

**COMP5046**

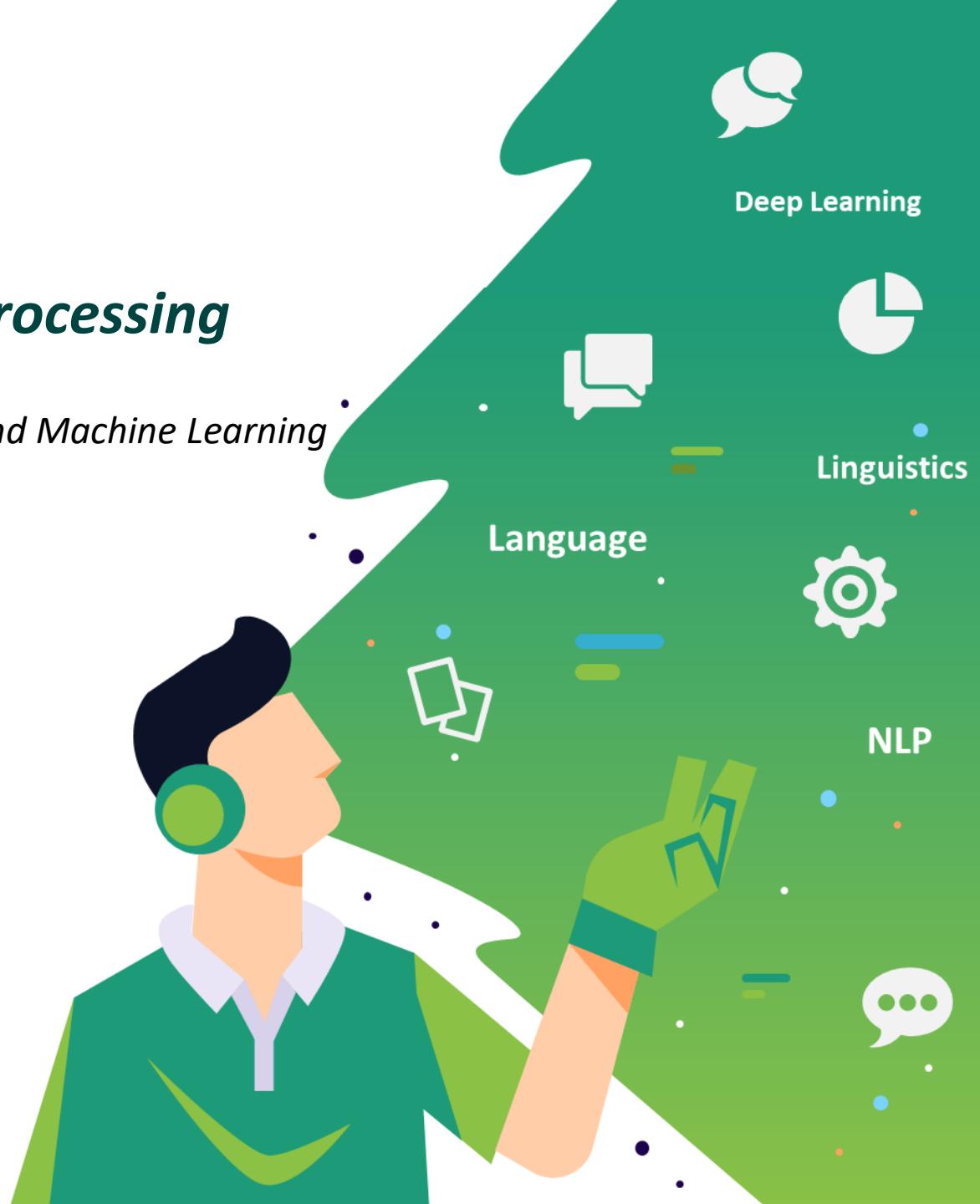
# *Natural Language Processing*

*Lecture 3: Word Classification and Machine Learning*

**Dr. Caren Han**

Semester 1, 2021

School of Computer Science,  
University of Sydney



# 0 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

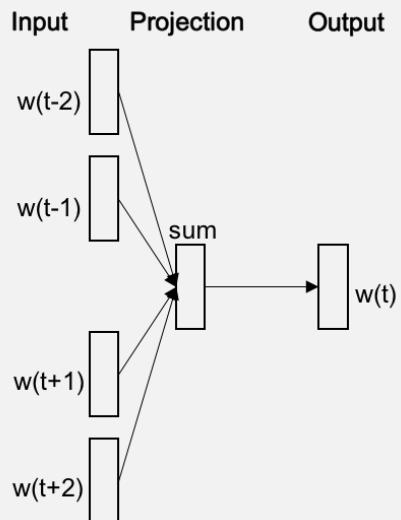
See how the Deep Learning can be used for NLP

  - Text Classification, etc.

## Word2Vec Models

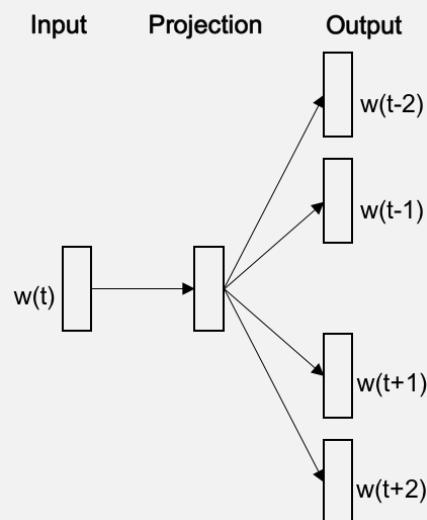
### CBOW

*Predict center word from (bag of) context words*



### Skip-gram

*Predict context words given center word*



# 1 Previous Lecture Review

## Word2Vec with Continuous Bag of Words (CBOW)

Predict center word from (bag of) context words

**Sentence: “Sydney is the state capital of NSW”**

*Using window slicing, develop the training data*

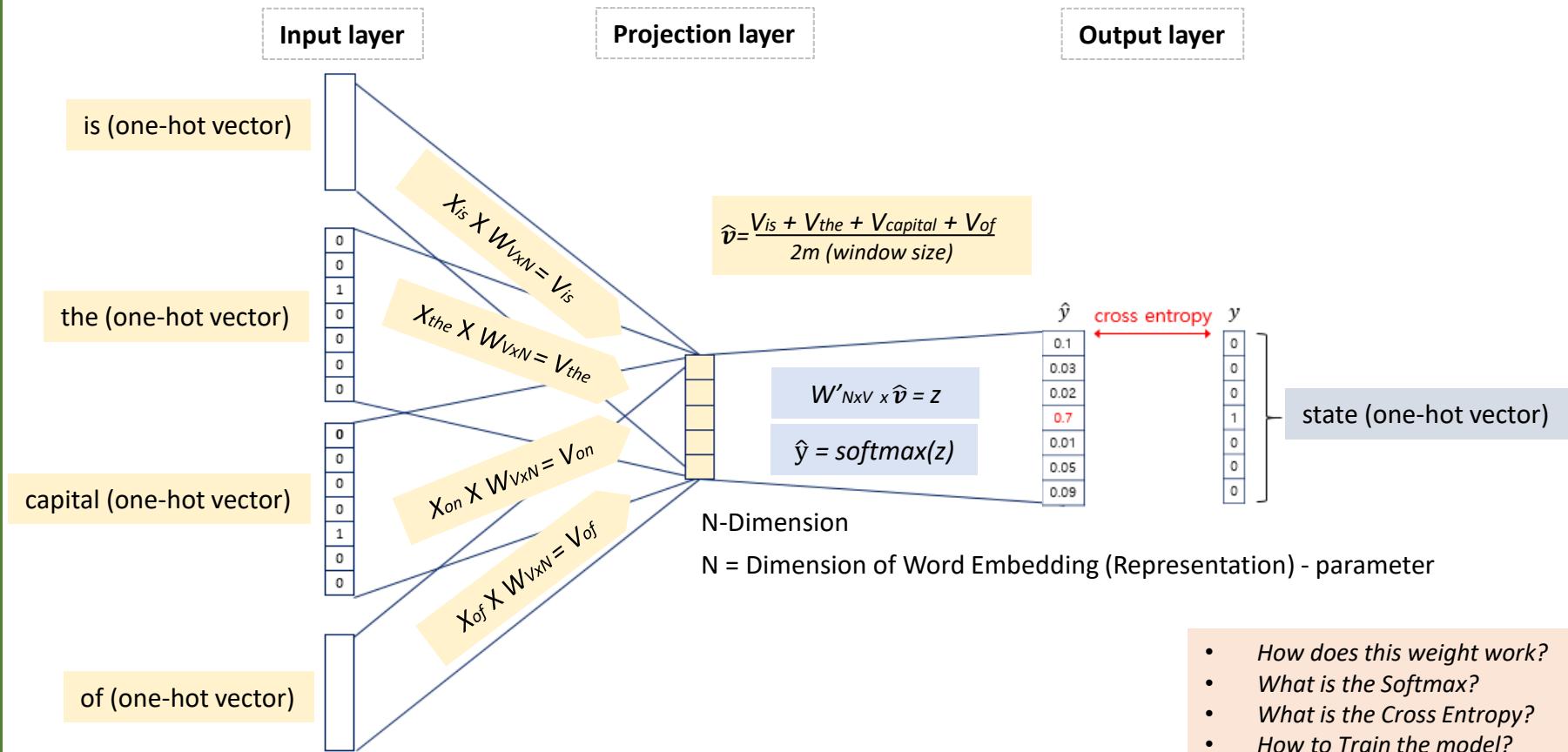
Center word	Context (“outside”) word	Sequence
[1,0,0,0,0,0,0]	[0,1,0,0,0,0,0], [0,0,1,0,0,0,0]	Sydney is the state capital of NSW
[0,1,0,0,0,0,0]	[1,0,0,0,0,0,0], [0,0,1,0,0,0,0], [0,0,0,1,0,0,0]	Sydney is the state capital of NSW
[0,0,1,0,0,0,0]	[1,0,0,0,0,0,0], [0,1,0,0,0,0,0] [0,0,0,1,0,0,0], [0,0,0,0,1,0,0]	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

## CBOW – Neural Network Architecture

Predict center word from (bag of) context words

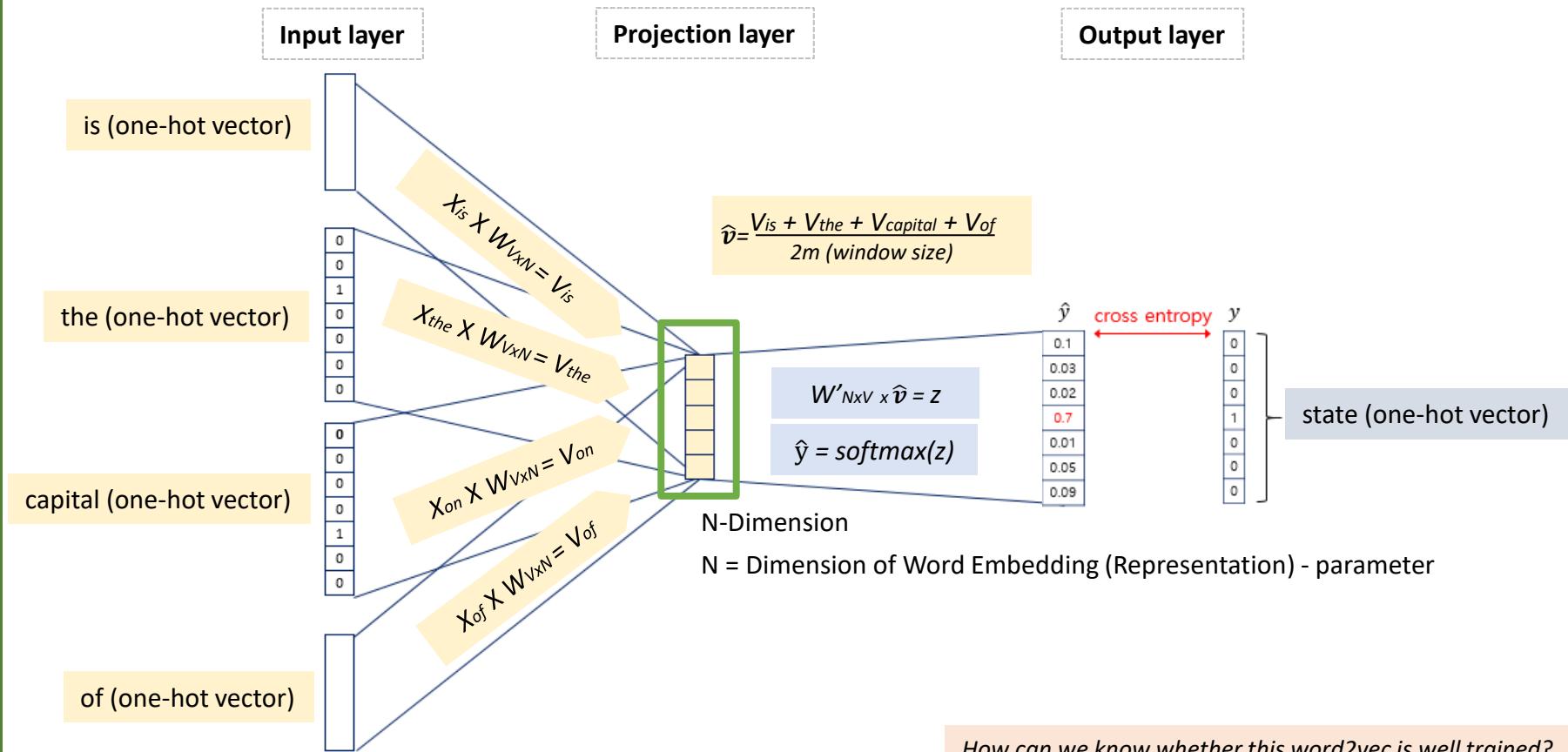
**Sentence: “Sydney is the state capital of NSW”**



## CBOW – Neural Network Architecture

Predict center word from (bag of) context words

**Sentence:** “Sydney is the state capital of NSW”



# 0 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

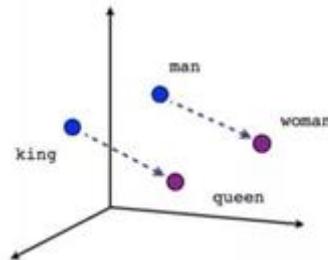
See how the Deep Learning can be used for NLP

  - Text Classification, etc.

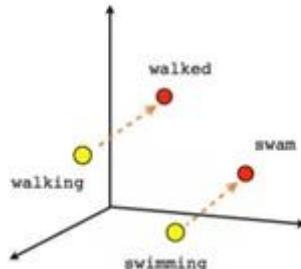
# Word Embedding Evaluation

## How to evaluate word vectors?

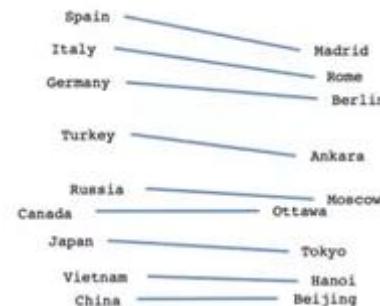
Type	How to work / Benefit
Intrinsic	Evaluation on a specific/intermediate subtask
Intrinsic	<ul style="list-style-type: none"> <li>Fast to compute</li> <li>Helps to understand that system</li> <li>Not clear if really helpful unless correlation to real task is established</li> </ul>
Extrinsic	Evaluation on a real task
Extrinsic	<ul style="list-style-type: none"> <li>Can take a long time to compute accuracy</li> <li>Unclear if the subsystem is the problem or its interaction or other subsystems</li> </ul>



Male-Female



Verb tense



Country-Capital

## 2

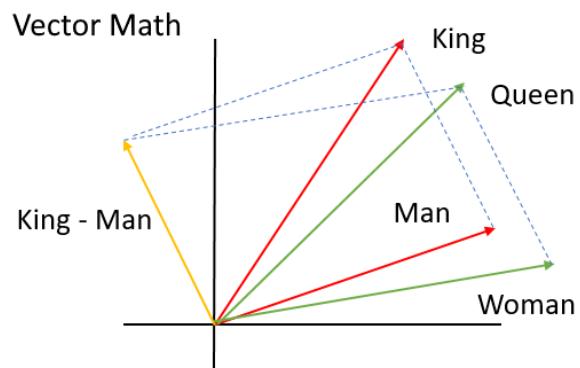
# Word Embedding Evaluation

## Intrinsic word vector evaluation

### Word Vector Analogies

$$\begin{array}{ccc} a \leftrightarrow b & :: & c \leftrightarrow ??? \\ man \leftrightarrow women & :: & king \leftrightarrow ??? \end{array}$$

- 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	Training 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-Syntactic 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% syntactic 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

two sets of datasets used to test the models

(Original Word2vec Paper - Mikolov et al.2013)

can see CBOW very good at Syntactic words but not so good at semantic relationships

# Word Embedding Evaluation

## Intrinsic word vector evaluation

### Evaluation Result Comparison

The Semantic-Syntactic word relationship tests for understanding of a wide variety of relationships as shown below.

Table 2: Results on the word analogy task, given as percent accuracy. Underlined scores are best within groups of similarly-sized models; bold scores are best overall. HPCA vectors are publicly available<sup>2</sup>; (i)vLBL results are from (Mnih et al., 2013); skip-gram (SG) and CBOW results are from (Mikolov et al., 2013a,b); we trained SG<sup>†</sup> and CBOW<sup>†</sup> using the word2vec tool<sup>3</sup>. See text for details and a description of the SVD models.

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<b>80.8</b>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW <sup>†</sup>	300	6B	63.6	<u>67.4</u>	65.7
SG <sup>†</sup>	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	71.7
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	<b>81.9</b>	<b>69.3</b>	<b>75.0</b>

(Original Glove Paper - Pennington et al.2014)

# Word Embedding Evaluation

## Intrinsic word vector evaluation

### Evaluation Result Comparison

The Semantic-Syntactic word relationship tests for understanding of a wide variety of relationships as shown below.

#### Window-Size (m) and Vector Dimension (N)

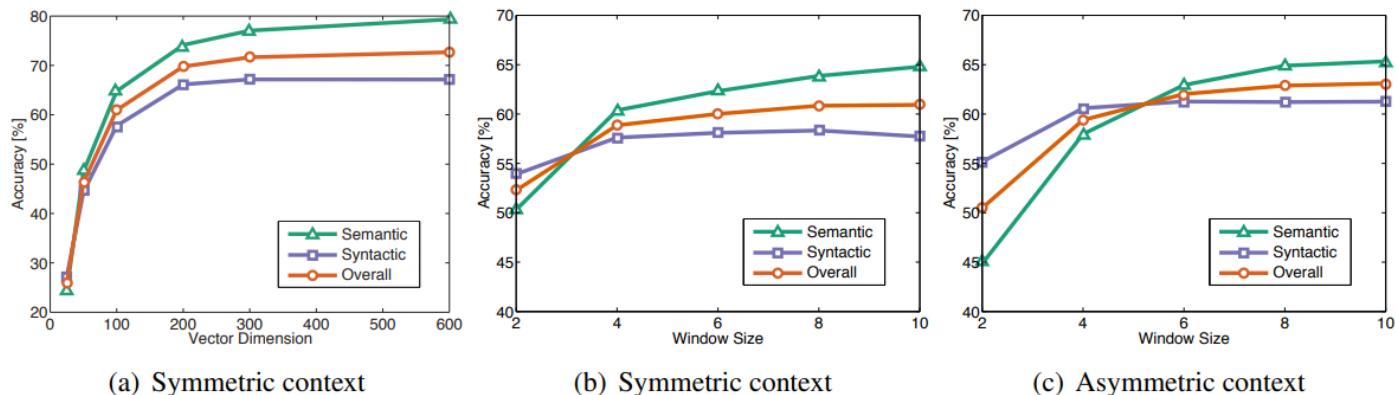


Figure 2: Accuracy on the analogy task as function of vector size and window size/type. All models are trained on the 6 billion token corpus. In (a), the window size is 10. In (b) and (c), the vector size is 100.

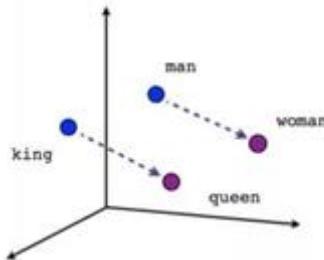
**these accuracy only for the semantic-syntactic dataset (*Original Glove Paper - Pennington et al.2014*) not a good representation for general text datasets**

# Word Embedding Evaluation

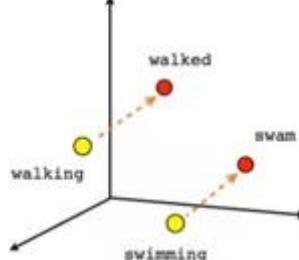
## How to evaluate word vectors?

