

COMP5046

Natural Language Processing

Lecture 3: Word Classification and Machine Learning

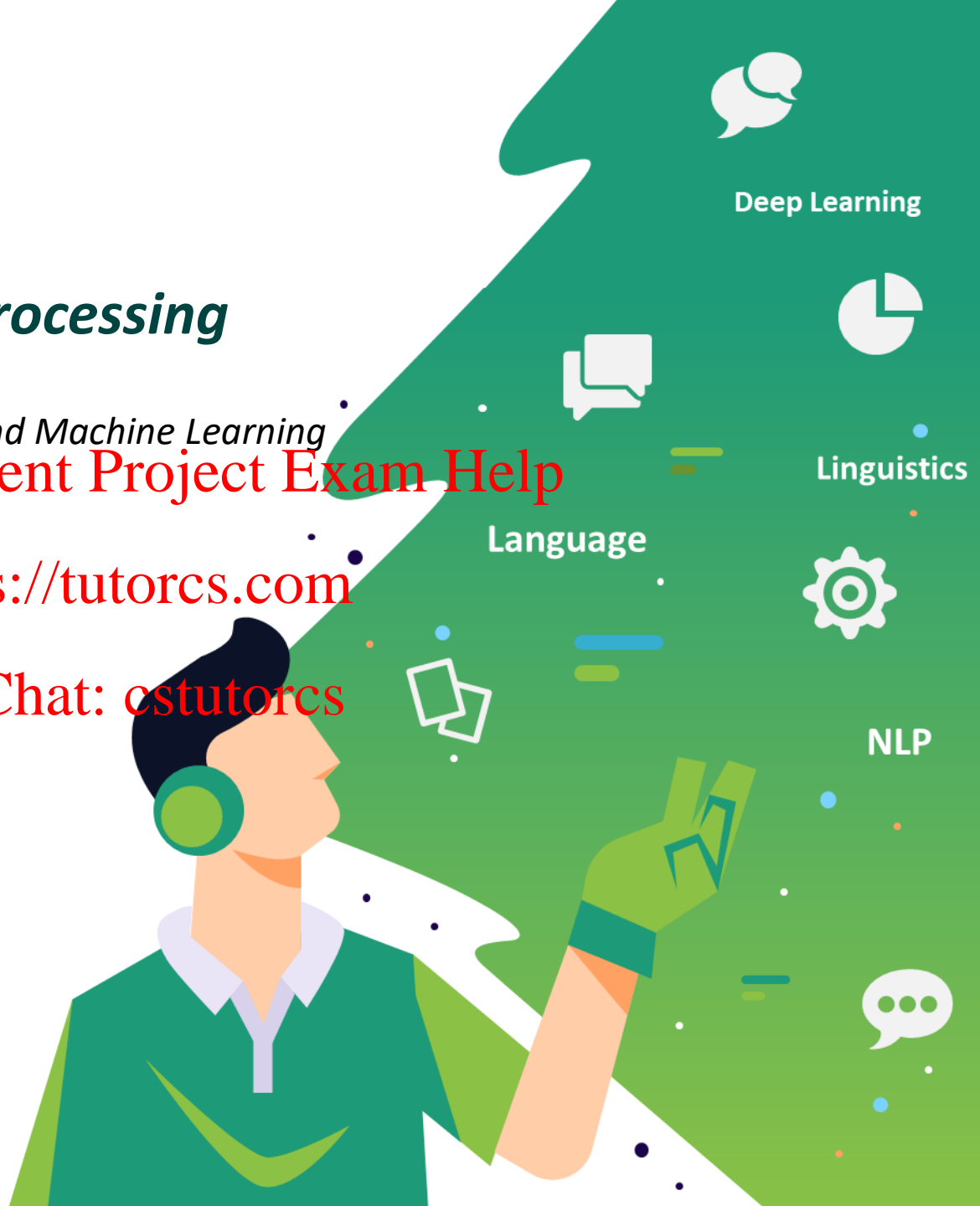
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Dr. Caren Han

Semester 1, 2021
School of Computer Science,
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.

Assignment Project Exam Help

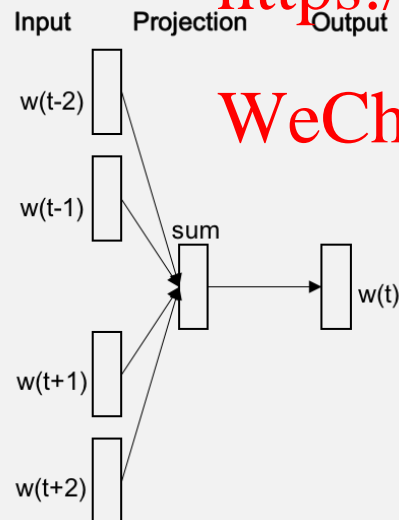
<https://tutorcs.com>

WeChat: cstutorcs

Word2Vec Models

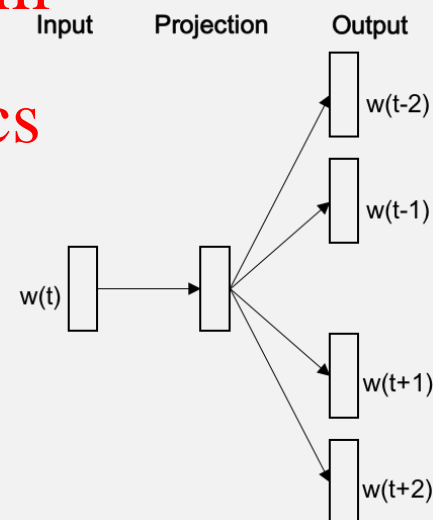
CBOW

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



Skip-gram

*Predict context words
given center word*



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

Assignment Project Exam Help

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]	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

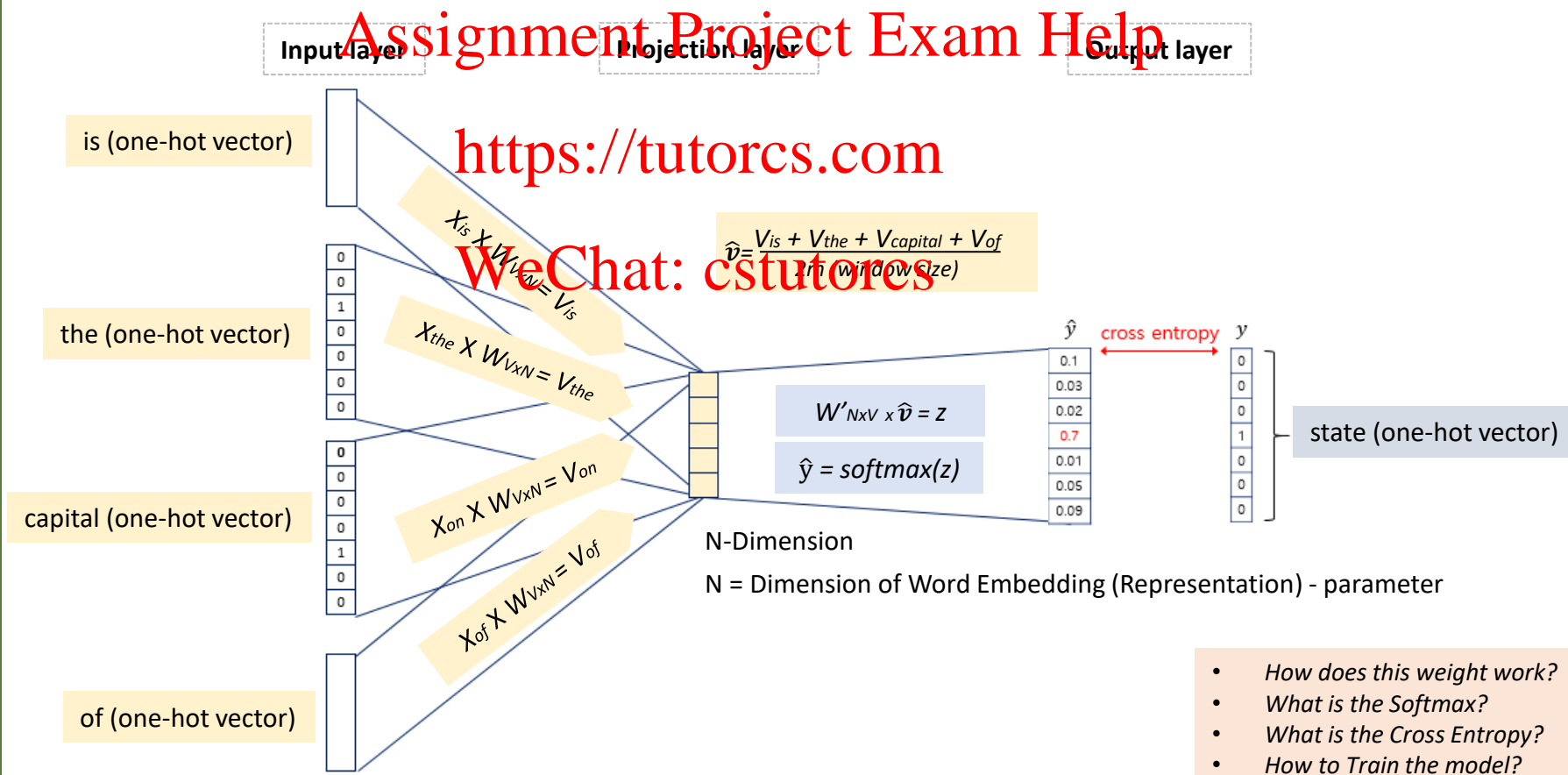
Center word

Context (“outside”) word

CBOW – Neural Network Architecture

Predict center word from (bag of) context words

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

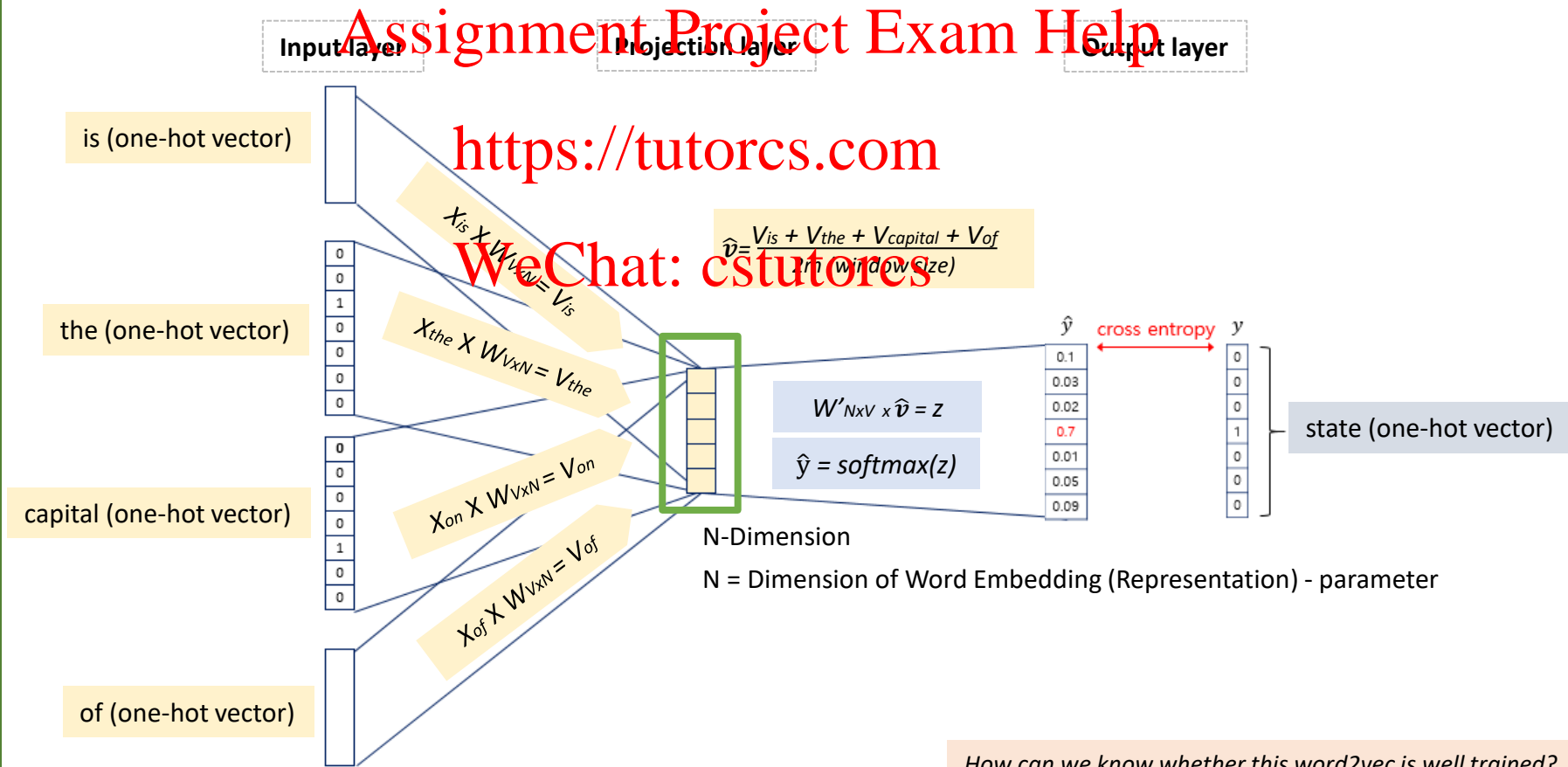


Previous Lecture Review

CBOW – Neural Network Architecture

Predict center word from (bag of) context words

Sentence: “Sydney is the state capital of NSW”



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.

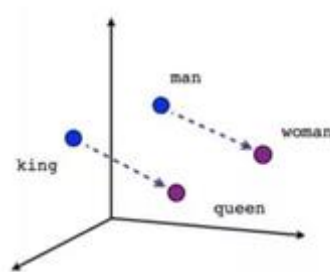
Assignment Project Exam Help

<https://tutorcs.com>

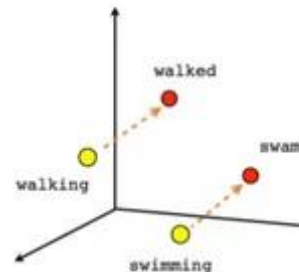
WeChat: cstutorcs

How to evaluate word vectors?

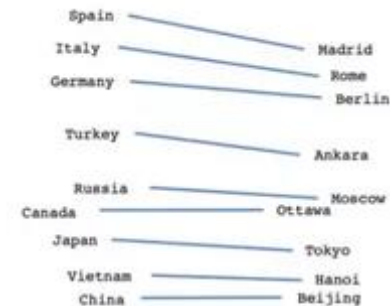
Type	How to work / Benefit
Intrinsic	Evaluation on a specific/intermediate subtask • Fast to compute • Helps to understand that system • Not clear if really helpful unless correlation to real task is established
Extrinsic	Evaluation on a real task • Can take a long time to compute accuracy • Unclear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



Male-Female



Verb tense



Country-Capital

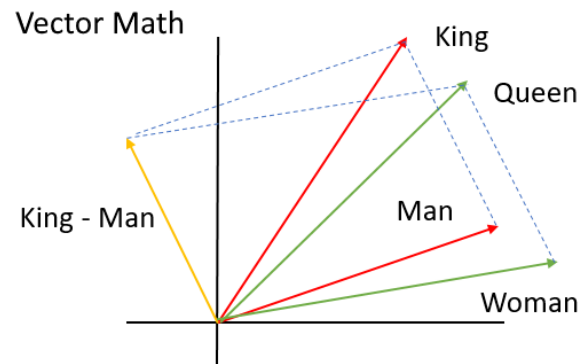
Intrinsic word vector evaluation

Word Vector Analogies

$a \leftrightarrow b$
 $man \leftrightarrow women :: king \leftrightarrow ???$

<https://tutorcs.com>

- 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 = ?

Assignment Project Exam Help

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

<https://tutorcs.com>

WeChat: cstutorcs

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.

<https://tutorcs.com>

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

(Original Word2vec Paper - Mikolov et al.2013)

Intrinsic word vector evaluation

Evaluation Result Comparison

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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 available²; (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[†] and CBOW[†] using the word2vec tool³. 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
ivLBL	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	<u>80.8</u>	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 [†]	300	6B	63.6	<u>67.4</u>	65.7
SG [†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
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	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

(Original Glove Paper - Pennington et al.2014)

Intrinsic word vector evaluation

Evaluation Result Comparison

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

Assignment Project Exam Help

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

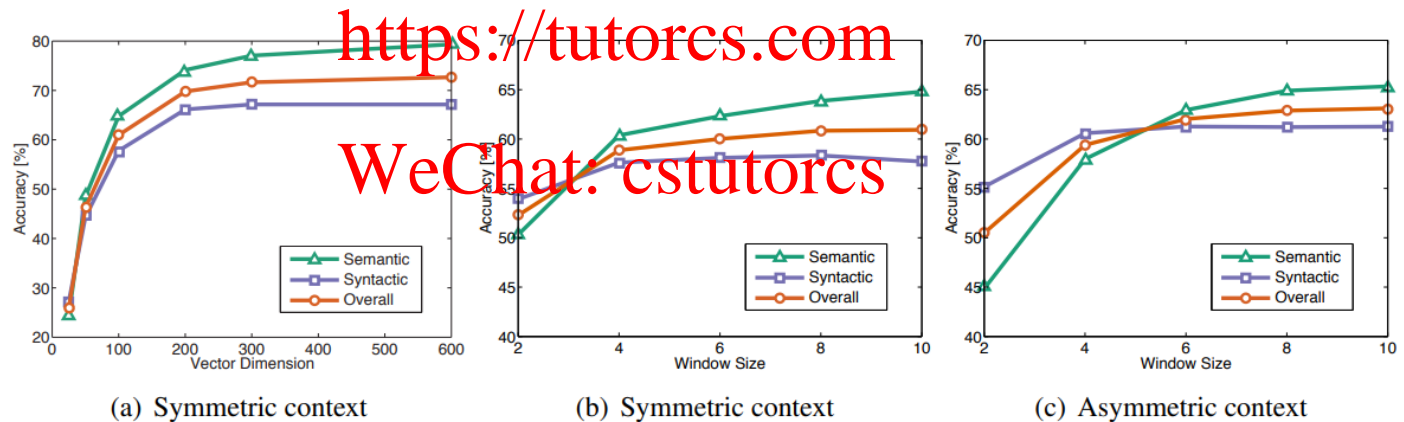
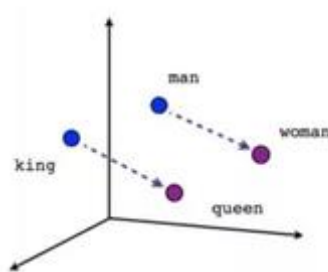


Figure 2: Accuracy on the analogy task as function of vector size and window size/type. All models are trained on the 6 billion token corpus. In (a), the window size is 10. In (b) and (c), the vector size is 100.

