

COMP5046

Natural Language Processing

Lecture 3: Word Classification and Machine Learning

Dr. Caren Han

Semester 1, 2022

School of Computer Science,
University of Sydney

Deep Learning



Language

Linguistics

NLP

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

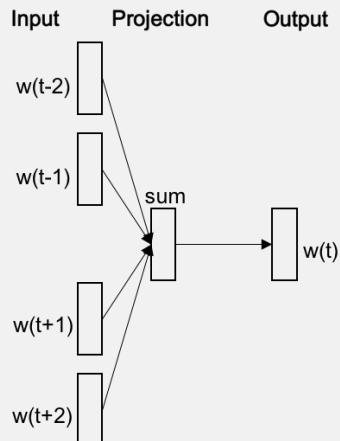
 - Text Classification, etc.

1 Previous Lecture Review

Word2Vec Models

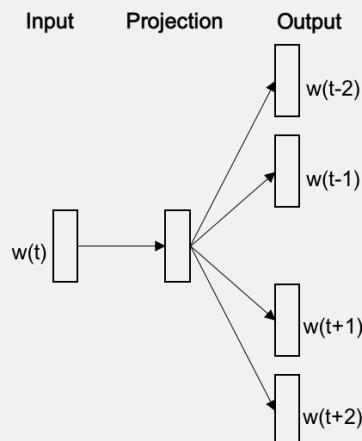
CBOW

*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 | |
|-----------------|--|------------------------------------|
| [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

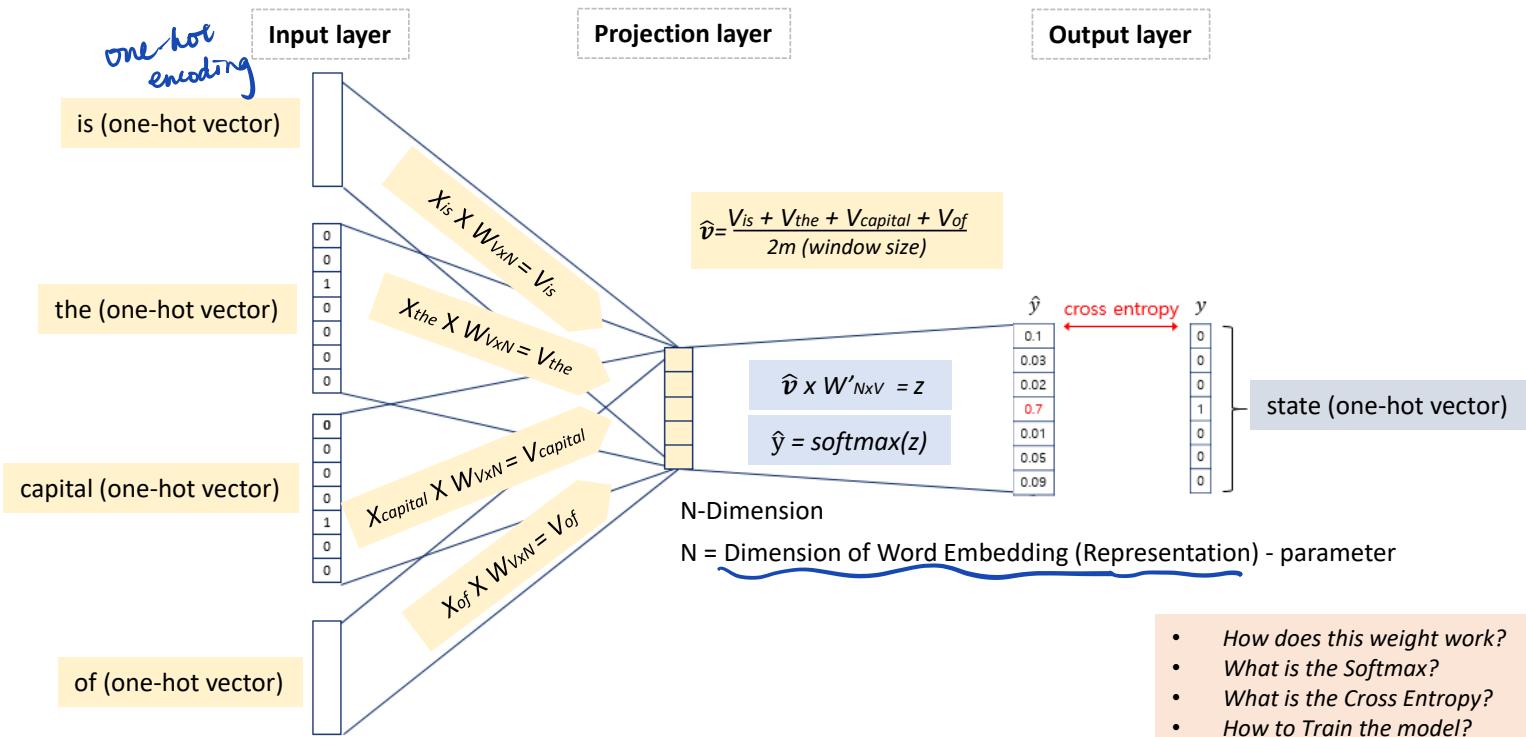
Context (“outside”) word

1 Previous Lecture Review

CBOW – Neural Network Architecture

Predict center word from (bag of) context words

Sentence: “Sydney is the state capital of NSW”

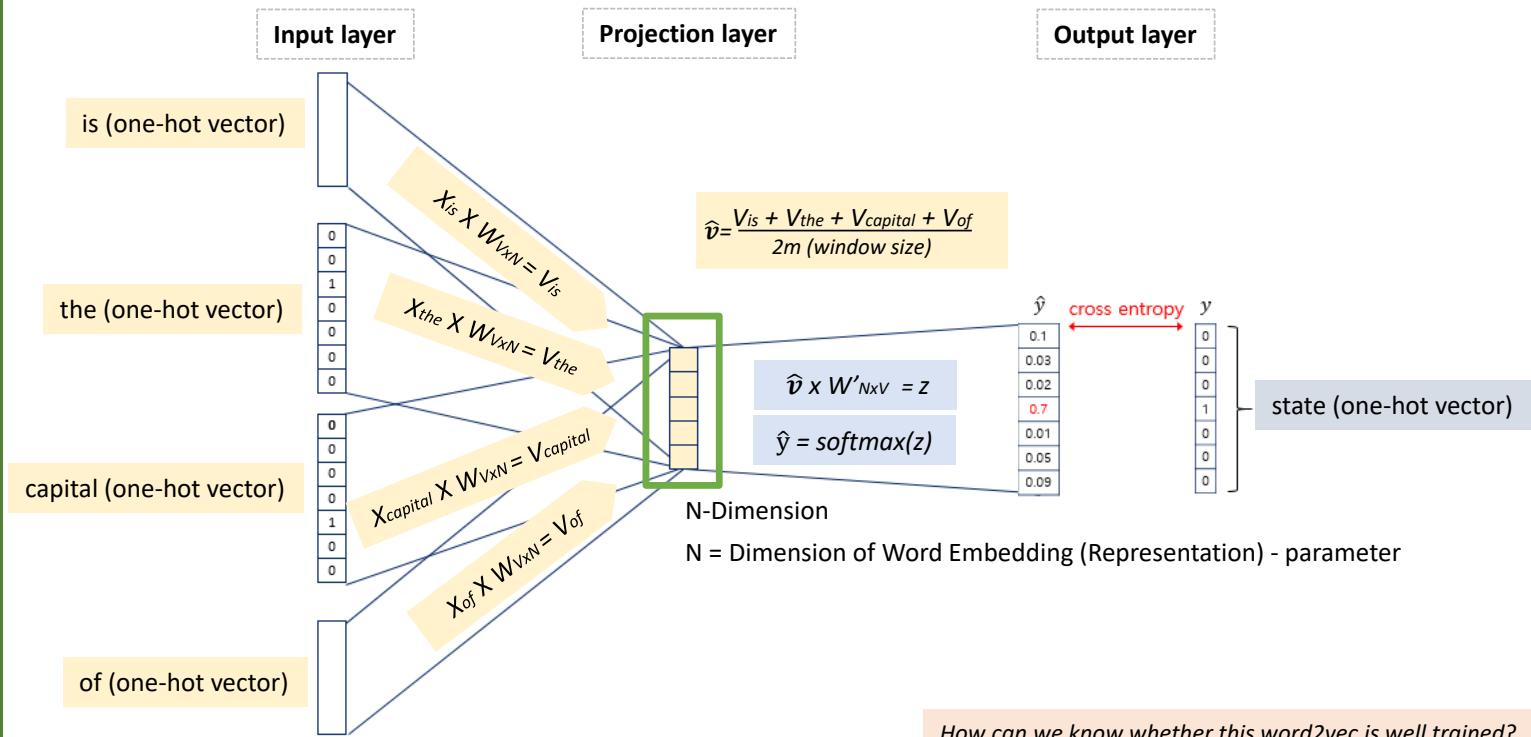


1 Previous Lecture Review

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

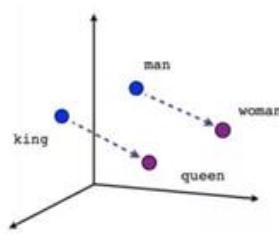
See how the Deep Learning can be used for NLP

 - Text Classification, etc.

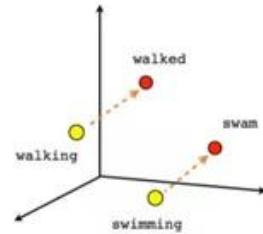
2 Word Embedding Evaluation

How to evaluate word vectors?

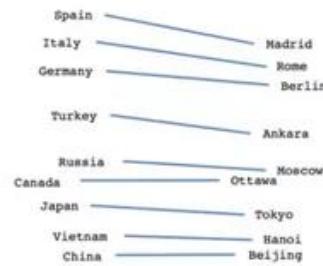
| Type | How to work / Benefit |
|-----------|--|
| Intrinsic | <p>Evaluation on a specific/intermediate subtask</p> <ul style="list-style-type: none"> • Fast to compute • Helps to understand that system • Not clear if really helpful unless correlation to real task is established |
| Extrinsic | <p>Evaluation on a real task</p> <ul style="list-style-type: none"> • Can take a long time to compute accuracy • Unclear if the subsystem is the problem or its interaction or other subsystems <p><i>"we try to find the similarity representation of actual word"</i></p> <p><i>"we on other tasks, eg task classification, question answering"</i></p> <p><i>"word embedding will be used as input to extra task"</i></p> |



Male-Female



Verb tense



Country-Capital

2 Word Embedding Evaluation

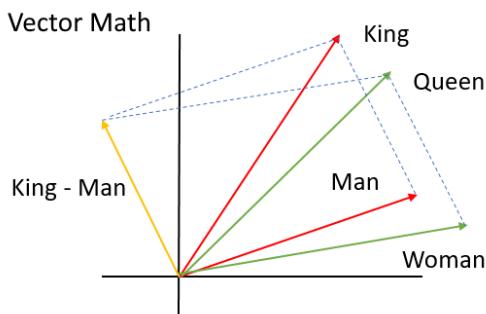
Intrinsic word vector evaluation

Word Vector Analogies

$$a \leftrightarrow b :: c \leftrightarrow ???$$

$$\text{man} \leftrightarrow \text{women} :: \text{king} \leftrightarrow ???$$

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



2 Word Embedding Evaluation

Intrinsic word vector evaluation

Word Vector Analogies

Depends on training set, type

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 |

2 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 |

(Original Word2vec Paper - Mikolov et al.2013)

2 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 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.

*CBOW: predict the actual center word.
skip-gram: predict context words*

| 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 | <u>16.1</u> | 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 | <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)

2 Word Embedding Evaluation

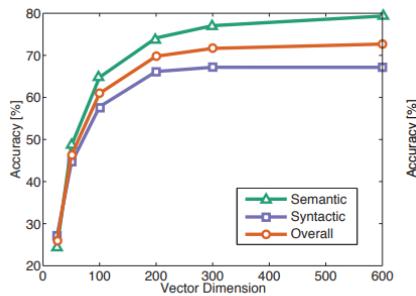
Intrinsic word vector evaluation

Evaluation Result Comparison

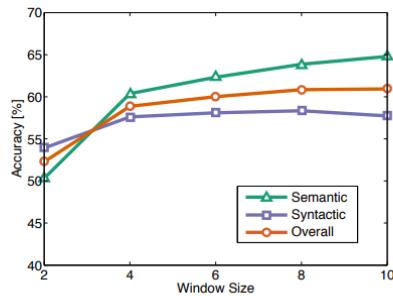
The Semantic-Syntactic word relationship tests for understanding of a wide variety of relationships as shown below.