Type	How to work / Benefit
Intrinsic	Evaluation on a specific/intermediate subtask <ul style="list-style-type: none"> <li>Fast to compute</li> <li>Helps to understand that system</li> <li>Not clear if really helpful unless correlation to real task is established</li> </ul>
Extrinsic	<b>Evaluation on a real task</b> <ul style="list-style-type: none"> <li>Can take a long time to compute accuracy</li> <li>Unclear if the subsystem is the problem or its interaction or other subsystems</li> </ul>

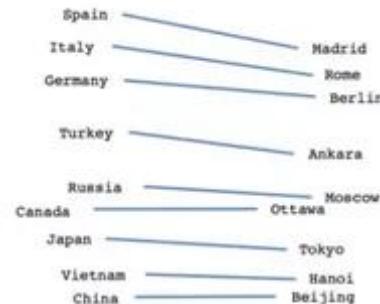
place word/text embedding as inputs into a text/document classification



Male-Female



Verb tense



Country-Capital

there is no best way to find best evaluator of a word vector

# 0 LECTURE PLAN

## Lecture 3: Word Classification and Machine Learning

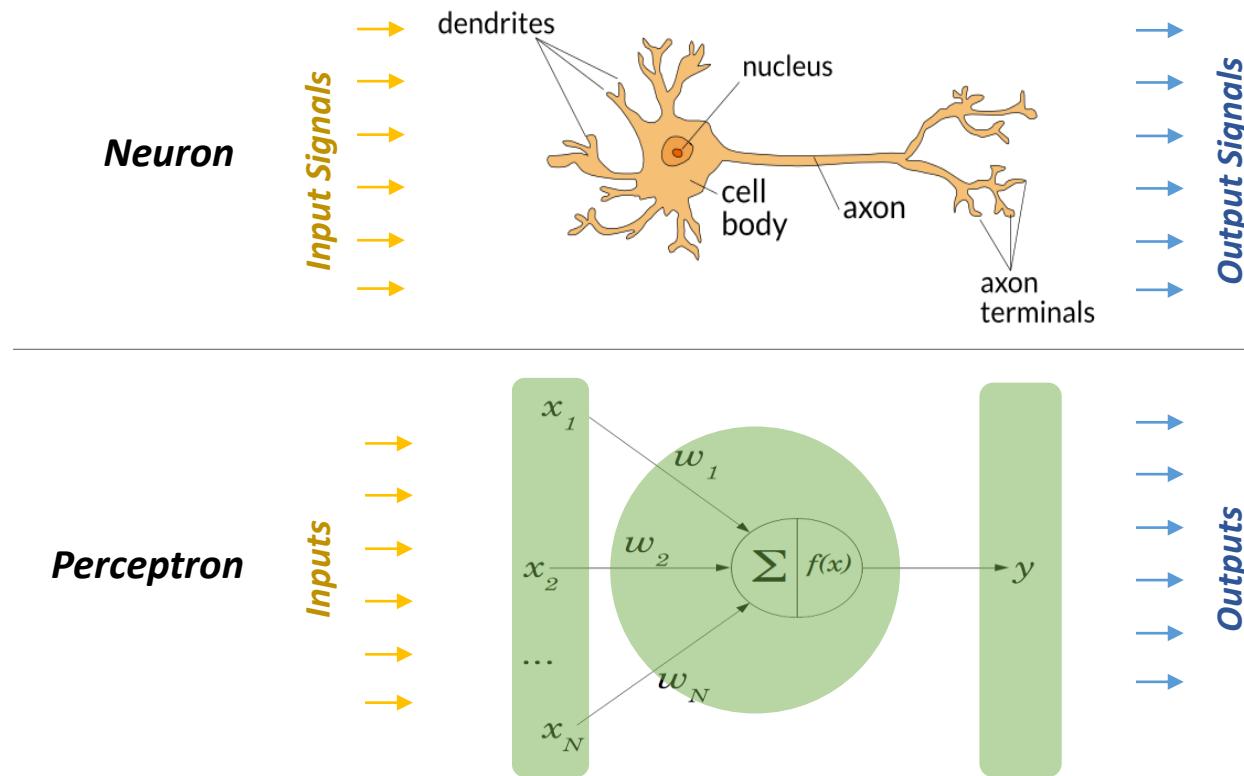
1. Previous Lecture: Word Embedding Review
2. Word Embedding Evaluation
3. **Deep Neural Network for Natural Language Processing**
  1. Perceptron and Neural Network (NN)
  2. Multilayer Perceptron
  3. Applications
4. Next Week Preview

See how the Deep Learning can be used for NLP

  - Text Classification, etc.

## Deep Learning with Neural Network

### Neuron and Perceptron

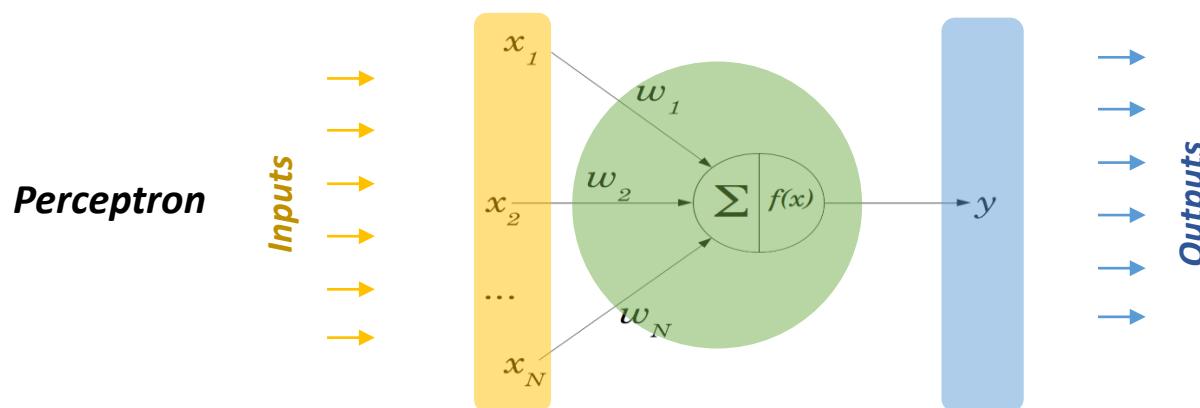


# 3 Deep Learning for NLP

## Deep Learning with Neural Network

Inputs and Outputs (Labels) for Natural Language Processing

$x_i$	Inputs	<b>Features</b> words (indices or vectors!), context windows, sentences, documents, etc.
$y_i$	Outputs (labels)	<b>What we try to predict/classify</b> <ul style="list-style-type: none"> <li>E.g. word meaning, sentiment, name entity</li> </ul>



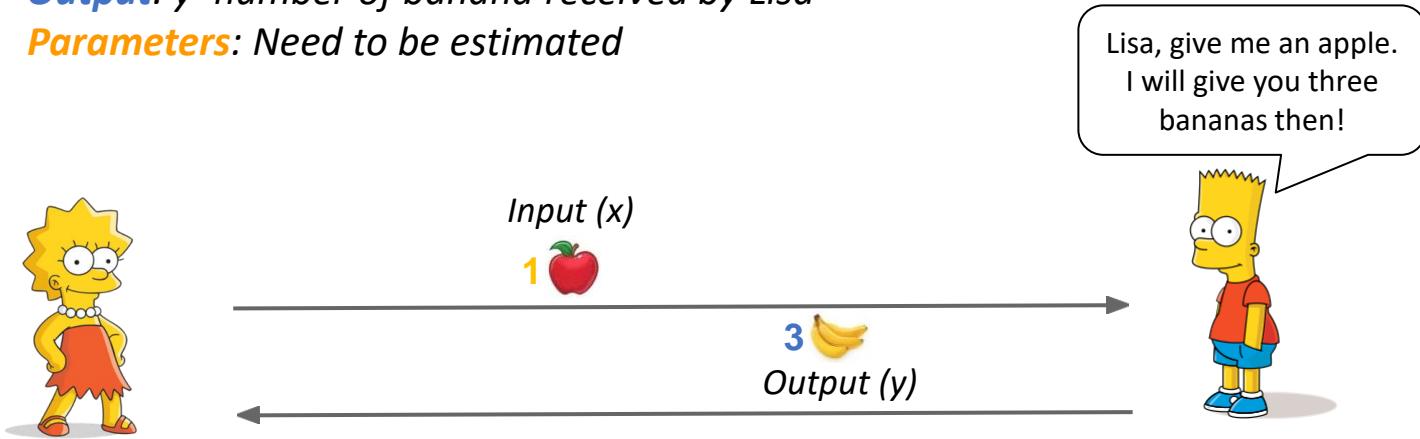
## 3 Deep Learning for NLP

## 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



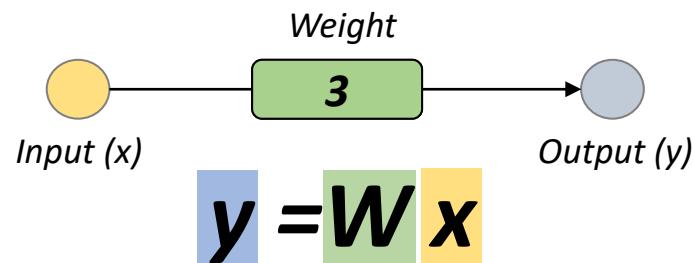
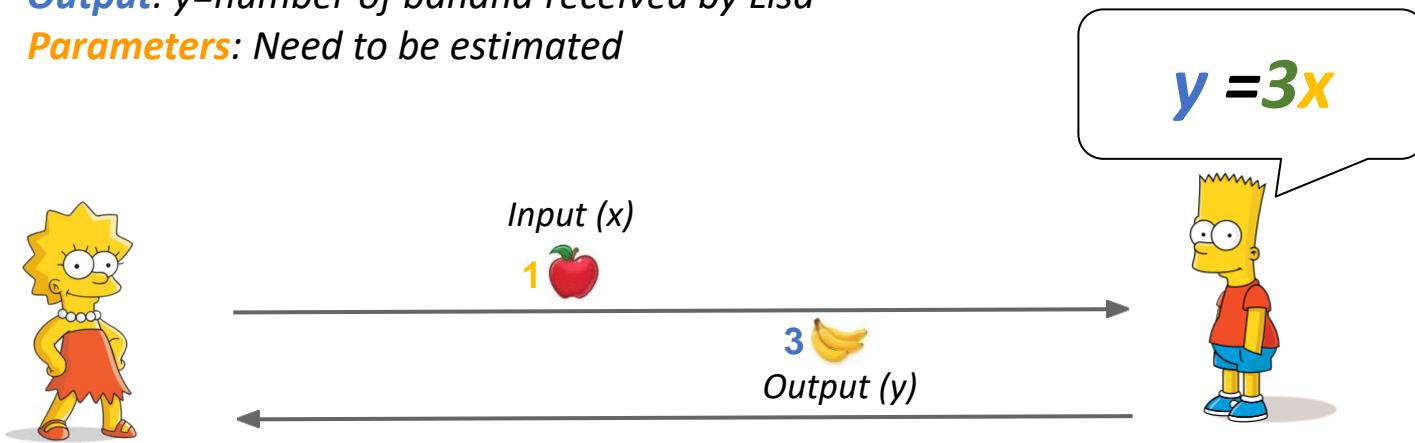
# 3 Deep Learning for NLP

## 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



## 3 Deep Learning for NLP



# 3 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!



# 3 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

initialise random weights then backpropagate to change weights



$$y = Wx$$

What is  $W$  then?

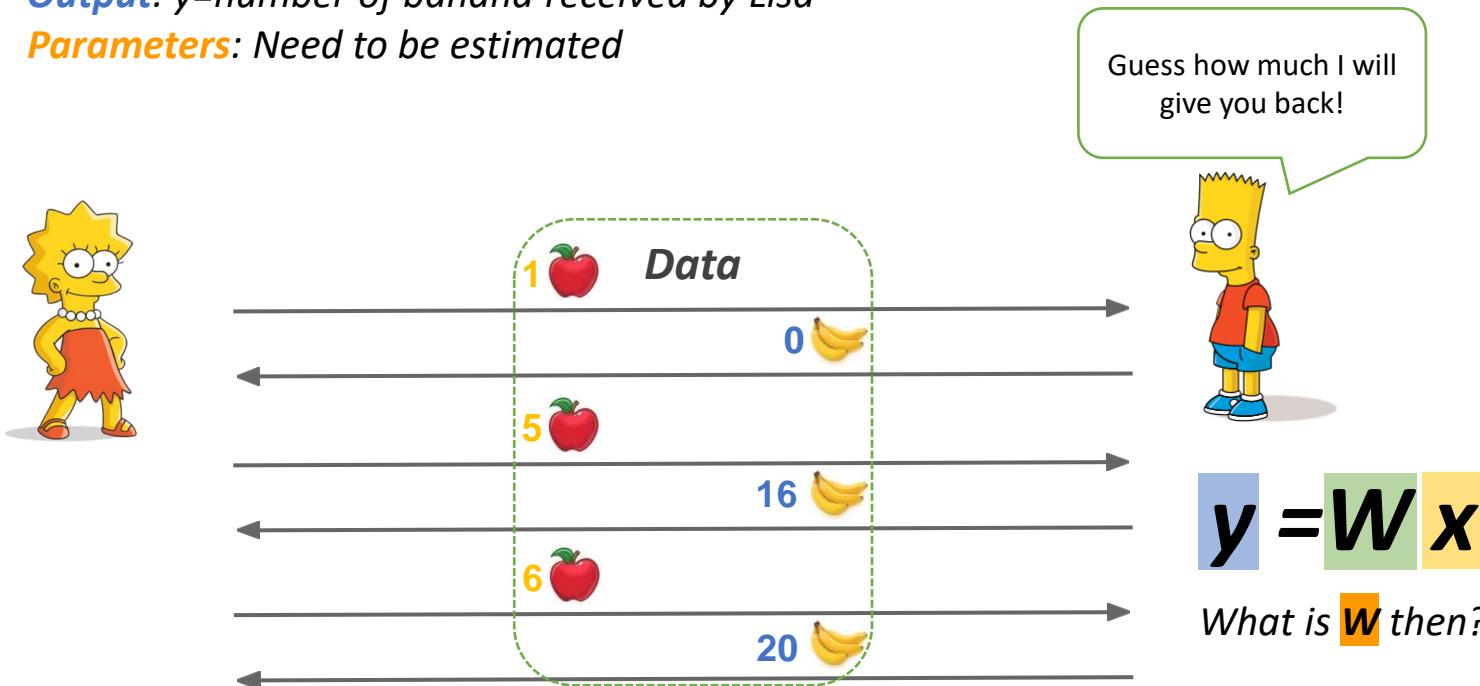
# 3 Deep Learning for NLP

## 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



## 3 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

$$\mathbf{y} = \mathbf{W} \mathbf{x}$$

What is  $\mathbf{W}$  then?

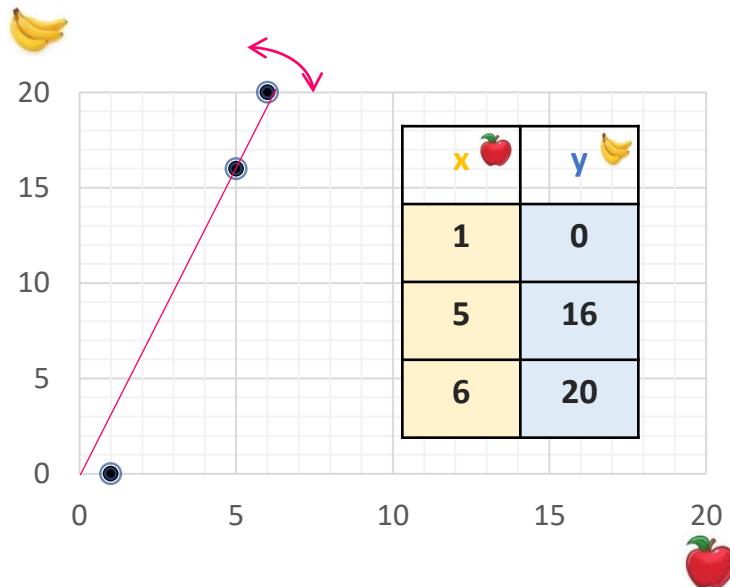
# 3 Deep Learning for NLP

## 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



$$y = Wx$$

What is  $W$  then?

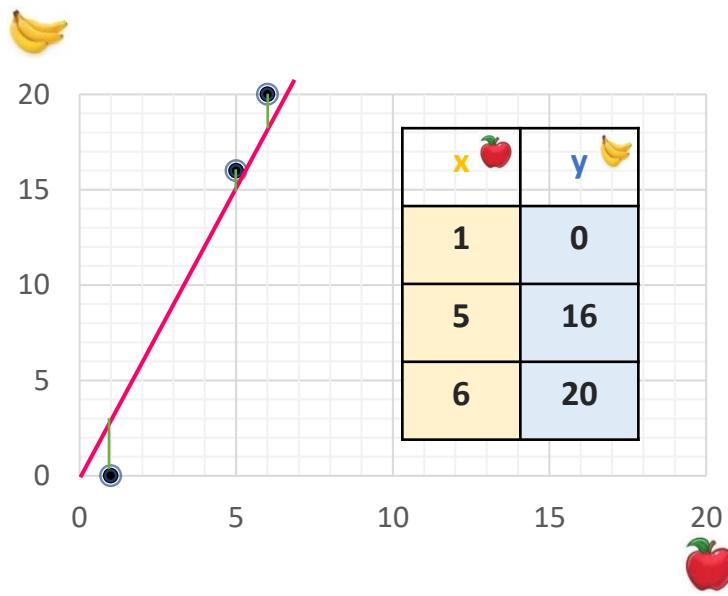
# 3 Deep Learning for NLP

## 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



$$\mathbf{y} = \mathbf{W} \mathbf{x}$$

What if  $W$  is 3?

$$3 = 3 \times 1$$

$$15 = 3 \times 5$$

$$20 = 3 \times 6$$

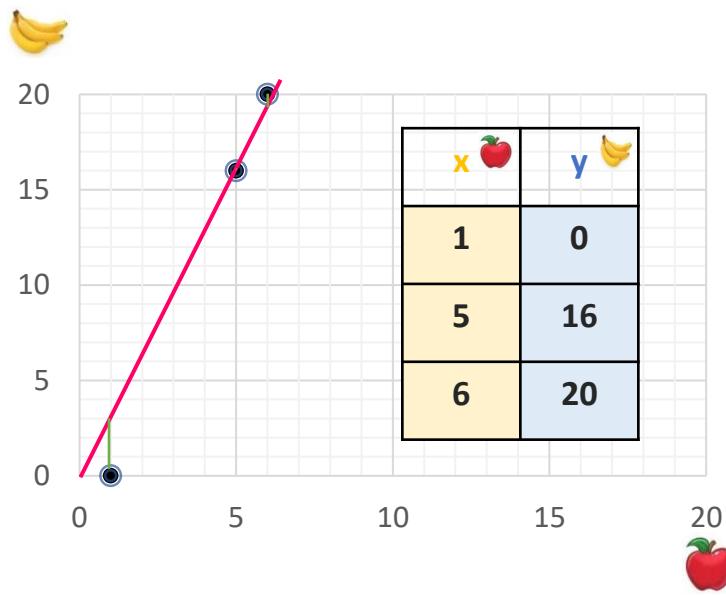
# 3 Deep Learning for NLP

## 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



$$y = Wx$$

What if  $W$  is 3.2?

$$3.2 = 3.2 \times 1$$

$$16 = 3.2 \times 5$$

$$19.2 = 3.2 \times 6$$

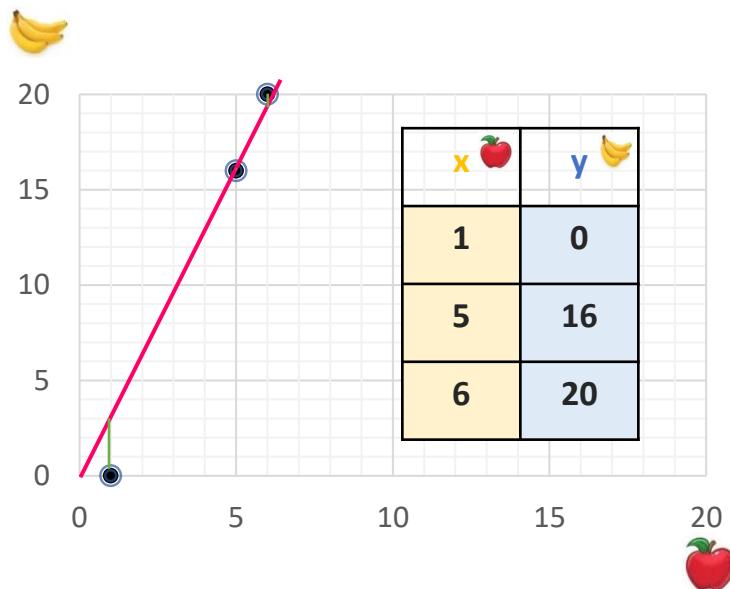
# 3 Deep Learning for NLP

## 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



$$y = \underset{\text{weight}}{W}x + \underset{\text{bias}}{b}$$

Weight is not enough...