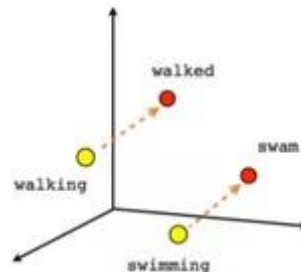
(Original Glove Paper - Pennington et al.2014)

How to evaluate word vectors?

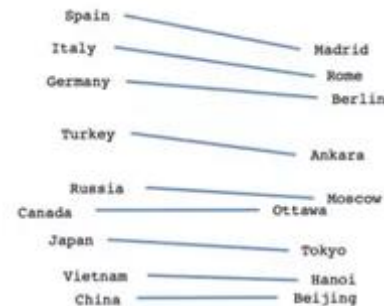
Type	How to work / Benefit
Intrinsic	Evaluation on a specific/intermediate subtask • Fast to compute • Helps to understand that system • Not clear if really helpful unless correlation to real task is established
Extrinsic	Evaluation on a real task • Can take a long time to compute accuracy • Unclear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



Male-Female



Verb tense



Country-Capital

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.

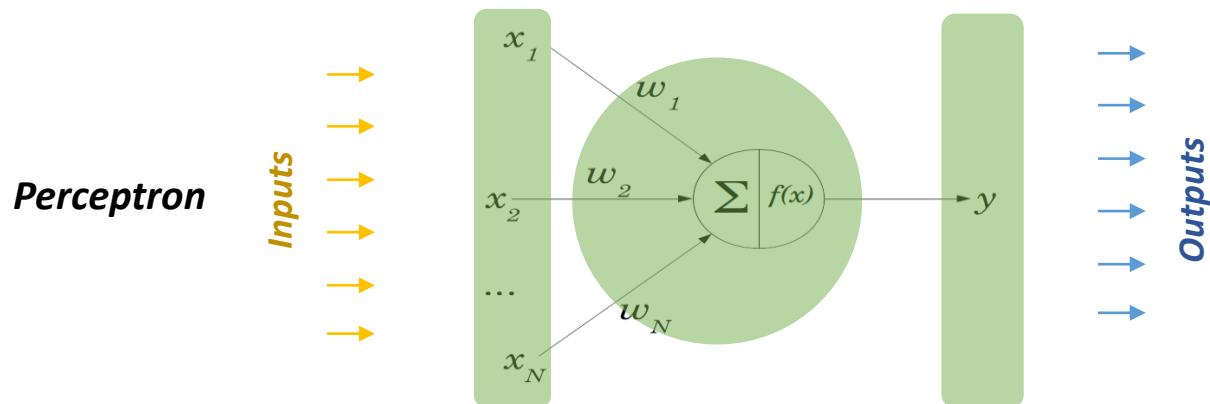
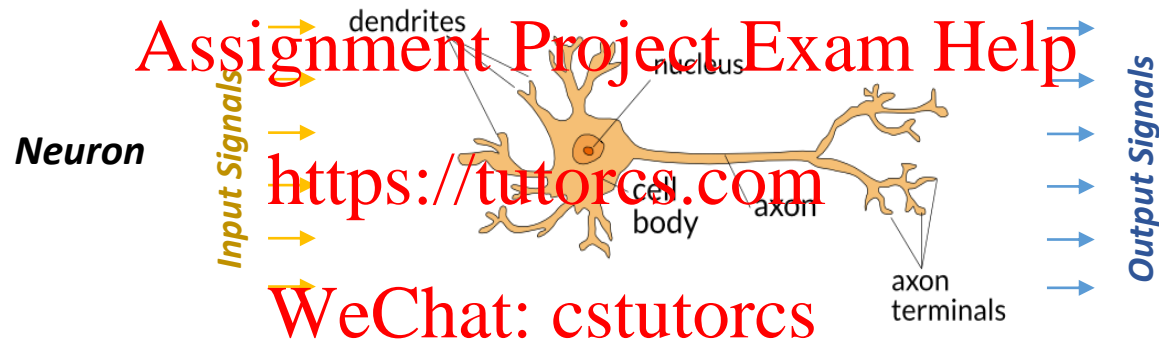
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Deep Learning with Neural Network

Neuron and Perceptron

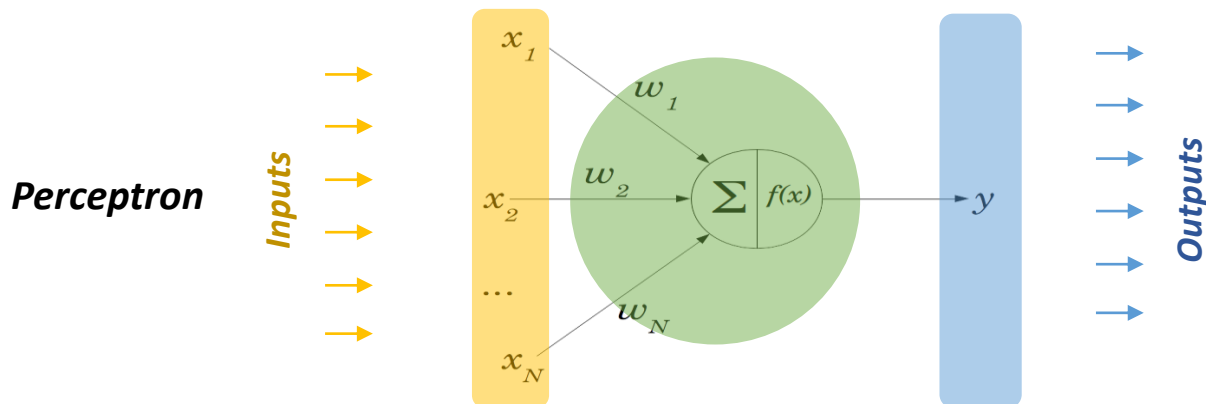


Deep Learning with Neural Network

Inputs and Outputs (Labels) for Natural Language Processing

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

WeChat: cstutorcs



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

Assignment Project Exam Help



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



Deep Learning with Neural Network - Model

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

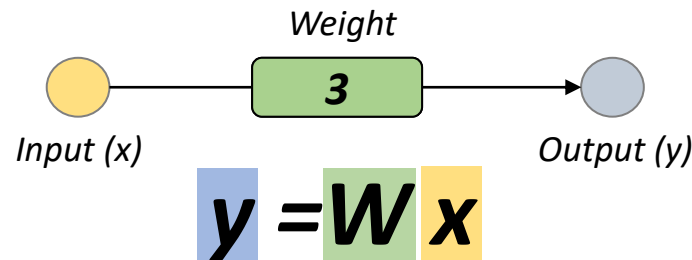
Parameters: Need to be estimated

Assignment Project Exam Help $y = 3x$



$\text{Input } (x)$
1 

$\text{Output } (y)$
3 
WeChat: cstutorcs



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

IS ANYBODY REALLY READY?

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

Assignment Project Exam Help



<https://tutorcs.com>

WeChat: cstutorcs



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

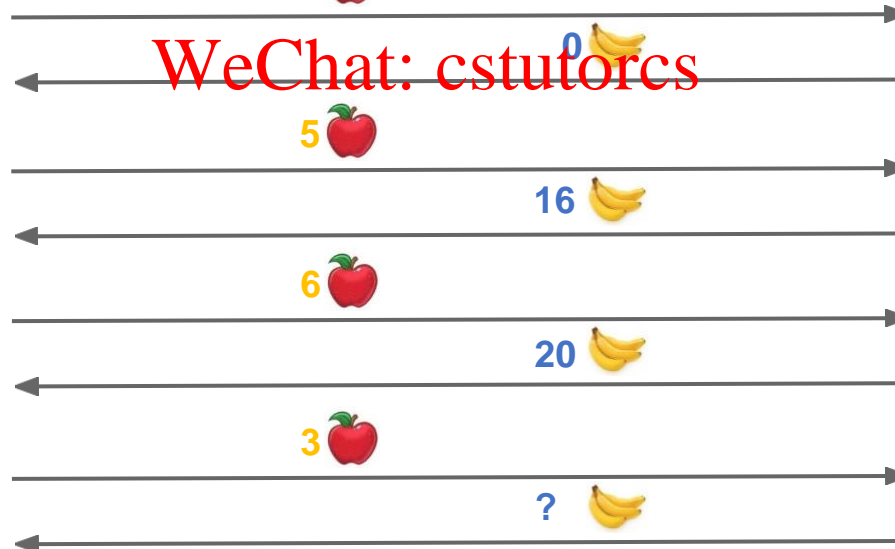
Assignment Project Exam Help

Guess how much I will give you back!



<https://tutorcs.com>

WeChat: cstutorcs



$$y = Wx$$

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

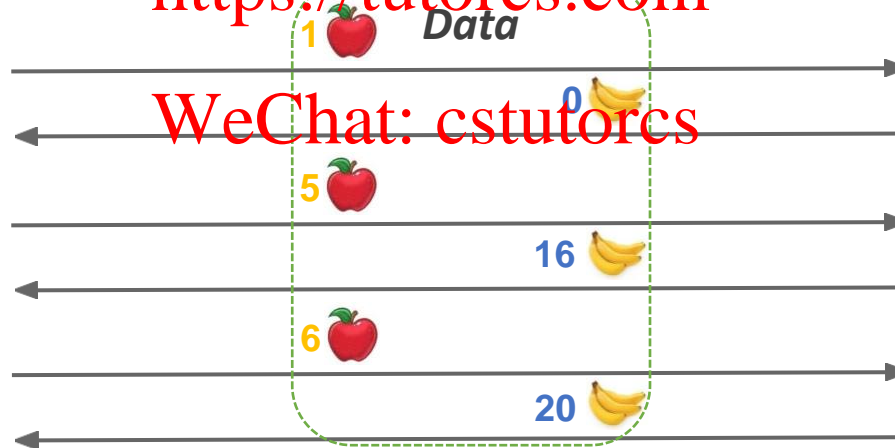
Assignment Project Exam Help

Guess how much I will give you back!



<https://tutorcs.com>

WeChat: cstutorcs



$$y = Wx$$

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat:  cstutorcs

Data

	
1	0
5	16
6	20

$$y = Wx$$

What is W then?

Deep Learning for NLP

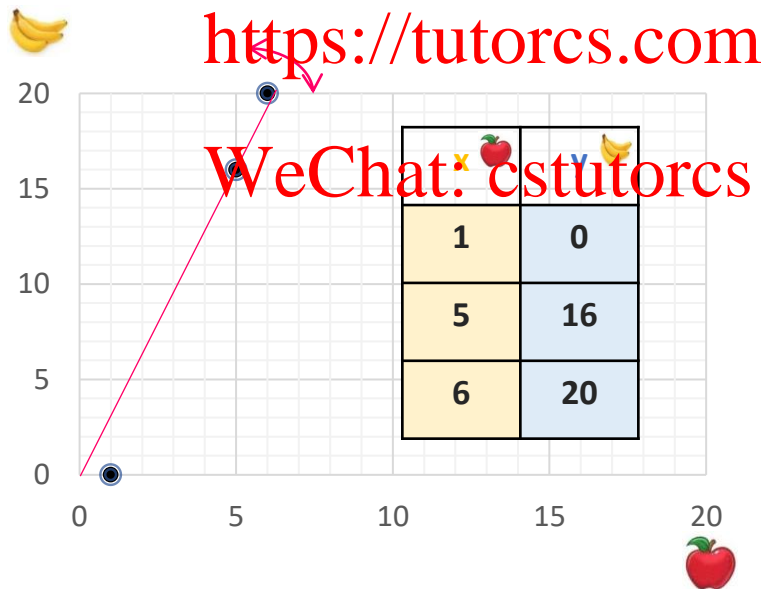
Deep Learning with Neural Network - Parameter

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx$$

What is W then?

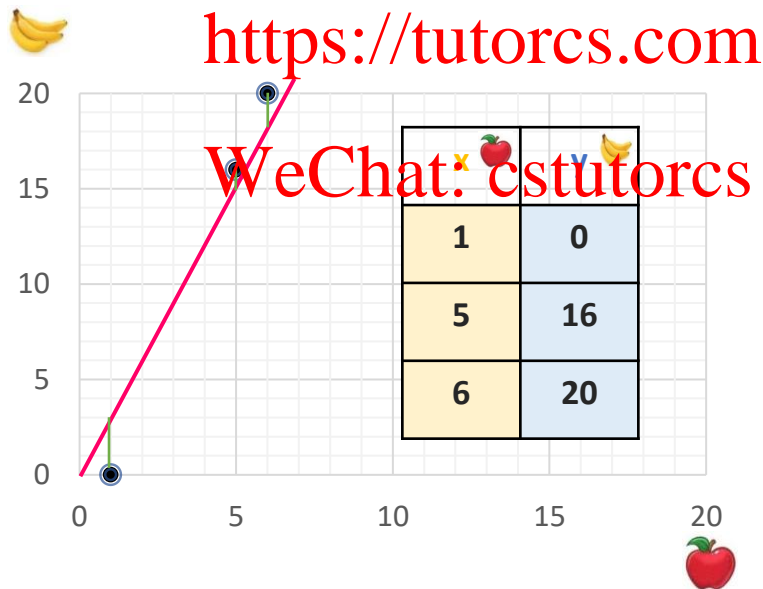
Deep Learning with Neural Network - Parameter

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx$$

What if W is 3?

$$3 = 3 \times 1$$

$$15 = 3 \times 5$$

$$20 = 3 \times 6$$

Deep Learning for NLP

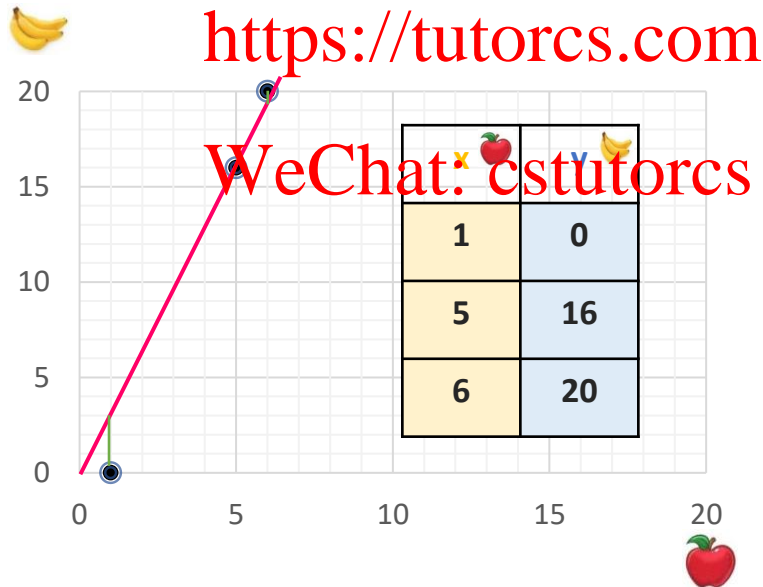
Deep Learning with Neural Network - Parameter

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx$$

What if W is 3.2?

$$3.2 = 3.2 \times 1$$

$$16 = 3.2 \times 5$$

$$19.2 = 3.2 \times 6$$

Deep Learning for NLP

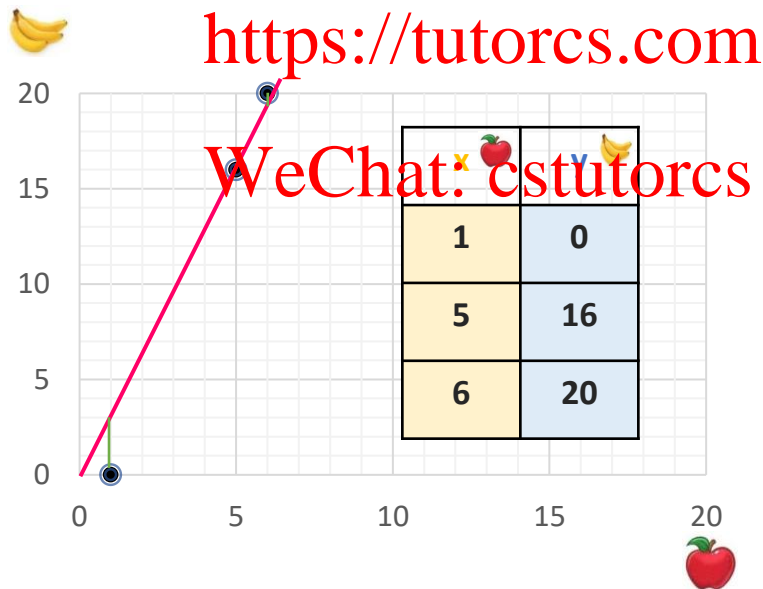
Deep Learning with Neural Network - Parameter

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx + b$$

weight bias

Weight is not enough...