Other parameters affects accuracy

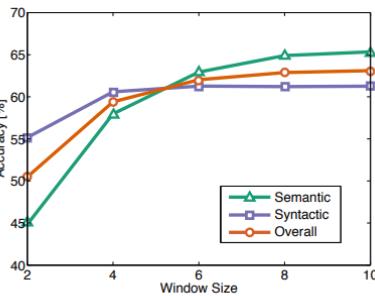
Window-Size (m) and Vector Dimension (N)



(a) Symmetric context



(b) Symmetric context



(c) Asymmetric context

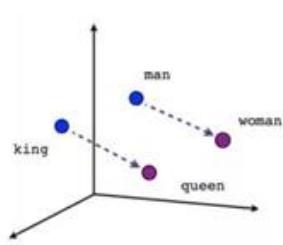
Figure 2: Accuracy on the analogy task as function of vector size and window size/type. All models are trained on the 6 billion token corpus. In (a), the window size is 10. In (b) and (c), the vector size is 100.

Also: dataset is general / representative, not specific
sometimes / only works well in this small dataset (Original Glove Paper - Pennington et al. 2014)
the algorithms

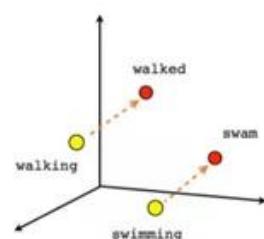
2 Word Embedding Evaluation

How to evaluate word vectors?

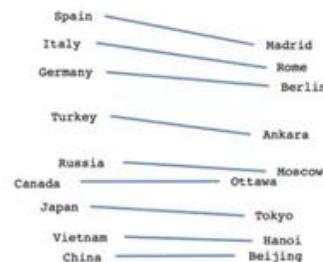
| Type | How to work / Benefit |
|-----------|--|
| Intrinsic | Evaluation on a specific/intermediate subtask <ul style="list-style-type: none"> • Fast to compute • Helps to understand that system • Not clear if really helpful unless correlation to real task is established |
| Extrinsic | Evaluation on a real task <ul style="list-style-type: none"> • 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

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

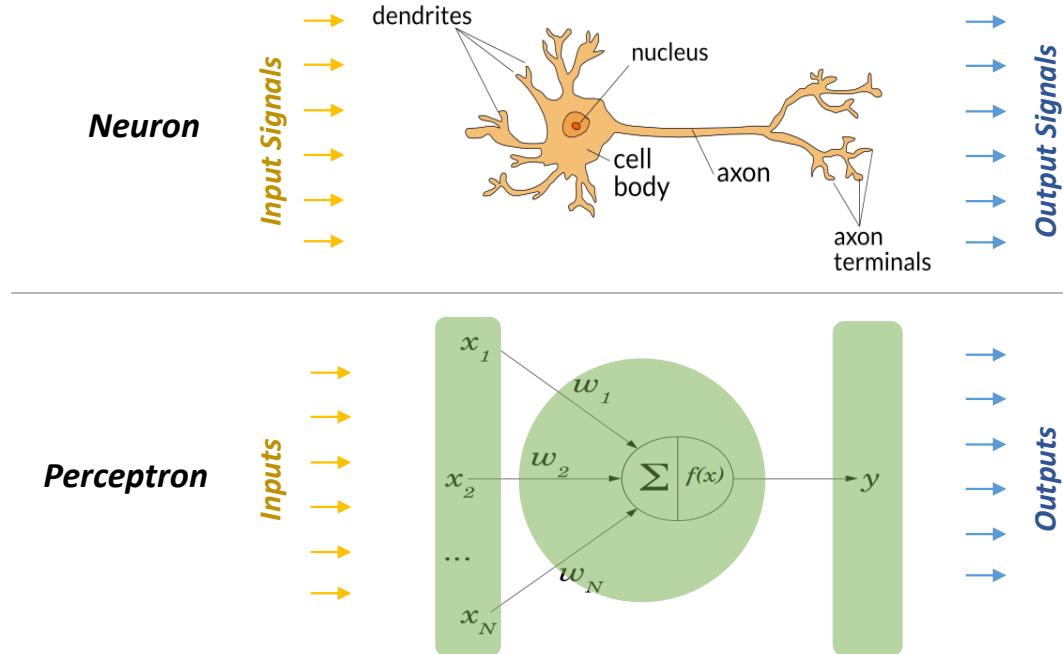
See how the Deep Learning can be used for NLP

 - Text Classification, etc.

3 Deep Learning for NLP

Deep Learning with Neural Network

Neuron and Perceptron

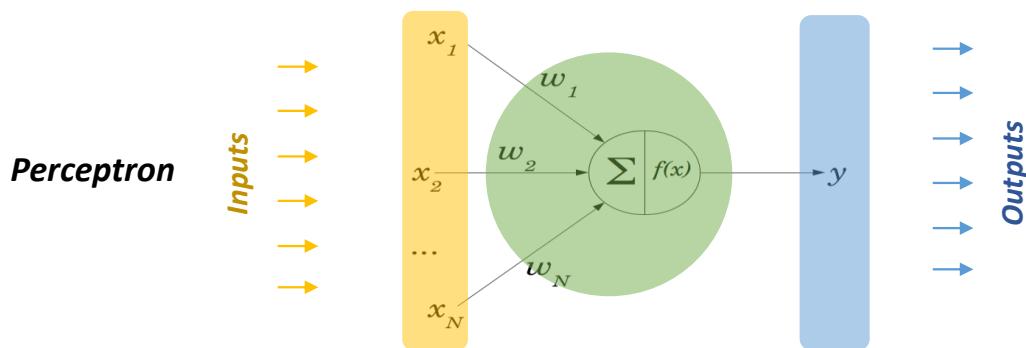


3 Deep Learning for NLP

Deep Learning with Neural Network

Inputs and Outputs (Labels) for Natural Language Processing

| | | |
|-------|------------------|--|
| x_i | Inputs | 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 |



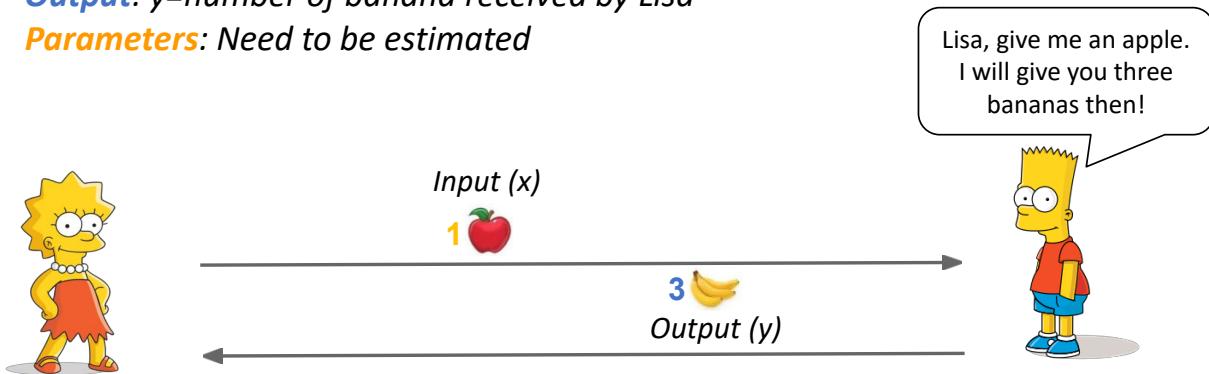
3 Deep Learning for NLP

Deep Learning with Neural Network

Input: $x=\text{number of apple given by Lisa}$

Output: $y=\text{number of banana received by Lisa}$

Parameters: Need to be estimated



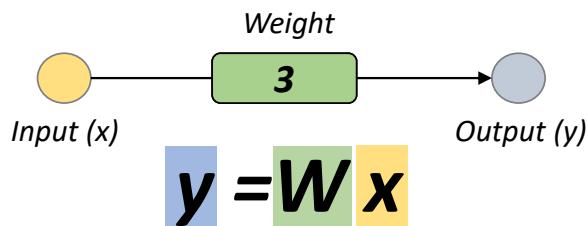
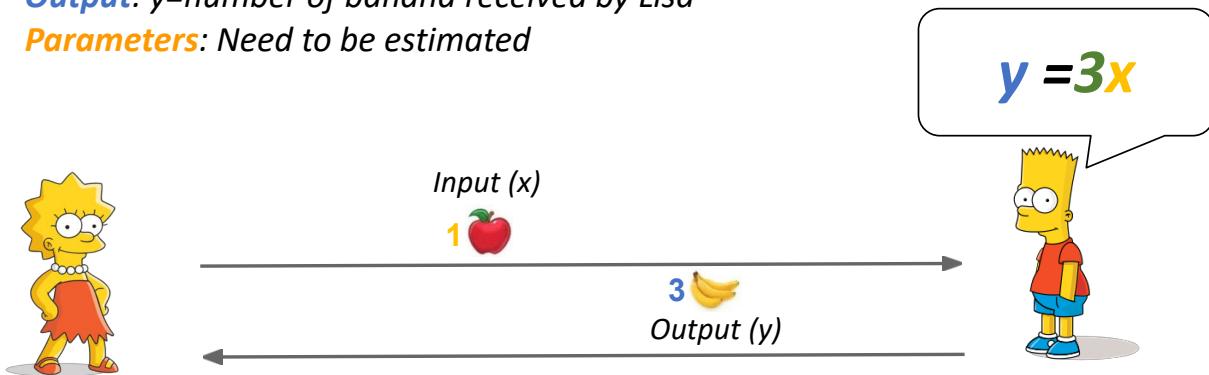
3 Deep Learning for NLP

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

Deep Learning with Neural Network - Model

Input: $x=\text{number of apple given by Lisa}$

Output: $y=\text{number of banana received by Lisa}$

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



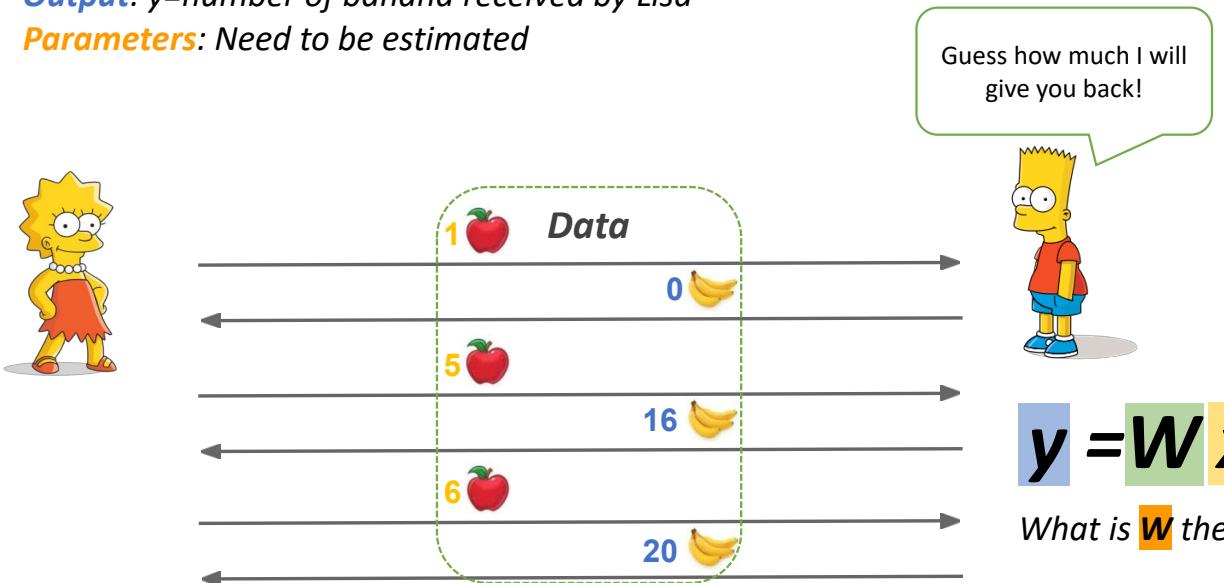
3 Deep Learning for NLP

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 |

$$y = Wx$$

What is W then?

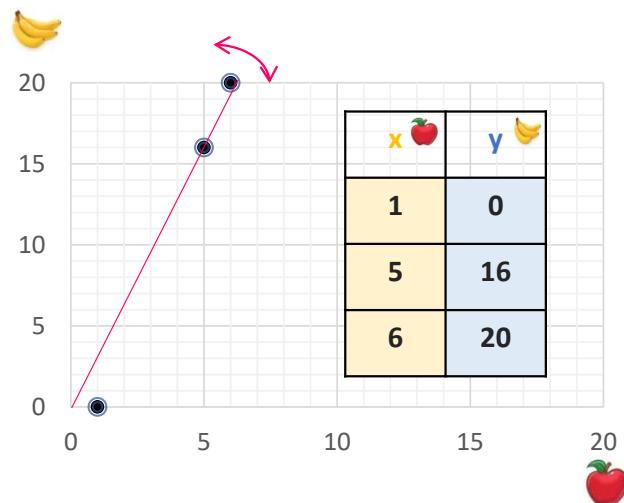
3 Deep Learning for NLP

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



$$y = Wx$$

What is W then?