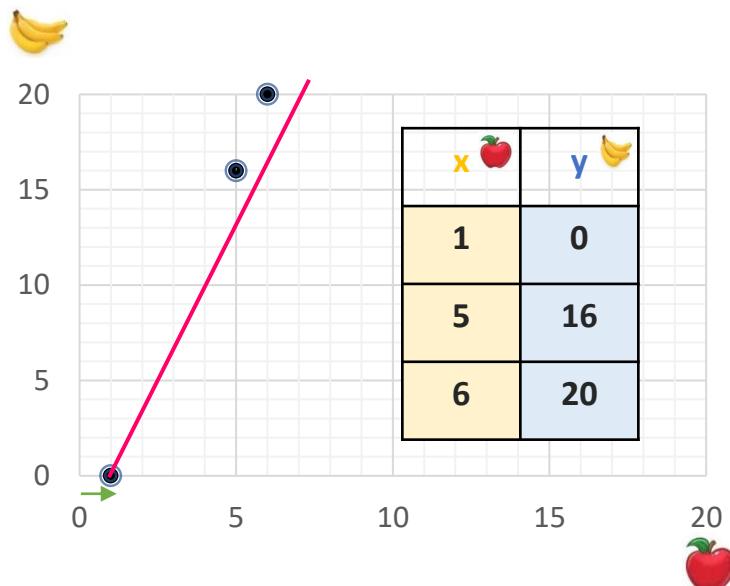
# 3 Deep Learning for NLP

## 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



push the line from 0,0 to the first data point

$$y = \textcolor{blue}{W}x + \textcolor{orange}{b}$$

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

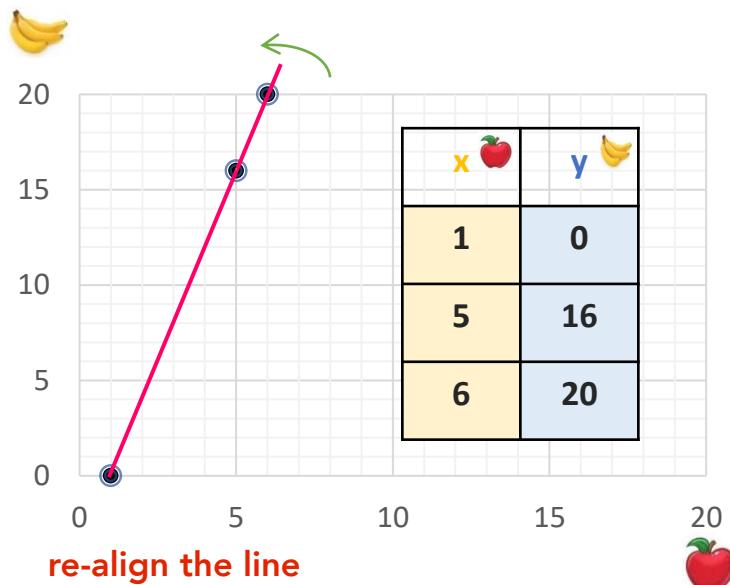
# 3 Deep Learning for NLP

## 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



$$y = \textcolor{blue}{W}x + \textcolor{orange}{b}$$

weight    bias

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

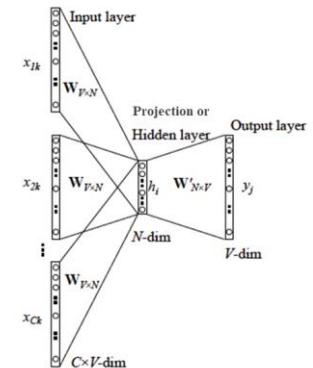
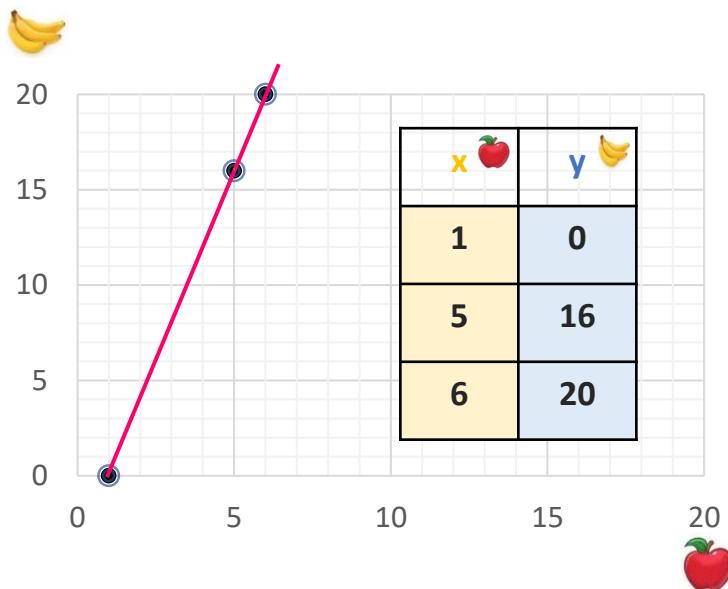
# 3 Deep Learning for NLP

## 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



Model  $\textcolor{blue}{y} = \textcolor{green}{W} \textcolor{orange}{x} + \textcolor{green}{b}$

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

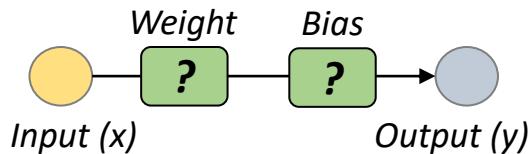
# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Cost

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



*Model Ex#1*

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

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

*Model Ex#2*

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

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



Which one is closer?

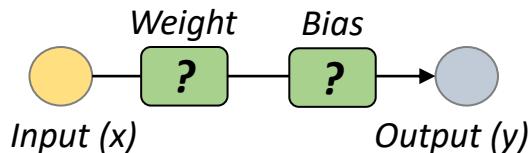
# 3 Deep Learning for NLP

## Deep Learning with Neural Network – Cost (loss)

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



*Model Ex#1*

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

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

*Model Ex#2*

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

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

😎 Let's calculate the cost(loss)!

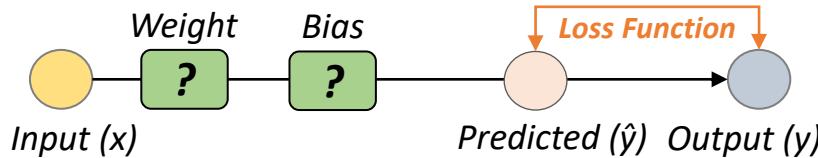
Mean Squared Error  
(MSE)

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

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

# 3 Deep Learning for NLP

## WAIT! Loss Function? Cost Calculation?



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

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

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

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

The diagram shows two probability distributions,  $\hat{\mathbf{y}}$  and  $\mathbf{y}$ , represented as vertical vectors of probabilities for six categories. A red double-headed arrow labeled "cross entropy" connects the two distributions.

$\hat{\mathbf{y}}$	$\mathbf{y}$
0.1	0
0.03	0
0.02	0
0.7	1
0.01	0
0.05	0
0.09	0

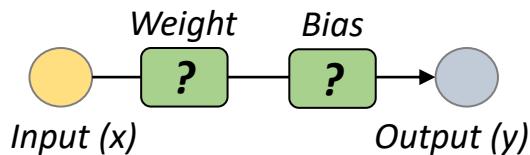
# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Cost (loss)

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



*Model Ex#1*

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

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

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

*Model Ex#2*

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

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

$$C(2,2) = 68$$



😎 Let's calculate the cost!

$$C(w,b) = \sum_{n \in \{0,1,2\}} (y_n - \hat{y}_n)$$

# 3 Deep Learning for NLP

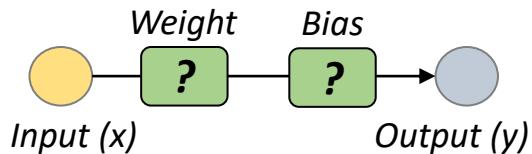
## Deep Learning with Neural Network - Cost (loss)

*Actual Data*

$$y = ? \ x + ?$$

weight      bias

x 🍎	y 🍌
1	0
5	16
6	20



*Model Ex#1*

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

weight      bias

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*

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

weight      bias

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

$$C(2,2) = 68$$



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

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

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

# 3 Deep Learning for NLP

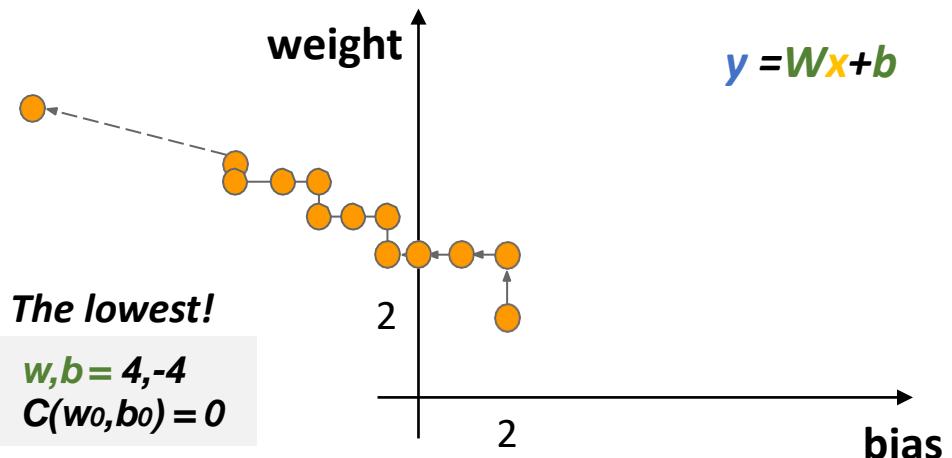
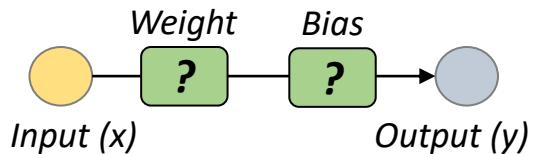
## Deep Learning with Neural Network - Optimizer

*Actual Data*

$$y = ? \times x + ?$$

weight                    bias

x 🍎	y 🍌
1	0
5	16
6	20



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

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

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

# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

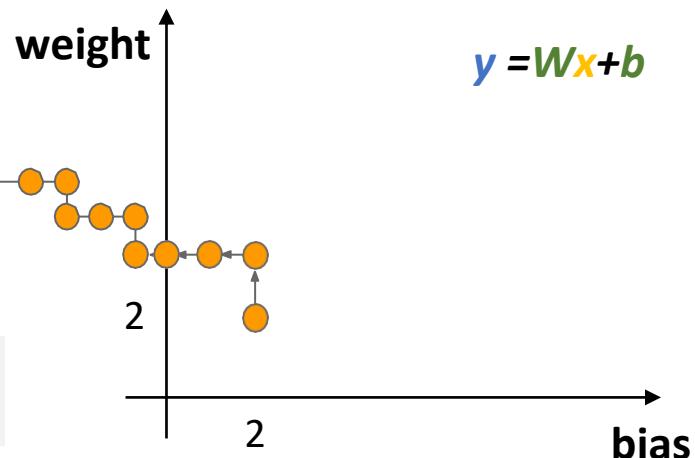
**Backpropagation (weight update)**

$$y = 4x - 4$$

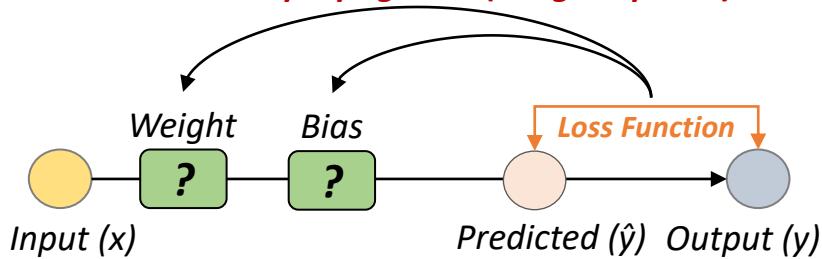
weight                    bias

x 🍎	y 🍌
1	0
5	16
6	20

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



**Backpropagation (weight update)**



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

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

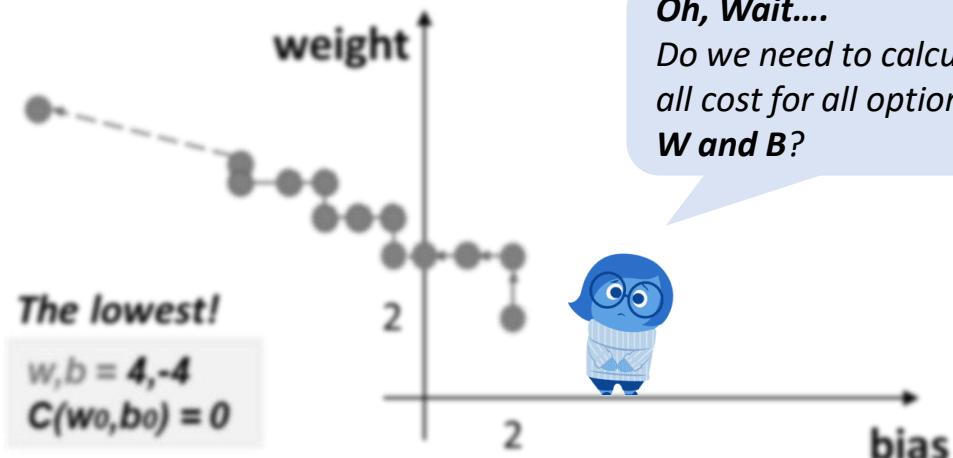
# 3 Deep Learning for NLP

*Backpropagation (weight update)*

## Deep Learning with Neural Network - Optimizer

$$y = 4 \text{ weight} | x - 4 \text{ bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Expensive to compute  
*(hours or days)*

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

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

## 3 Deep Learning for NLP

## Finding the Optimal weight and bias – Gradient Descent



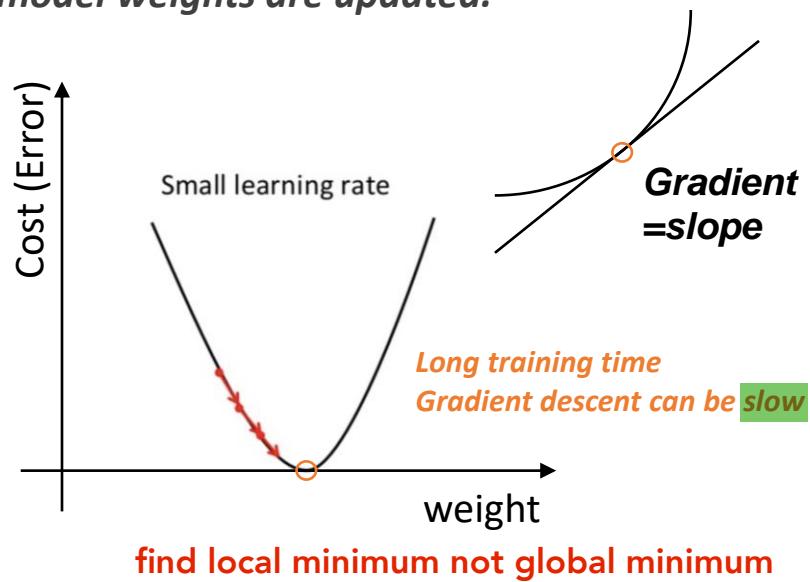
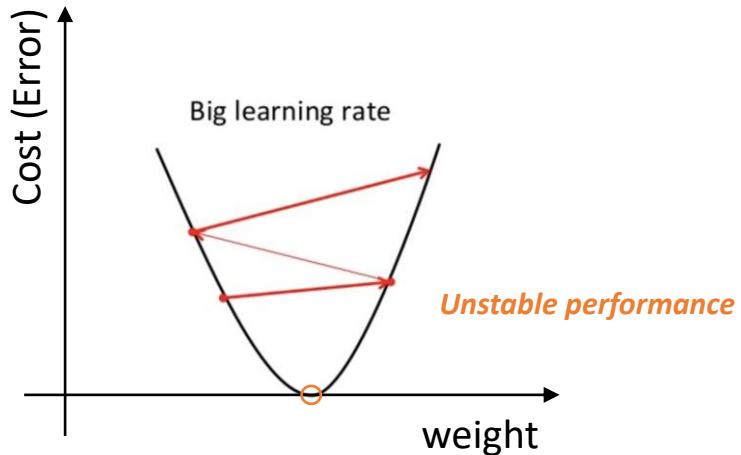
***There are different types of Gradient descent optimization algorithms:***

***Batch Gradient Descent, Stochastic Gradient Descent, Momentum, Adam, etc.***

# 3 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.



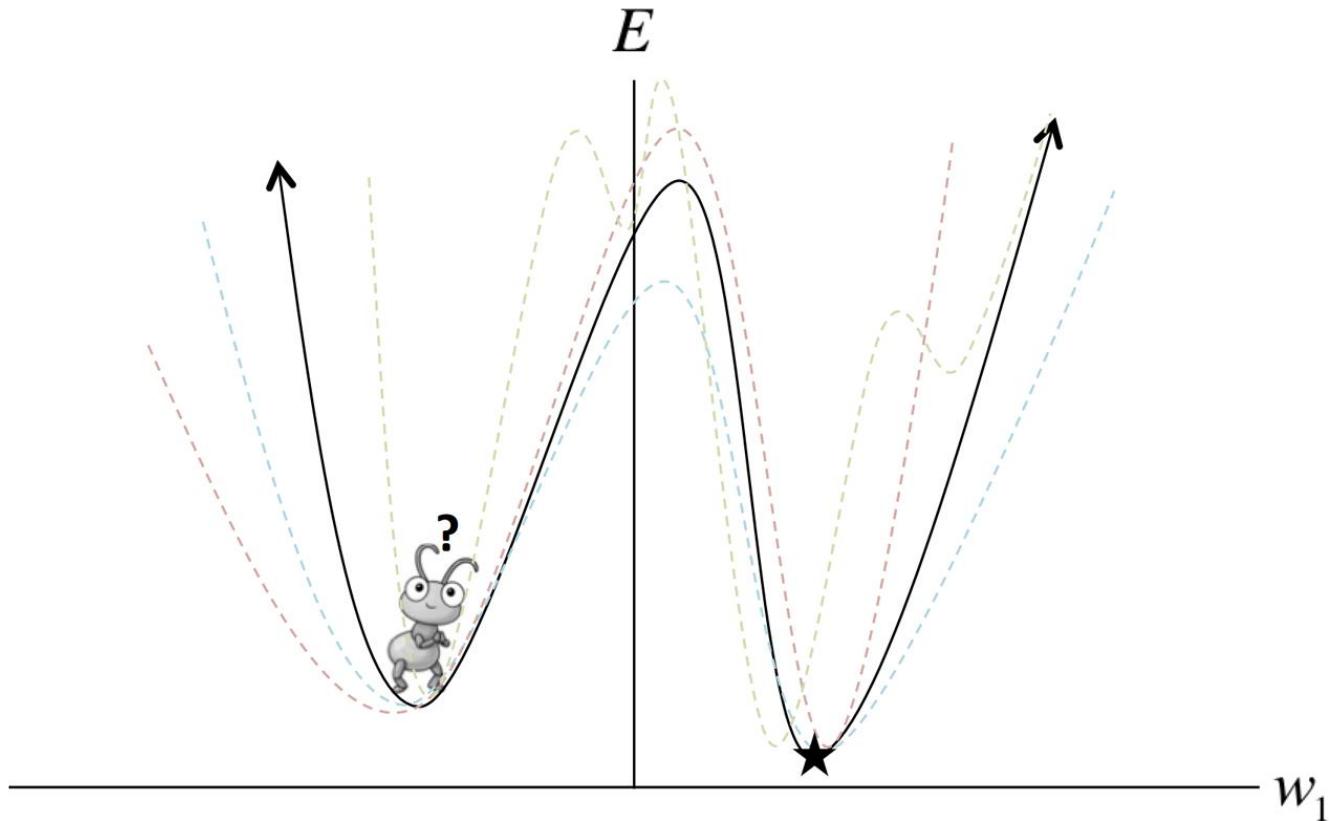
`new_weight = existing_weight - learning_rate * gradient`

`new_weight = existing_weight - learning_rate * (current_output - desired_output)`

`*gradient(current output) * existing_input`

## 3 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.*

# 3 Deep Learning for NLP

## Stochastic Gradient Descent

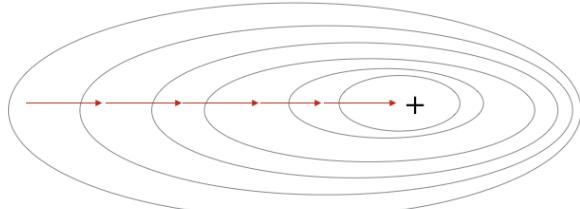
**The cost would be very expensive if we calculate it for all windows in the corpus!**  
You would wait a very long time before making a single update!