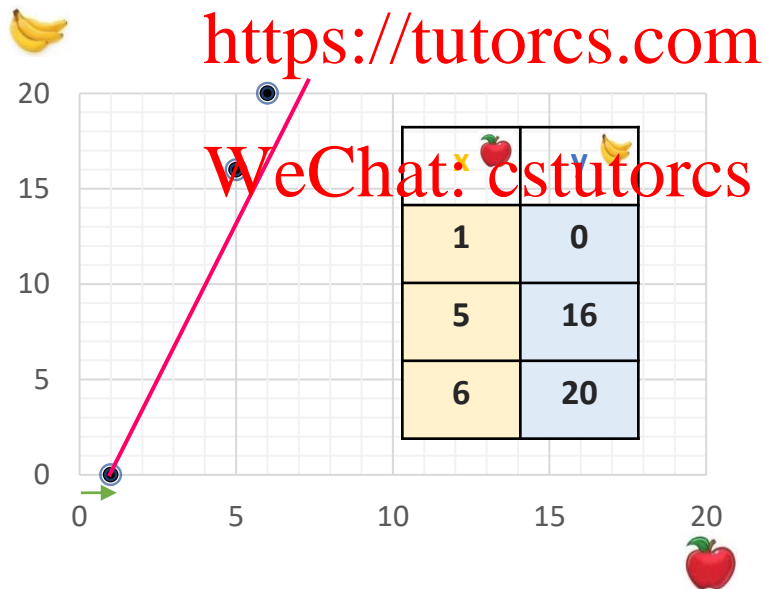
Deep Learning with Neural Network - Parameter

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx + b$$

weight
bias

How can we find the parameters, w and b ?

Deep Learning for NLP

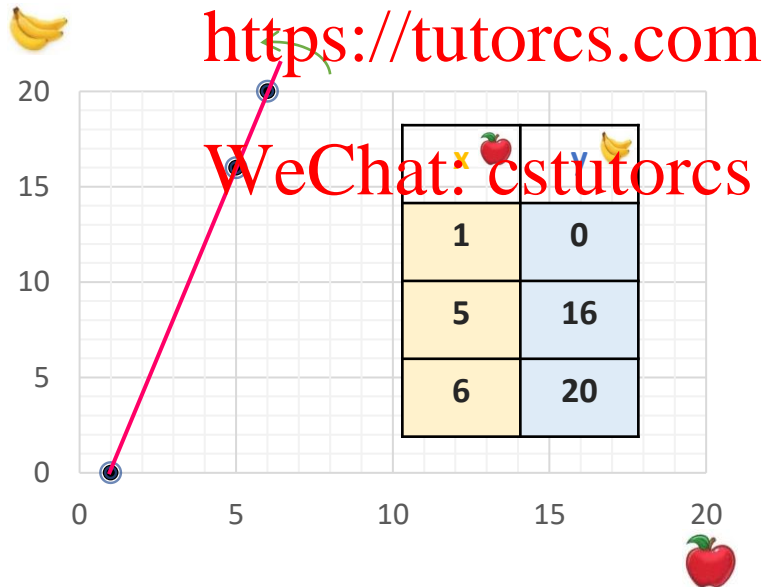
Deep Learning with Neural Network - Parameter

Input: x =number of apple given by Lisa

Output: y =number of banana received by Lisa

Parameters: Need to be estimated

Assignment Project Exam Help



$$y = Wx + b$$

weight bias

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

Assignment Project Exam Help



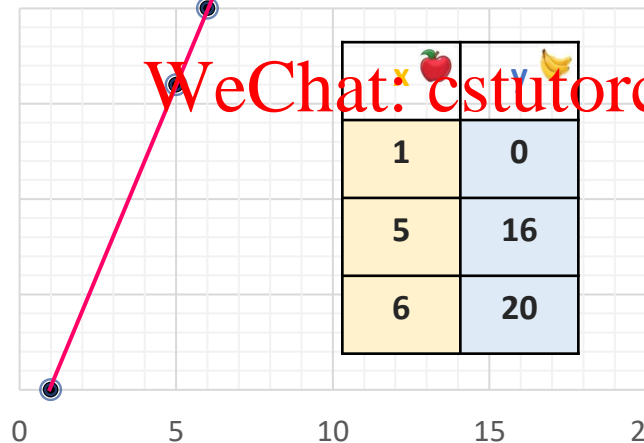
20

15

10

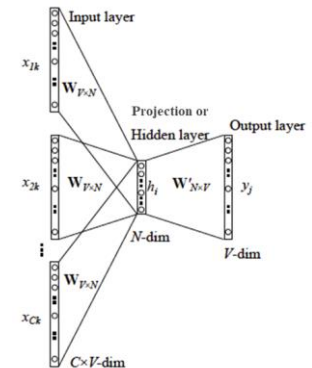
5

0



<https://tutorcs.com>

WeChat: cstutorcs



Model

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

weight bias

How can we find the parameters, \mathbf{w} and \mathbf{b} ?

Deep Learning with Neural Network - Cost

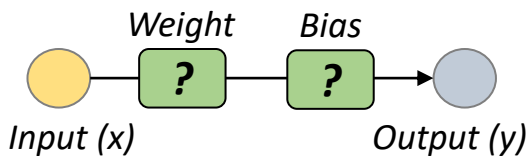
Actual Data

weight

bias

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Model Ex#1

weight

bias

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

x 🍎	predicted \hat{y} 🍌	actual y 🍌
1	1	0
5	5	16
6	6	20

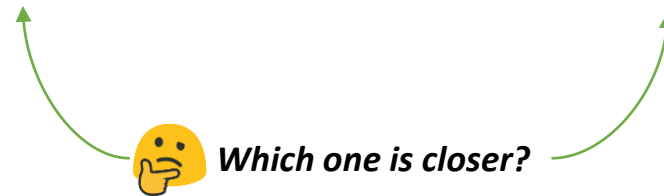
Model Ex#2

weight

bias

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

x 🍎	predicted \hat{y} 🍌	actual y 🍌
1	4	0
5	12	16
6	14	20



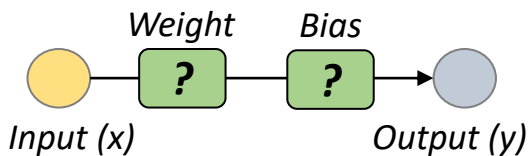
Assignment Project Exam Help
<https://tutorcs.com>
 WeChat: cstutorcs

Deep Learning with Neural Network – Cost (loss)

Actual Data

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Model Ex#1

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

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

Model Ex#2

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

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

Assignment Project Exam Help
<https://tutorcs.com>
 WeChat: cstutorcs



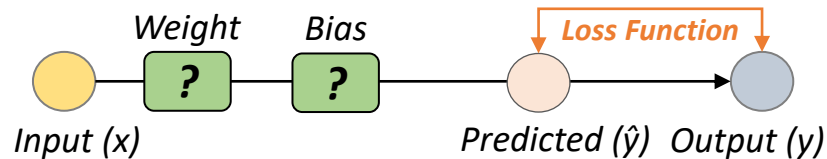
Let's calculate the cost(loss)!

Mean Squared Error (MSE)

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

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

WAIT! Loss Function? Cost Calculation?



Assignment Project Exam Help

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

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

<https://tutorcs.com>

WeChat: cstutorcs

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

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

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

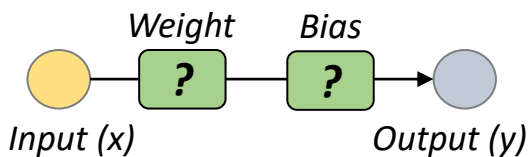
3 Deep Learning for NLP

Deep Learning with Neural Network - Cost (loss)

Actual Data

$$y = ?x + ?$$

weight	bias
x 🍎	y 🍌
1	0
5	16
6	20



Model Ex#1

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

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

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

Model Ex#2

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

predicted	actual	cost
\hat{y} 🍌	y 🍌	$(y - \hat{y})^2$
4	0	16
12	16	16
14	20	36

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



Let's calculate the cost!

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

Assignment Project Exam Help
<https://tutorcs.com>
 WeChat: cstutorcs

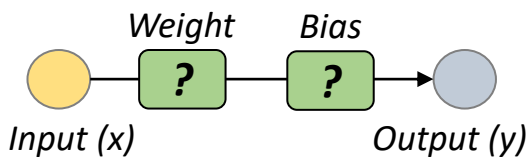
Deep Learning for NLP

Deep Learning with Neural Network - Cost (loss)

Actual Data

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Model Ex#1

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

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

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

Model Ex#2

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

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

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



Assignment Project Exam Help
<https://tutorcs.com>
 WeChat: cstutorcs



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

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

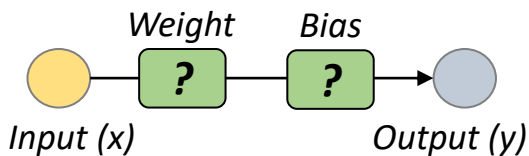
$$w,b \in [-\infty, \infty]$$

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

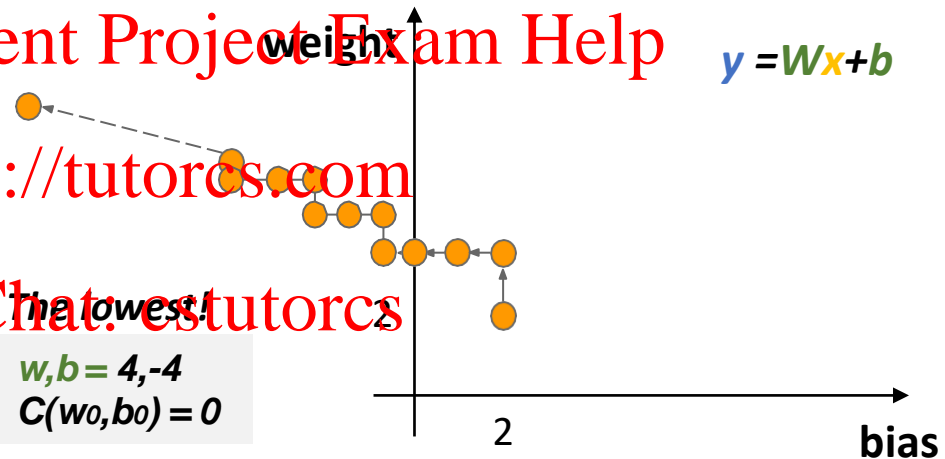
WeChat: estutorcs

The lowest!

$$w, b = 4, -4$$

$$C(w_0, b_0) = 0$$

$$y = Wx + b$$



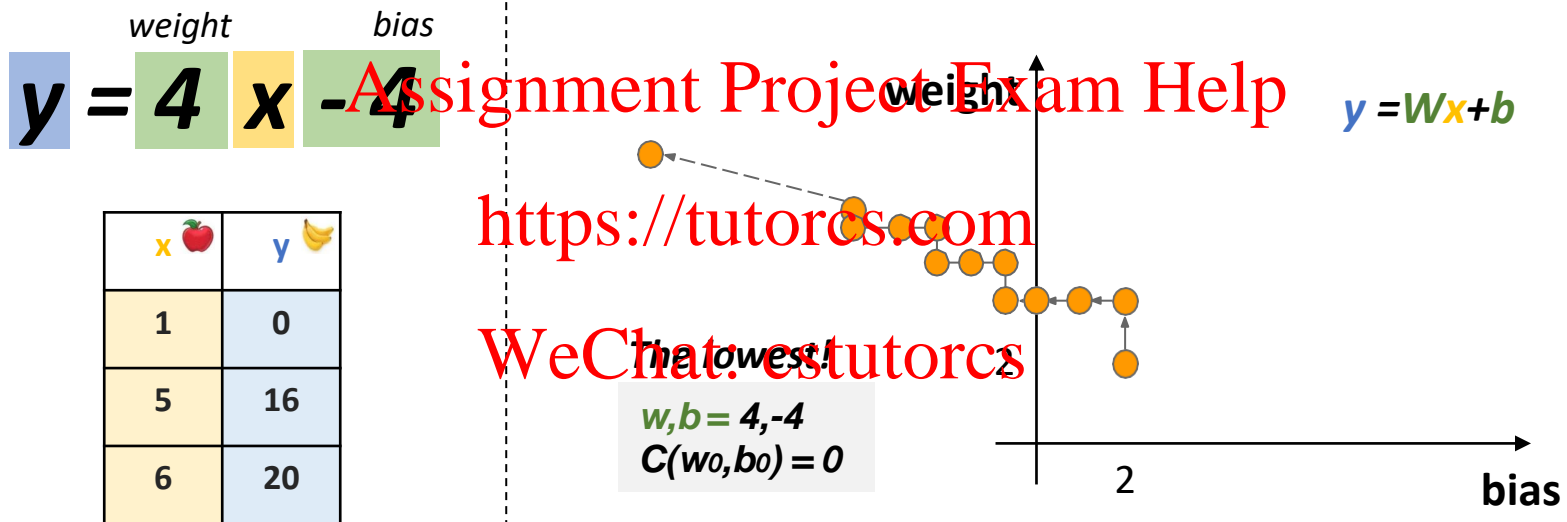
Let's calculate the costs and get the lowest one!

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

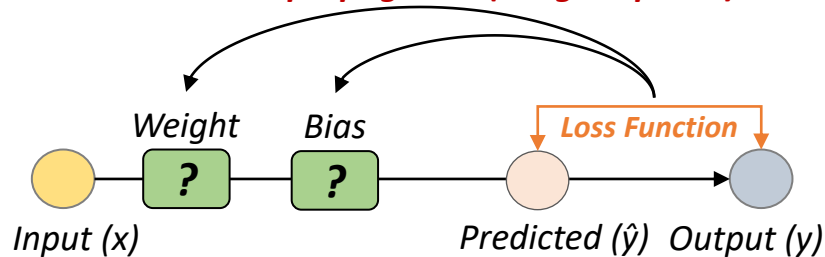
$$w, b \in [-\infty, \infty]$$

Deep Learning with Neural Network - Optimizer

Backpropagation (weight update)



Backpropagation (weight update)



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

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

Backpropagation (weight update)

Deep Learning with Neural Network - Optimizer

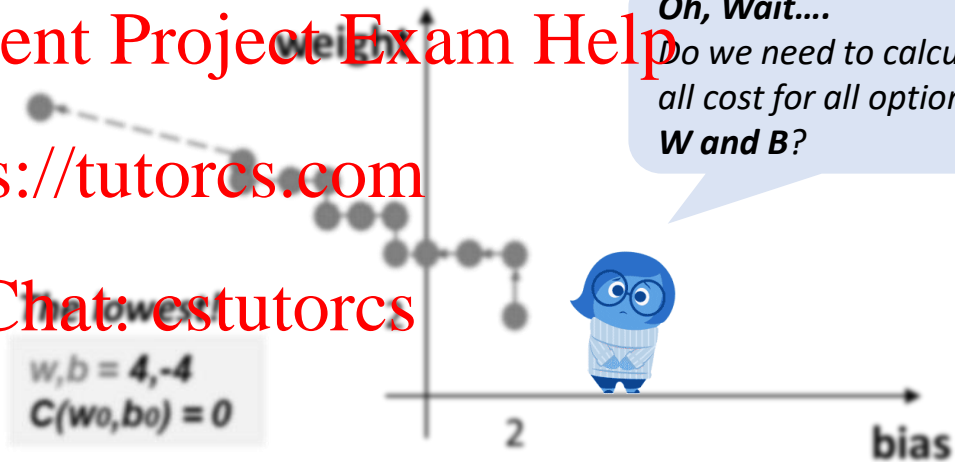
$$y = \overset{\text{weight}}{4} x \overset{\text{bias}}{-4}$$

x 🍎	y 🍌
1	0
5	16
6	20

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: estutorcs



Expensive to compute
(hours or days)

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

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

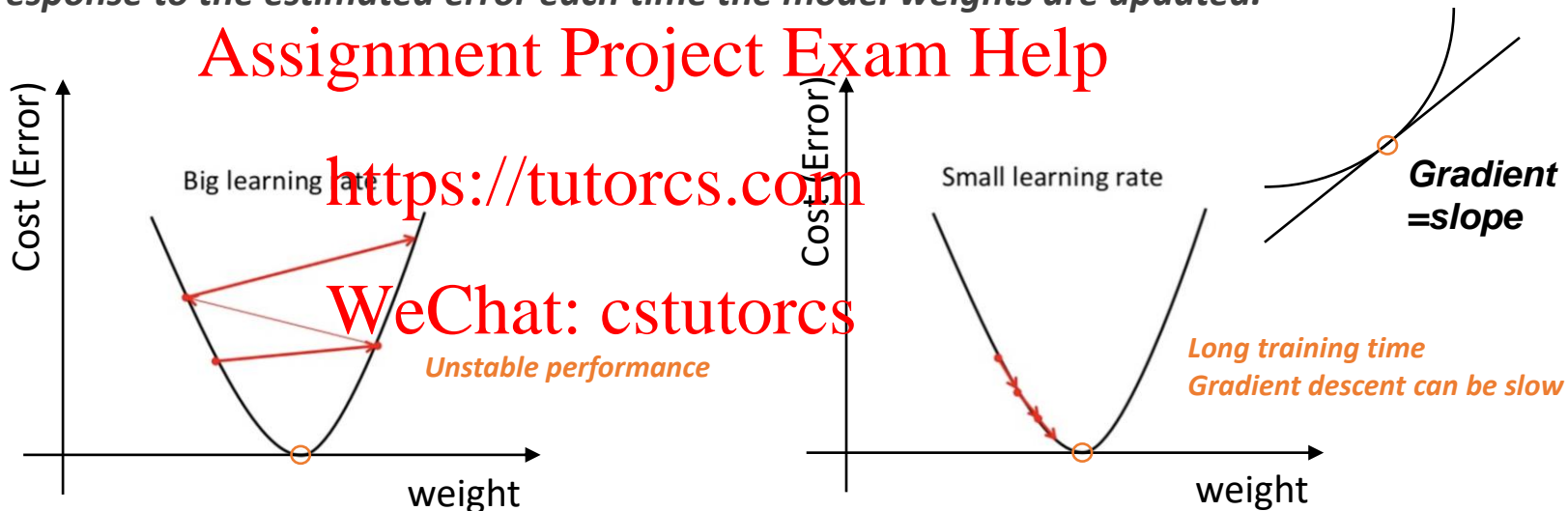
Finding the Optimal weight and bias – Gradient Descent



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

Choose the optimal Learning Rate!