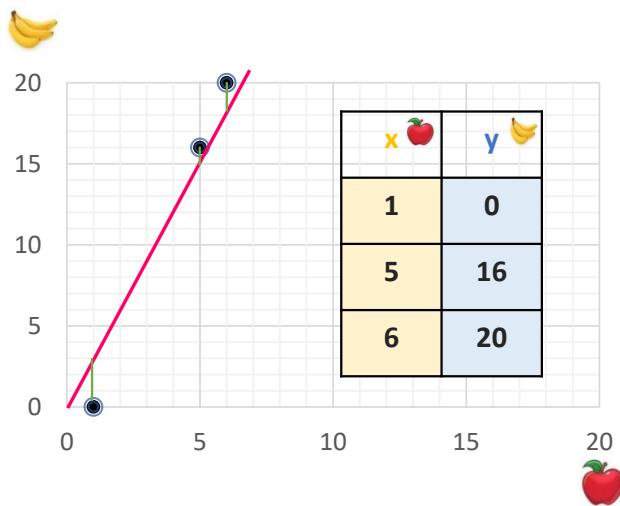
3 Deep Learning for NLP

Deep Learning with Neural Network - Parameter

Input: $x = \text{number of apple given by Lisa}$

Output: $y = \text{number of banana received by Lisa}$

Parameters: Need to be estimated



$$y = Wx$$

What if W is 3?

$$3 = 3 \times 1$$

$$15 = 3 \times 5$$

$$20 = 3 \times 6$$

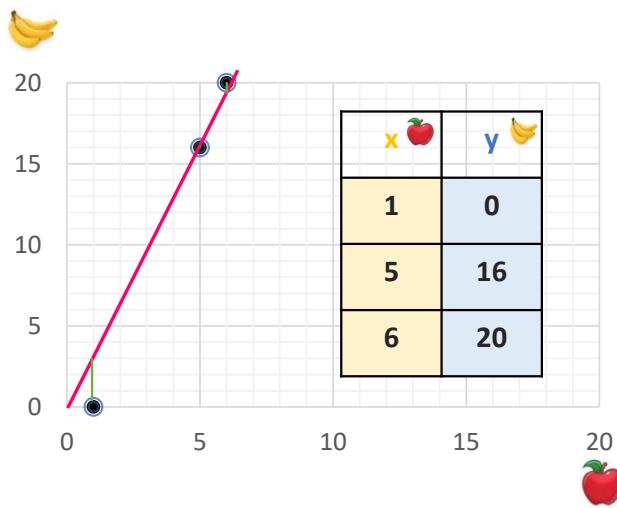
3 Deep Learning for NLP

Deep Learning with Neural Network - Parameter

Input: $x = \text{number of apple given by Lisa}$

Output: $y = \text{number of banana received by Lisa}$

Parameters: Need to be estimated



$$y = Wx$$

What if W is 3.2?

$$3.2 = 3.2 \times 1$$

$$16 = 3.2 \times 5$$

$$19.2 = 3.2 \times 6$$

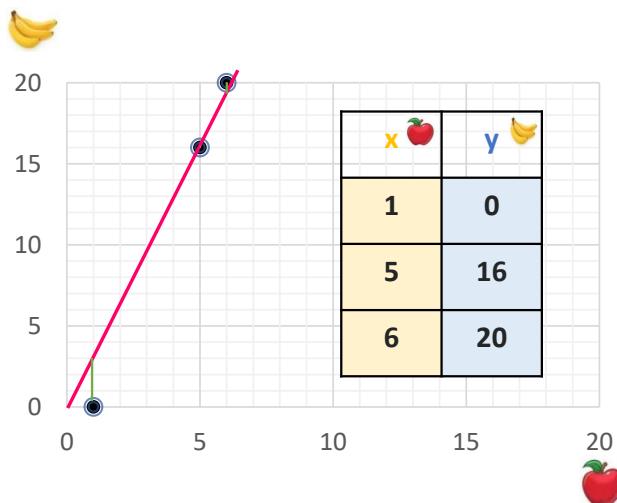
3 Deep Learning for NLP

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



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

Weight is not enough...

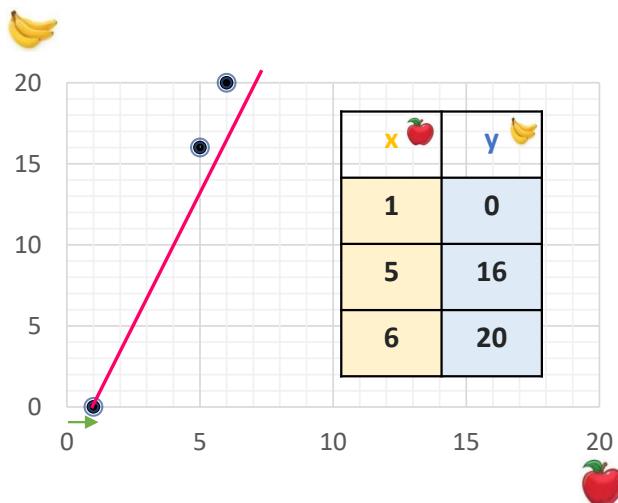
3 Deep Learning for NLP

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



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

How can we find the parameters, w and b ?

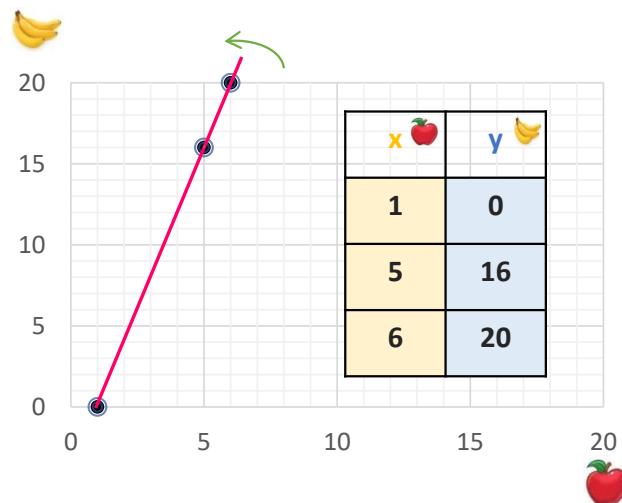
3 Deep Learning for NLP

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



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

How can we find the parameters, W and b ?

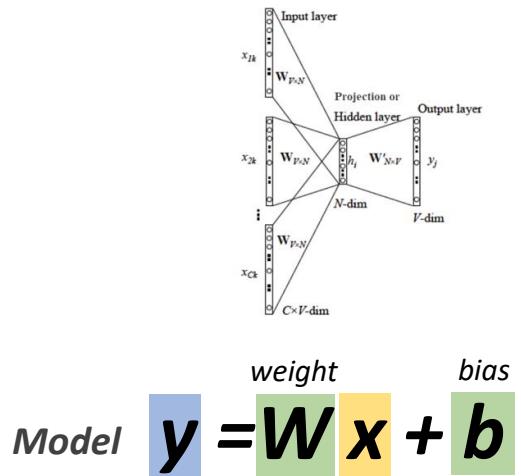
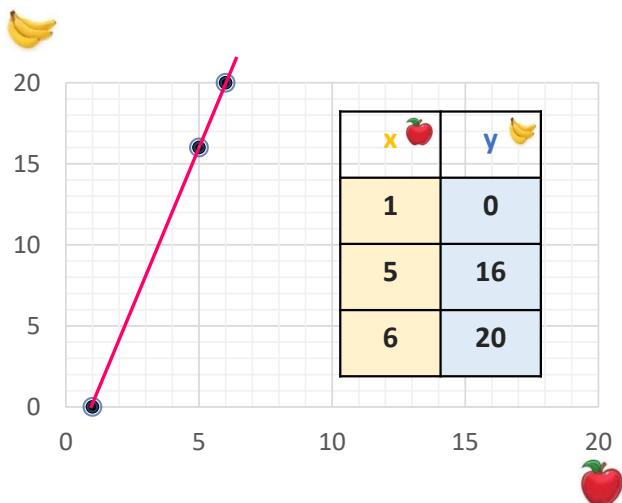
3 Deep Learning for NLP

Deep Learning with Neural Network - Parameter

Input: $x = \text{number of apple given by Lisa}$

Output: $y = \text{number of banana received by Lisa}$

Parameters: Need to be estimated



How can we find the parameters, w and b ?

3 Deep Learning for NLP

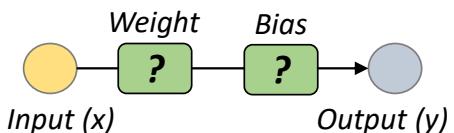
Deep Learning with Neural Network - Cost

Actual Data

$$y = ? \ x + ?$$

weight bias

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



Model Ex#1

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

weight bias

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

Model Ex#2

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

weight bias

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



Which one is closer?

3 Deep Learning for NLP

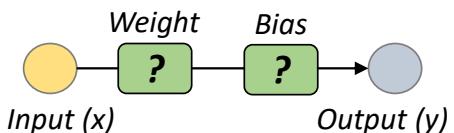
Deep Learning with Neural Network – Cost (loss)

Actual Data

$$y = ? \ x + ?$$

weight bias

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



Model Ex#1

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

weight bias

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

$$(y - \hat{y})^2$$

Model Ex#2

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

weight bias

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

$$(y - \hat{y})^2$$

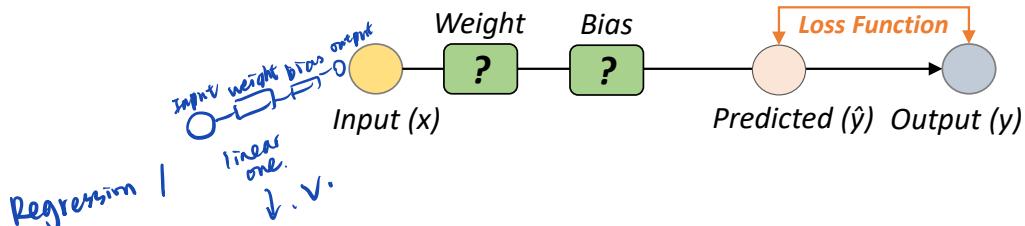

Let's calculate the cost(loss)!

Mean Squared Error (MSE) $C(w, b) = \sum (y_n - \hat{y}_n)^2$

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

3 Deep Learning for NLP

WAIT! Loss Function? Cost Calculation?



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

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

Classification ↗

Decision boundary for classification is large
why cross entropy is better for classification?
• MSE doesn't punish misclassification enough

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

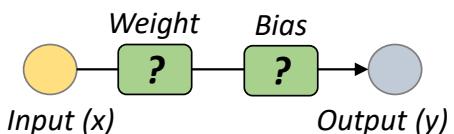
Deep Learning with Neural Network - Cost (loss)

Actual Data

$$y = ? \ x + ?$$

weight bias

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



Model Ex#1

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

weight bias

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

$$(y - \hat{y})^2$$

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

Model Ex#2

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

weight bias

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

$$(y - \hat{y})^2$$

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



Let's calculate the cost!

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

3 Deep Learning for NLP

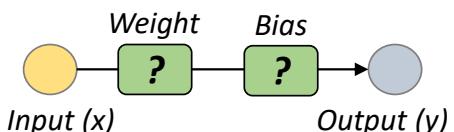
Deep Learning with Neural Network - Cost (loss)

Actual Data

weight bias

$$y = ? \ x + ?$$

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



Model Ex#1

weight bias

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

| 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

weight bias

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

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

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



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

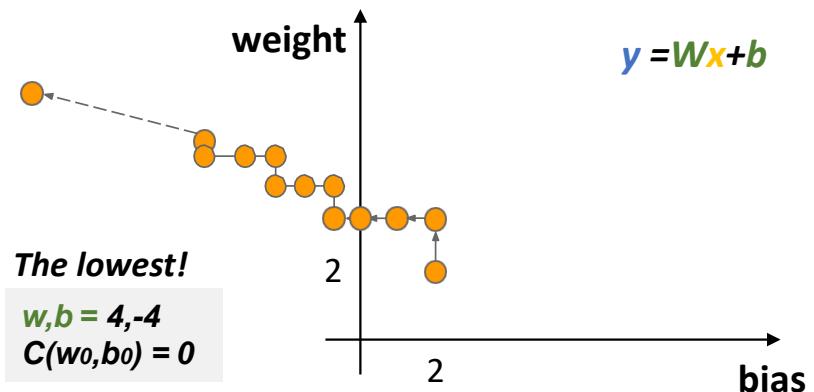
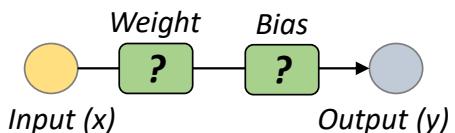
Deep Learning with Neural Network - Optimizer