*The Solution can be used different Gradient Descent Method.  
The most common – “**Stochastic Gradient Descent (SGD)**”*

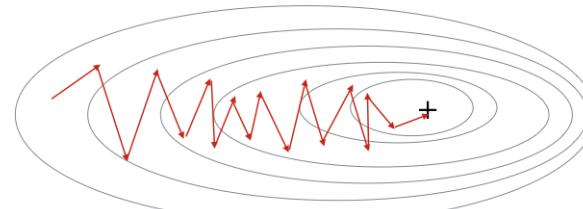
*Vanilla (Batch) gradient descent performs redundant computations for large datasets, as it recomputes gradients for similar examples before each parameter update. SGD does away with this redundancy by performing one update at a time. It is therefore usually much faster and can also be used to learn online.*

**recomputes gradient one by one to get to current gradient.**

Gradient Descent

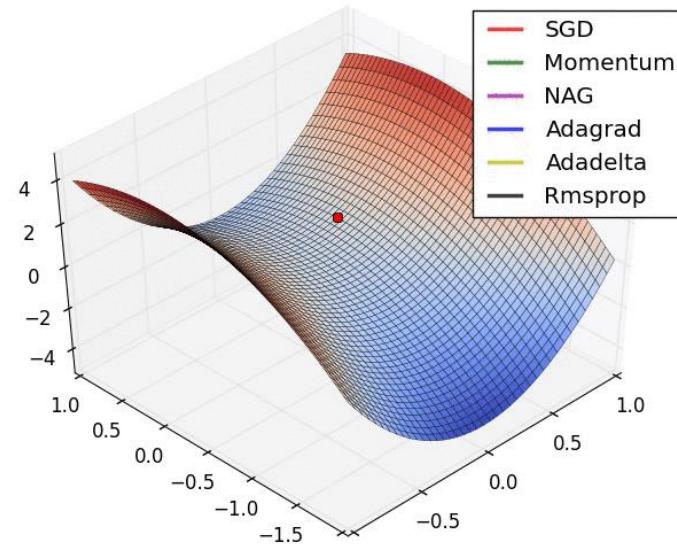
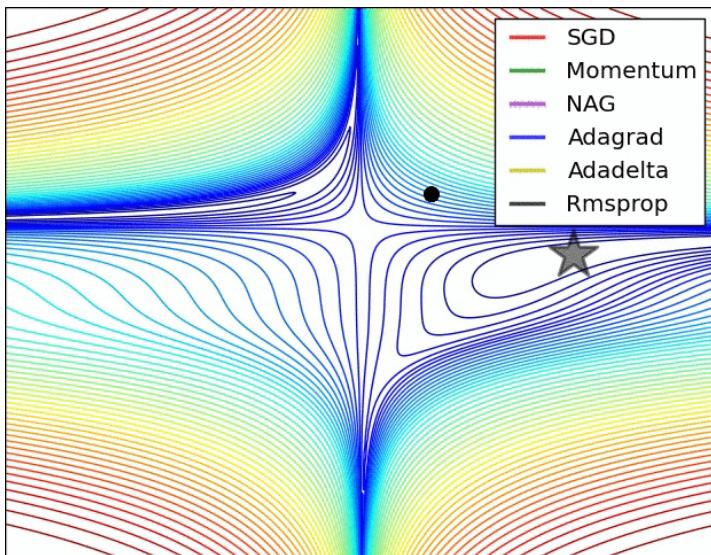


Stochastic Gradient Descent



## 3 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.*

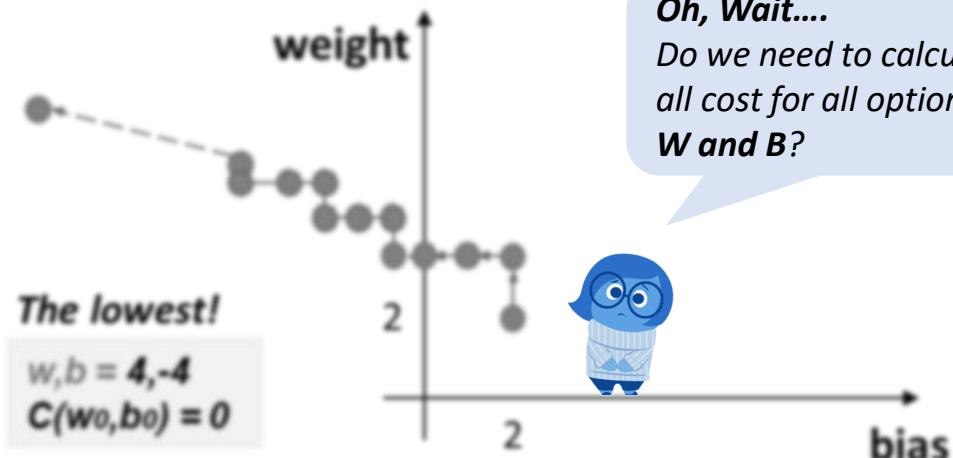
# 3 Deep Learning for NLP

*Backpropagation (weight update)*

## Deep Learning with Neural Network - Optimizer

$$y = 4 \text{ weight} | x - 4 \text{ bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Expensive to compute  
*(hours or days)*

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

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

# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

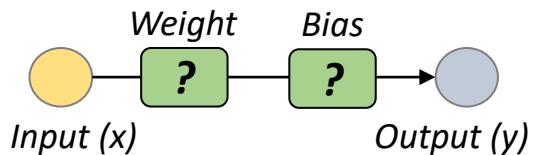
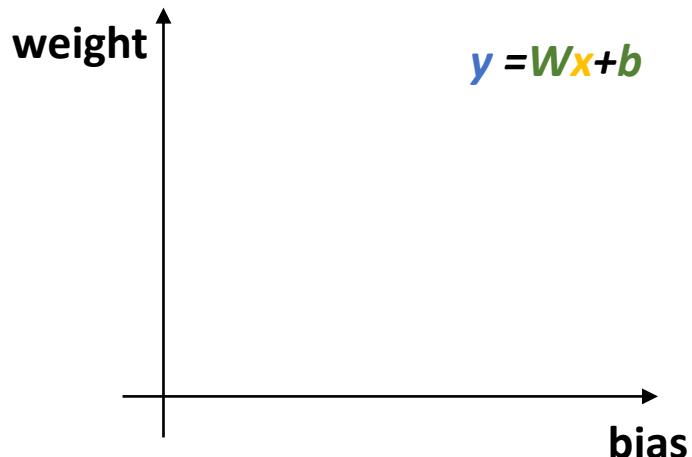
*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

weight                    bias

x	y
1	0
5	16
6	20



Should be used sparingly

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

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

# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

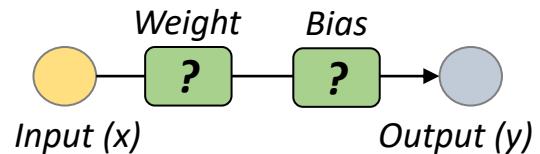
*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

weight      bias

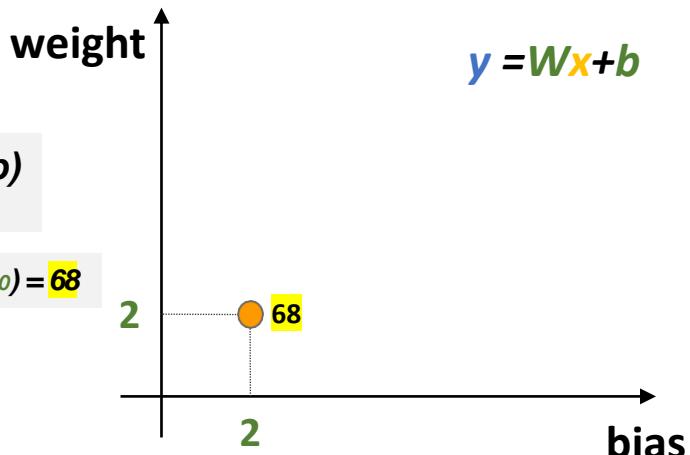
x	y
1	0
5	16
6	20



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

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

$w,b = 2,2 \quad C(w_0,b_0) = 68$



# 3 Deep Learning for NLP

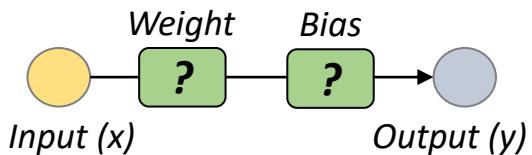
## Deep Learning with Neural Network - Optimizer

*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



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

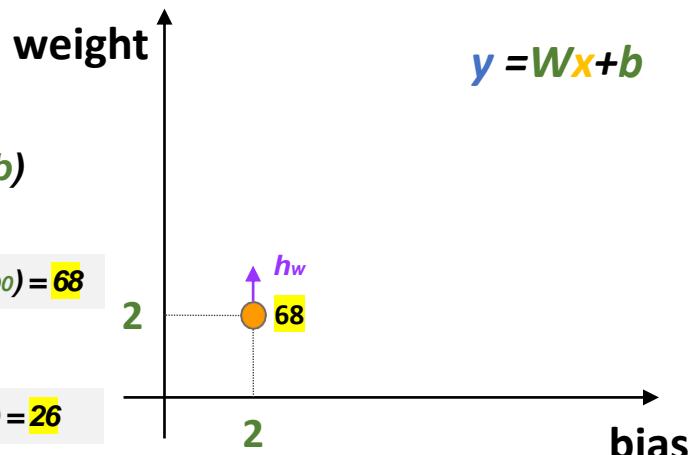
$$w, b \in [-\infty, \infty]$$

$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$h_w = 1$$

$$C(w_0 + h_w, b_0) = C(3, 2) = 26$$

eg. learning rate is 1



# 3 Deep Learning for NLP

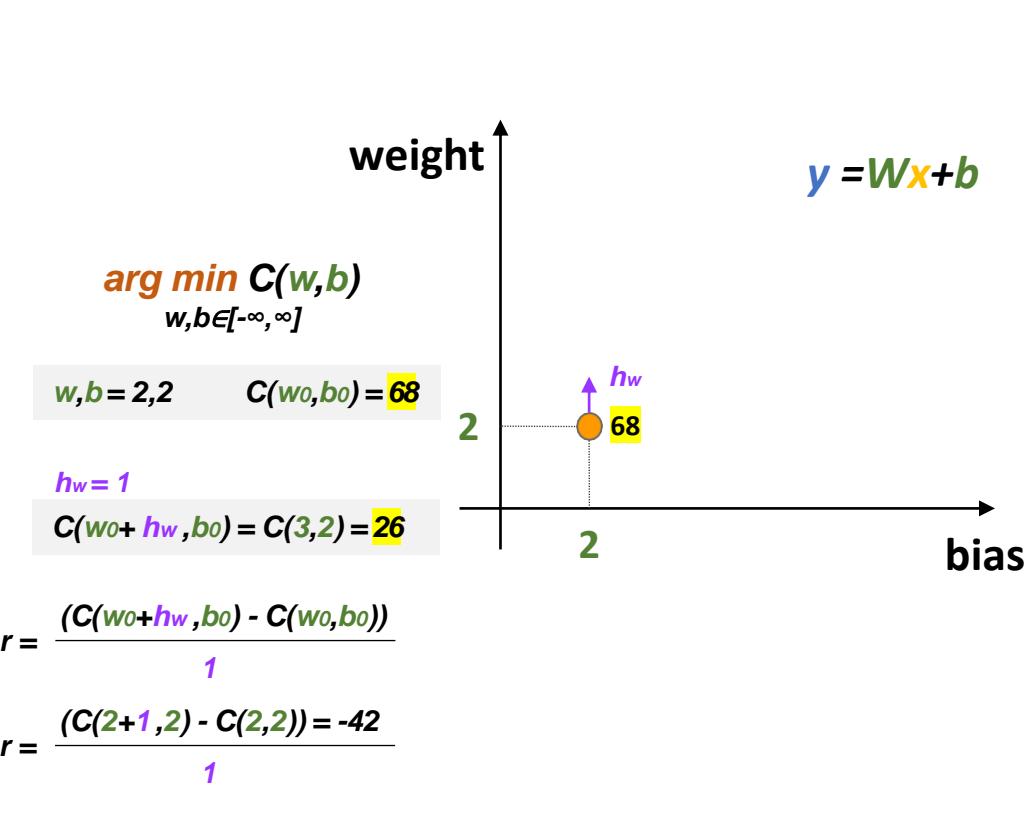
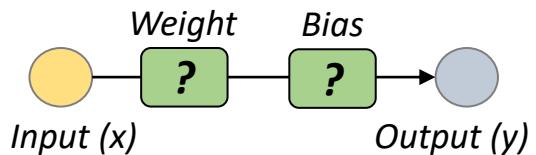
## Deep Learning with Neural Network - Optimizer

*Gradient!*

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



# 3 Deep Learning for NLP

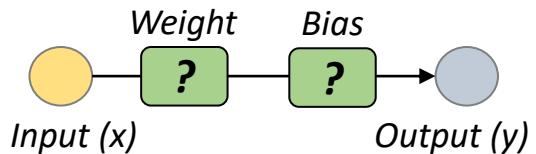
## Deep Learning with Neural Network - Optimizer

*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20

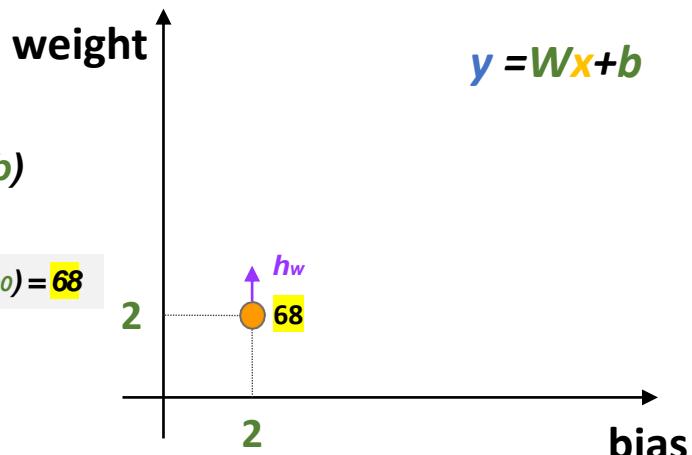


$$\arg \min_{w,b} C(w,b)$$

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

$$w,b = 2,2 \quad C(w_0,b_0) = 68$$

$$\begin{aligned}
 h_w &= 1, r = -42 \\
 h_w &= 0.1, r = -98 \\
 h_w &= 0.01, r = -104 \\
 h_w &= 0.001, r = -104
 \end{aligned}$$



# 3 Deep Learning for NLP

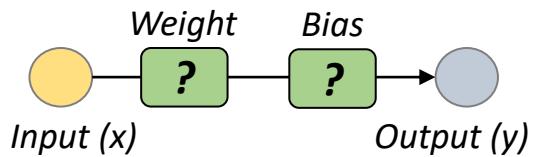
## Deep Learning with Neural Network - Optimizer

*Gradient!*

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



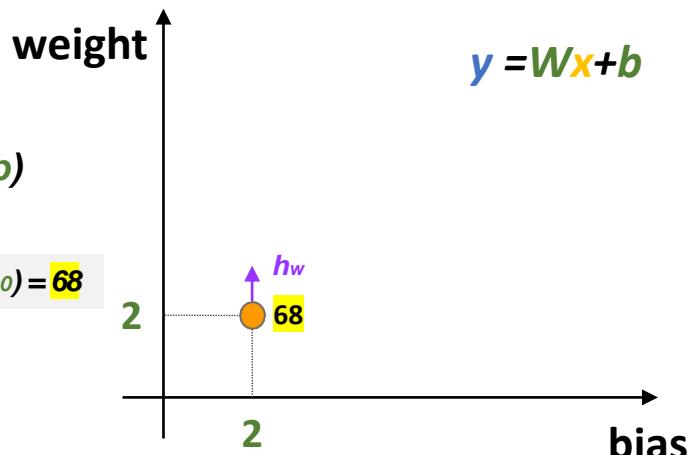
$$\arg \min_{w,b \in [-\infty, \infty]} C(w,b)$$

$$w,b=2,2 \quad C(w_0,b_0)=68$$

$$\begin{aligned} h_w &= 1, r = -42 \\ h_w &= 0.1, r = -98 \\ h_w &= 0.01, r = -104 \\ h_w &= 0.001, r = -104 \end{aligned}$$

*r = reward?*

$$h_w \rightarrow 0, \quad r = \frac{\partial C}{\partial w}(w_0,b_0)$$



# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

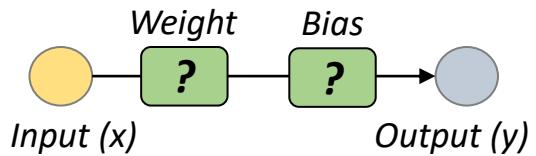
*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

weight      bias

x	y
1	0
5	16
6	20

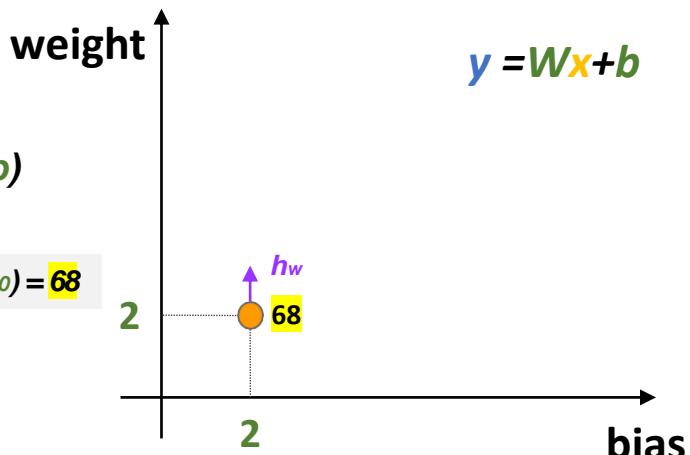


$$\arg \min_{w,b} C(w,b)$$

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

$$w,b = 2,2 \quad C(w_0,b_0) = 68$$

$$\frac{\partial C}{\partial w} = \frac{\partial \sum_n (y_n - \hat{y}_n)^2}{\partial w}$$



# 3 Deep Learning for NLP

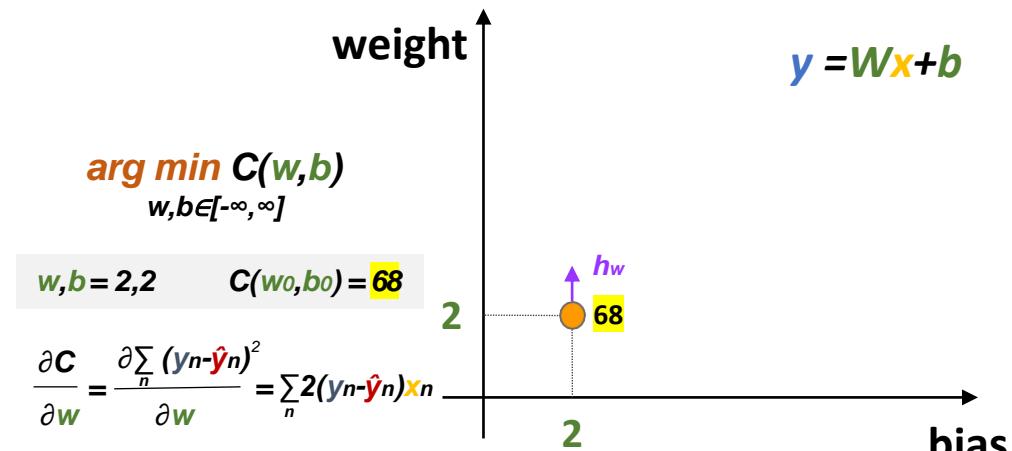
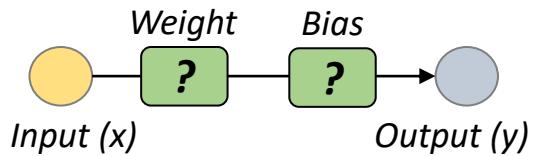
## Deep Learning with Neural Network - Optimizer

*Gradient!*

*Actual Data*

$$y = ? \times x + ?$$

x	y
1	0
5	16
6	20



# 3 Deep Learning for NLP

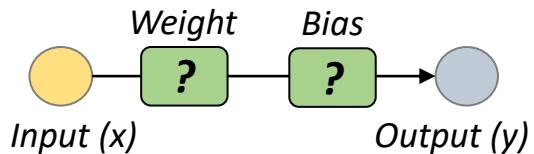
*Gradient!*

## Deep Learning with Neural Network - Optimizer

*Actual Data*

$$y = ? \times x + ?$$

x	y
1	0
5	16
6	20



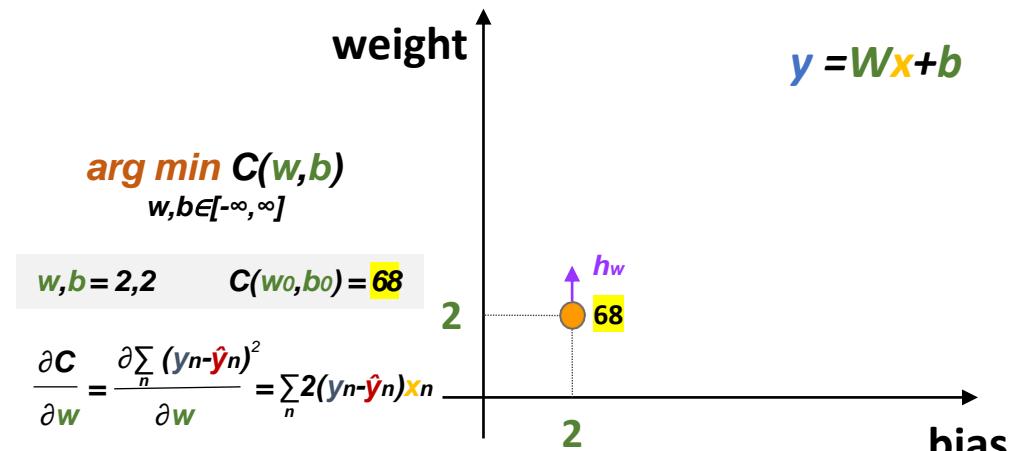
$$\arg \min_{w,b} C(w,b)$$

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

$$w,b = 2,2 \quad C(w_0,b_0) = 68$$

$$\frac{\partial C}{\partial w} = \frac{\partial \sum_n (y_n - \hat{y}_n)^2}{\partial w} = \sum_n 2(y_n - \hat{y}_n)x_n$$

$$h_w \rightarrow 0, r = \frac{\partial C}{\partial w}(w_0, b_0)$$



# 3 Deep Learning for NLP

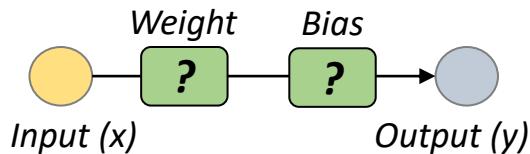
## Deep Learning with Neural Network - Optimizer

**Gradient!**

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



$$\arg \min_{w,b \in [-\infty, \infty]} C(w, b)$$

$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$\frac{\partial C}{\partial w} = \frac{\partial \sum_n (y_n - \hat{y}_n)^2}{\partial w} = \sum_n 2(y_n - \hat{y}_n)x_n$$

$$h_w \rightarrow 0, r = \frac{\partial C}{\partial w}(w_0, b_0) = 104$$

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

	<i>predicted</i>	<i>actual</i>		
x 🍎	$\hat{y}$ 🍚	y 🍌	$(y - \hat{y})$	$2(y - \hat{y})x$
1	4	0	-4	-8
5	12	16	4	40
6	14	20	6	72

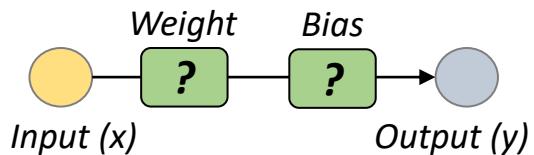
# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