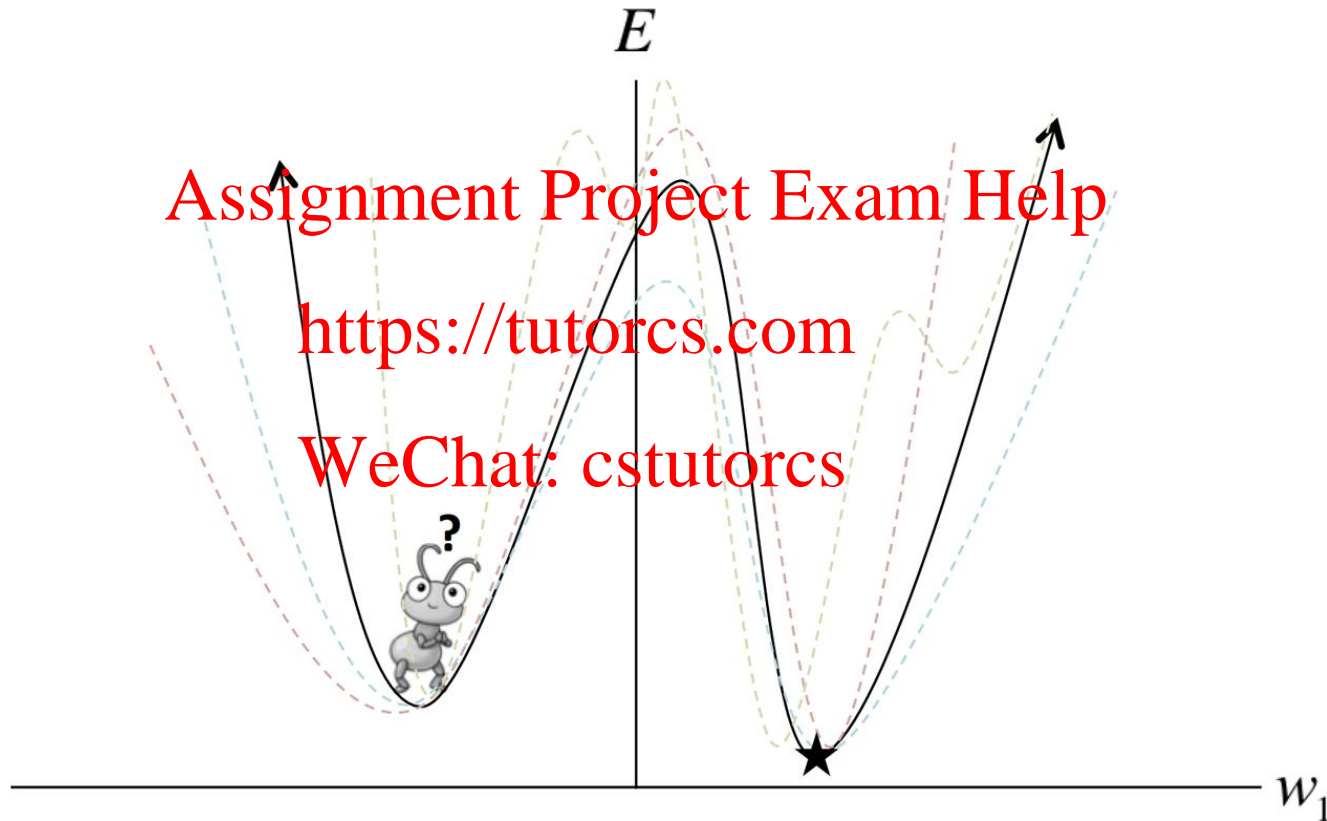
Learning Rate: a hyperparameter that controls how much to change the model in response to the estimated error each time the model weights are updated.



$$\text{new_weight} = \text{existing_weight} - \text{learning_rate} * \text{gradient}$$

$$\text{new_weight} = \text{existing_weight} - \text{learning_rate} * (\text{current_output} - \text{desired output}) * \text{gradient}(\text{current output}) * \text{existing_input}$$

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.

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!*

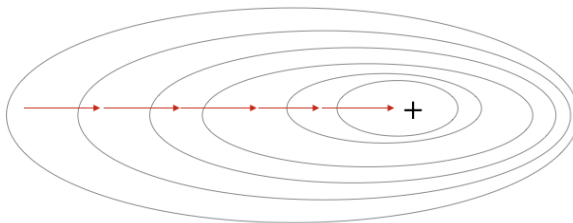
The Solution can be used different Gradient Descent Method.

The most common – “Stochastic Gradient Descent (SGD)”

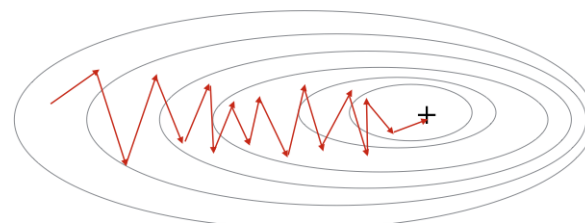
<https://tutorcs.com>

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.

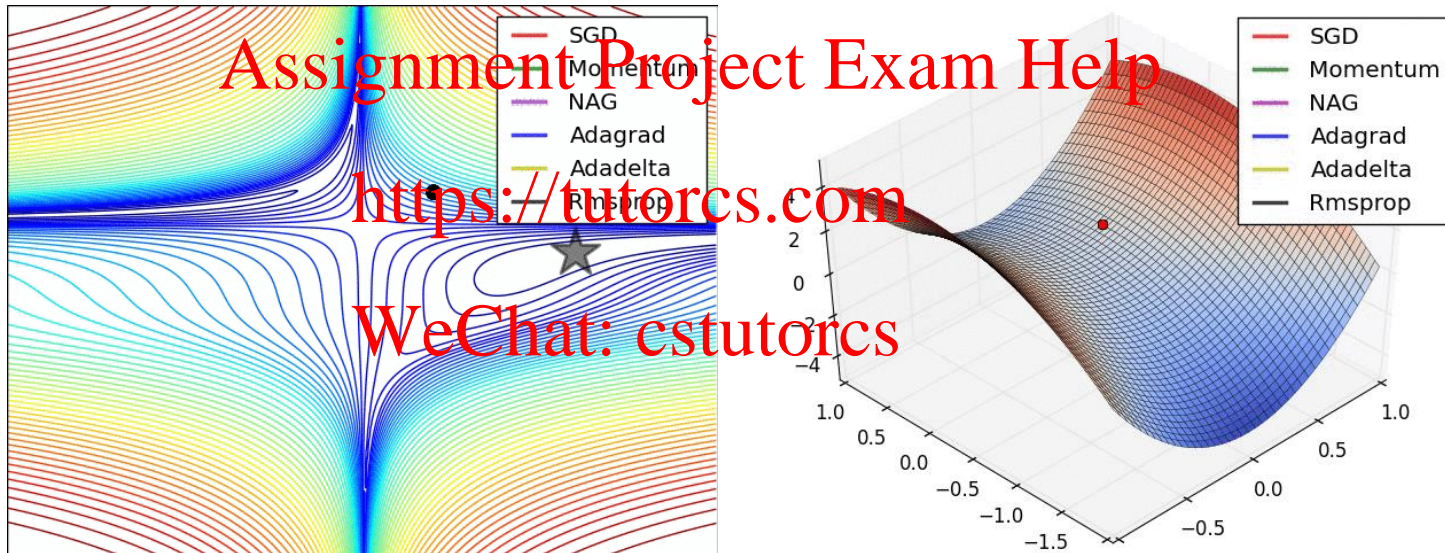
Gradient Descent



Stochastic Gradient Descent



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 with Neural Network - Optimizer

Backpropagation (weight update)

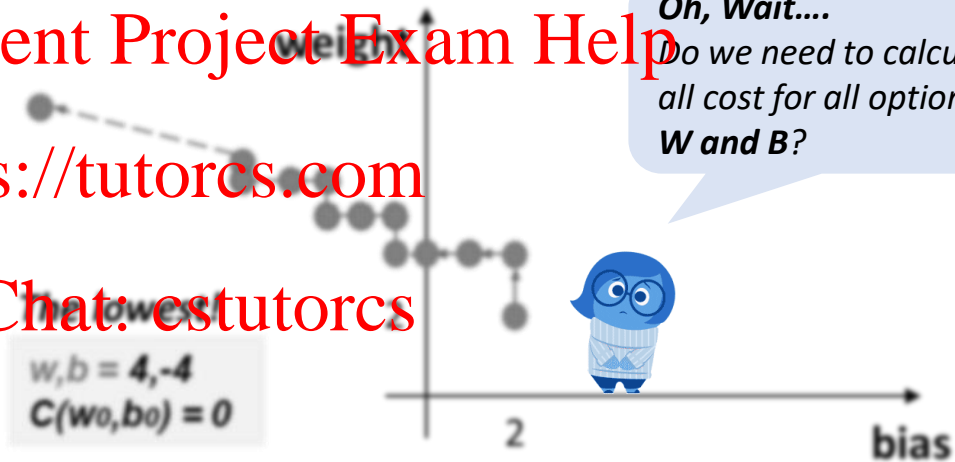
$$y = \overset{\text{weight}}{4} x + \overset{\text{bias}}{-4}$$

x 🍎	y 🍌
1	0
5	16
6	20

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: estutorcs



Expensive to compute
(hours or days)

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

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

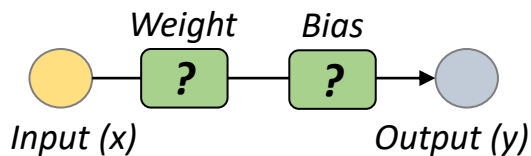
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20

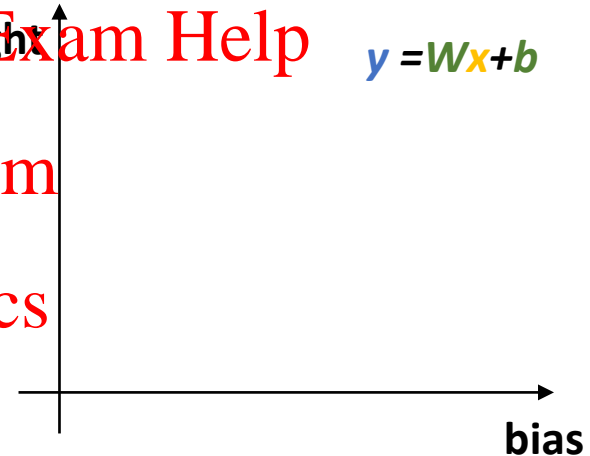


Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = Wx + b$$



Should be used sparingly

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

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

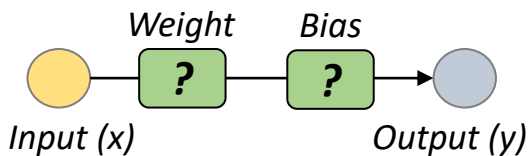
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20

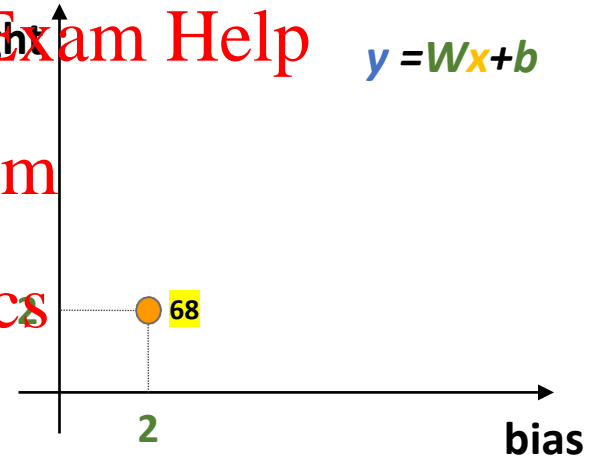


Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = Wx + b$$



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

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

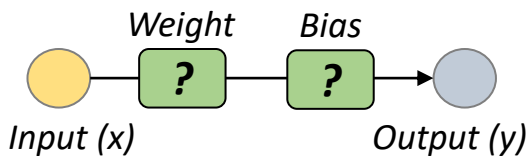
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

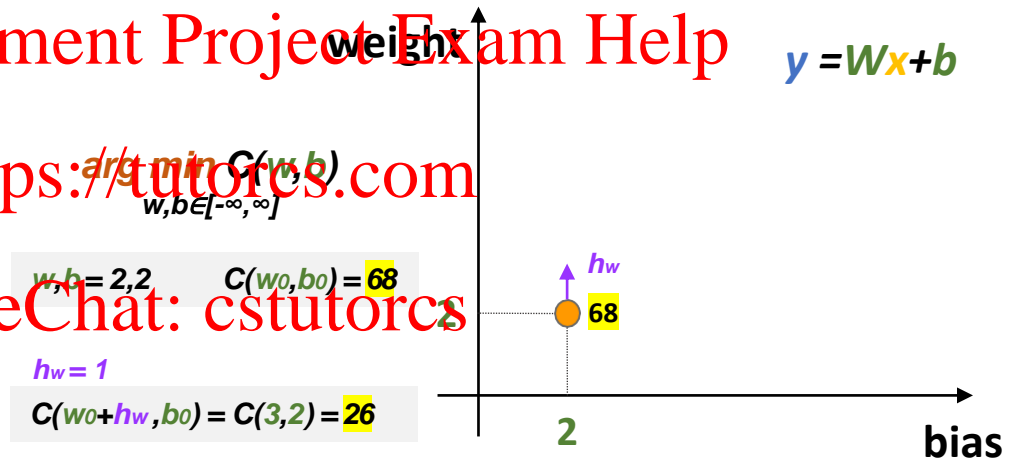
x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



Gradient!

Deep Learning with Neural Network - Optimizer

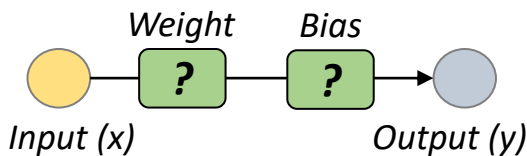
Actual Data

weight

bias

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = Wx + b$$

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

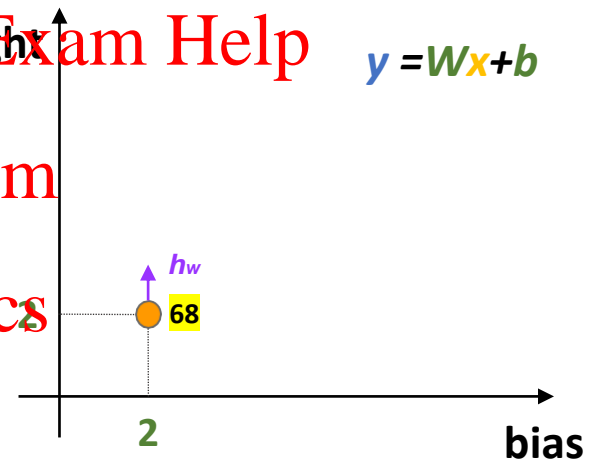
$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$h_w = 1$$

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

$$r = \frac{C(w_0 + h_w, b_0) - C(w_0, b_0)}{1}$$

$$r = \frac{C(2 + 1, 2) - C(2, 2)}{1} = -42$$



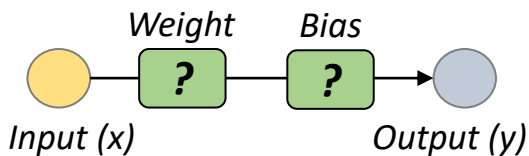
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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

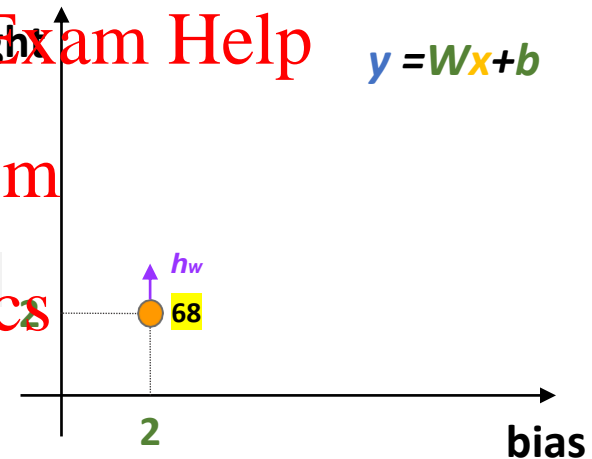
$$w, b = 2, 2 \quad C(w_0, b_0) = 68$$

$$h_w = 1, r = -42$$

$$h_w = 0.1, r = -98$$

$$h_w = 0.01, r = -104$$

$$h_w = 0.001, r = -104$$



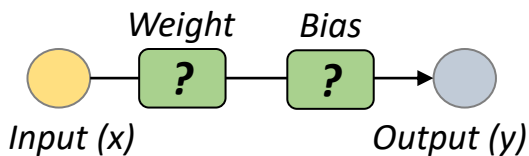
Deep Learning with Neural Network - Optimizer

Gradient!