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

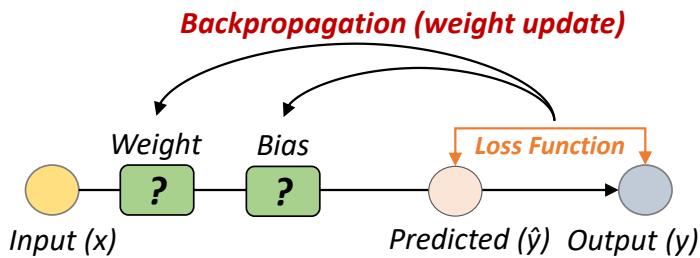
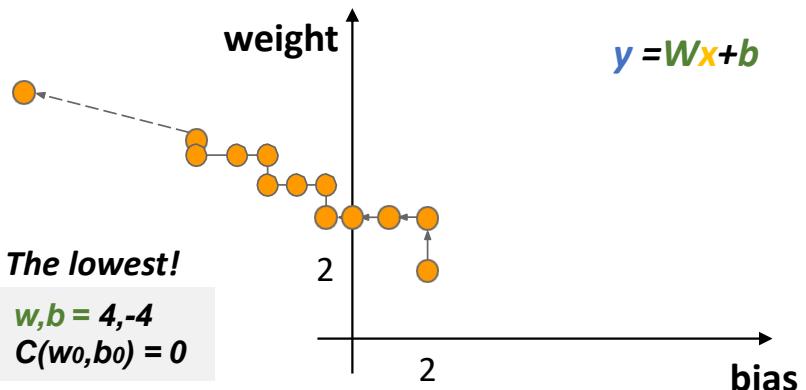
Deep Learning with Neural Network - Optimizer

Backpropagation (weight update)

$$y = 4x - 4$$

weight bias

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



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

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

3 Deep Learning for NLP

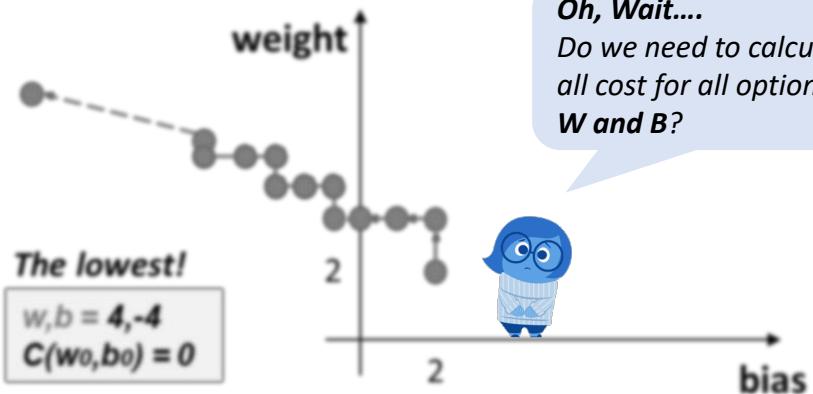
Backpropagation (weight update)

Deep Learning with Neural Network - Optimizer

$$y = 4x - 4$$

weight 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



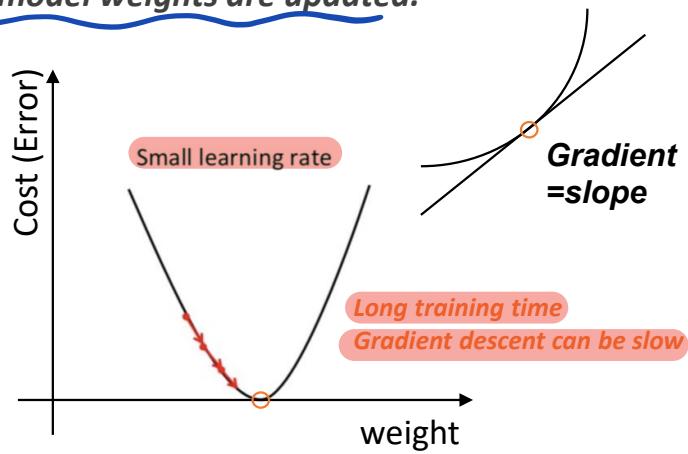
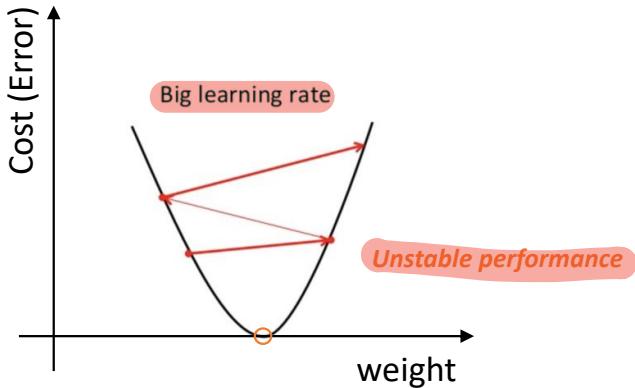
There are different types of Gradient descent optimization algorithms:

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.



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

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

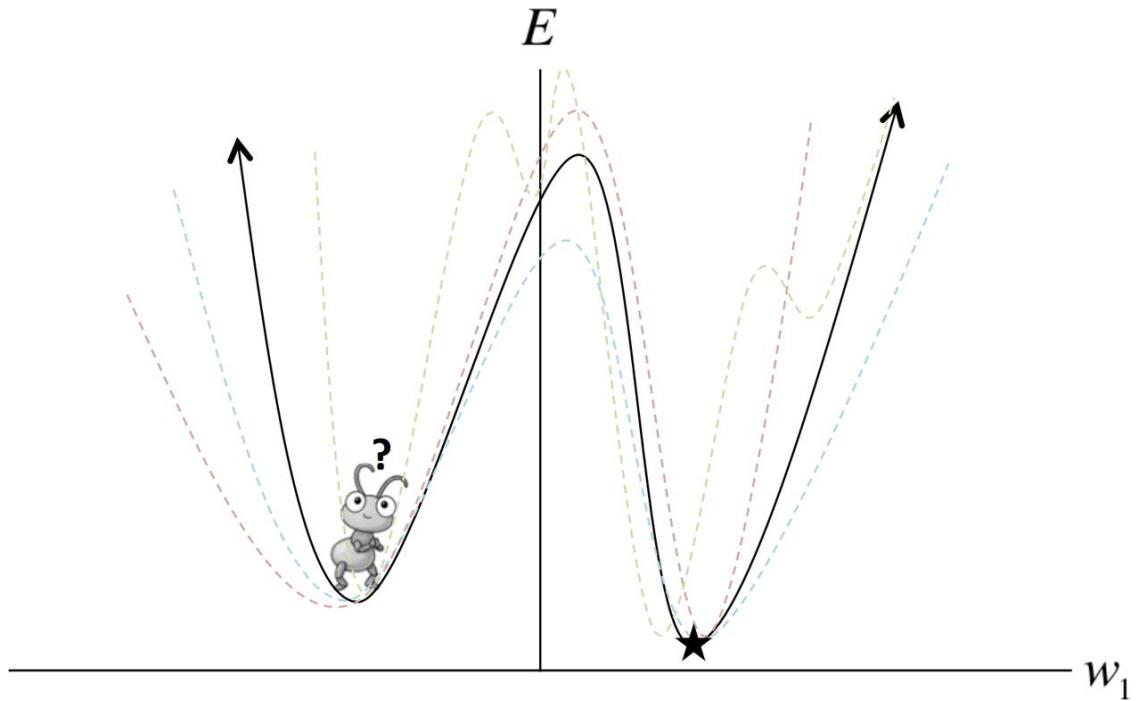
$$L = \frac{1}{N} \sum_{i=1}^N (c_i - d_i)^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial w_i} = \frac{\partial c_i}{\partial w_i} = \frac{1}{N} \sum_{i=1}^N (c_i - d_i) \cdot \frac{\partial (c_i - d_i)}{\partial w_i} \rightarrow \text{grad (current output)}$$

$$\frac{\partial c_i}{\partial w_i} = \frac{\partial (w_i * x_i + b)}{\partial w_i} \rightarrow \text{existing input}$$

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

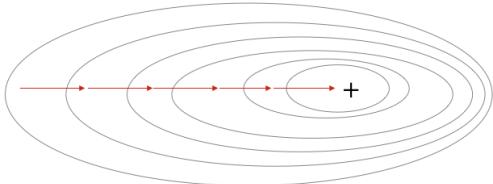
Stochastic Gradient Descent

**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!**

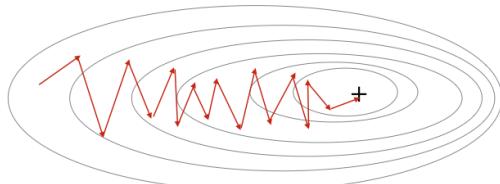
*The Solution can be used different Gradient Descent Method.
The most common – “**Stochastic Gradient Descent (SGD)**”*

all training data
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

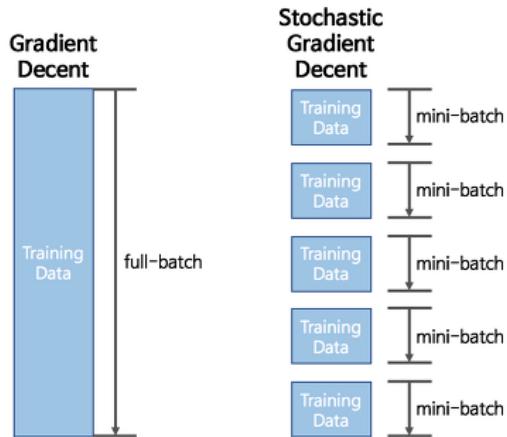


Stochastic Gradient Descent



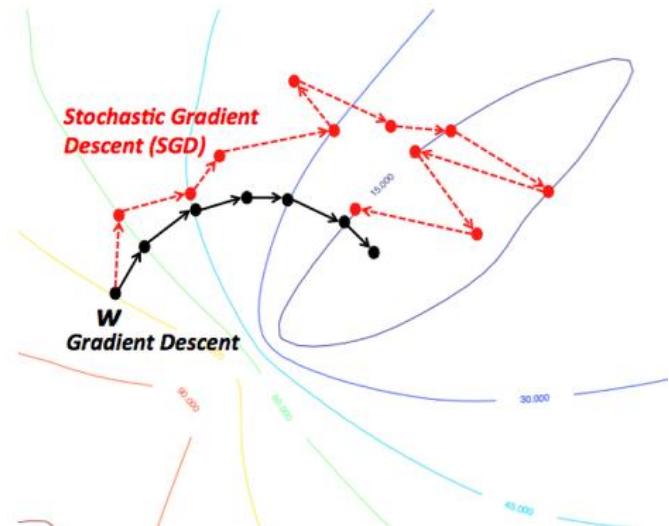
3 Deep Learning for NLP

Stochastic Gradient Descent



Gradient Descent

- Calculate with all training data
- Forward the optimal one step (at a time)
- Accurate but very slow
- e.g. 6 steps * 1 hour = 6 hours

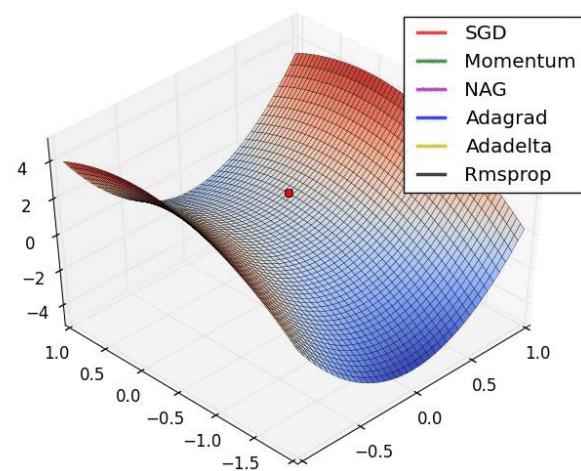
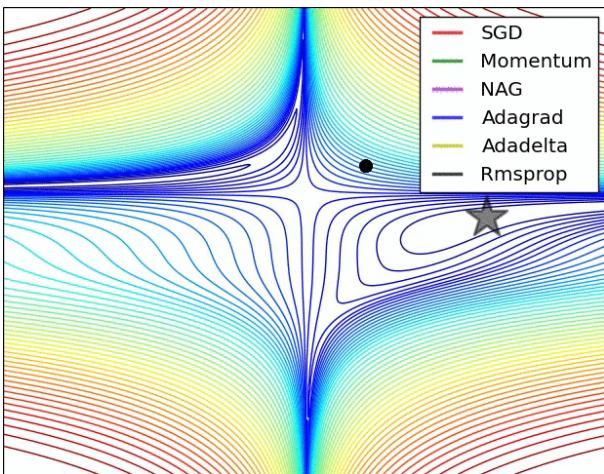


Stochastic Gradient Descent

- Calculate with mini-batch data
- Forward Faster
- A bit fluctuated but faster
- e.g. 10 steps * 5 mins = 50 mins

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

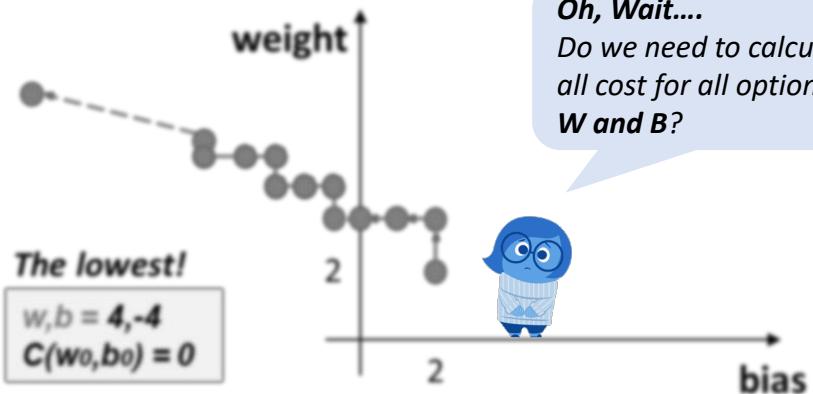
Backpropagation (weight update)

Deep Learning with Neural Network - Optimizer

$$y = 4x - 4$$

weight 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

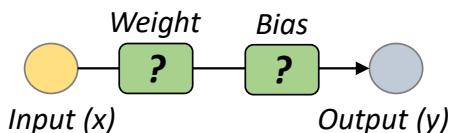
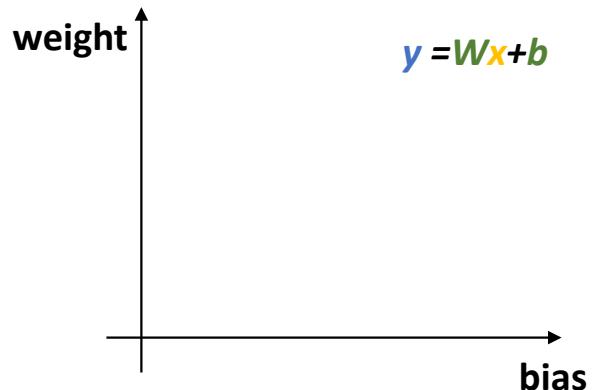
Actual Data

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

weight bias

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

Gradient!