*Actual Data*

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

x	y
1	0
5	16
6	20



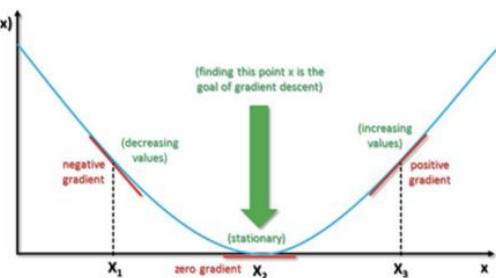
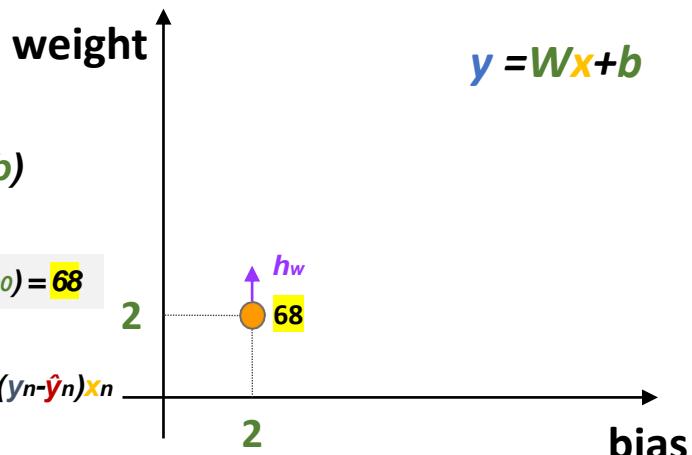
$$\arg \min C(w, b)$$

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

$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$\frac{\partial C}{\partial w} = \frac{\partial \sum_n (y_n - \hat{y}_n)^2}{\partial w} = \sum_n 2(y_n - \hat{y}_n) x_n$$

$$\frac{\partial C}{\partial b} = \frac{\partial \sum_n (y_n - \hat{y}_n)^2}{\partial b} = \sum_n 2(y_n - \hat{y}_n)$$



# 3 Deep Learning for NLP

## Deep Learning with Neural Network - Optimizer

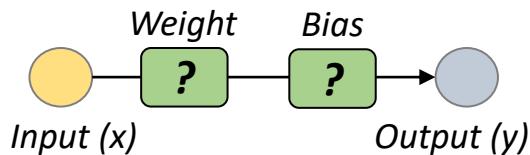
**Gradient!**

*Actual Data*

$$y = ? \ x + ?$$

weight      bias

x 🍎	y 🍌
1	0
5	16
6	20



$$\arg \min_{w,b \in [-\infty, \infty]} C(w,b)$$

$$w,b = 2,2 \quad C(w_0, b_0) = 68$$

$$h_w \rightarrow 0, r = \frac{\partial C}{\partial w} (w_0, b_0) = 104$$

$$h_b \rightarrow 0, r = \frac{\partial C}{\partial b} (w_0, b_0) = 12$$

$$y = 2 \ x + 2$$

weight      bias

	<i>predicted</i>	<i>actual</i>		
x 🍎	ŷ 🍎	y 🍌	(y - ŷ)	2(y - ŷ)
1	4	0	-4	-8
5	12	16	4	8
6	14	20	6	12

r = 12 because it is not multiplied by x like previous

# 3 Deep Learning for NLP

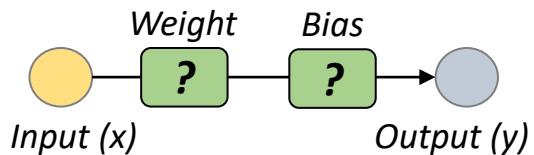
## Deep Learning with Neural Network - Optimizer

**Gradient!**

*Actual Data*

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

x 🍎	y 🍌
1	0
5	16
6	20



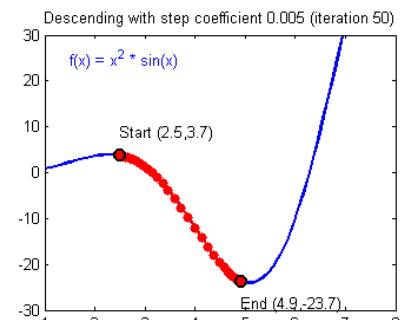
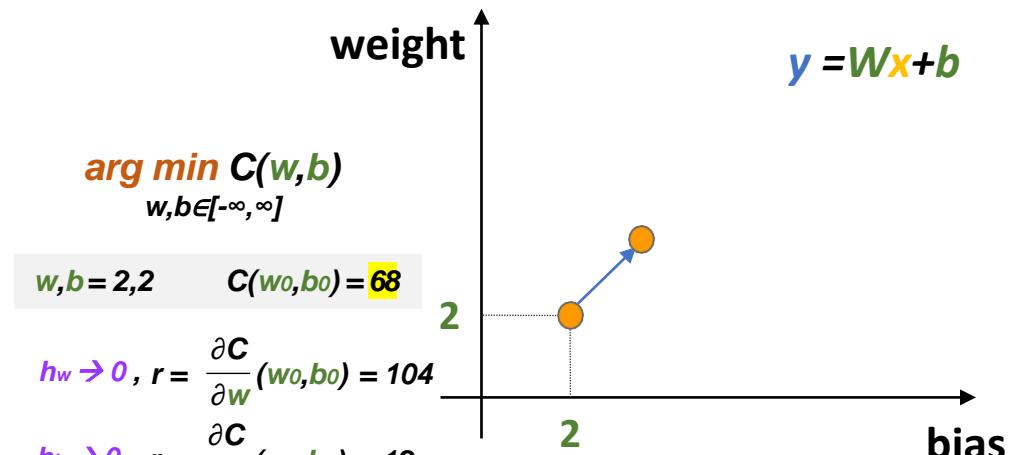
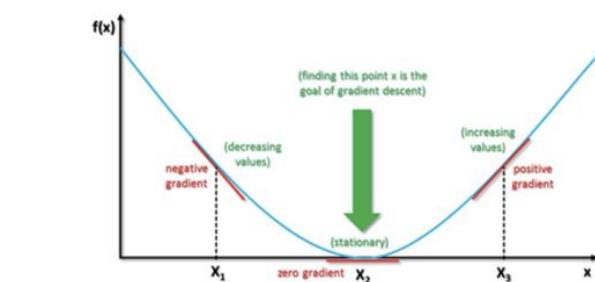
$$\arg \min C(w, b)$$

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

$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$h_w \rightarrow 0, r = \frac{\partial C}{\partial w} (w_0, b_0) = 104$$

$$h_b \rightarrow 0, r = \frac{\partial C}{\partial b} (w_0, b_0) = 12$$



## Deep Learning with Neural Network

### Data

x	y
1	0
5	16
6	20

### Model

$$y = ? \underset{\text{weight}}{x} + ? \underset{\text{bias}}$$

### Cost

$$C(w, b) = \sum_{n \in \{0, 1, 2\}} (y_n - \hat{y}_n)$$

### Optimizer

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

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



### System

$$y = \underset{\text{weight}}{4} \underset{x}{x} - \underset{\text{bias}}{4}$$

## 3 Deep Learning for NLP



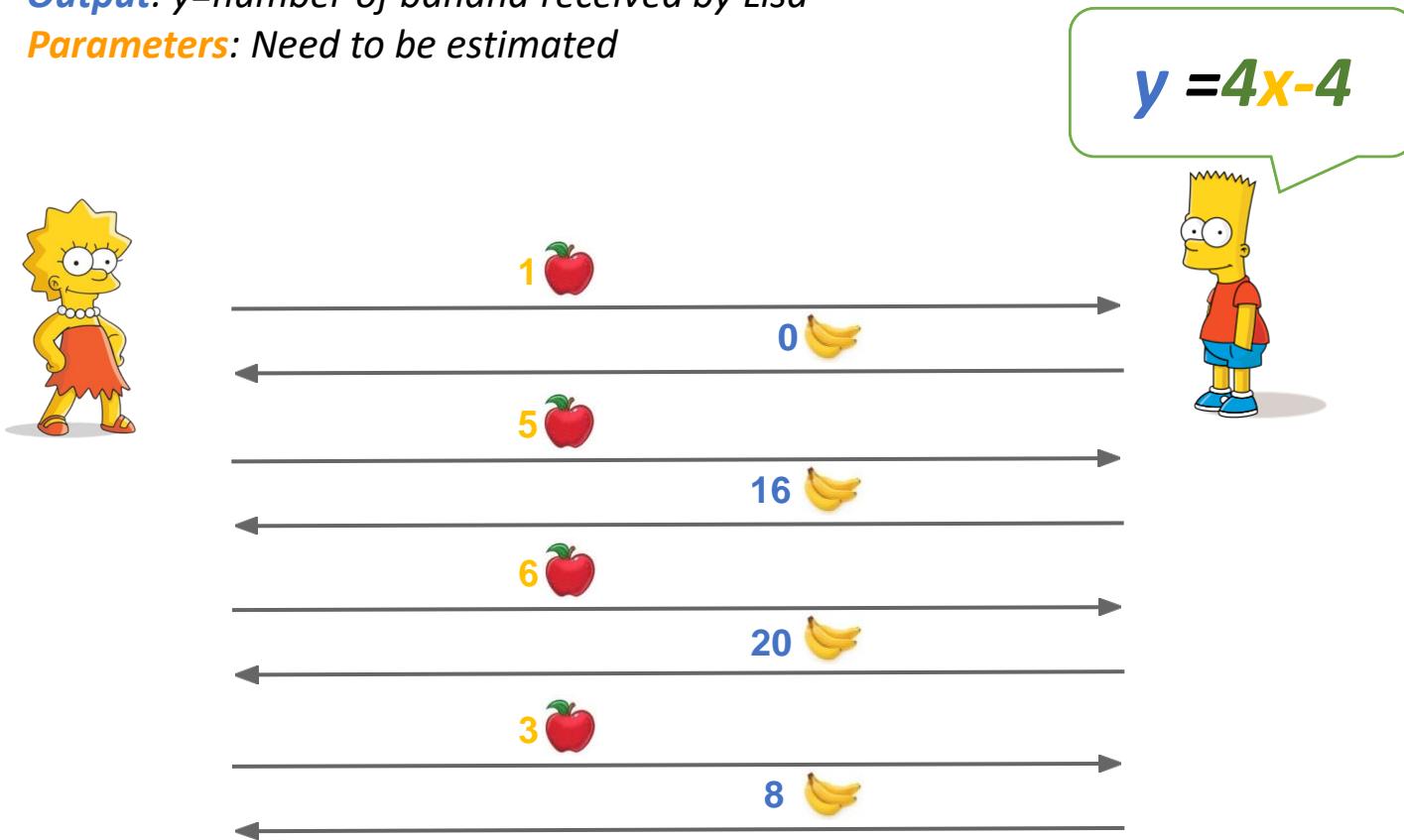
## 3 Deep Learning for NLP

## Deep Learning with Neural Network

**Input:**  $x$ =number of apple given by Lisa

**Output:**  $y$ =number of banana received by Lisa

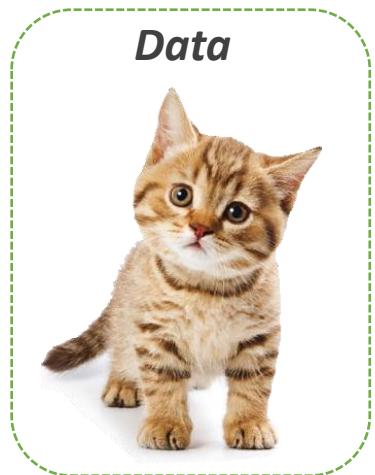
**Parameters:** Need to be estimated



## 3 Deep Learning for NLP

## Deep Learning with Neural Network

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



*Millions of Parameters*  
*Millions of Samples*

# 3 Deep Learning for NLP

## Deep Learning with Neural Network

$$y = W_1 \overset{\text{Vector1}}{X_1} + W_2 \overset{\text{Vector2}}{X_2} + W_3 X_3 + W_4 X_4 + \dots + W_n X_n + b$$



*Millions of Parameters*  
*Millions of Samples*

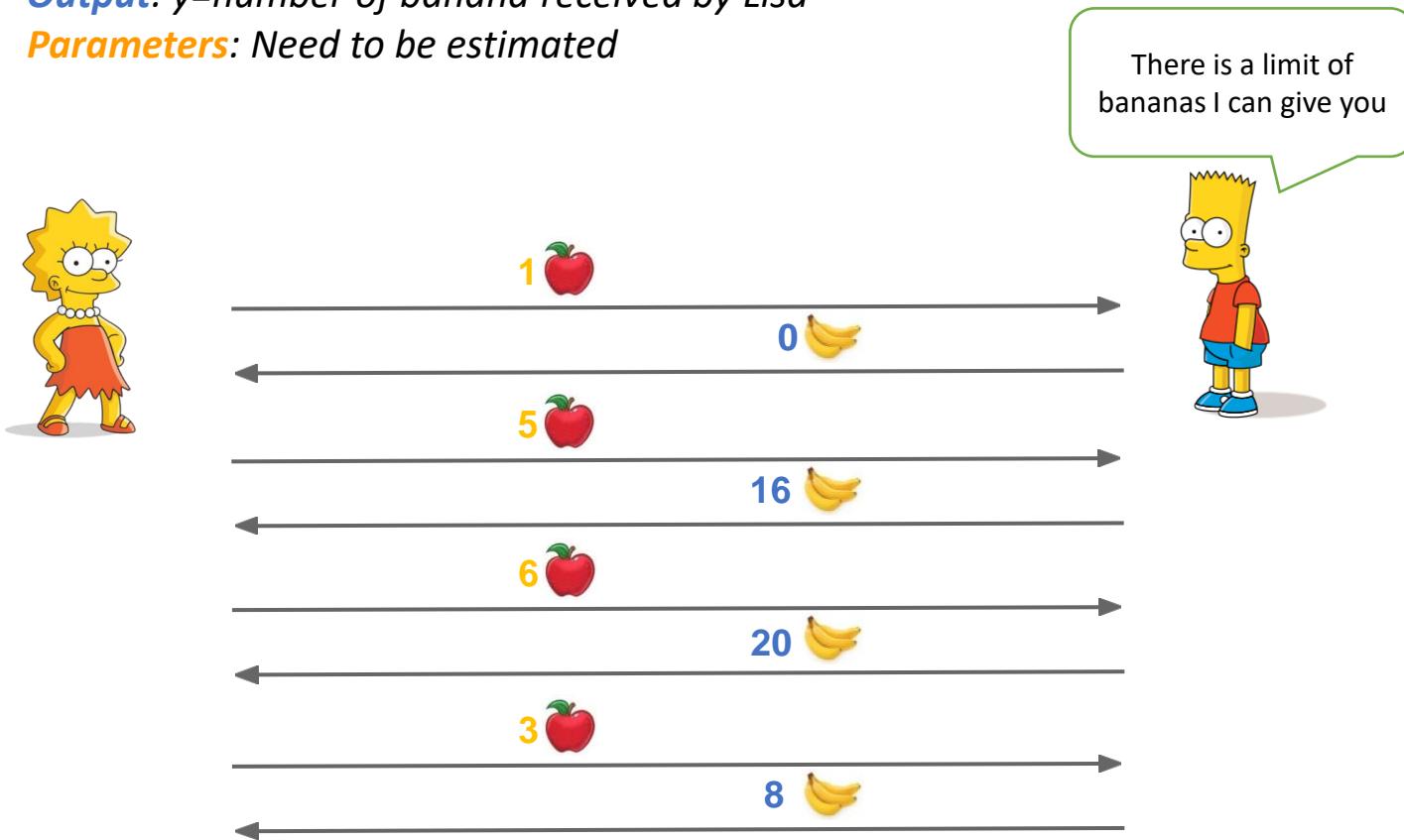
# 3 Deep Learning for NLP

## 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



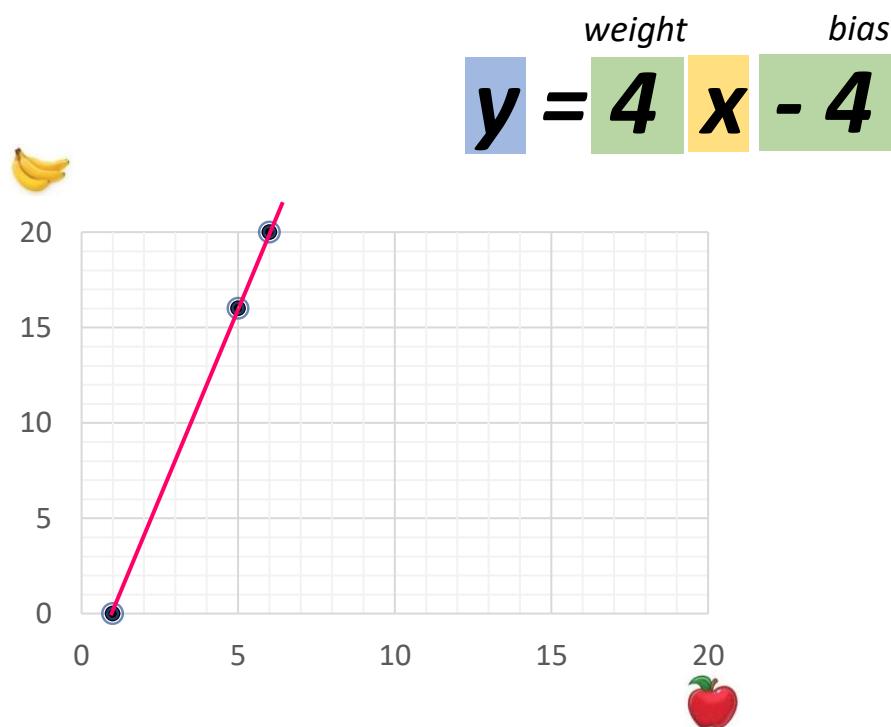
## 3 Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

*Data*

x 🍎	y 🍌
1	0
5	16
6	20



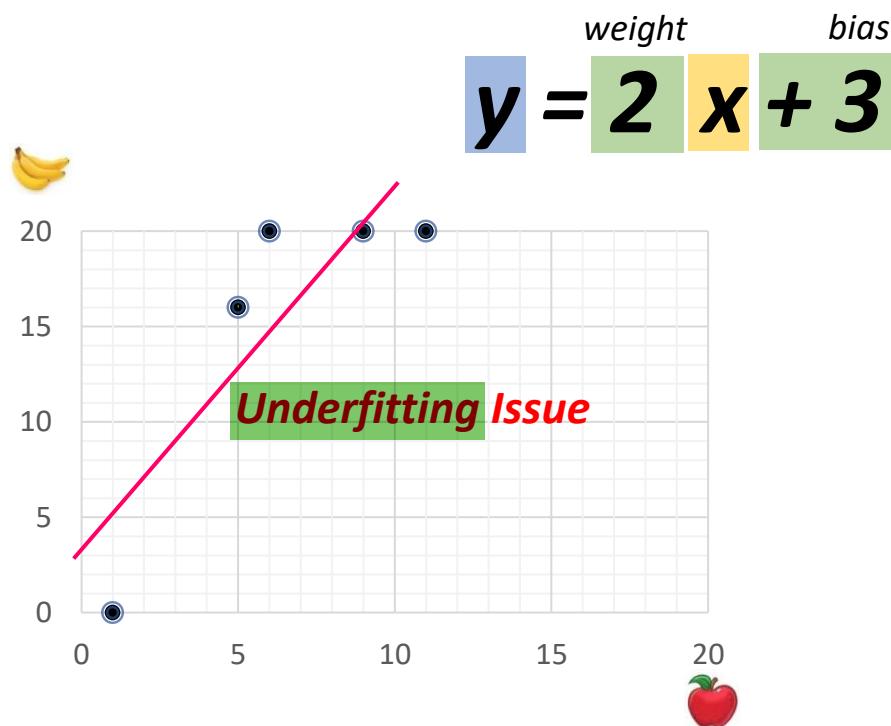
# 3 Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

*Data*

x 🍎	y 🍌
1	0
5	16
6	20
9	20
11	20



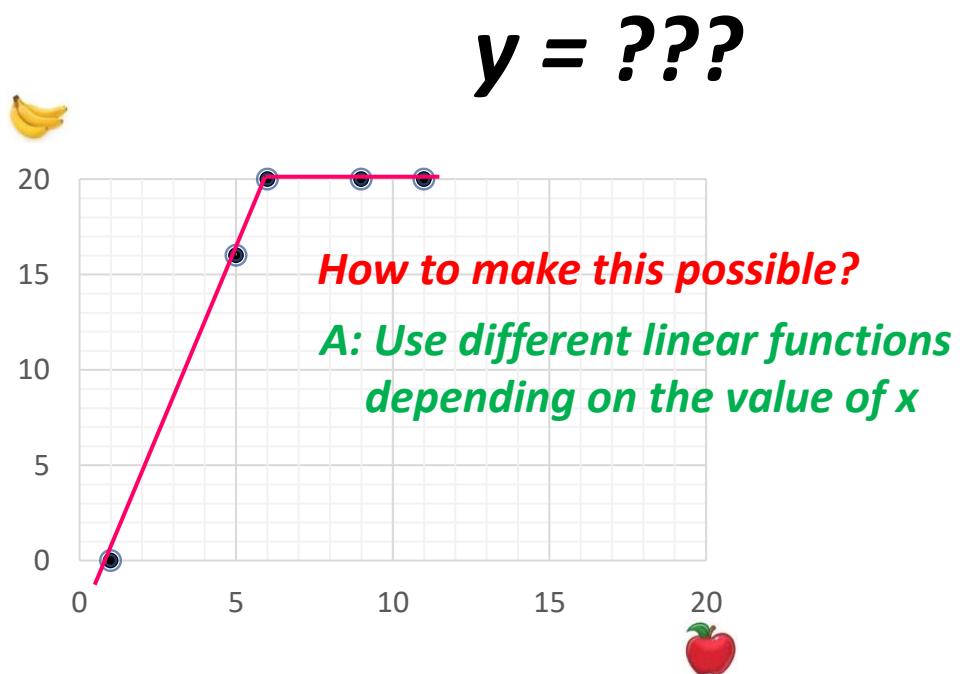
# 3 Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