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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

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

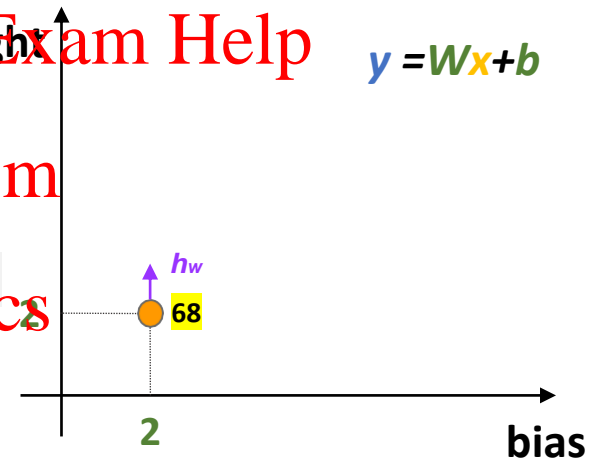
$$h_w = 1, r = -42$$

$$h_w = 0.1, r = -98$$

$$h_w = 0.01, r = -104$$

$$h_w = 0.001, r = -104$$

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



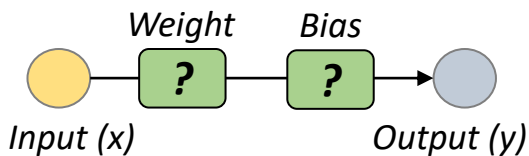
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

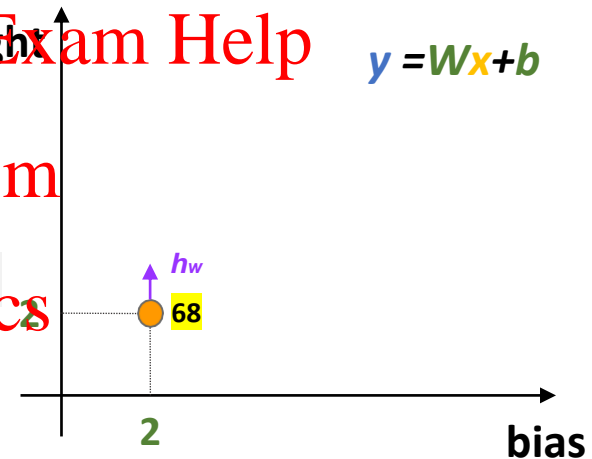
<https://tutorcs.com>

WeChat: cstutorcs

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

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

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



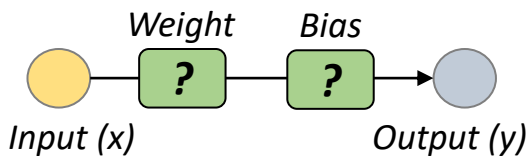
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



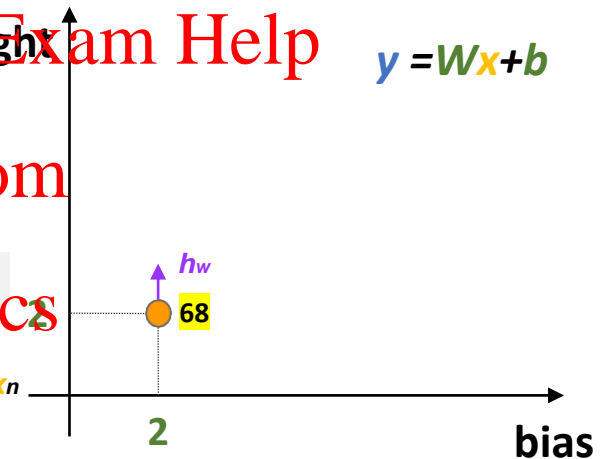
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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

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



Gradient!

Deep Learning with Neural Network - Optimizer

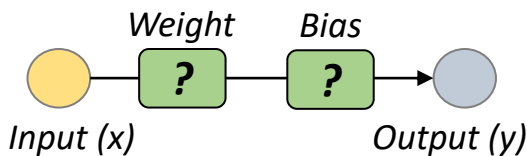
Actual Data

weight

bias

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

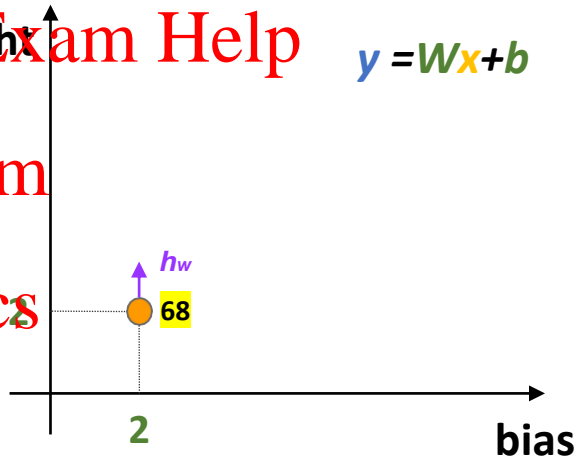
$$y = Wx + b$$

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

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

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

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



3 Deep Learning for NLP

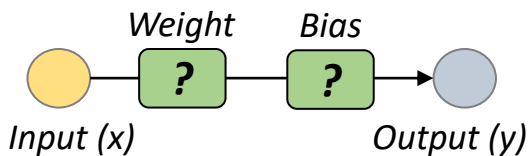
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



$$y = 2x + 2$$

	predicted	actual		
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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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

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

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

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

Gradient!

Deep Learning with Neural Network - Optimizer

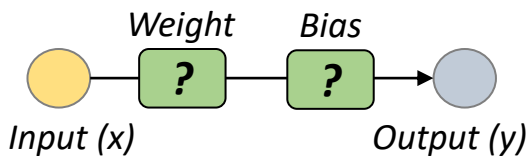
Actual Data

weight

bias

$$y = ?x + ?$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

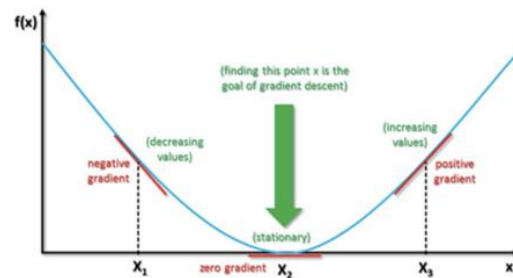
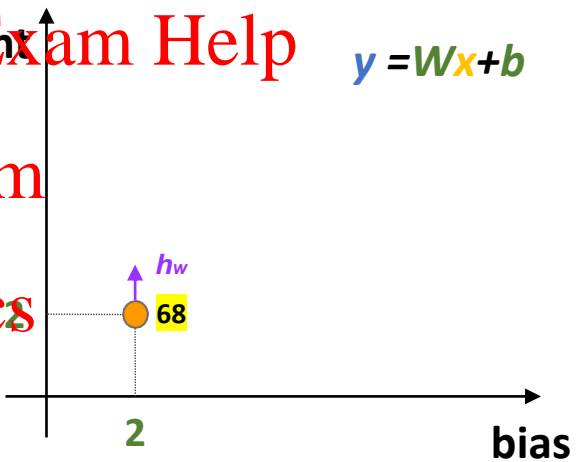
WeChat: cstutorcs

$$y = Wx + b$$

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

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

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



Gradient!

Deep Learning with Neural Network - Optimizer

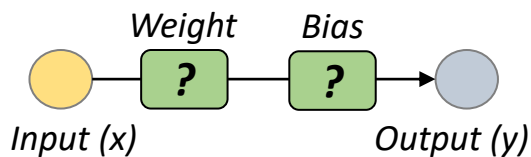
Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

weight bias

$y = ? \cdot x + ?$

x 🍎	y 🍌
1	0
5	16
6	20



$$y = \text{weight} \cdot x + \text{bias}$$

weight bias

$y = 2 \cdot x + 2$

x 🍎	predicted \hat{y} 🍌	actual y 🍌	$(y - \hat{y})$	$2(y - \hat{y})$
1	4	0	-4	-8
5	12	16	4	8
6	14	20	6	12

<https://tutorcs.com>

WeChat: cstutorcs

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

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

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

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

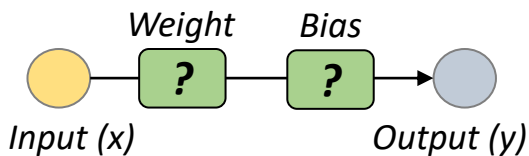
3 Deep Learning for NLP

Deep Learning with Neural Network - Optimizer

Actual Data

$$y = \text{weight} \cdot x + \text{bias}$$

x 🍎	y 🍌
1	0
5	16
6	20



Assignment Project Exam Help

<https://tutorcs.com>

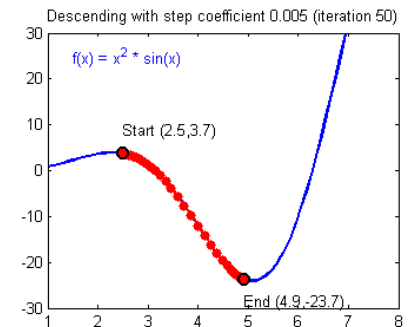
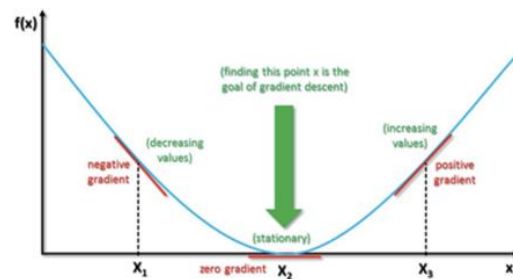
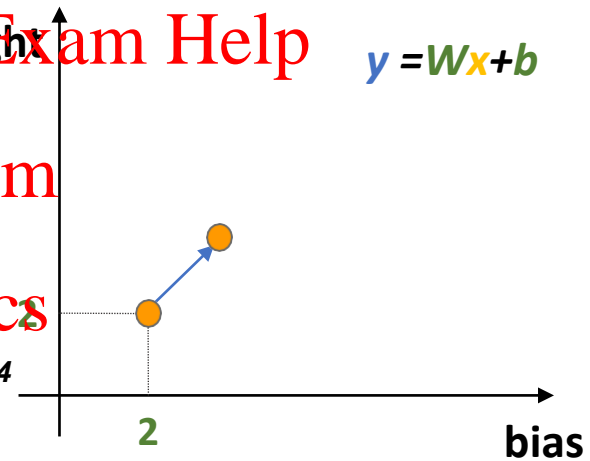
WeChat: cstutorcs

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

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

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

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



3 Deep Learning for NLP

Deep Learning with Neural Network

<u>Data</u>	<u>Model</u>	<u>Cost</u>	<u>Optimizer</u>								
<table><tr><th>x 🍎</th><th>y 🍌</th></tr><tr><td>1</td><td>0</td></tr><tr><td>5</td><td>16</td></tr><tr><td>6</td><td>20</td></tr></table>	x 🍎	y 🍌	1	0	5	16	6	20	$\underset{\text{weight}}{y} = \underset{\text{weight}}{?} \underset{\text{bias}}{x} + \underset{\text{bias}}{?}$	$C(w,b) = \sum_{n \in \{0,1,2\}} (y_n - \hat{y}_n)$	$\arg \min C(w,b)$ $w, b \in [-\infty, \infty]$
x 🍎	y 🍌										
1	0										
5	16										
6	20										

Assignment Project Exam Help

<https://tutorcs.com>

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



System

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



Deep Learning with Neural Network

Input: x = number of apple given by Lisa

Output: y = number of banana received by Lisa

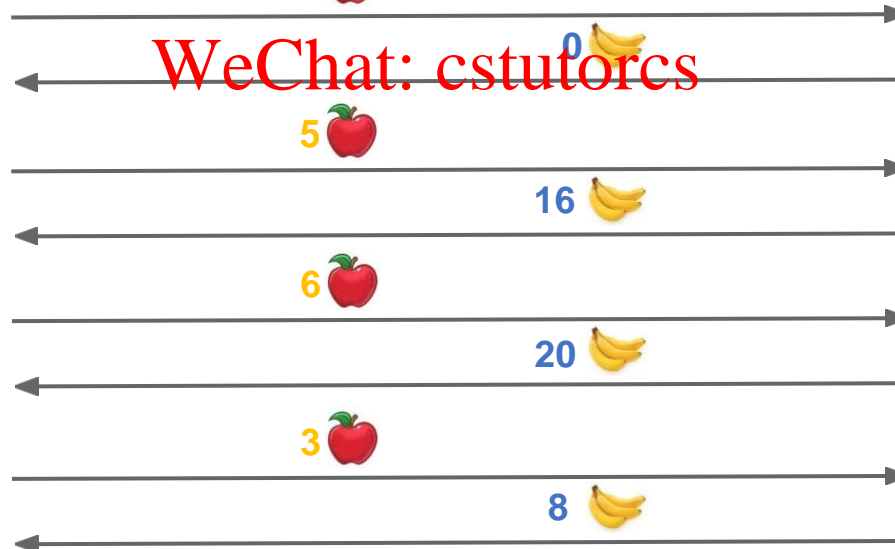
Parameters: Need to be estimated

Assignment Project Exam Help $y = 4x - 4$



<https://tutorcs.com>

WeChat: cstutorcs



Deep Learning with Neural Network

$$y = \overset{\text{pixel}(1,1)}{W_1} \overset{\text{pixel}(1,2)}{X_1} + W_2 X_2 + W_3 X_3 + W_4 X_4 + \dots + W_n X_n + b$$

Assignment Project Exam Help

Data

<https://tutorcs.com>

WeChat: cstutorcs



Millions of Parameters
Millions of Samples

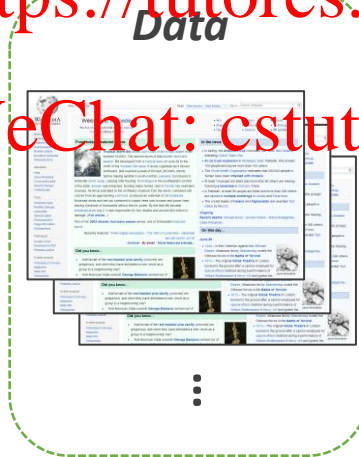
Deep Learning with Neural Network

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



Millions of *Parameters*

Millions of *Samples*

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

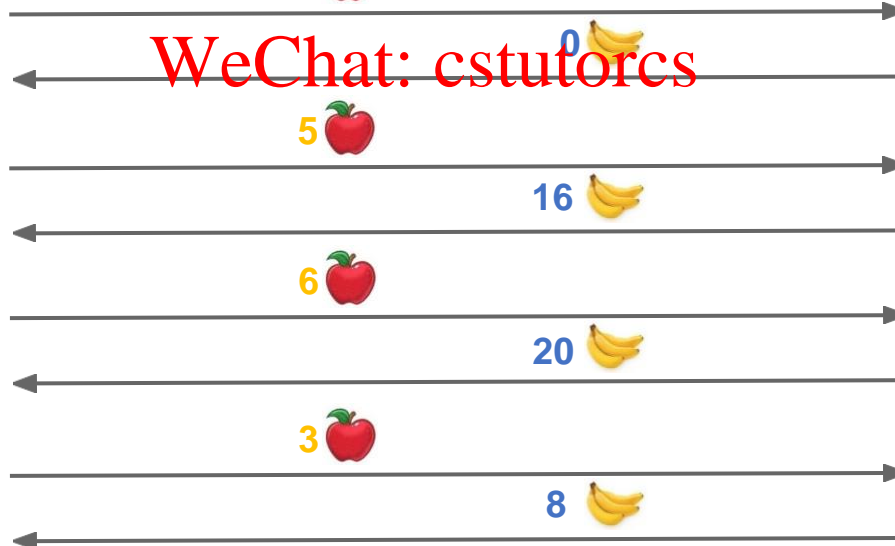
Assignment Project Exam Help

There is a limit of bananas I can give you



<https://tutorcs.com>

WeChat: cstutorcs



Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

Data

x 🍎	y 🍌
1	0
5	16
6	20

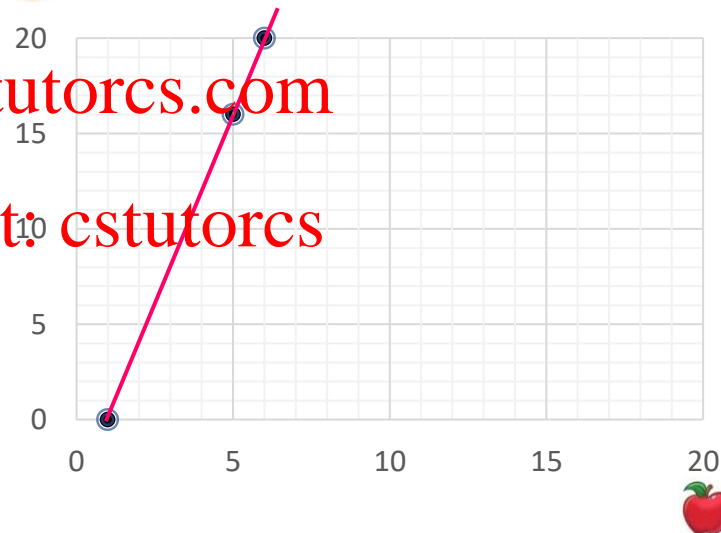
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = 4x - 4$$

weight bias



Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

Data

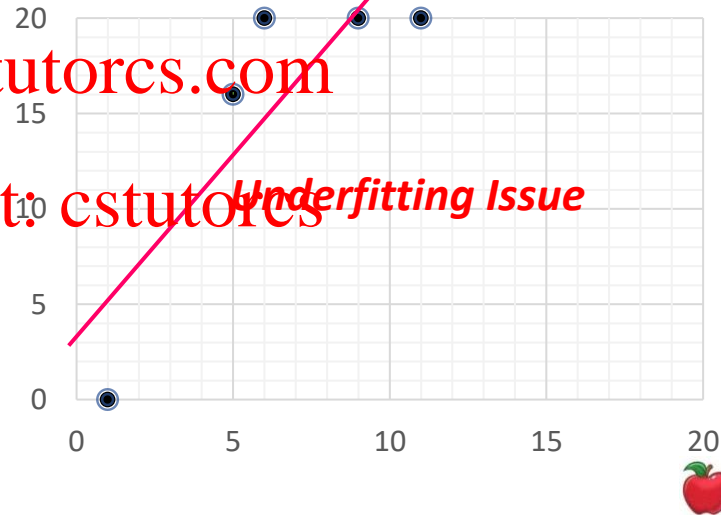
x 🍎	y 🍌
1	0
5	16
6	20
9	20
11	20

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = \overset{\text{weight}}{2} \overset{\text{bias}}{x} + 3$$



Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

Data

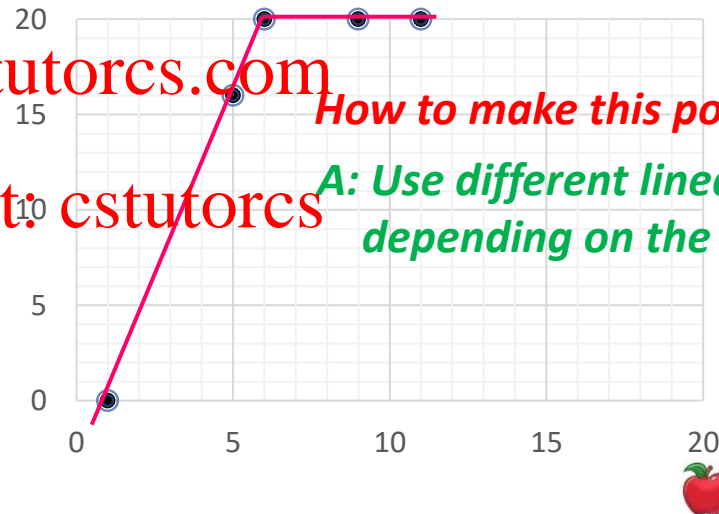
x 🍎	y 🍌
1	0
5	16
6	20
9	20
11	20

Assignment Project Exam Help

$y = ???$

<https://tutorcs.com>

WeChat: cstutorcs



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

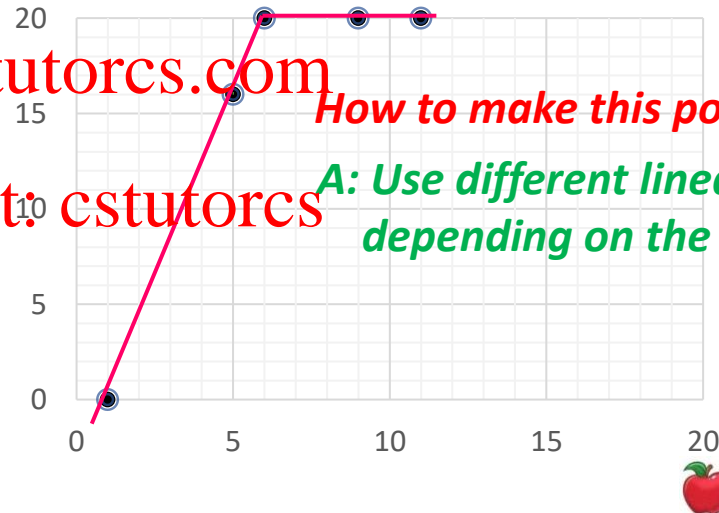
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

$$y = (\overset{\text{weight}}{w_1}x + \overset{\text{bias}}{b_1})s_1 + (\overset{\text{weight}}{w_2}x + \overset{\text{bias}}{b_2})s_2$$

If $x < 6$ and 0 If $x \geq 6$ and 0



How to make this possible?