Should be used sparingly

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

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

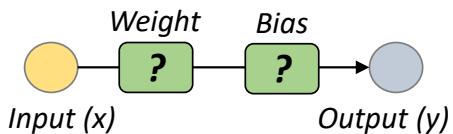
3 Deep Learning for NLP

Deep Learning with Neural Network - Optimizer

Actual Data

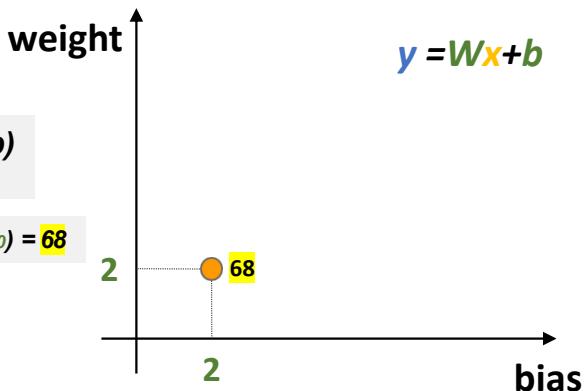
$$y = ? \times x + ?$$

| | |
|---|--|
|  | y  |
| 1 | 0 |
| 5 | 16 |
| 6 | 20 |



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

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



3 Deep Learning for NLP

Gradient!

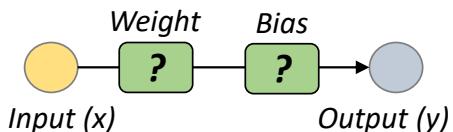
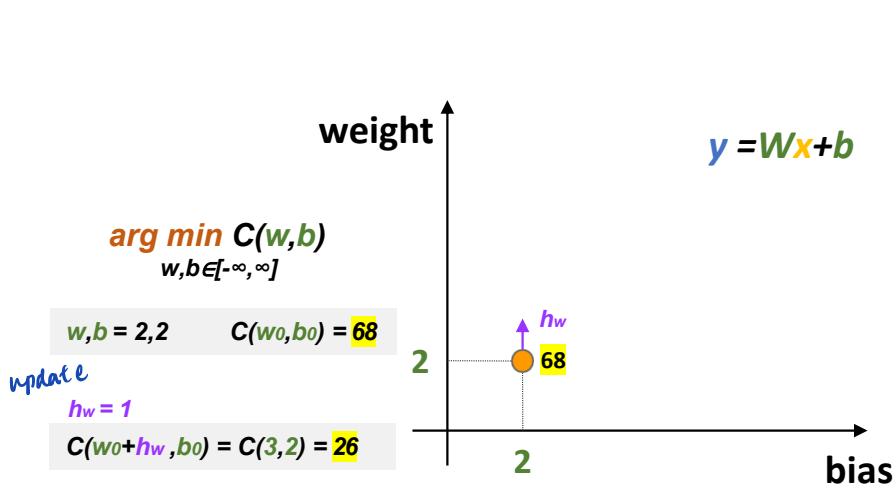
Deep Learning with Neural Network - Optimizer

Actual Data

$$y = ? \quad x + ?$$

weight bias

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



3 Deep Learning for NLP

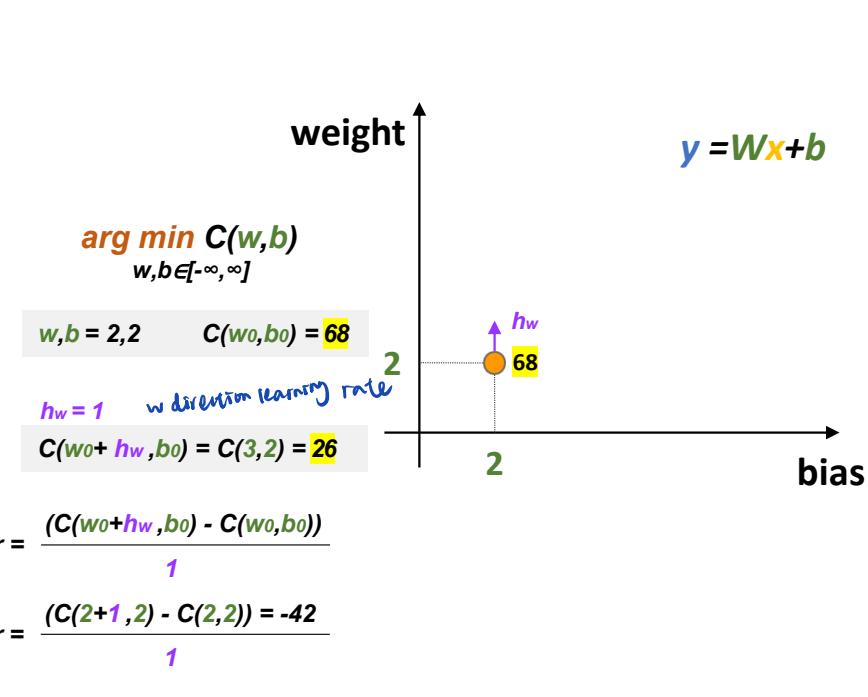
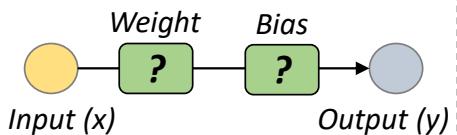
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

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

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



3 Deep Learning for NLP

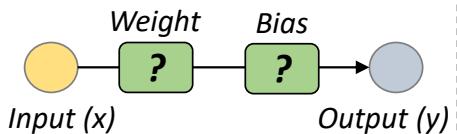
Deep Learning with Neural Network - Optimizer

Actual Data

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

weight bias

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



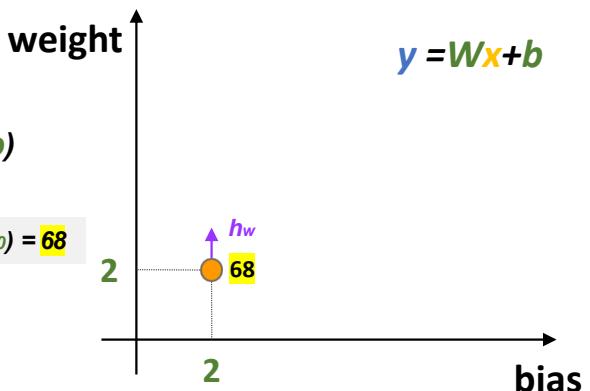
Gradient!

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

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

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

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



3 Deep Learning for NLP

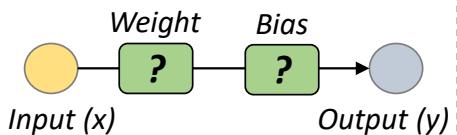
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

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

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



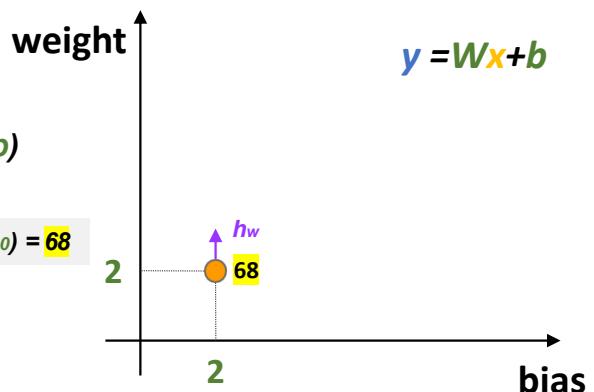
$$\arg \min_{w,b \in [-\infty, \infty]} C(w,b)$$

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

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

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

reward
partial differentiation



3 Deep Learning for NLP

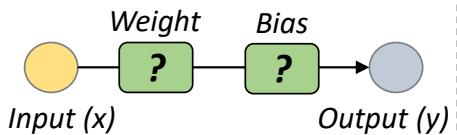
Deep Learning with Neural Network - Optimizer

Gradient!

Actual Data

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

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

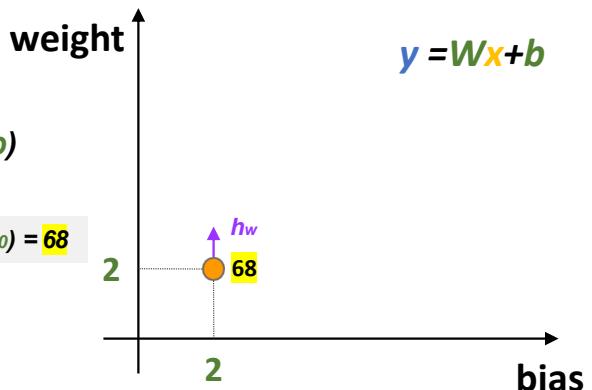


$$\arg \min_{w,b} C(w,b)$$

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

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

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



3 Deep Learning for NLP

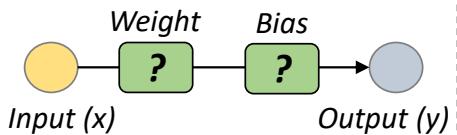
Deep Learning with Neural Network - Optimizer

Actual Data

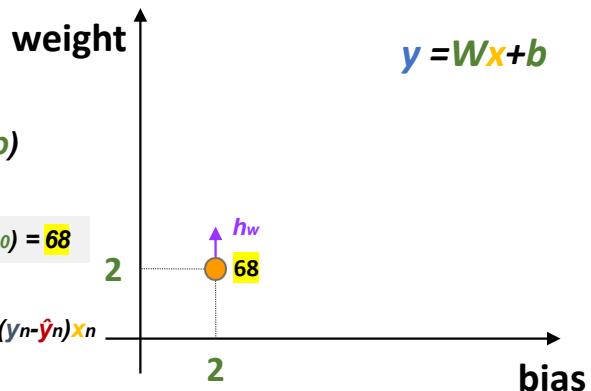
$$y = ? \times x + ?$$

weight bias

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



Gradient!



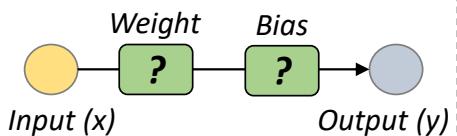
3 Deep Learning for NLP

Deep Learning with Neural Network - Optimizer

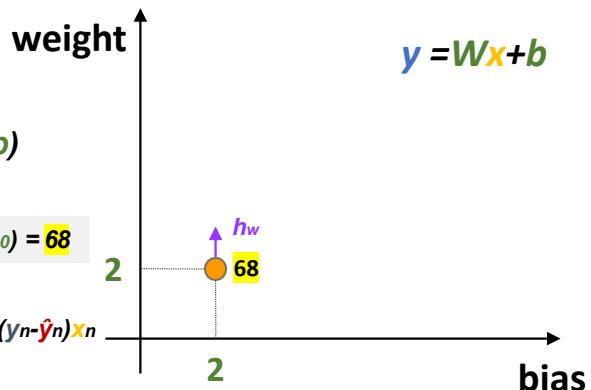
Actual Data

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

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



Gradient!



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

3 Deep Learning for NLP

Gradient!