*Data*

x 🍎	y 🍌
1	0
5	16
6	20
9	20
11	20



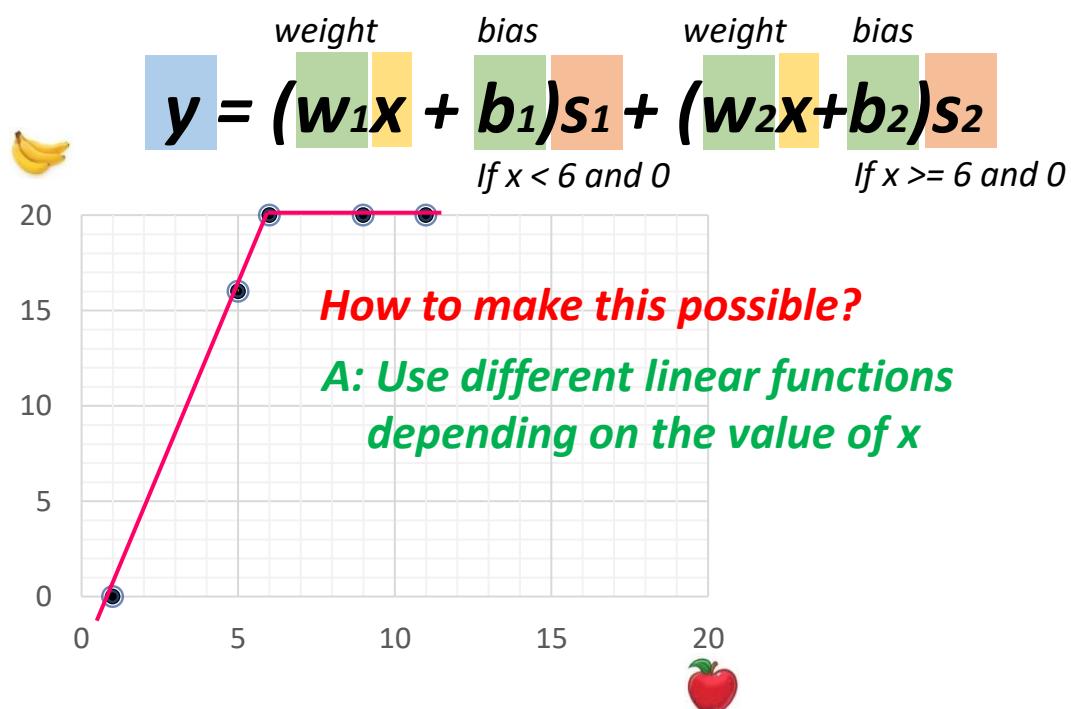
# 3 Deep Learning for NLP

## Deep Learning with Neural Network

Nonlinear Neural Network

*Data*

x	y
1	0
5	16
6	20
9	20
11	20



# 3 Deep Learning for NLP

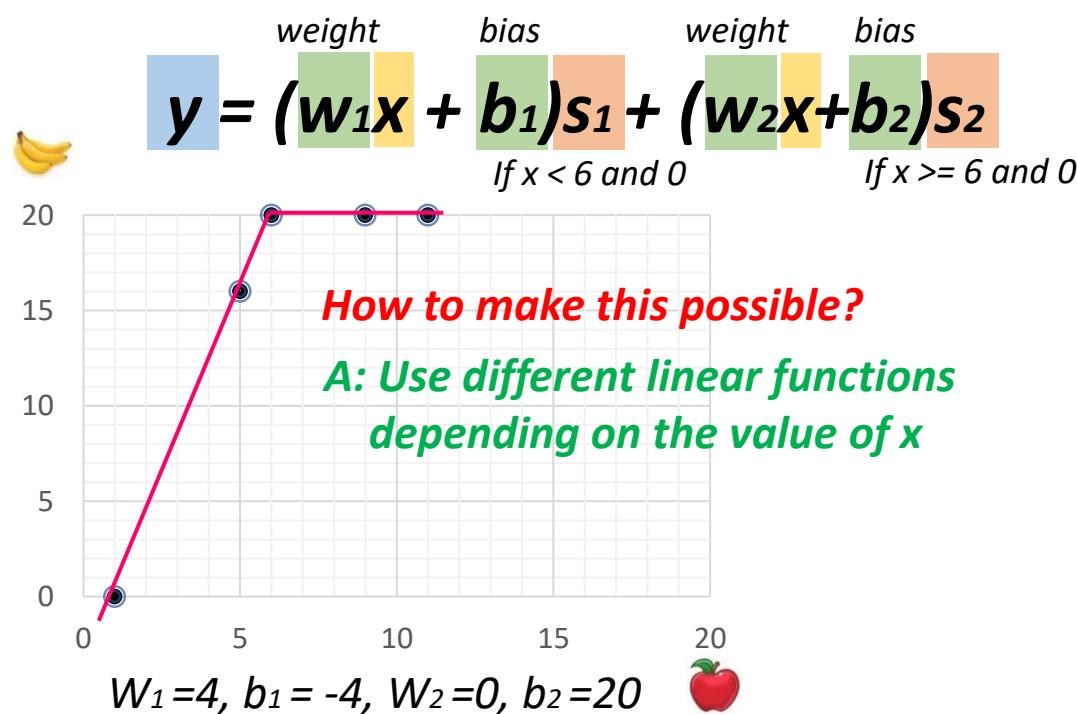
## Deep Learning with Neural Network

multiple linear functions

Nonlinear Neural Network

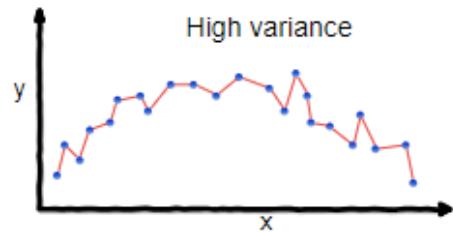
*Data*

x	y
1	0
5	16
6	20
9	20
11	20

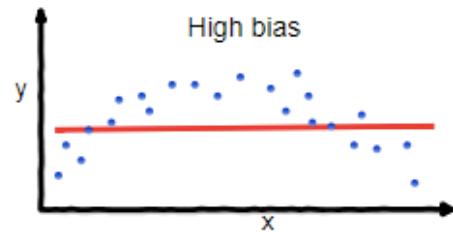


## 3 Deep Learning for NLP

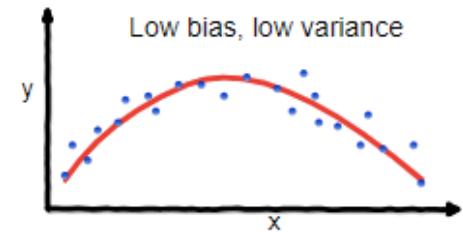
## Deep Learning with Neural Network



overfitting

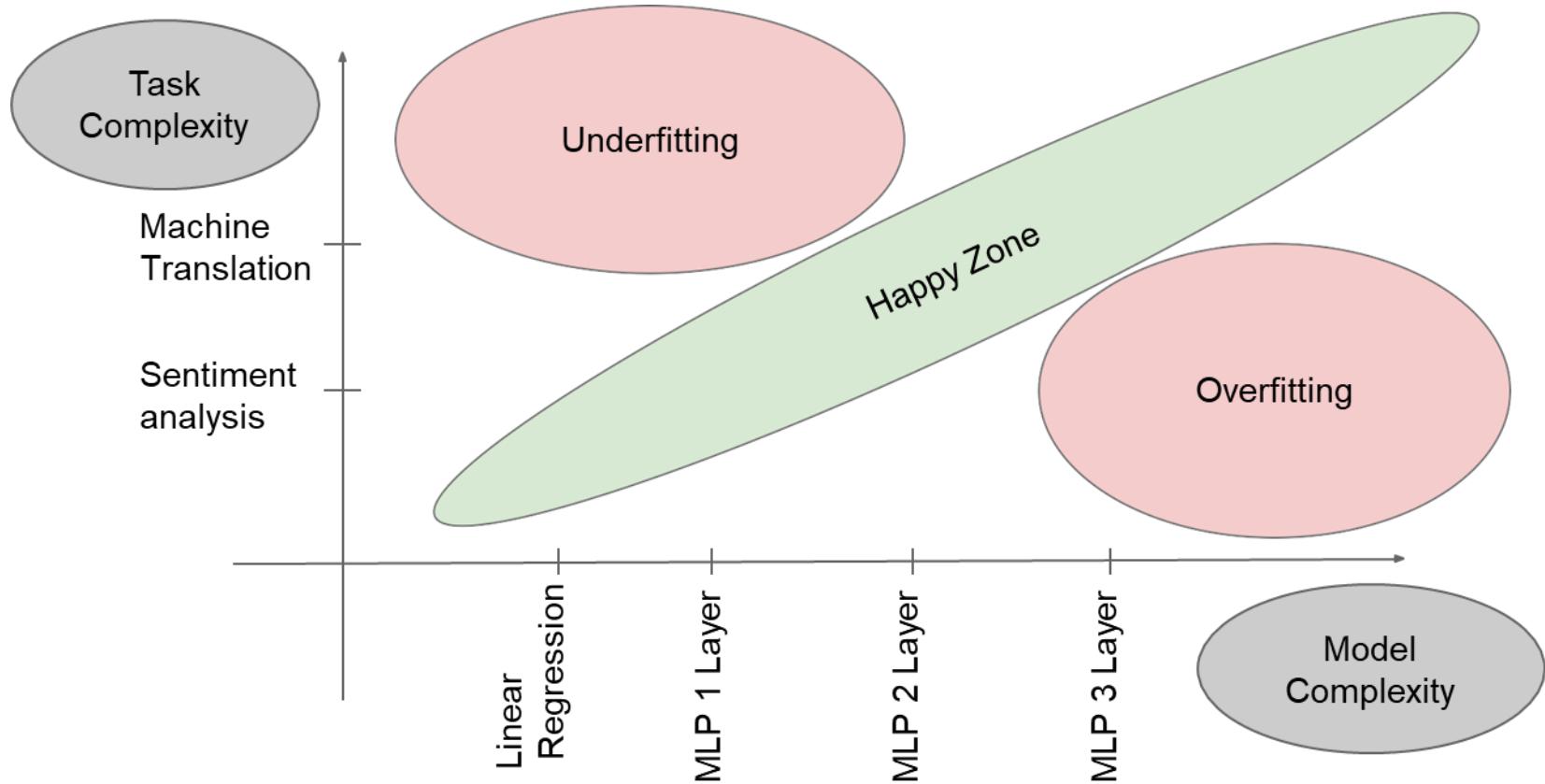


underfitting

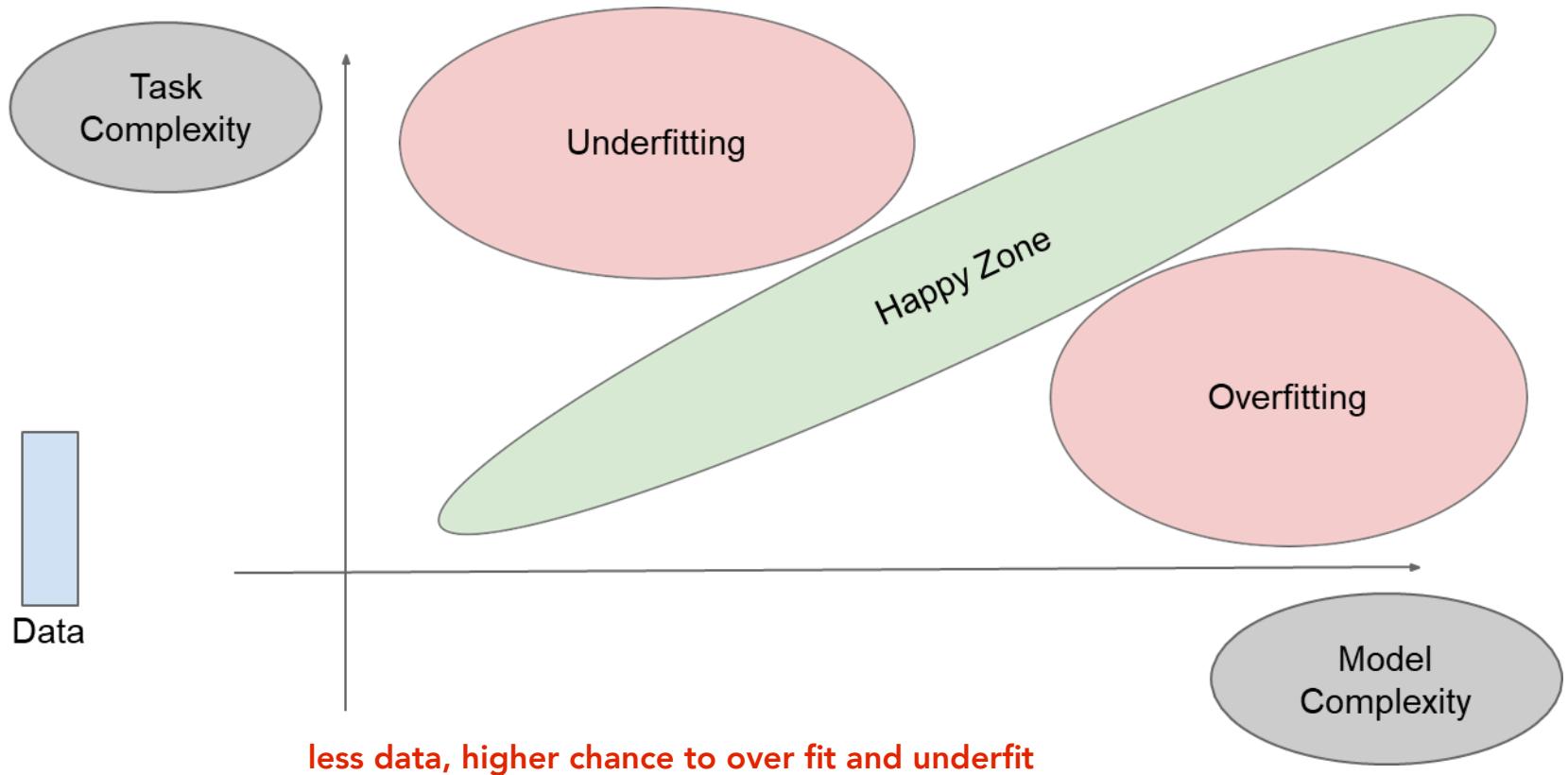


Good balance

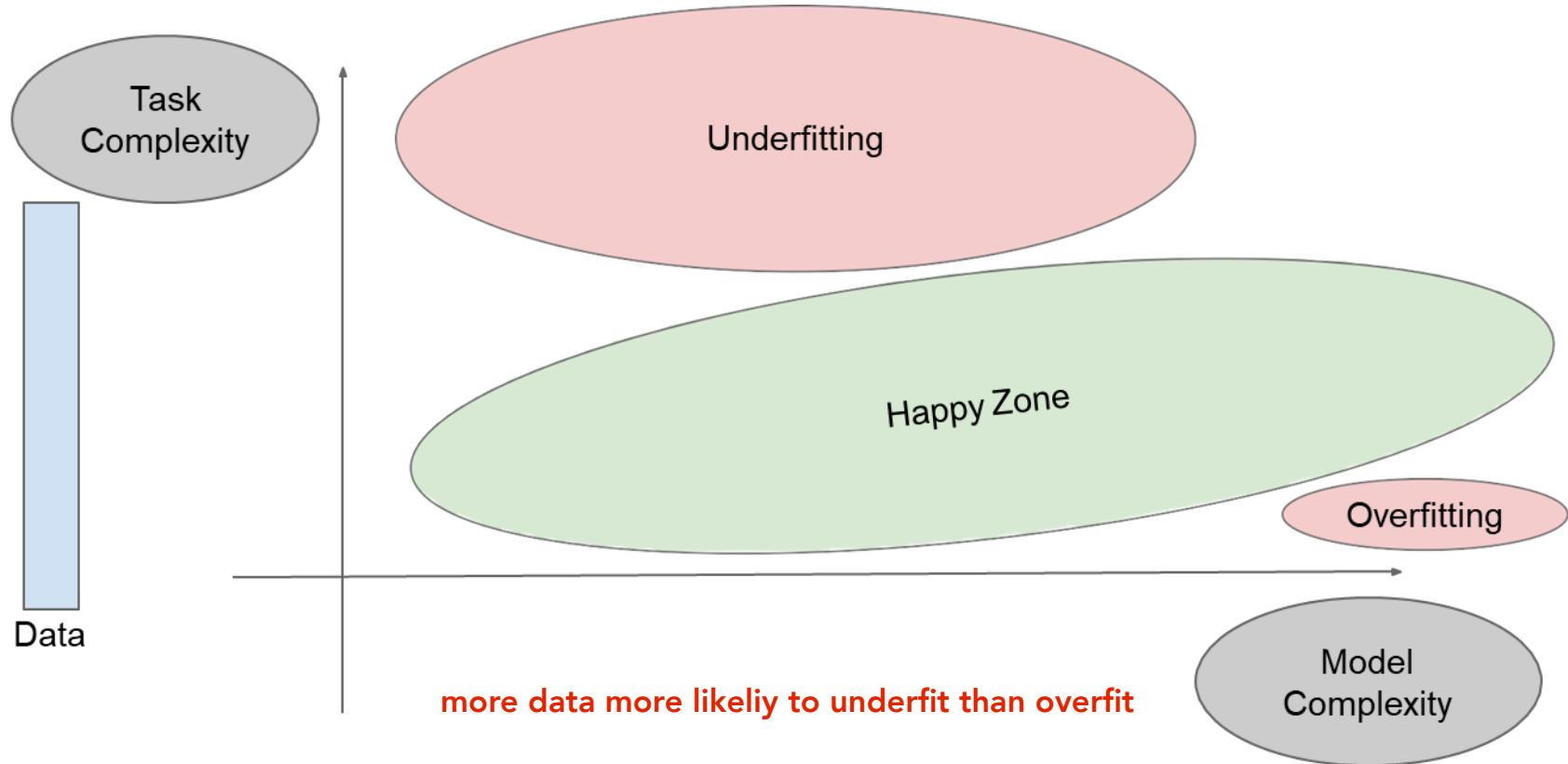
## 3 Deep Learning for NLP



# 3 Deep Learning for NLP

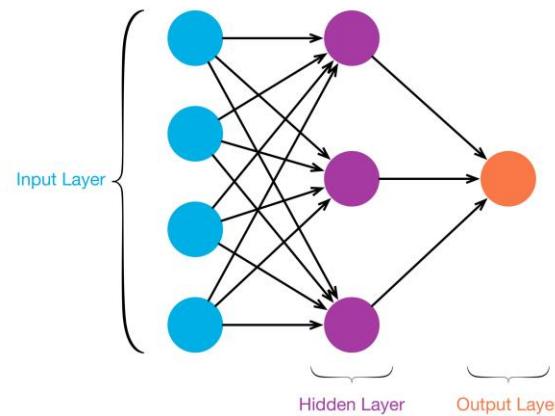
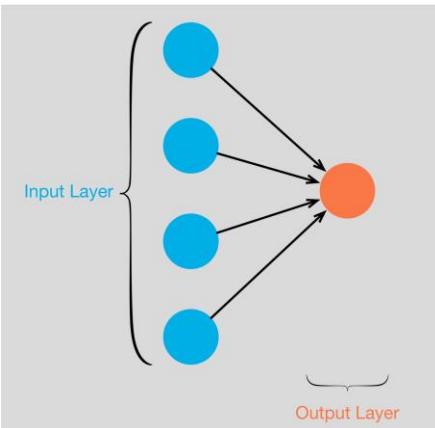
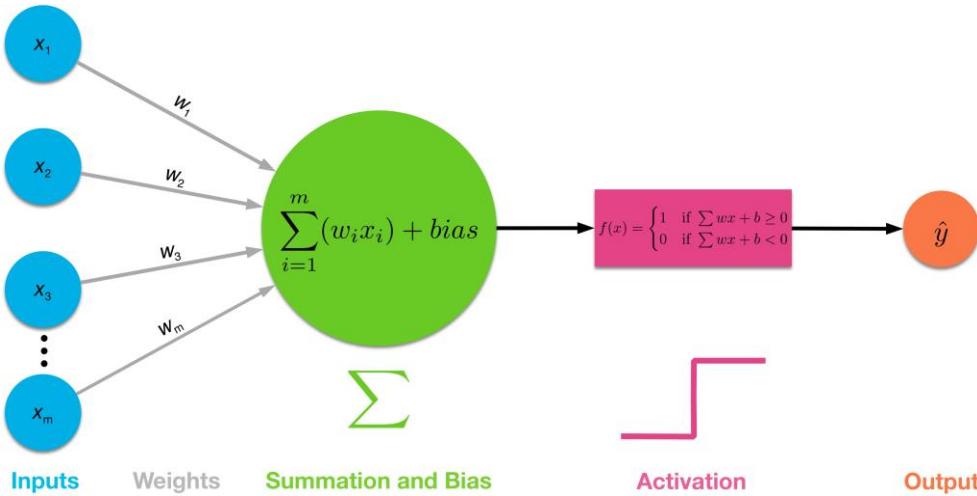