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

Deep Learning with Neural Network

Nonlinear Neural Network

Data

x 🍎	y 🍌
1	0
5	16
6	20
9	20
11	20

Assignment Project Exam Help

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

weight
bias
weight
bias

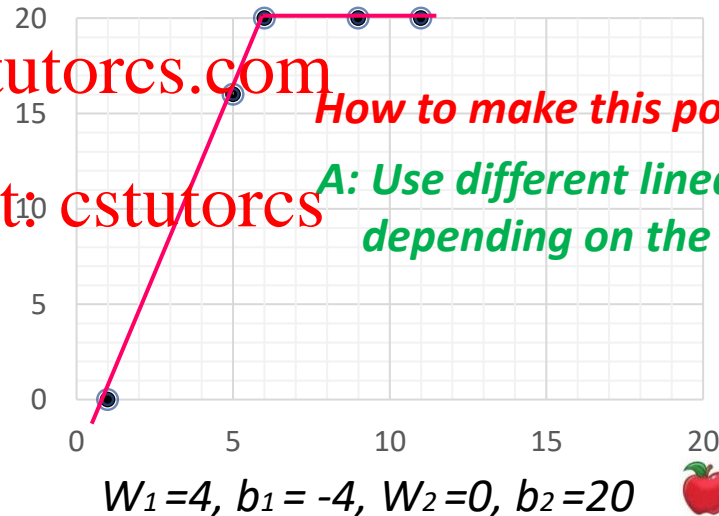
If $x < 6$ and 0
If $x \geq 6$ and 0

<https://tutorcs.com>

WeChat: cstutorcs

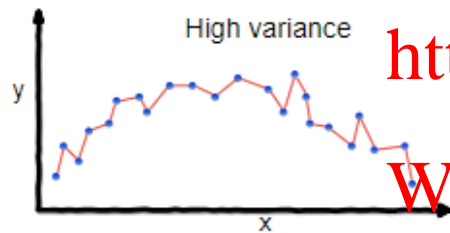
How to make this possible?

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

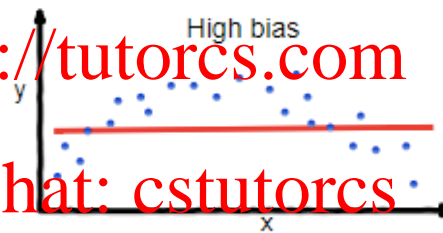


Deep Learning with Neural Network

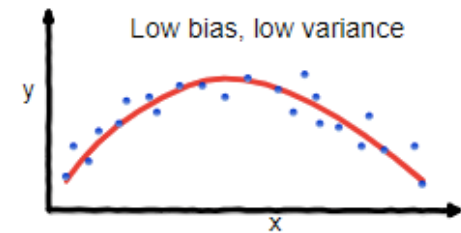
Assignment Project Exam Help



overfitting

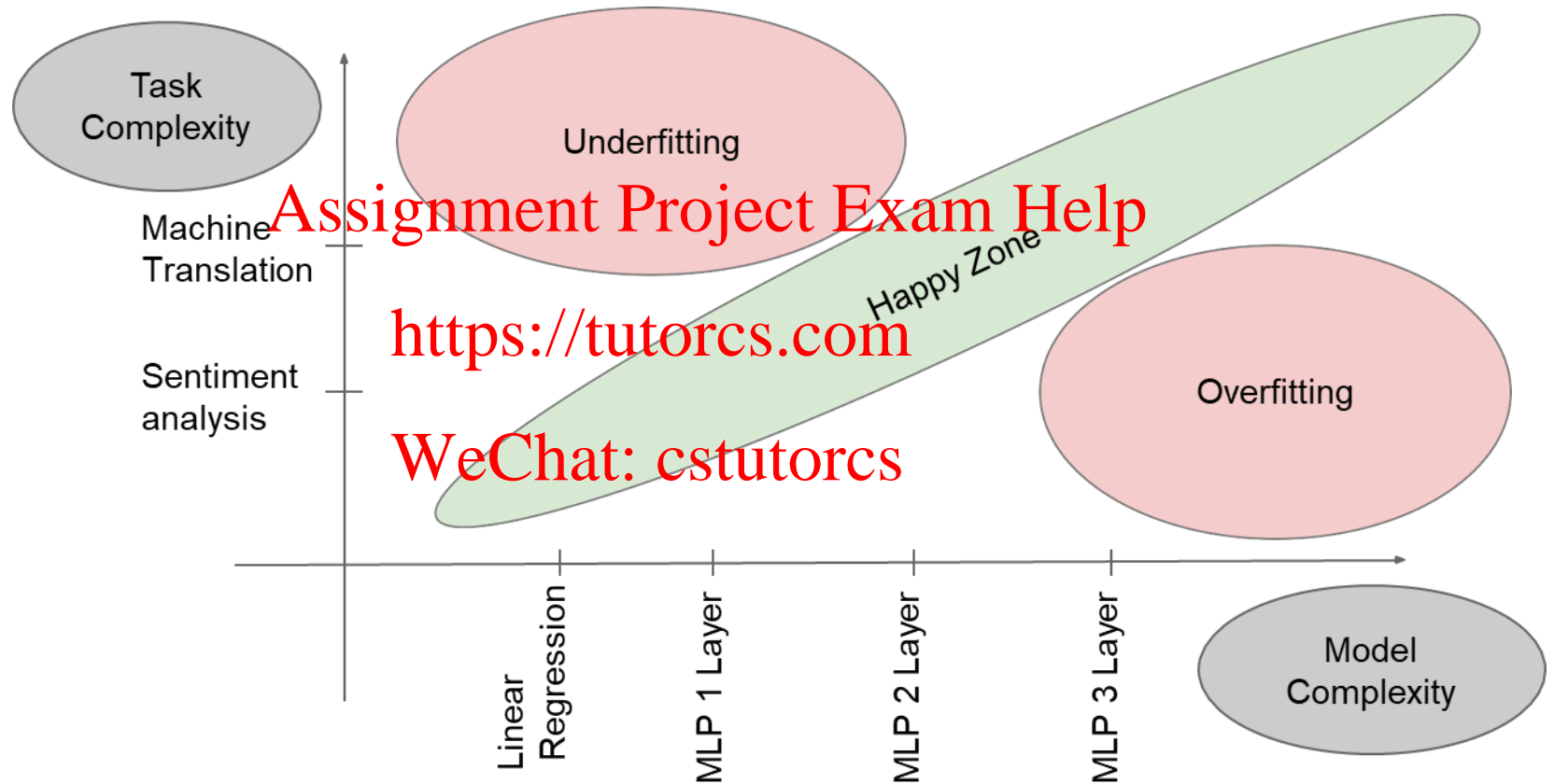


underfitting



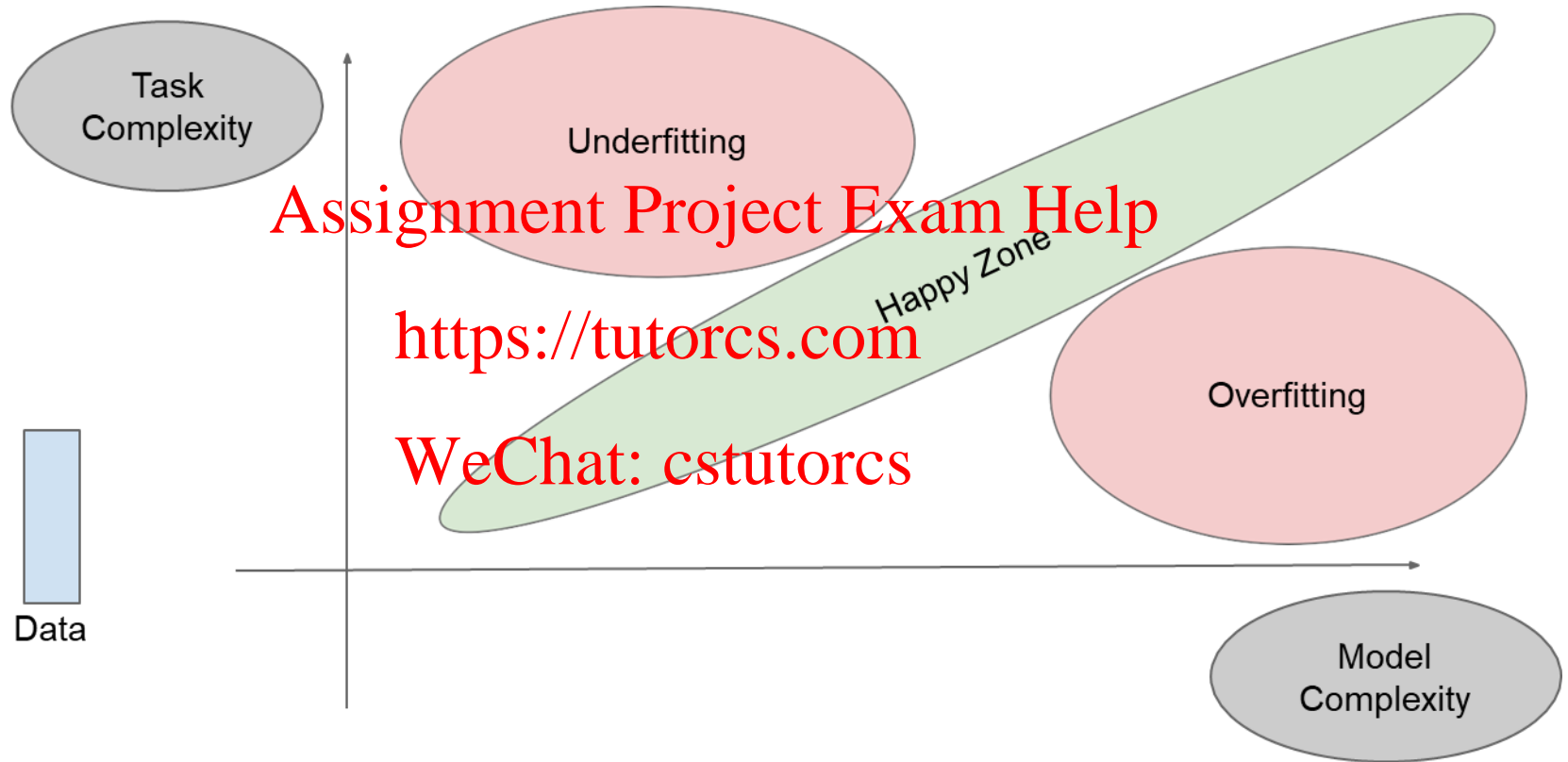
Good balance

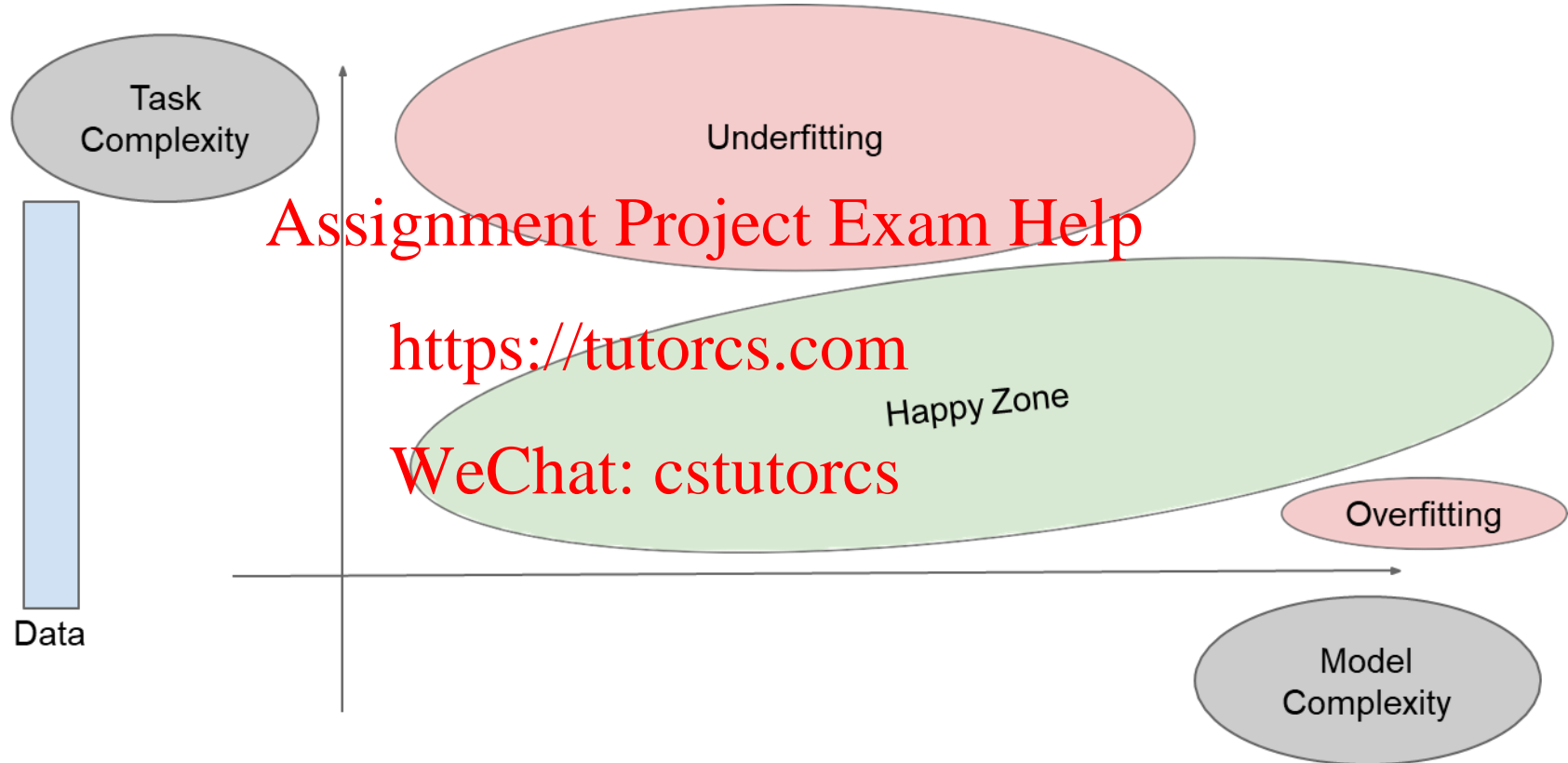
<https://tutorcs.com>
WeChat: cstutorcs



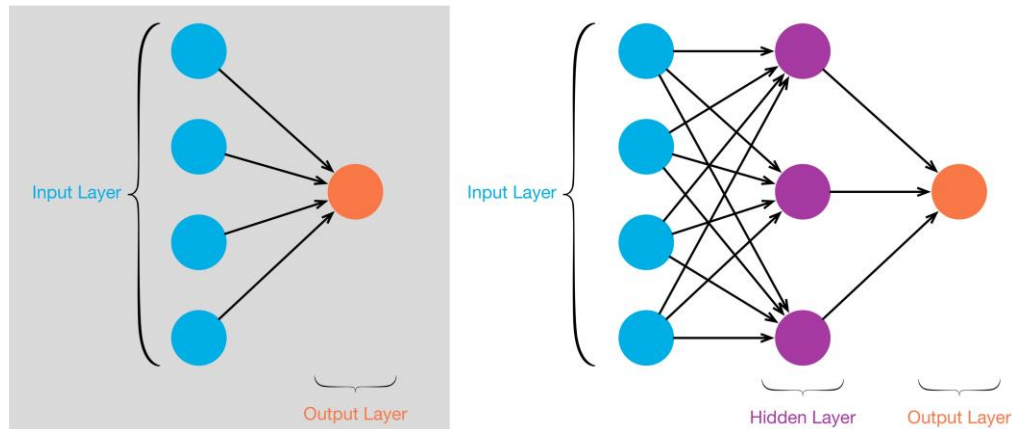
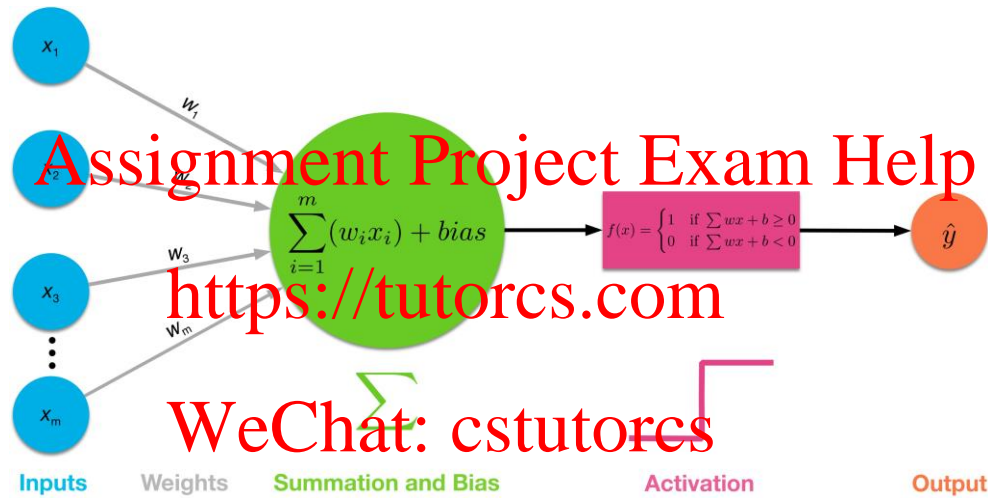
3