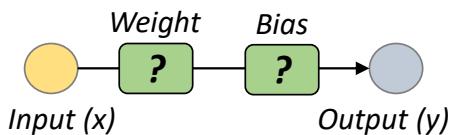
Deep Learning with Neural Network - Optimizer

Actual Data

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

weight bias

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



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

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

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

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

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

weight bias

| | 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 |

3 Deep Learning for NLP

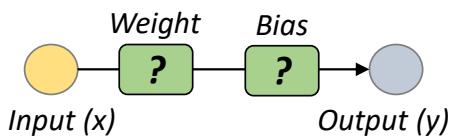
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

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

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

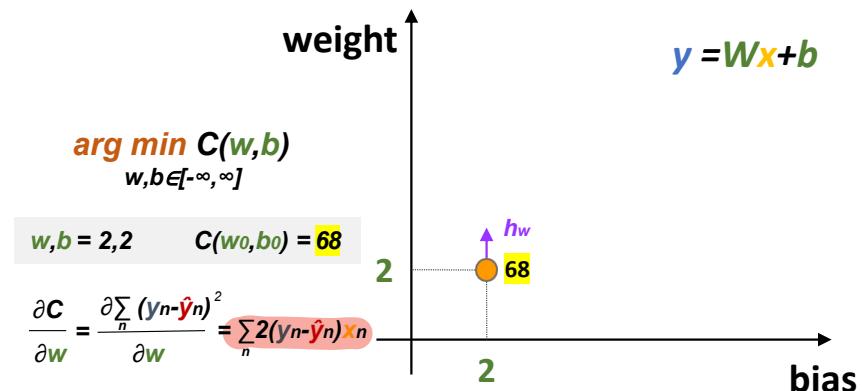


$$\arg \min_{w,b \in [-\infty, \infty]} C(w,b)$$

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

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

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



3 Deep Learning for NLP

Gradient!

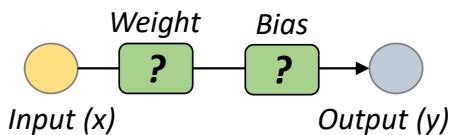
Deep Learning with Neural Network - Optimizer

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

| | predicted | actual | | |
|-----|-----------|--------|-------|--------|
| x 🍎 | ŷ 🍔 | y 🍌 | (y-ŷ) | 2(y-ŷ) |
| 1 | 4 | 0 | -4 | -8 |
| 5 | 12 | 16 | 4 | 8 |
| 6 | 14 | 20 | 6 | 12 |

3 Deep Learning for NLP

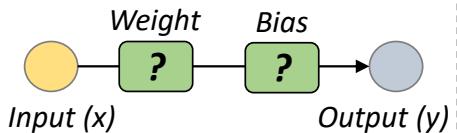
Gradient!

Deep Learning with Neural Network - Optimizer

Actual Data

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

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

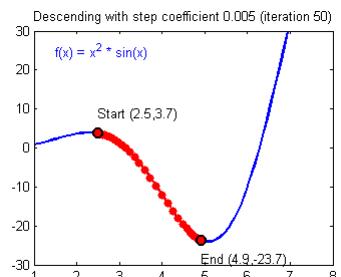
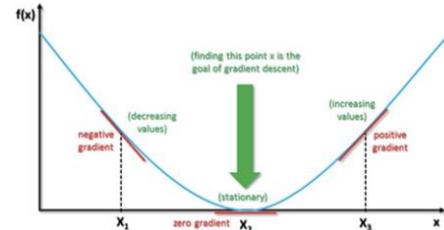
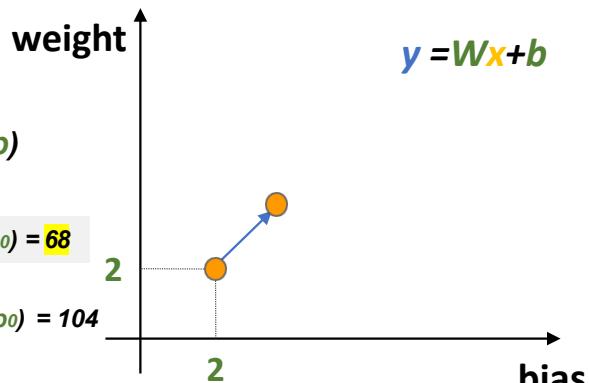


$$\arg \min_{w,b \in [-\infty, \infty]} C(w,b)$$

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

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

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



3 Deep Learning for NLP

Deep Learning with Neural Network

Data

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

Model

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

Cost

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

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

Optimizer

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

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



System

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

3 Deep Learning for NLP

Deep Learning with Neural Network

Input: $x = \text{number of apple given by Lisa}$

Output: $y = \text{number of banana received by Lisa}$

Parameters: Need to be estimated

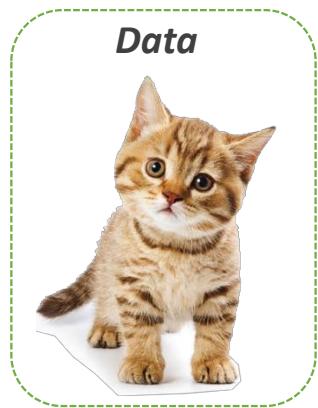
$$y = 4x - 4$$



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

$$y = \underset{\text{Vector1}}{W_1}x_1 + \underset{\text{Vector2}}{W_2}x_2 + W_3x_3 + W_4x_4 + \dots + W_nx_n + b$$



Millions of Parameters
Millions of Samples

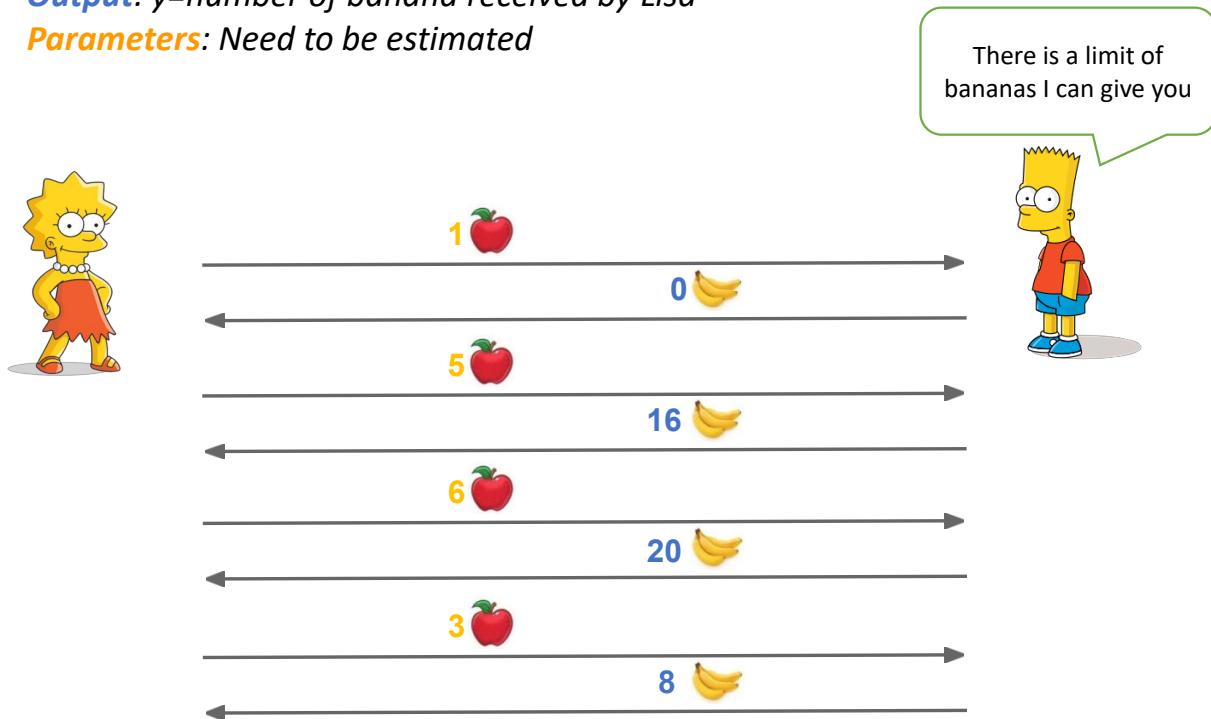
3 Deep Learning for NLP

Deep Learning with Neural Network

Input: $x = \text{number of apple given by Lisa}$

Output: $y = \text{number of banana received by Lisa}$

Parameters: Need to be estimated



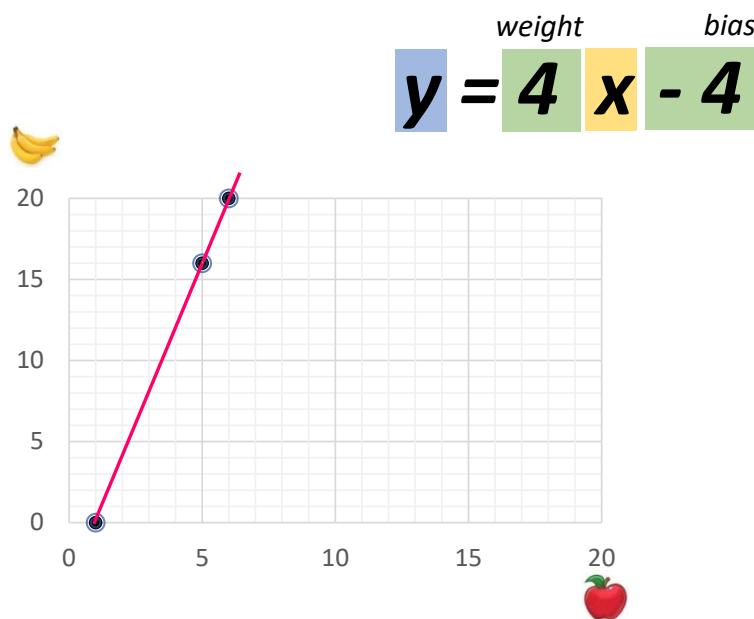
3 Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

Data

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



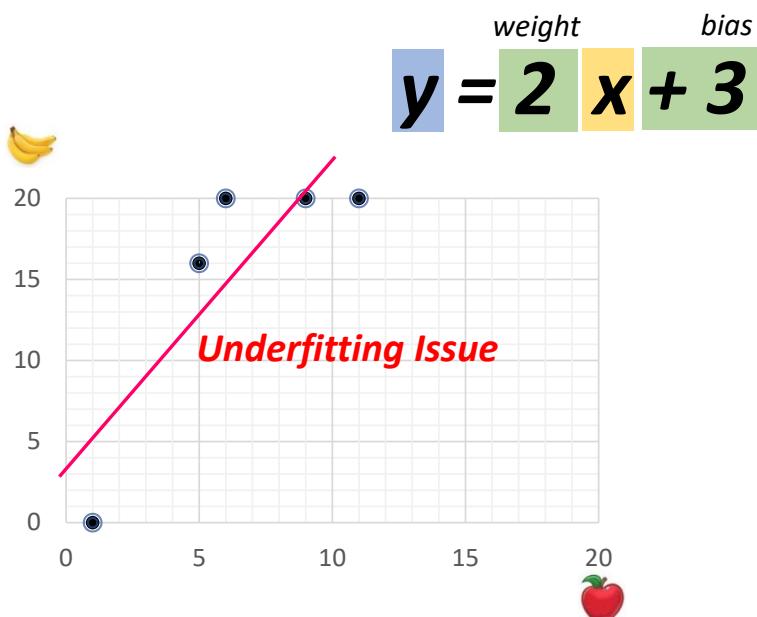
3 Deep Learning for NLP

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

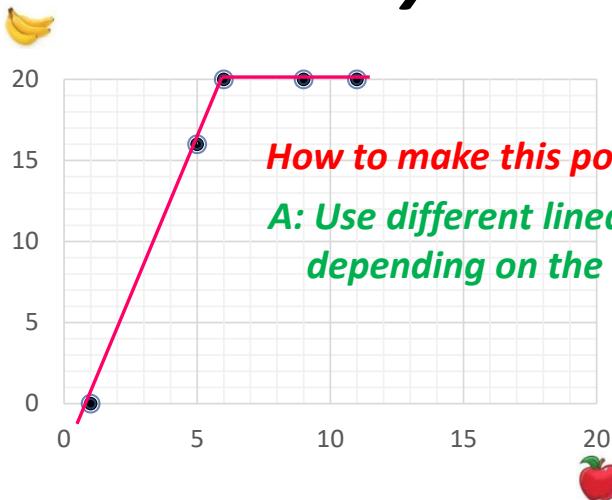
Deep Learning with Neural Network

Nonlinear Neural Network

Data

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

$$y = ???$$



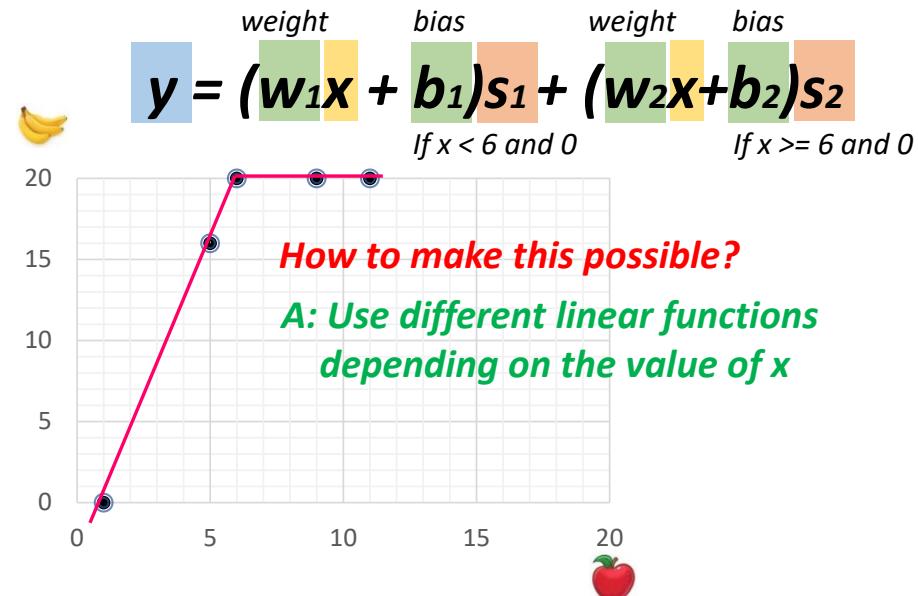
3 Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