## 3 Deep Learning for NLP



# 3 Deep Learning for NLP

## Single Neuron VS Multilayer

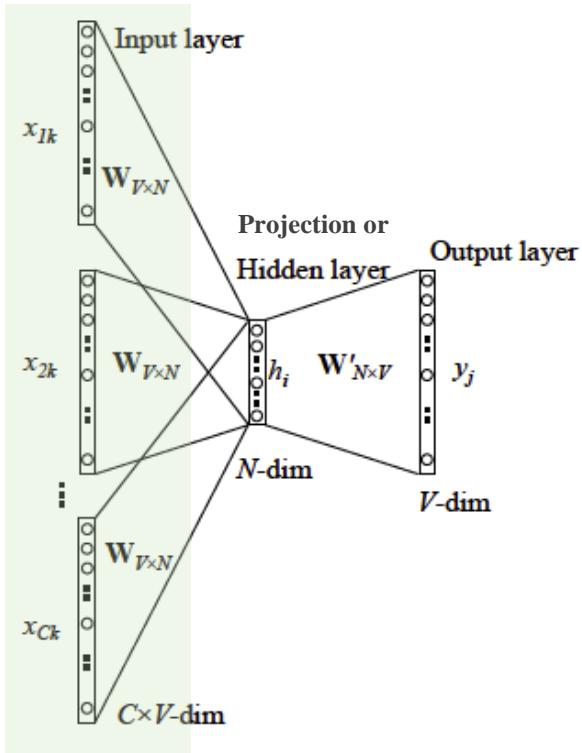


# 3 Deep Learning for NLP

## CBOB – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOW Training** (Review your understanding with equations)



1. Initialise each word in a one-hot vector form.

$$x_k = [0, \dots, 0, 1, 0, \dots, 0]$$

2. Use context words ( $2m$ , based on window size =  $m$ ) as input of the Word2Vec-CBOW model.

$$(x^{c-m}, x^{c-m+1}, \dots, x^{c-1}, x^{c+1}, \dots, x^{c+m-1}, x^{c+m}) \in \mathbb{R}^{|V|}$$

3. Has two Parameter Matrices:

1) Parameter Matrix (from Input Layer to Hidden/Projection Layer)  
 $W \in \mathbb{R}^{V \times N}$

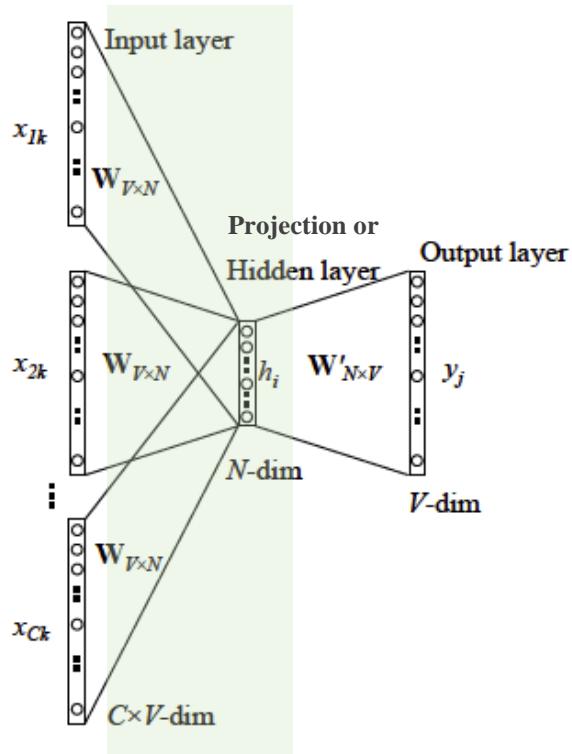
2) Parameter Matrix (to Output Layer)  
 $W' \in \mathbb{R}^{N \times V}$

# 3 Deep Learning for NLP

## CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOW Training** (Review your understanding with equations)



4. Initial words are represented in one hot vector so multiplying a **one hot vector** with  $W_{V \times N}$  will give you a  $1 \times N$  (embedded word) vector.

$$\text{e.g. } [0 \ 1 \ 0 \ 0] \times \begin{bmatrix} 10 & 2 & 18 \\ 15 & 22 & 3 \\ 25 & 11 & 19 \\ 4 & 7 & 22 \end{bmatrix} = [15 \ 22 \ 3]$$

$$(\mathbf{v}_{c-m} = Wx^{c-m}, \dots, \mathbf{v}_{c+m} = Wx^{c+m}) \in \mathbb{R}^n$$

5. **Average** those  $2m$  embedded vectors to calculate the value of the Hidden Layer.

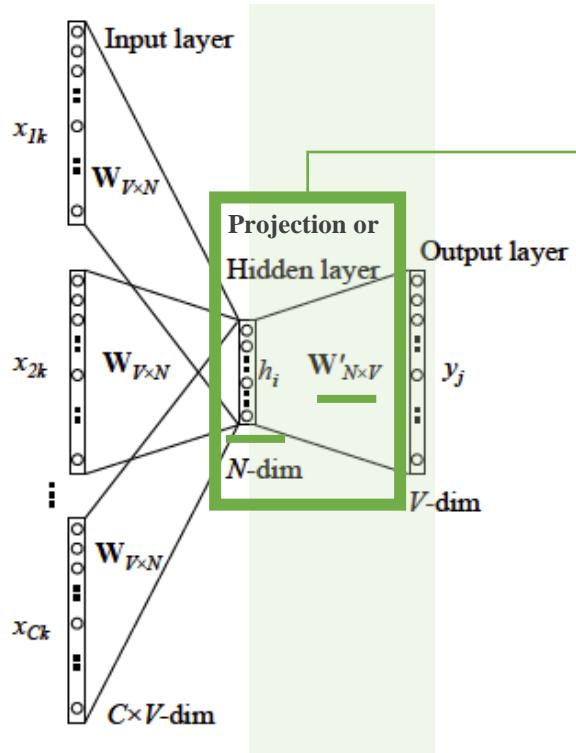
$$\hat{v} = \frac{\mathbf{v}_{c-m} + \mathbf{v}_{c-m+1} + \dots + \mathbf{v}_{c+m}}{2m}$$

# 3 Deep Learning for NLP

## CBOV – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOV Training** (Review your understanding with equations)



6. Calculate the score value for the output layer. The higher score is produced when words are closer.  
 $\mathbf{z} = \mathbf{W}' \times \hat{\mathbf{v}} \in \mathbb{R}^{|V|}$

7. Calculate the probability using softmax  
 $\hat{y} = \text{softmax}(\mathbf{z}) \in \mathbb{R}^{|V|}$

8. Train the parameter matrix using objective function.

$$H(\hat{y}, y) = - \sum_{j=1}^{|V|} y_j \log(\hat{y}_j)$$

\* Focus on minimising the value

We use an one-hot vector (one 1, the rest 0) so it will be calculated in only one.

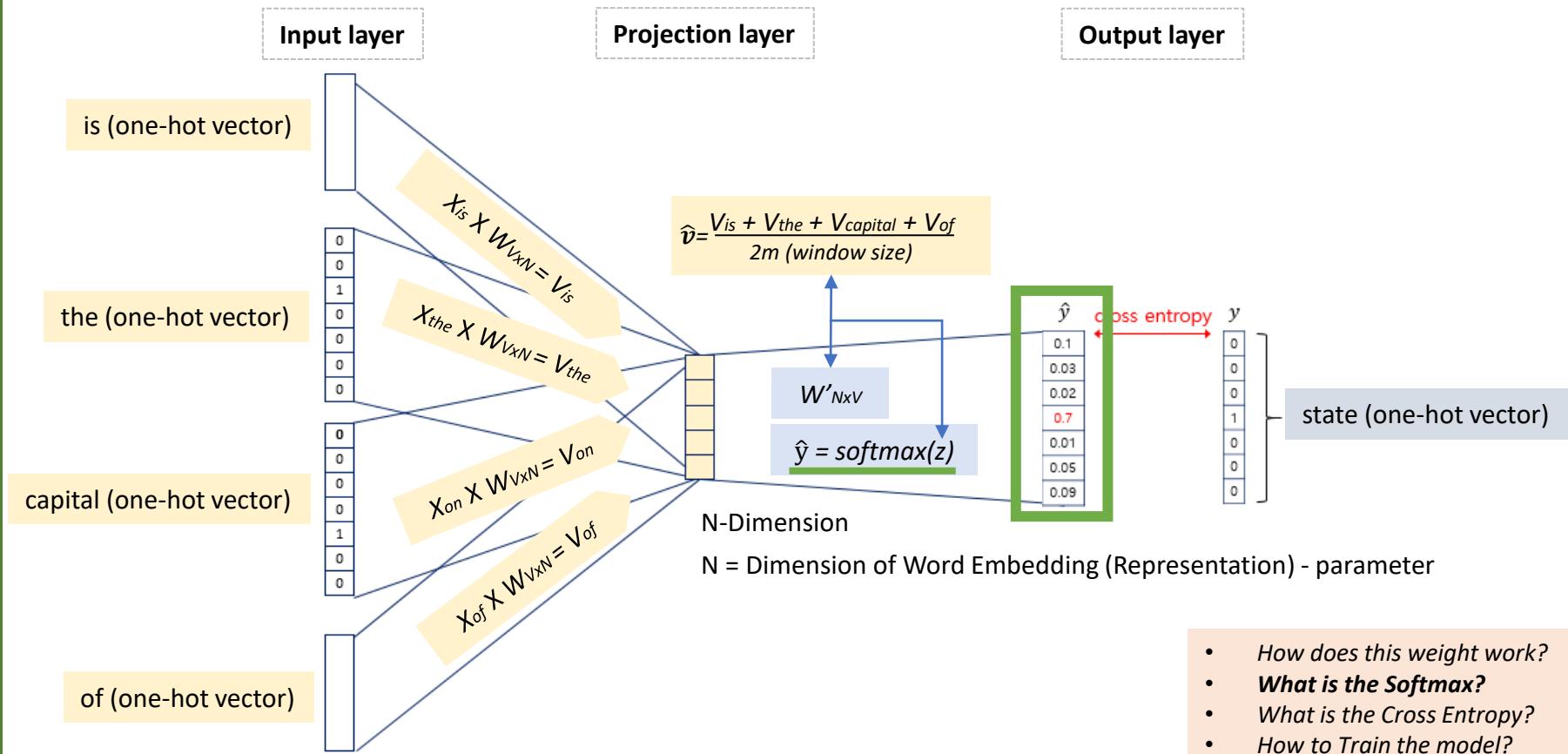
$$H(\hat{y}, y) = -y_j \log(\hat{y}_j)$$

# 3 Deep Learning for NLP

## CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words

**Sentence:** “Sydney is the state capital of NSW”

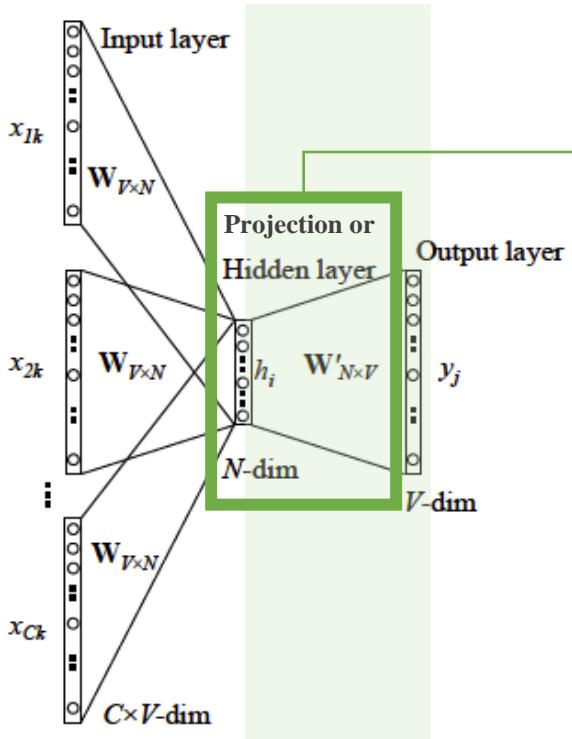


# 3 Deep Learning for NLP

## CBOV – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOV Training** (Review your understanding with equations)



6. Calculate the score value for the output layer. The higher score is produced when words are closer.

$$\mathbf{z} = \mathbf{W}' \times \hat{\mathbf{v}} \in \mathbb{R}^{|V|}$$

7. Calculate the probability using softmax

$$\hat{y} = \text{softmax}(\mathbf{z}) \in \mathbb{R}^{|V|}$$

*The softmax is an operator that will be used frequently. It transforms a vector into a vector whose  $i$ -th component is:*

$$\frac{e^{\hat{y}_i}}{\sum_{j=1}^{|V|} e^{\hat{y}_j}}$$

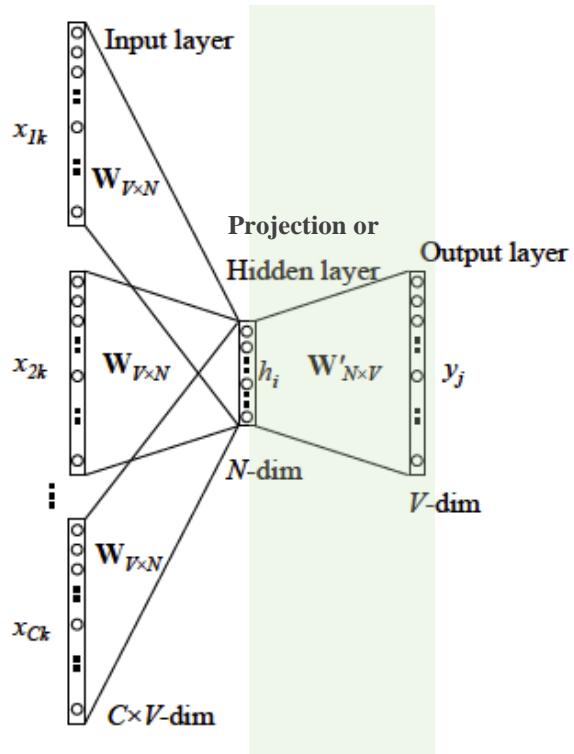
- Exponentiate to make positive
- Dividing by  $\sum_{j=1}^{|V|} e^{\hat{y}_j}$  normalizes the vector ( $\sum_{j=1}^n \hat{y}_j = 1$ ) to give probability

# 3 Deep Learning for NLP

## CBOV – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOV Training** (Review your understanding with equations)



6. Calculate the score value for the output layer. The higher score is produced when words are closer.

$$\mathbf{z} = \mathbf{W}' \times \hat{\mathbf{v}} \in \mathbb{R}^{|V|}$$

7. Calculate the probability using softmax

$$\hat{y} = \text{softmax}(\mathbf{z}) \in \mathbb{R}^{|V|}$$

8. Train the parameter matrix using **objective function**.

$$H(\hat{y}, y) = - \sum_{j=1}^{|V|} y_j \log(\hat{y}_j) \quad \boxed{\text{Cross Entropy}}$$

\* Focus on **minimising** the value

We use an one-hot vector (one 1, the rest 0) so it will be calculated in only one.

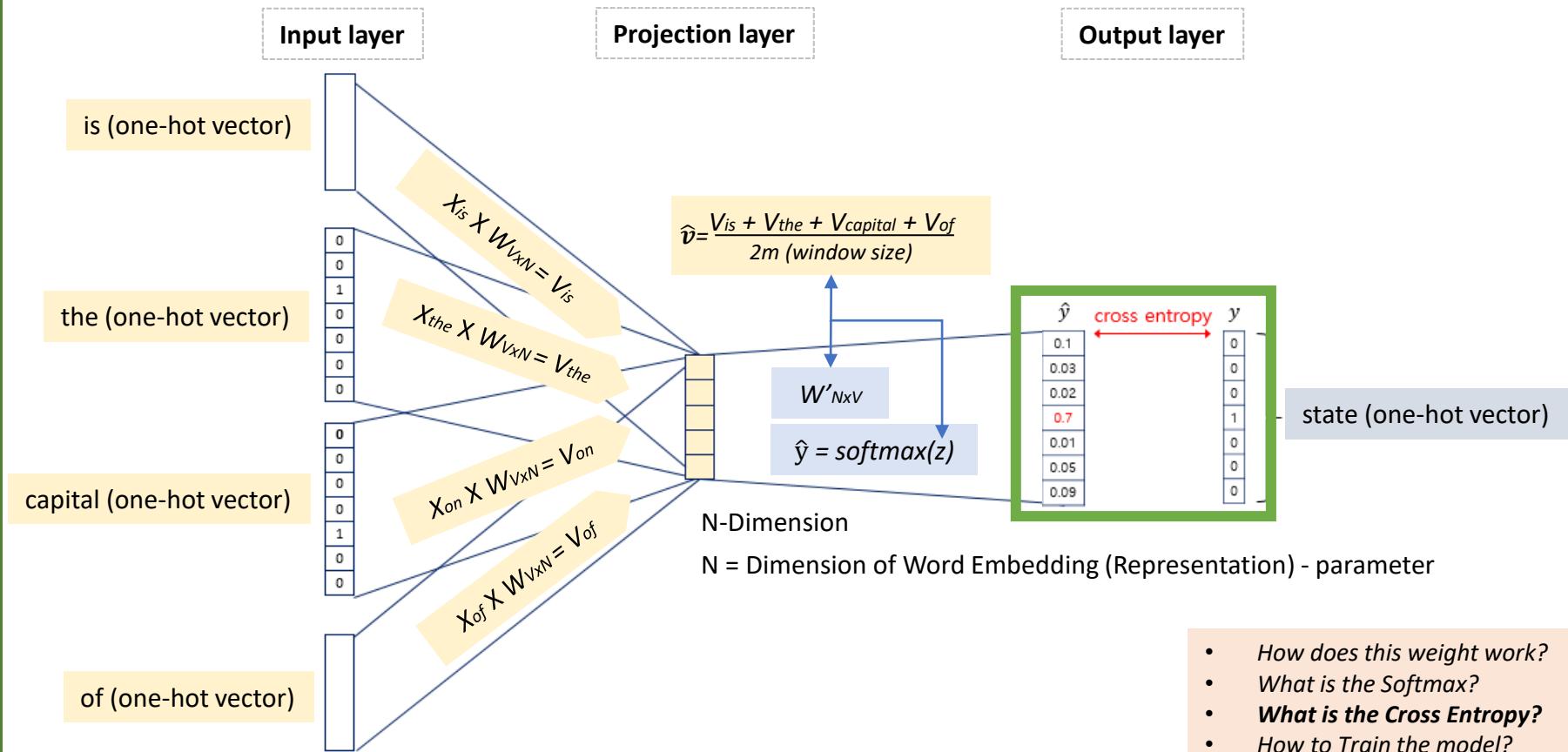
$$H(\hat{y}, y) = -y_j \log(\hat{y}_j)$$

# 3 Deep Learning for NLP

## CBOV – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words

**Sentence:** “Sydney is the state capital of NSW”

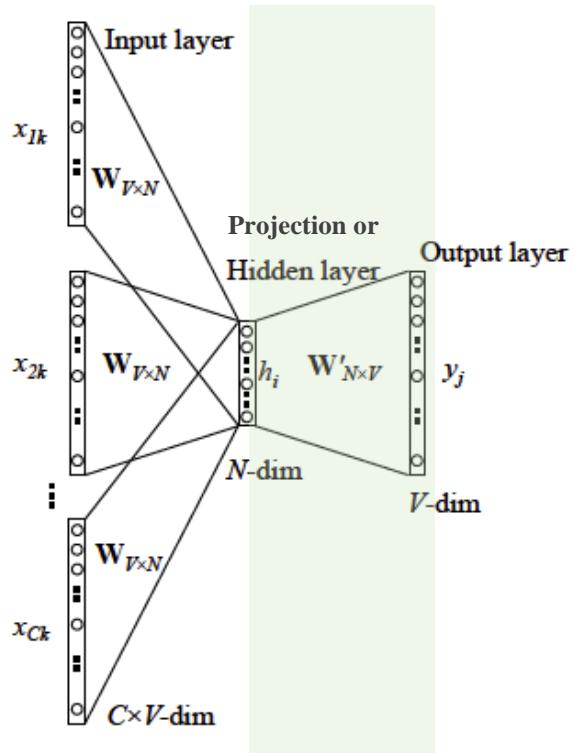


# 3 Deep Learning for NLP

## CBOB – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

**Summary of CBOB Training** (Review your understanding with equations)



8-1. Optimization Objective Function can be presented:

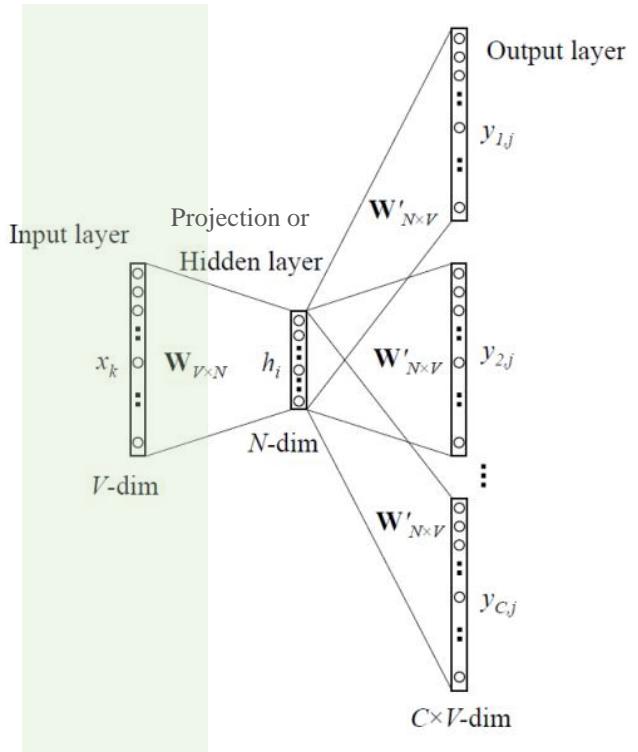
$$\begin{aligned}
 \text{minimize } J &= -\log P(w_c | w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m}) \\
 &= -\log P(u_c | \hat{v}) \\
 &= -\log \frac{\exp(u_c^T \hat{v})}{\sum_{j=1}^{|V|} \exp(u_j^T \hat{v})} \\
 &= -u_c^T \hat{v} + \log \sum_{j=1}^{|V|} \exp(u_j^T \hat{v})
 \end{aligned}$$