Deep Learning for NLP





Single Neuron VS Multilayer



CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)

Assignment Project Exam Help

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

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

<https://tutorcs.com>

WeChat: cstutorcs

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

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

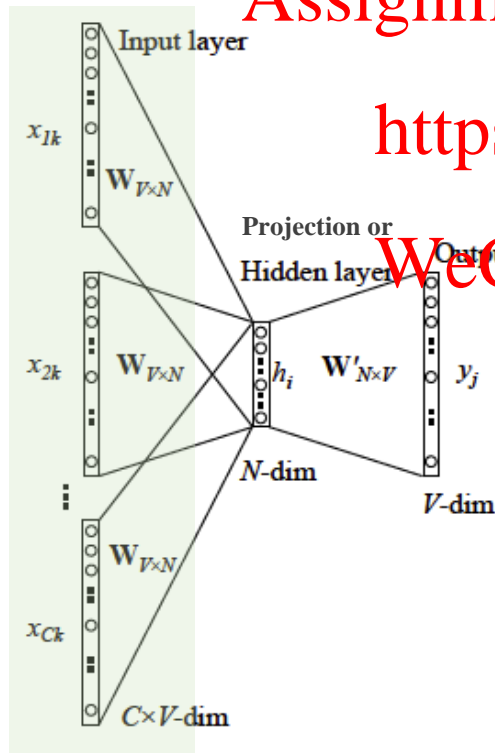
3. Has two Parameter Matrices:

1) Parameter Matrix (from Input Layer to Hidden/Projection Layer)

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

2) Parameter Matrix (to Output Layer)

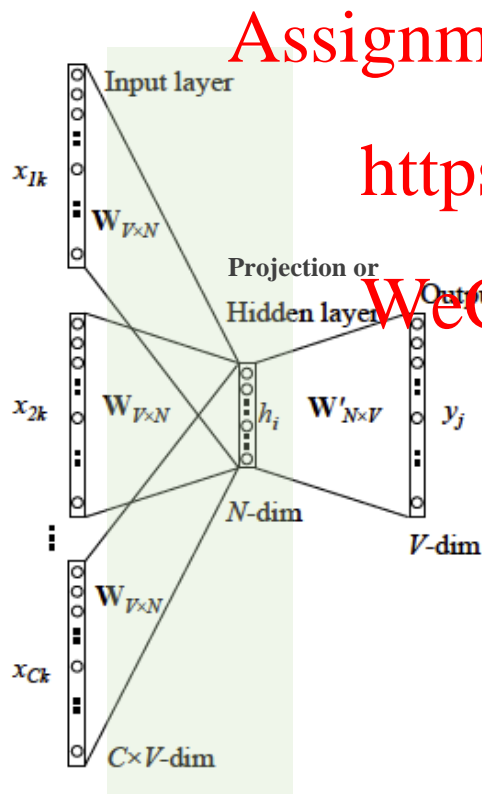
$$\mathbf{W}' \in \mathbb{R}^{N \times V}$$



CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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.

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

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

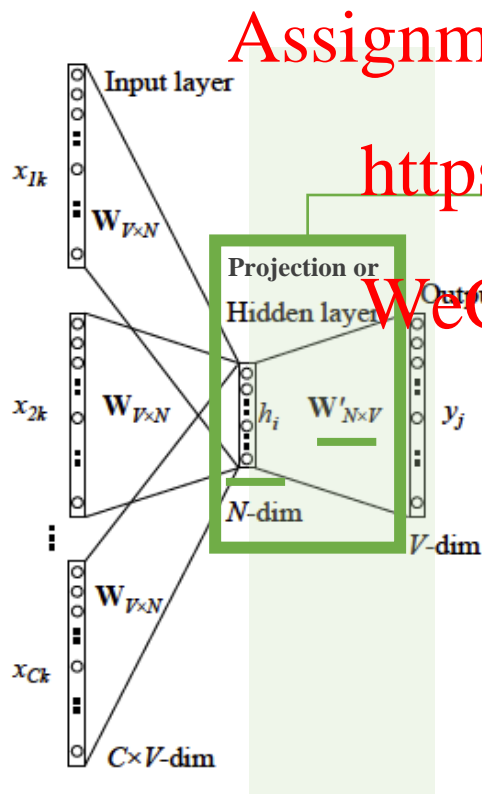
5. Average those $2m$ embedded vectors to calculate the value of the Hidden Layer.

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

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



Assignment Project Exam Help

<https://tutores.com>

WeChat: cstutores

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

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

7. Calculate the probability using softmax

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

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

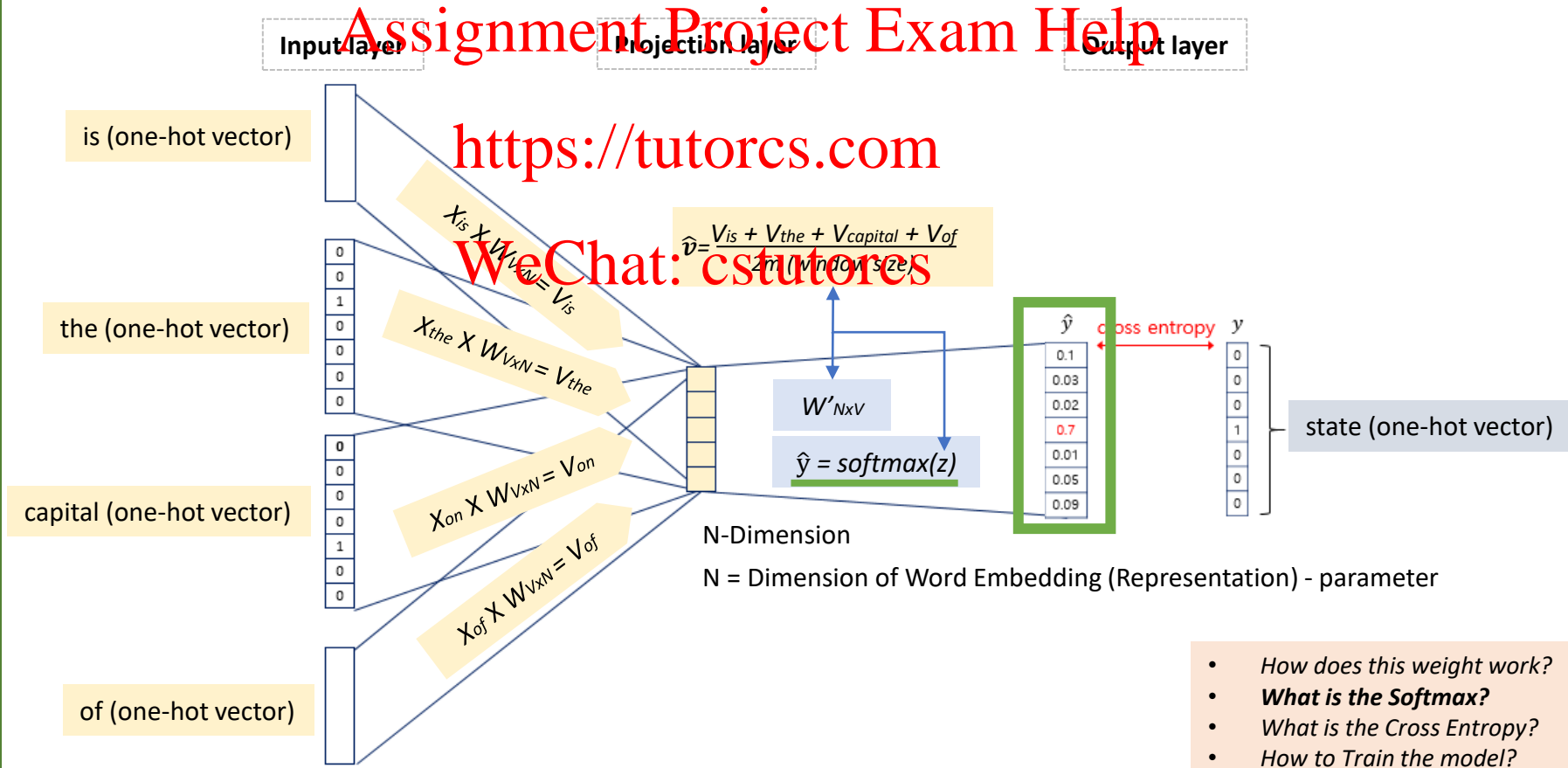
We use an one-hot vector (one 1, the rest 0) so it will be calculated in only one.

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

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words

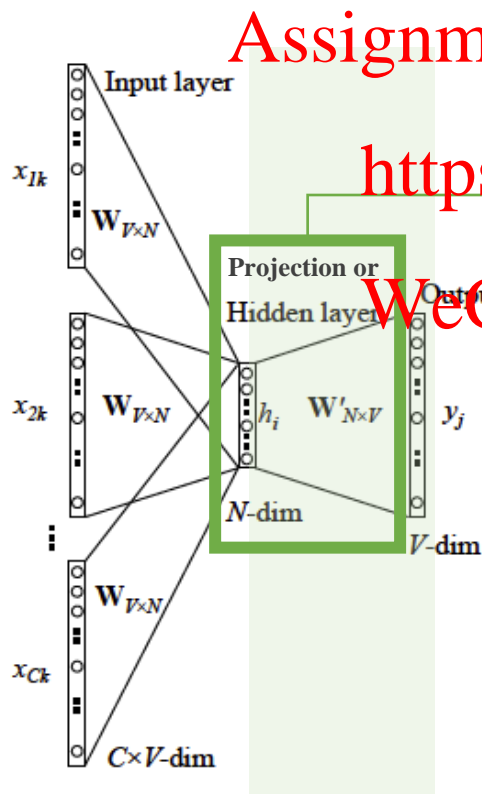
Sentence: “Sydney is the state capital of NSW”



CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



Assignment Project Exam Help

<https://tutores.com>

WeChat: cstutores

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

$$\mathbf{z} = \mathbf{W} \times \mathbf{h} \in \mathbb{R}^{|V|}$$

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:

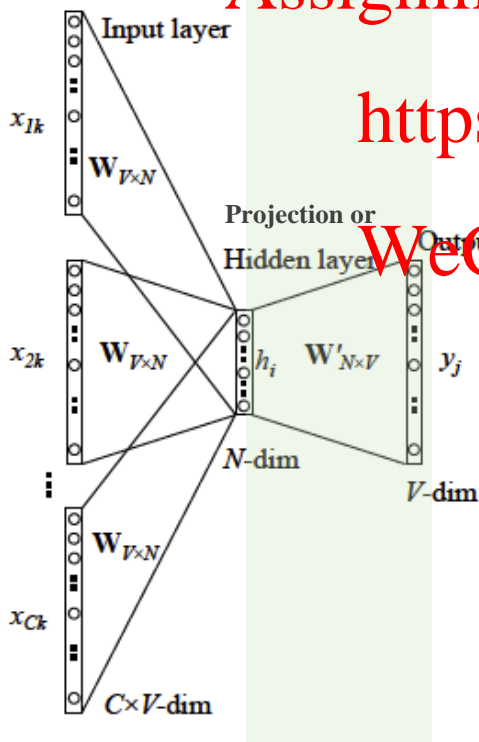
$$\frac{e^{\hat{y}_i}}{\sum_{j=1}^{|V|} e^{\hat{y}_j}}$$

- 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

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

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

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

7. Calculate the probability using softmax

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

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

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

Cross Entropy

* Focus on minimising the value

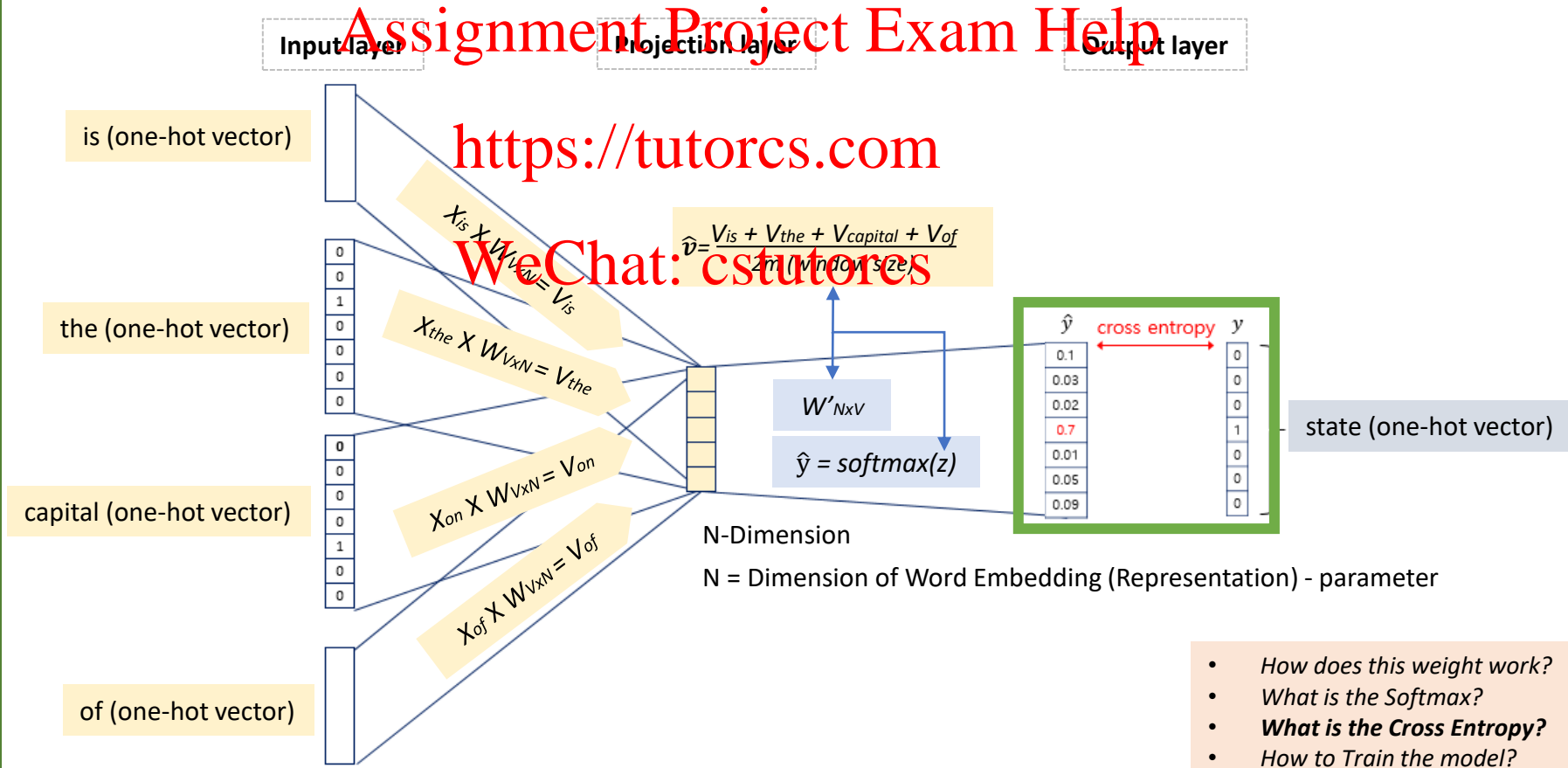
We use an one-hot vector (one 1, the rest 0) so it will be calculated in only one.