Data

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



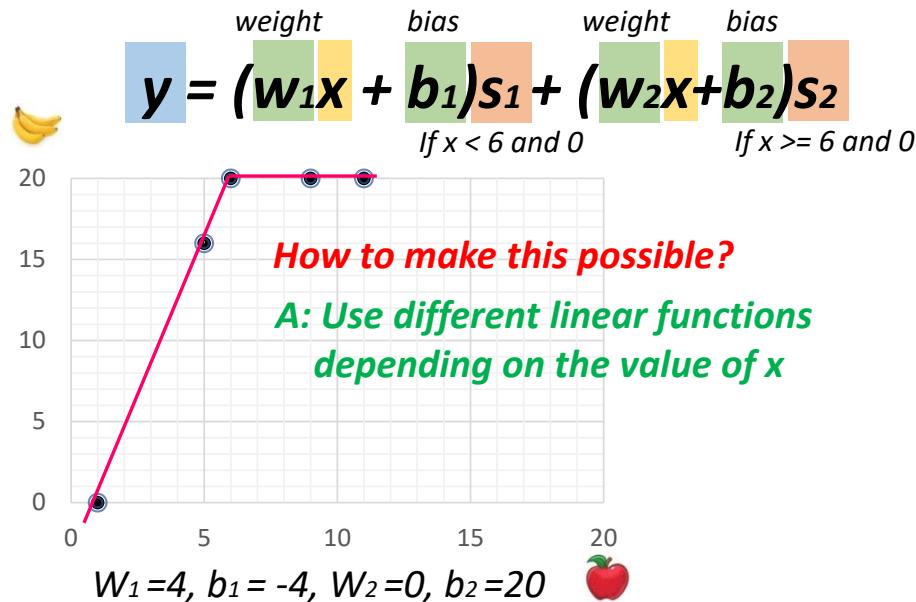
3 Deep Learning for NLP

Deep Learning with Neural Network

Nonlinear Neural Network

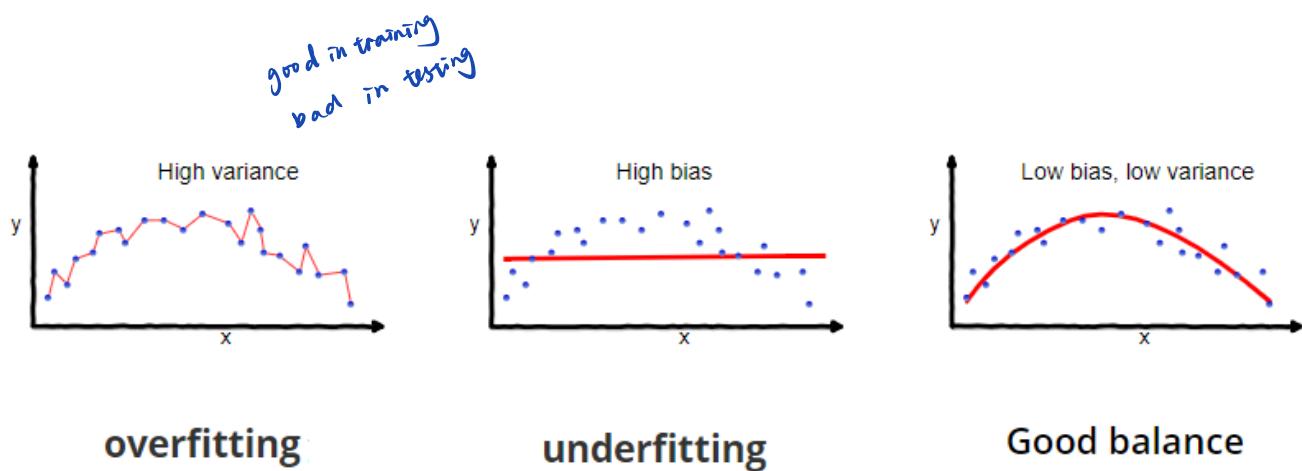
Data

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

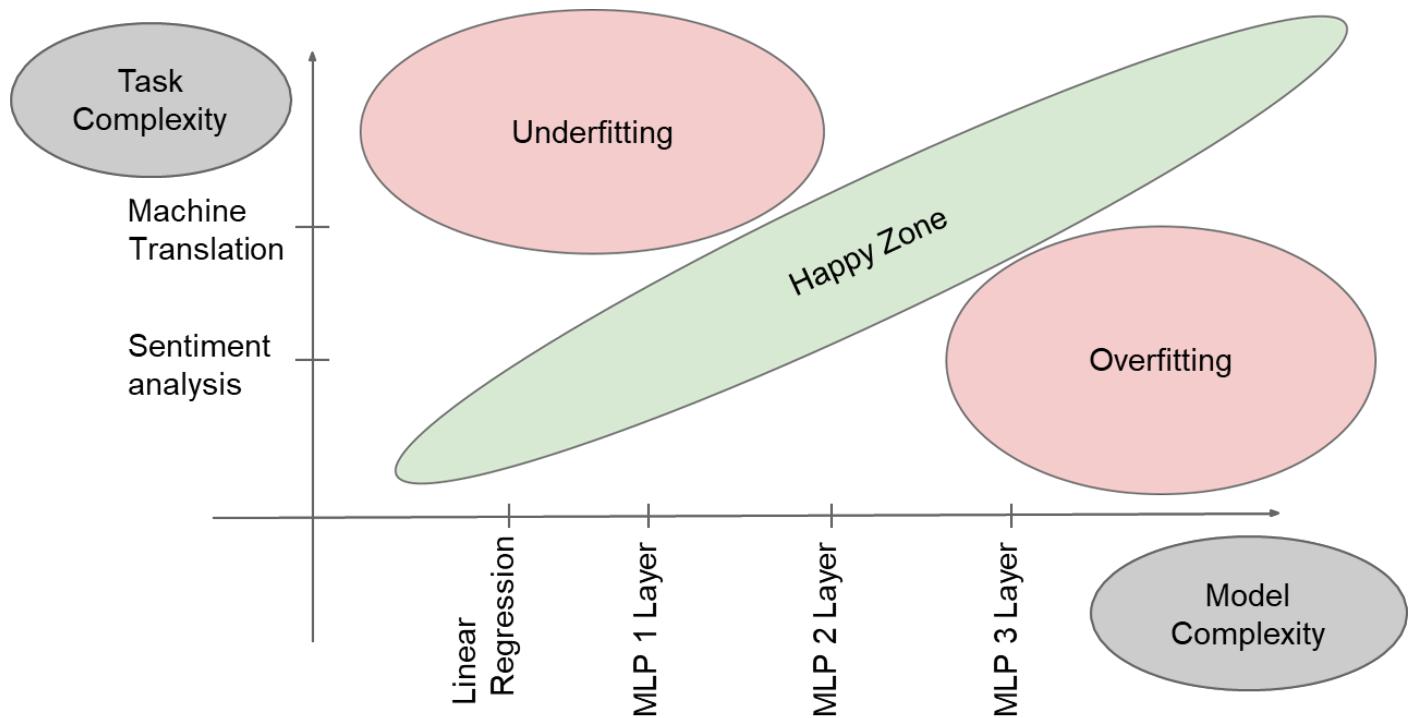


3 Deep Learning for NLP

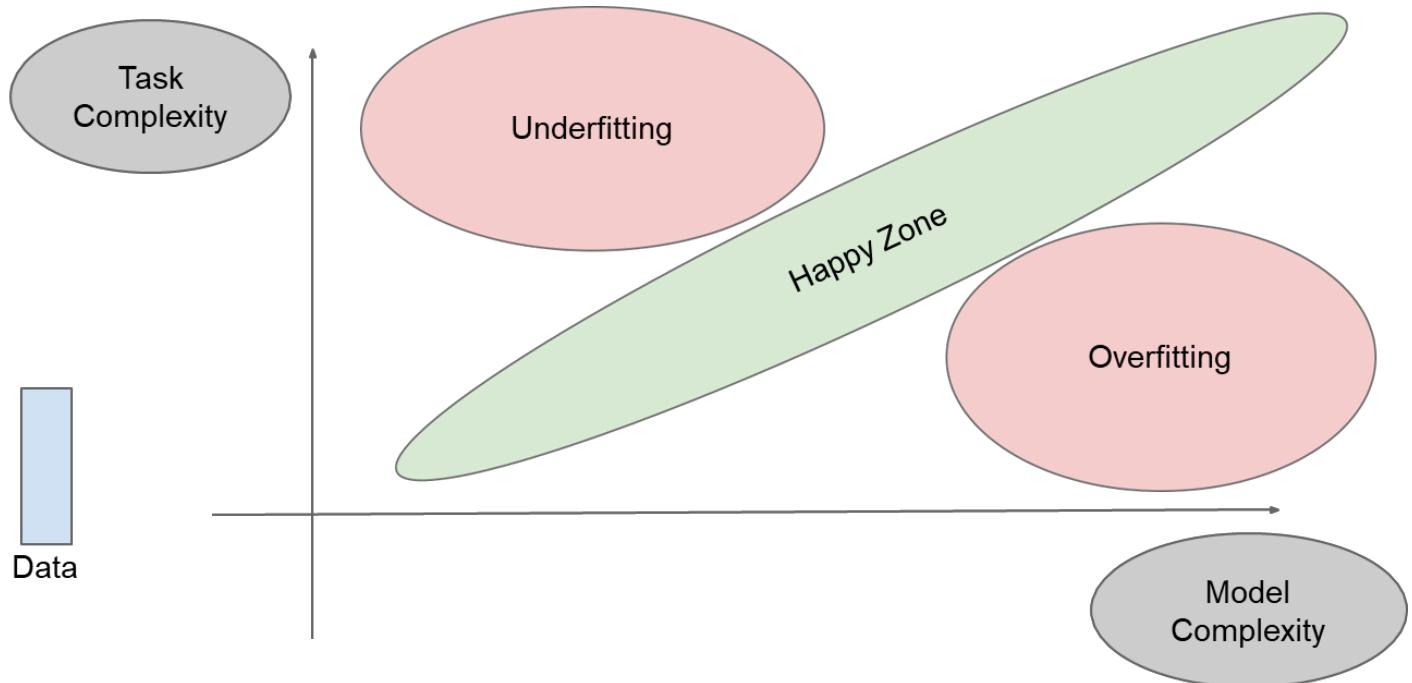
Deep Learning with Neural Network



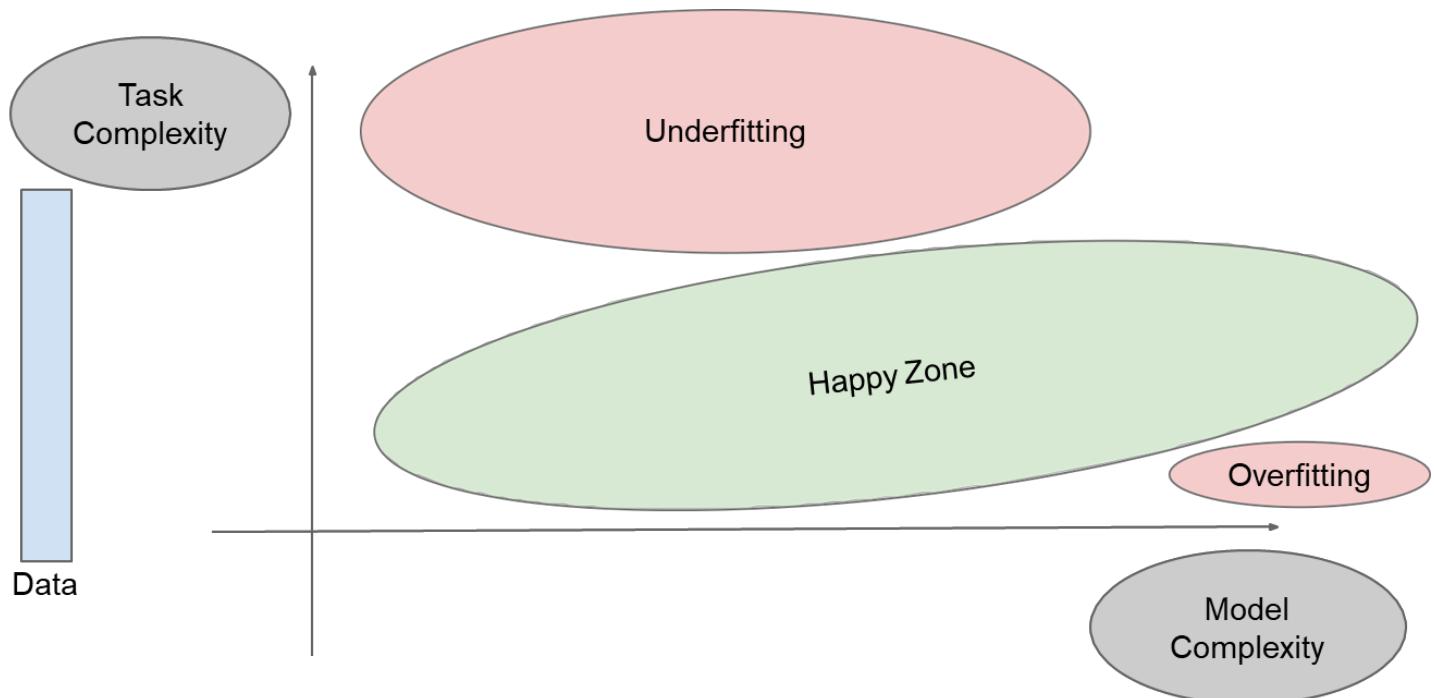
3 Deep Learning for NLP



3 Deep Learning for NLP

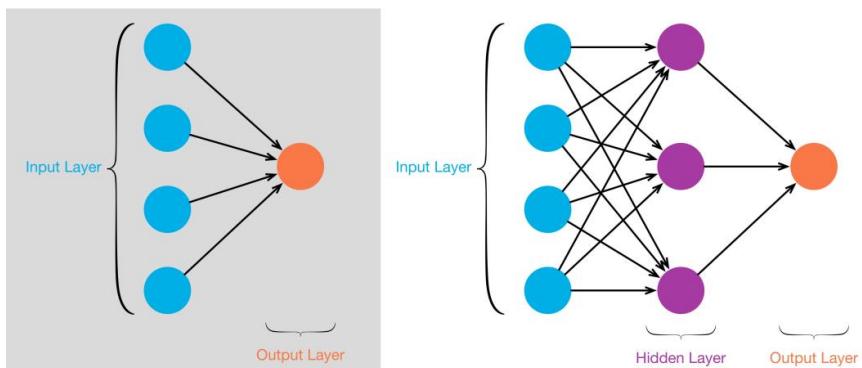
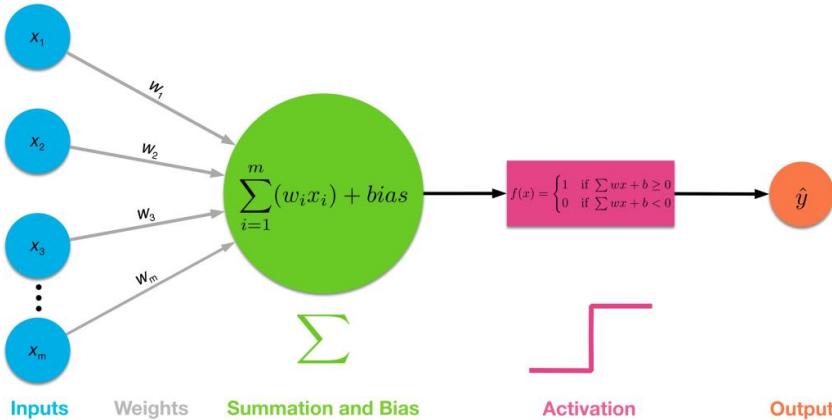


3 Deep Learning for NLP



3 Deep Learning for NLP

Single Neuron VS Multilayer

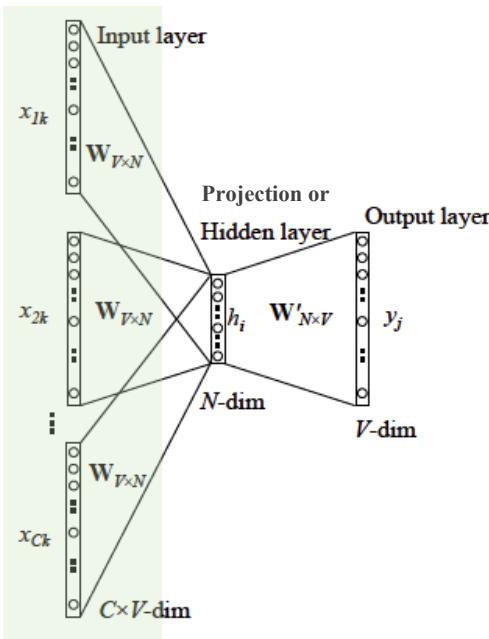


3 Deep Learning for NLP

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

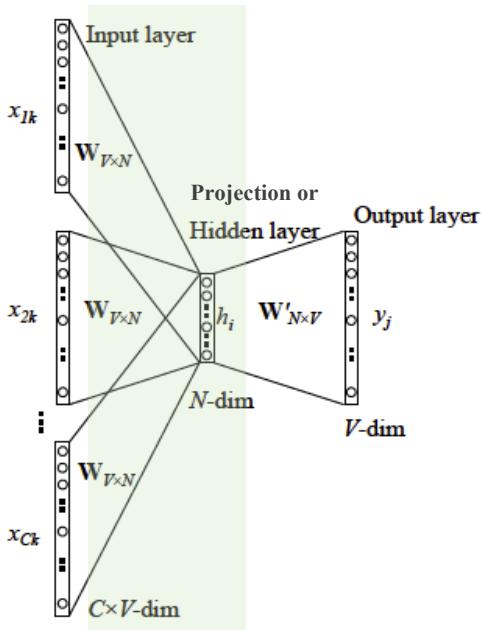
2) Parameter Matrix (to Output Layer)
 $W' \in \mathbb{R}^{N \times V}$

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.

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

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

