likes | to | eat | apples | and | bananas | *translation#2* |
| 0.2 | 0.1 | 0.3 |

# Deep Learning for NLP

## Application

Application #2: Translation Rescoring

| Lisa | gosta | de | comer | Macas | e | bananas | **source** |

Lisa    likes    to    eat    apples    and    bananas    *translation#2*

0.2         0.1      0.3

# Deep Learning for NLP

## Application

Application #2: Translation Rescoring

Lisa     gosta     de     comer     macas     e     bananas     **source**

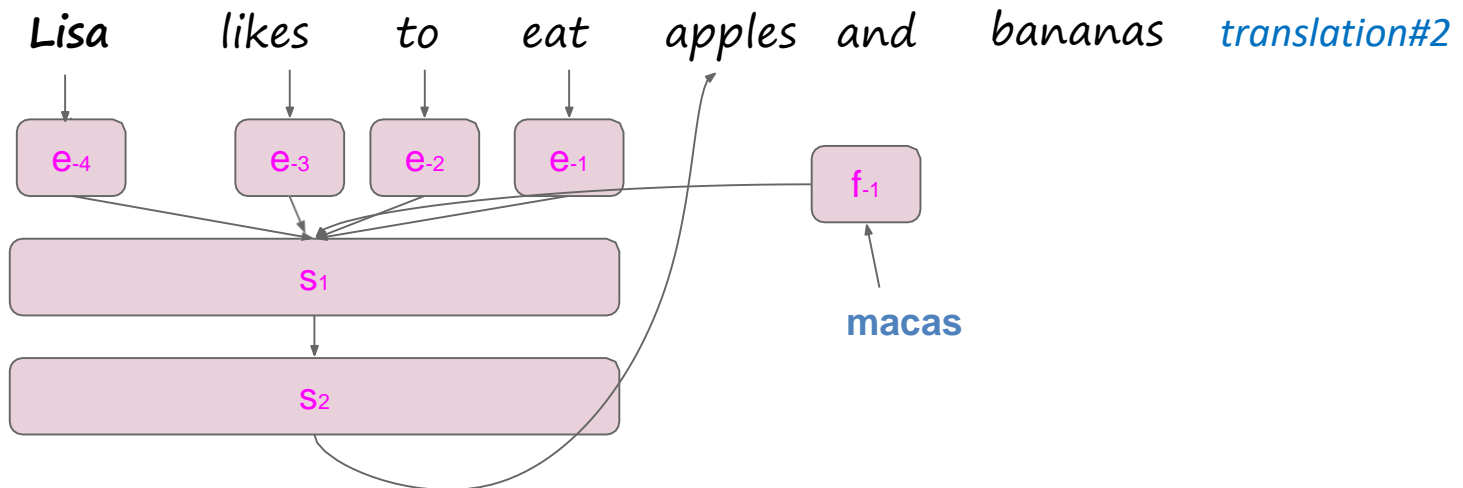Lisa     likes     to     eat     apples     and     bananas     *translation#2*

$e_{-4}$     $e_{-3}$     $e_{-2}$     $e_{-1}$     $f_{-1}$

$s_1$

$s_2$

**macas**

## Application

Application #2: Translation Rescoring

| Lisa | gosta | de | corner | Macas | e | bananas | **source** |

| **Bart** | does | to | eat | coconuts | and | bananas | *0.00003* |

| **Lisa** | likes | to | eat | apples | and | bananas | ***0.000378*** |

| **Lisa** | dislikes | to | drink | apples | and | bananas | *0.00012* |

**LECTURE PLAN**

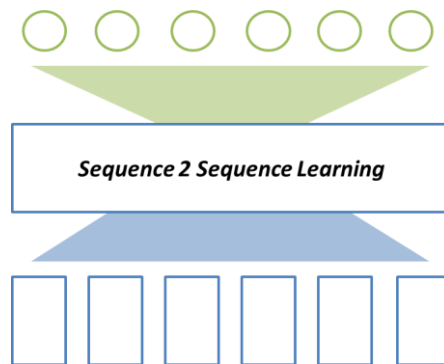**Lecture 3: Word Classification and Machine Learning**

1.  Previous Lecture: Word Embedding Review
2.  Word Embedding Evaluation
3.  Deep Neural Network for Natural Language Processing
    1.  Perceptron and Neural Network (NN)
    2.  Multilayer Perceptron
    3.  Applications
4.  **Next Week Preview**
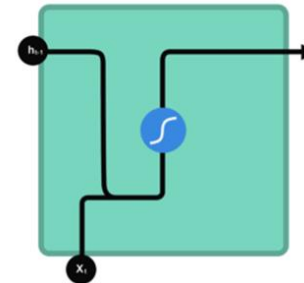    See how the Deep Learning can be used for NLP
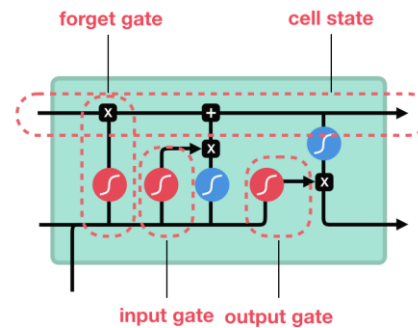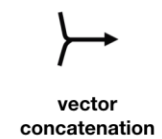    -  Text Classification, etc.

*RNN Cell*

Sequence 2 Sequence Learning

Tanh function
new hidden state
previous hidden state
input
concatenation

LSTM

forget gate
cell state
input gate
output gate

GRU

reset gate
update gate

sigmoid
tanh
pointwise multiplication
pointwise addition
vector concatenation

THE UNIVERSITY OF SYDNEY

## Reference for this lecture

- Deng, L., & Liu, Y. (Eds.). (2018). Deep Learning in Natural Language Processing. Springer.
- Rao, D., & McMahan, B. (2019). Natural Language Processing with PyTorch: Build Intelligent Language Applications Using Deep Learning. " O'Reilly Media, Inc.".
- Manning, C. D., Manning, C. D., & Schütze, H. (1999). Foundations of statistical natural language processing. MIT press.
- Blunsom, P 2017, Deep Natural Language Processing, lecture notes, Oxford University
- Manning, C 2017, Natural Language Processing with Deep Learning, lecture notes, Stanford University