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

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words

Sentence: “Sydney is the state capital of NSW”

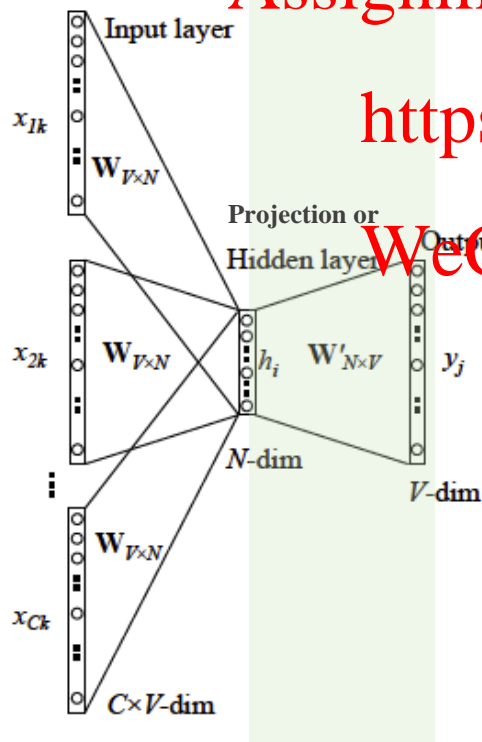


Deep Learning for NLP

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



Assignment Project Exam Help

8-1. Optimization Objective Function can be presented:

<https://tutorcs.com>

WeChat: cstutorcs

$$J = -\log P(w_c | w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m})$$

$$= -\log P(u_c | \hat{v})$$

$$= -\log \frac{\exp(u_c^T \hat{v})}{\sum_{j=1}^{|V|} \exp(u_j^T \hat{v})}$$

$$= -u_c^T \hat{v} + \log \sum_{j=1}^{|V|} \exp(u_j^T \hat{v})$$

u_i = the output vector representation of word w_i

Skip Gram – Neural Network Architecture

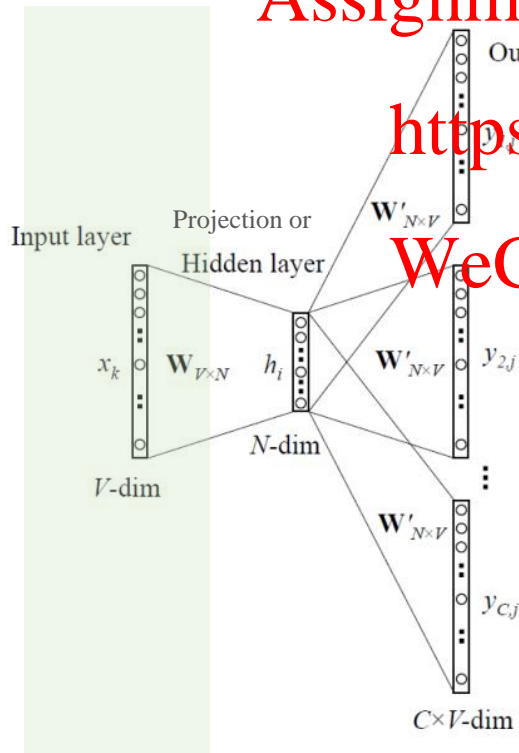
Predict context (“outside”) words (position independent) given center word

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



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

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

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

2. Initialise Parameter Matrices:

1) Parameter Matrix (from Input Layer to Hidden/Projection Layer)

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

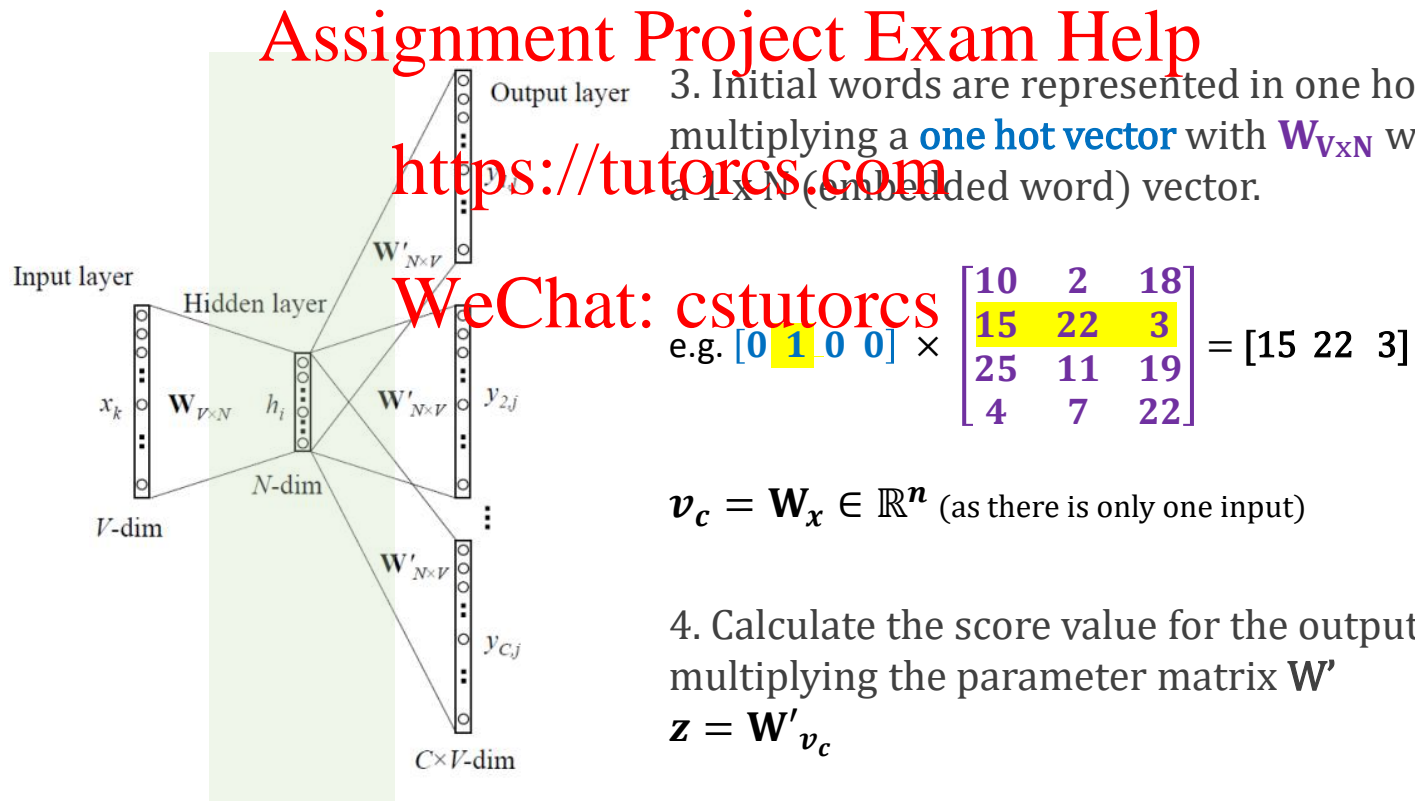
2) Parameter Matrix (to Output Layer)

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

Skip Gram – Neural Network Architecture

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

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



Skip Gram – Neural Network Architecture

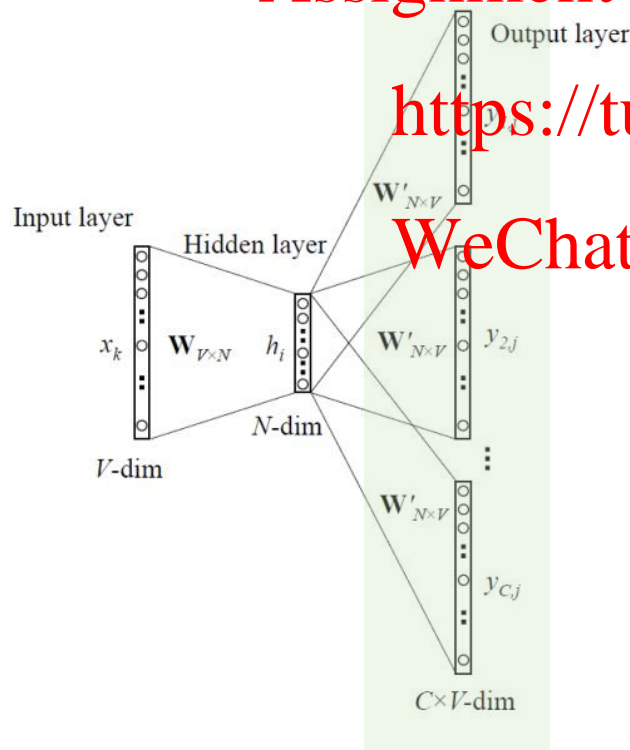
Predict context (“outside”) words (position independent) given center word

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



5. Calculate the probability using softmax

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

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)

$$y^{(c-m)}, \dots, y^{(c-1)}, y^{(c+1)}, \dots, y^{(c+m)}$$

Skip Gram – Neural Network Architecture

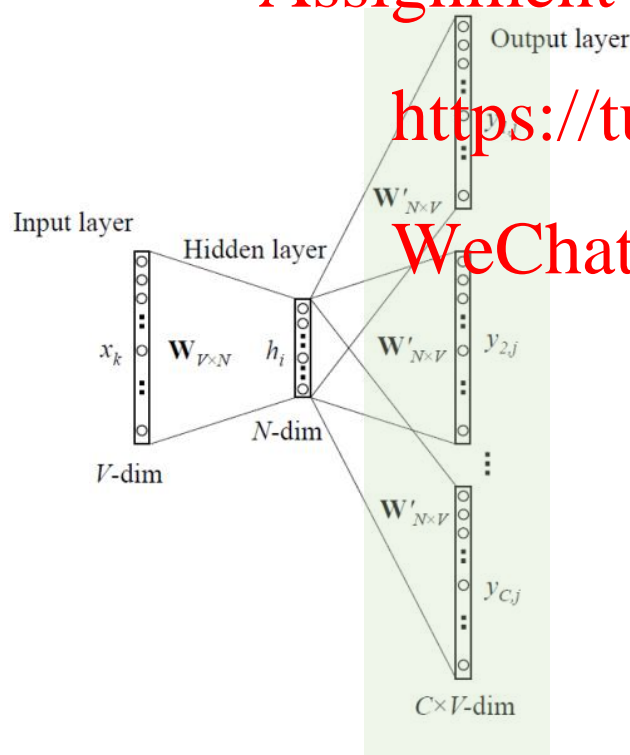
Predict context (“outside”) words (position independent) given center word

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



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.

$$\text{minimize } J = -\log P(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c)$$

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

u_i = the output vector representation of word w_i

Skip Gram – Neural Network Architecture

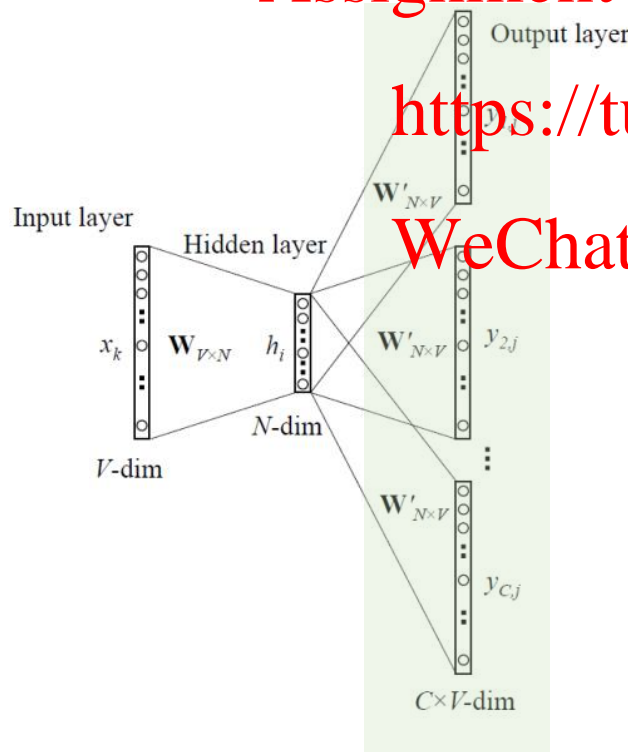
Predict context (“outside”) words (position independent) given center word

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

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



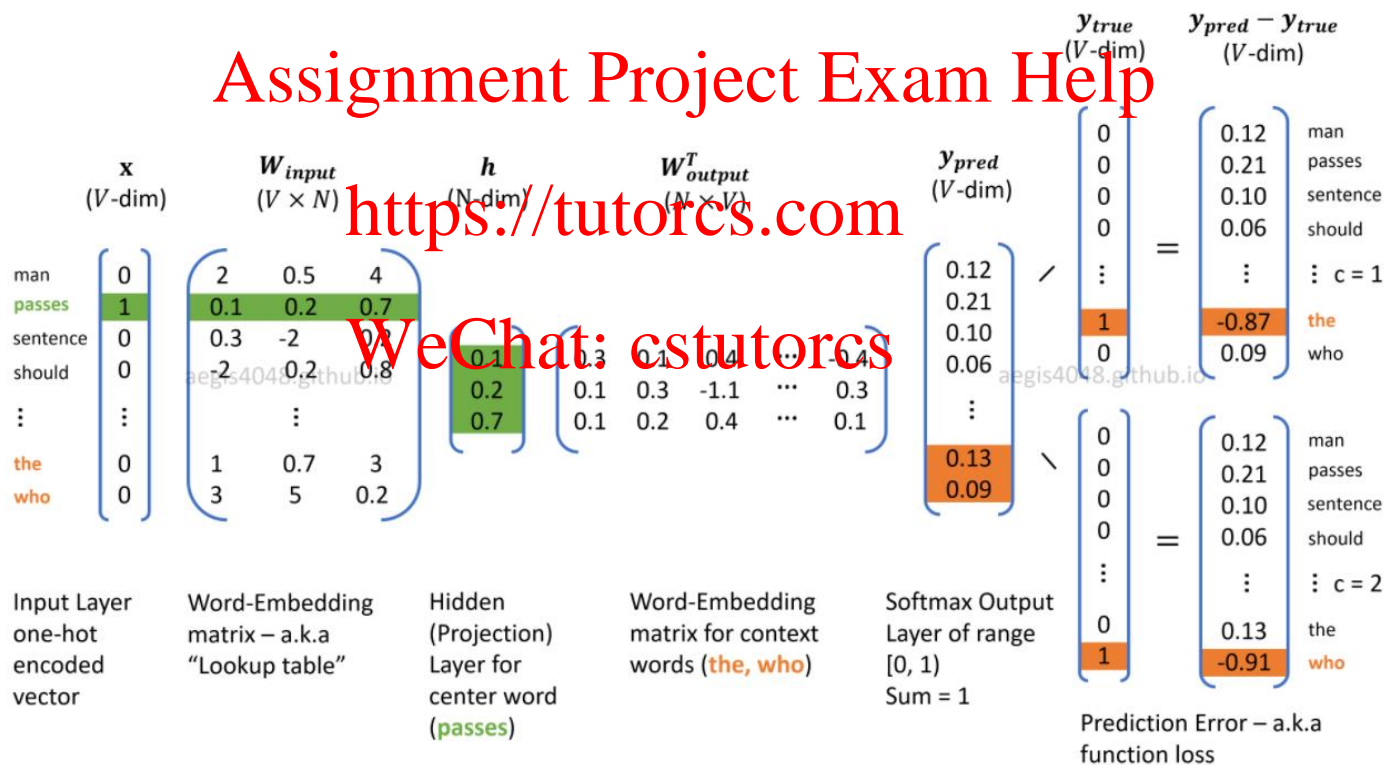
8-1. With this objective function, we can compute the gradients with respect to the unknown parameters and at each iteration update them via Stochastic Gradient Descent

$$J = - \sum_{j=0, j \neq m}^{2m} \log P(u_{c-m+j} | v_c)$$

$$= \sum_{j=0, j \neq m}^{2m} H(\hat{y}, y_{c-m+j})$$

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

Negative sample: 5

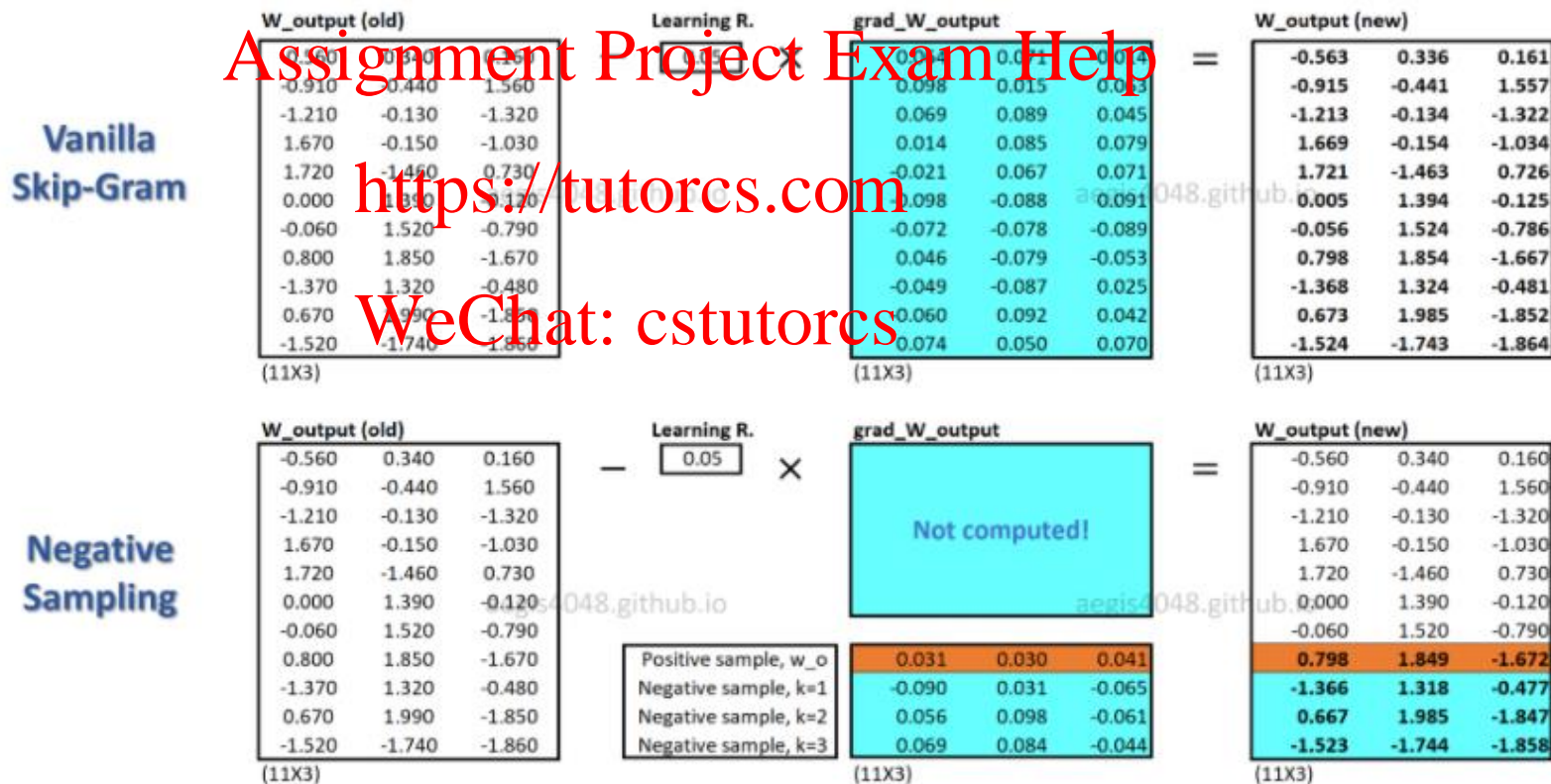
<i>Input word</i>	<i>Output word</i>	<i>Target</i>
eat	mango	1
eat	exam	0
eat	tobacco	0
eat	pool	0
eat	supervisor	0

**1= Appeared, 0=Not Appeared*

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.

Word2Vec-SkipGram Overview – negative sampling

With a simple diagram



Application

Application #1: Embedding Pretraining



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 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

Assignment Project Exam Help

<https://tutores.com>

WeChat: cstutorcs