$u_i$ =the output vector representation of word  $w_i$

## Skip Gram – Neural Network Architecture

Predict context (“outside”) words (position independent) given center word

**Summary of Skip Gram Training** (Review your understanding with equations)



1. Initialise the centre word in a one-hot vector form.

$$\mathbf{x}_k = [0, \dots, 0, 1, 0, \dots, 0]$$

$$\mathbf{x} \in \mathbb{R}^{|V|}$$

2. Has two Parameter Matrices:

- 1) Parameter Matrix (from Input Layer to Hidden/Projection Layer)  
 $\mathbf{W} \in \mathbb{R}^{V \times N}$

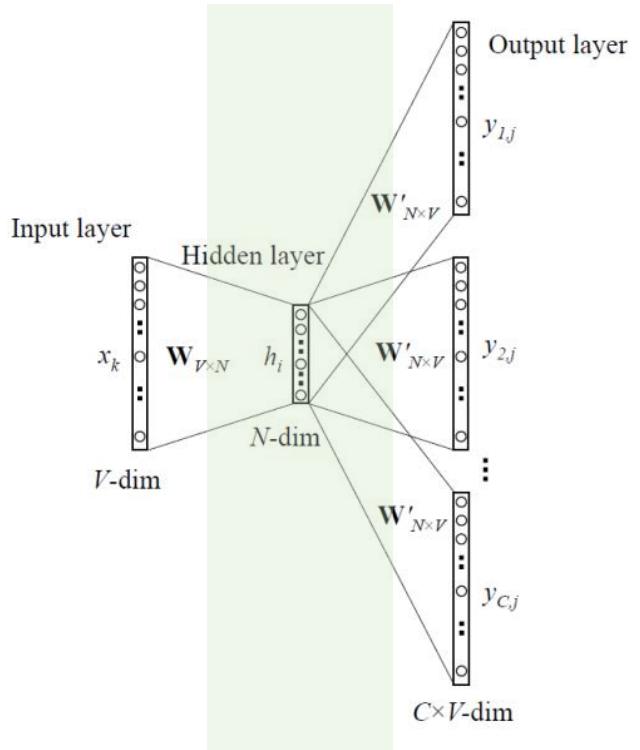
2) Parameter Matrix (to Output Layer)

$$\mathbf{W}' \in \mathbb{R}^{N \times V}$$

## Skip Gram – Neural Network Architecture

Predict context (“outside”) words (position independent) given center word

**Summary of Skip Gram Training** (Review your understanding with equations)



3. Initial words are represented in one hot vector so multiplying a **one hot vector** with  $\mathbf{W}_{V \times N}$  will give you a  $1 \times N$  (embedded word) vector.

$$\text{e.g. } [0 \ 1 \ 0 \ 0] \times \begin{bmatrix} 10 & 2 & 18 \\ 15 & 22 & 3 \\ 25 & 11 & 19 \\ 4 & 7 & 22 \end{bmatrix} = [15 \ 22 \ 3]$$

$$\mathbf{v}_c = \mathbf{W}_x \in \mathbb{R}^n \text{ (as there is only one input)}$$

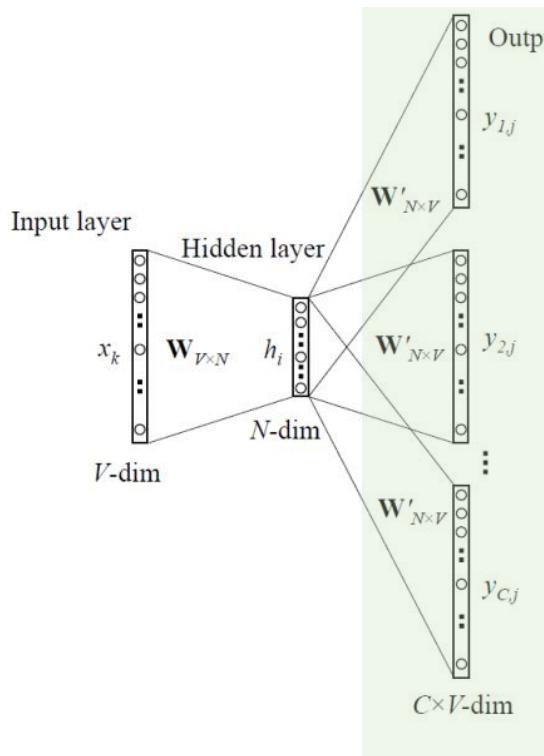
4. Calculate the score value for the output layer by multiplying the parameter matrix  $\mathbf{W}'$   
 $\mathbf{z} = \mathbf{W}' \mathbf{v}_c$

# Prediction based Word representation

## Skip Gram – Neural Network Architecture

Predict context (“outside”) words (position independent) given center word

**Summary of Skip Gram Training** (Review your understanding with equations)



5. Calculate the probability using softmax  
 $\hat{y} = \text{softmax}(\mathbf{z})$

6. Calculate  $2m$  probabilities as we need to predict  **$2m$  context words.**

$$\hat{y}_{c-m}, \dots, \hat{y}_{c-1}, \hat{y}_{c+1}, \dots, \hat{y}_{c+m}$$

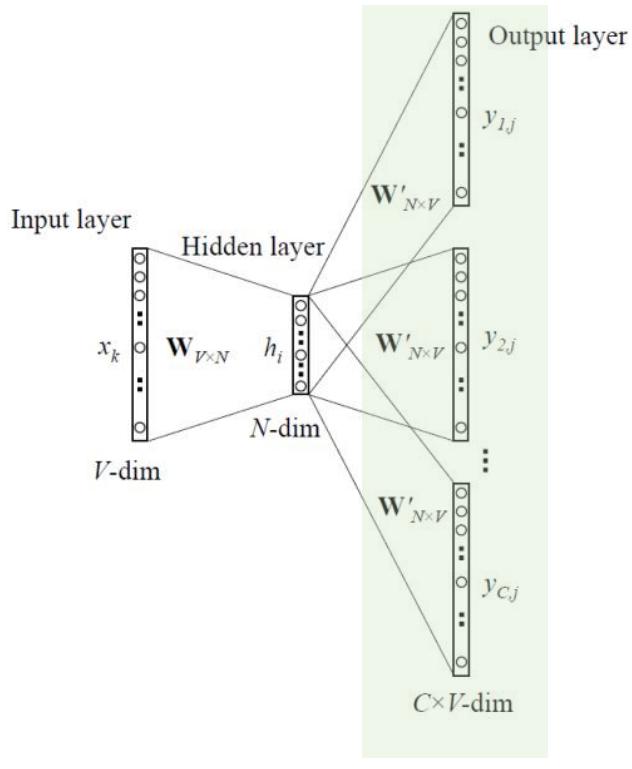
and compare with the ground truth (one-hot vector)  
 $y^{(c-m)}, \dots, y^{(c-1)}, y^{(c+1)}, \dots, y^{(c+m)}$

# Prediction based Word representation

## Skip Gram – Neural Network Architecture

Predict context (“outside”) words (position independent) given center word

**Summary of Skip Gram Training** (Review your understanding with equations)



8. As in CBOW, use an objective function for us to evaluate the model. A key difference here is that we invoke a **Naïve Bayes assumption to break out the probabilities**. It is a strong naïve conditional independence assumption. Given the centre word, all output words are completely **independent**.

$$\begin{aligned}
 \text{minimize } J &= -\log P(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c) \\
 &= -\log \prod_{j=0, j \neq m}^{2m} P(w_{c-m+j} | w_c) \\
 &= -\log \prod_{j=0, j \neq m}^{2m} \frac{\exp(u_{c-m+j}^\top v_c)}{\sum_{k=1}^{|V|} \exp(u_k^\top v_c)} \\
 &= - \sum_{j=0, j \neq m}^{2m} u_{c-m+j}^\top v_c + 2m \log \sum_{k=1}^{|V|} \exp(u_k^\top v_c)
 \end{aligned}$$

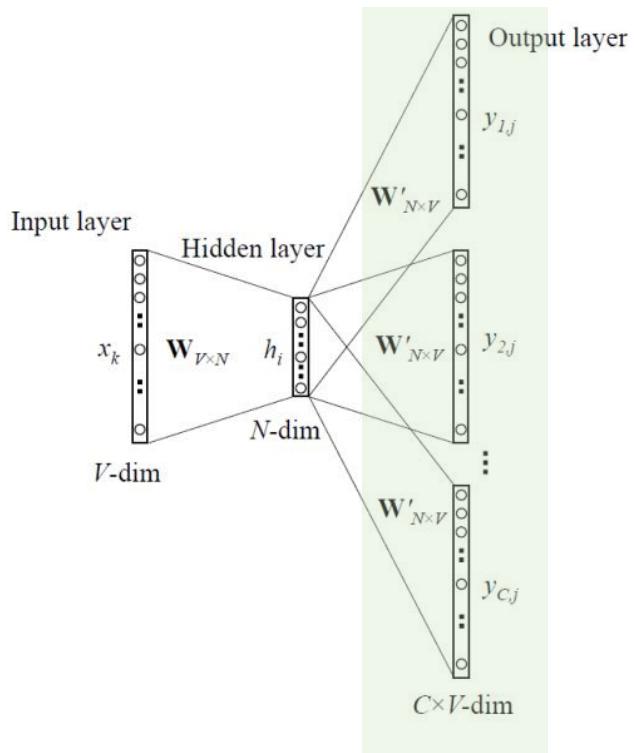
$u_i$ =the output vector representation of word  $w_i$

# Prediction based Word representation

## Skip Gram – Neural Network Architecture

Predict context (“outside”) words (position independent) given center word

**Summary of Skip Gram Training** (Review your understanding with equations)



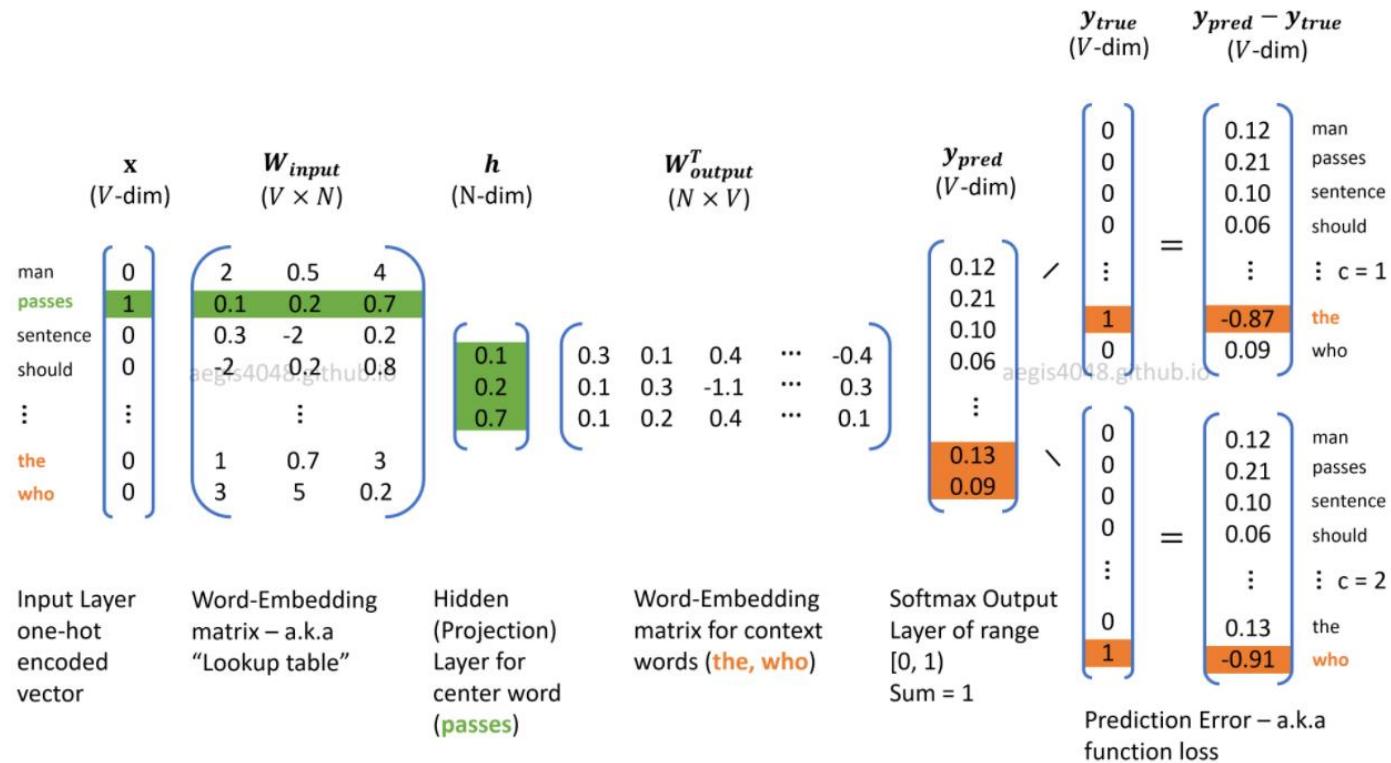
8-1. With this objective function, we can compute the gradients with respect to the unknown parameters and at each iteration update them via Stochastic Gradient Descent

$$\begin{aligned}
 J &= - \sum_{j=0, j \neq m}^{2m} \log P(u_{c-m+j} | v_c) \\
 &= \sum_{j=0, j \neq m}^{2m} H(\hat{y}, y_{c-m+j})
 \end{aligned}$$

# 4 Deep Learning for NLP

## Word2Vec-SkipGram Overview

With a simple diagram



## Key Parameter (2) for Training methods: Negative Samples (From lecture 2)

The number of negative samples is another factor of the training process.

Negative samples to our dataset – samples of words that are not neighbors

*Negative sample: 2*

<i>Input word</i>	<i>Output word</i>	<i>Target</i>
eat	mango	1
eat	exam	0
eat	tobacco	0

\*1=Appeared, 0=Not Appeared

*Negative sample: 5*

<i>Input word</i>	<i>Output word</i>	<i>Target</i>
eat	mango	1
eat	exam	0
eat	tobacco	0
eat	pool	0
eat	supervisor	0

The original paper prescribes **5-20 as being a good number of negative samples**. It also states that **2-5 seems to be enough when you have a large enough dataset**.

# 4 Deep Learning for NLP

## Word2Vec-SkipGram Overview – negative sampling

With a simple diagram

**don't compute words that are not neighbours.**

**Vanilla  
Skip-Gram**

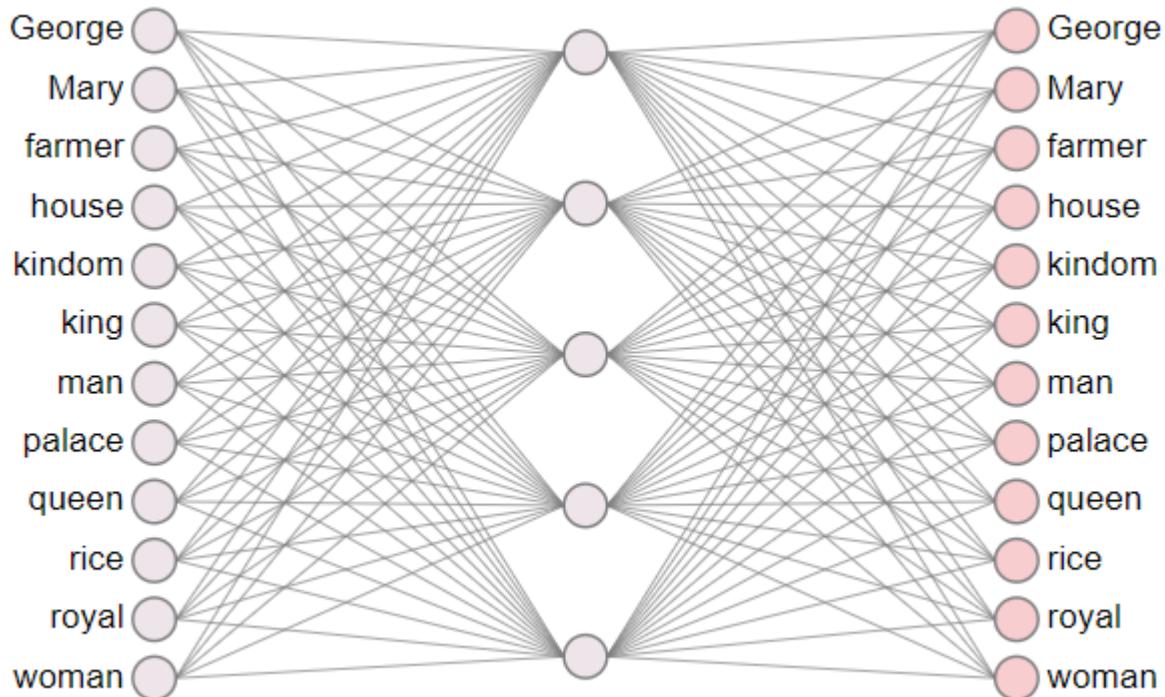
$$\begin{array}{c}
 \text{W\_output (old)} \\
 \begin{array}{ccc} -0.560 & 0.340 & 0.160 \\ -0.910 & -0.440 & 1.560 \\ -1.210 & -0.130 & -1.320 \\ 1.670 & -0.150 & -1.030 \\ 1.720 & -1.460 & 0.730 \\ 0.000 & 1.390 & -0.120 \\ -0.060 & 1.520 & -0.790 \\ 0.800 & 1.850 & -1.670 \\ -1.370 & 1.320 & -0.480 \\ 0.670 & 1.990 & -1.850 \\ -1.520 & -1.740 & -1.860 \end{array} \\
 (11 \times 3)
 \end{array}
 \quad - \quad \boxed{0.05} \quad \times \quad
 \begin{array}{c}
 \text{grad\_W\_output} \\
 \begin{array}{ccc} 0.064 & 0.071 & -0.014 \\ 0.098 & 0.015 & 0.063 \\ 0.069 & 0.089 & 0.045 \\ 0.014 & 0.085 & 0.079 \\ -0.021 & 0.067 & 0.071 \\ -0.098 & -0.088 & 0.091 \\ -0.072 & -0.078 & -0.089 \\ 0.046 & -0.079 & -0.053 \\ -0.049 & -0.087 & 0.025 \\ -0.060 & 0.092 & 0.042 \\ 0.074 & 0.050 & 0.070 \end{array} \\
 (11 \times 3)
 \end{array}
 \quad = \quad
 \begin{array}{c}
 \text{W\_output (new)} \\
 \begin{array}{ccc} -0.563 & 0.336 & 0.161 \\ -0.915 & -0.441 & 1.557 \\ -1.213 & -0.134 & -1.322 \\ 1.669 & -0.154 & -1.034 \\ 1.721 & -1.463 & 0.726 \\ 0.005 & 1.394 & -0.125 \\ -0.056 & 1.524 & -0.786 \\ 0.798 & 1.854 & -1.667 \\ -1.368 & 1.324 & -0.481 \\ 0.673 & 1.985 & -1.852 \\ -1.524 & -1.743 & -1.864 \end{array} \\
 (11 \times 3)
 \end{array}$$

**Negative  
Sampling**

$$\begin{array}{c}
 \text{W\_output (old)} \\
 \begin{array}{ccc} -0.560 & 0.340 & 0.160 \\ -0.910 & -0.440 & 1.560 \\ -1.210 & -0.130 & -1.320 \\ 1.670 & -0.150 & -1.030 \\ 1.720 & -1.460 & 0.730 \\ 0.000 & 1.390 & -0.120 \\ -0.060 & 1.520 & -0.790 \\ 0.800 & 1.850 & -1.670 \\ -1.370 & 1.320 & -0.480 \\ 0.670 & 1.990 & -1.850 \\ -1.520 & -1.740 & -1.860 \end{array} \\
 (11 \times 3)
 \end{array}
 \quad - \quad \boxed{0.05} \quad \times \quad
 \begin{array}{c}
 \text{grad\_W\_output} \\
 \begin{array}{c} \text{Not computed!} \end{array}
 \end{array}
 \quad = \quad
 \begin{array}{c}
 \text{W\_output (new)} \\
 \begin{array}{ccc} -0.560 & 0.340 & 0.160 \\ -0.910 & -0.440 & 1.560 \\ -1.210 & -0.130 & -1.320 \\ 1.670 & -0.150 & -1.030 \\ 1.720 & -1.460 & 0.730 \\ 0.000 & 1.390 & -0.120 \\ -0.060 & 1.520 & -0.790 \\ \text{Positive sample, } w_o & 0.031 & 0.030 & 0.041 \\ \text{Negative sample, k=1} & -0.090 & 0.031 & -0.065 \\ \text{Negative sample, k=2} & 0.056 & 0.098 & -0.061 \\ \text{Negative sample, k=3} & 0.069 & 0.084 & -0.044 \end{array} \\
 (11 \times 3)
 \end{array}
 \quad (11 \times 3)$$

## Application

### Application #1: Embedding Pretraining



# 0 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.

# / Reference

## 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