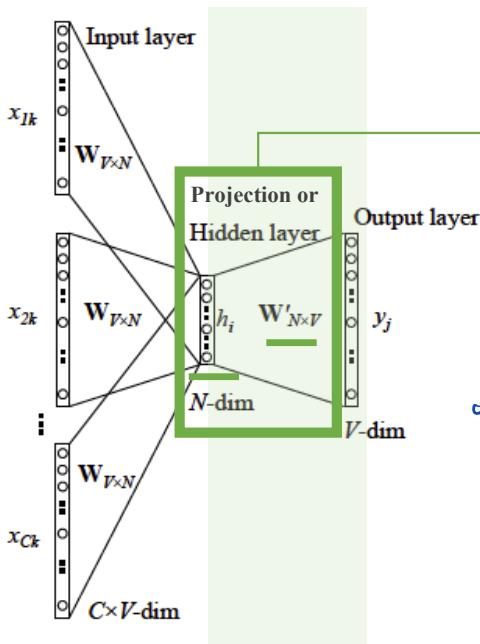
$$\hat{v} = \frac{\mathbf{v}_{c-m} + \mathbf{v}_{c-m+1} + \dots + \mathbf{v}_{c+m}}{2m}$$

3 Deep Learning for NLP

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



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

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

8. Train the parameter matrix using objective function.
 cross entropy

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

* Focus on minimising the value

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

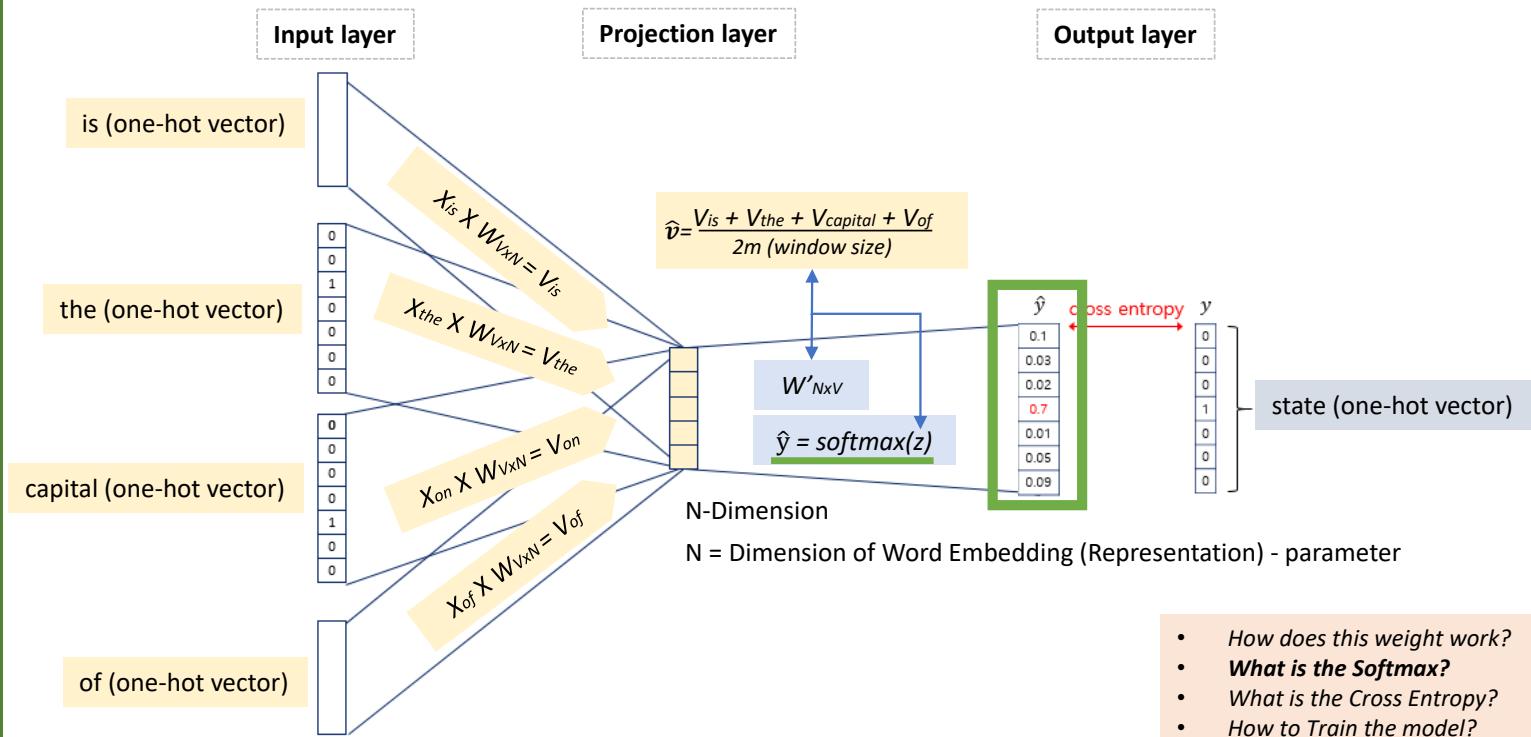
$$H(\hat{y}, y) = -y_j \log(\hat{y}_j)$$

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”

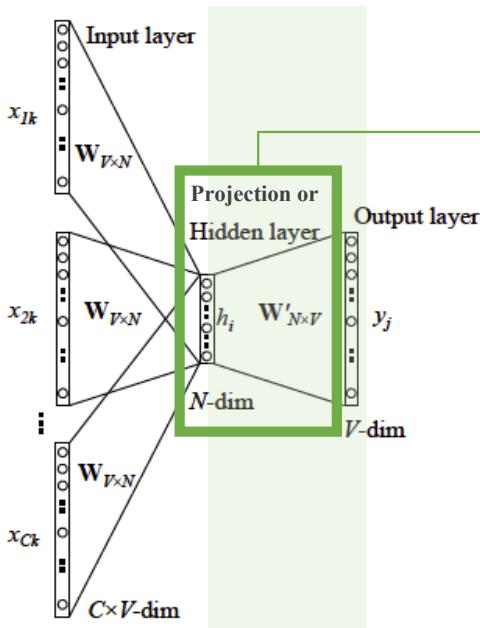


3 Deep Learning for NLP

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



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

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:

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

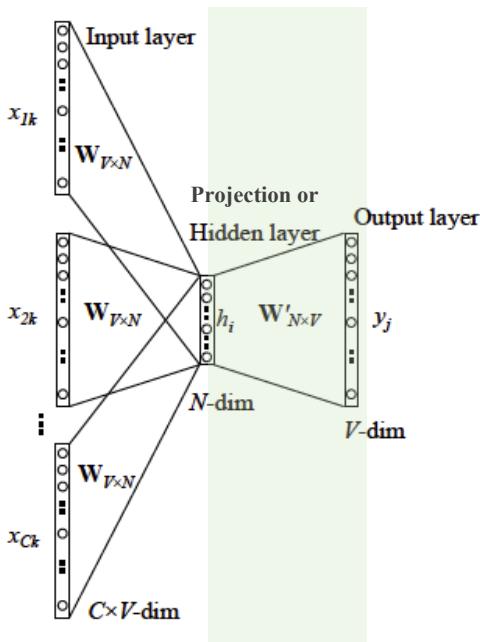
- Exponentiate to make positive
- Dividing by $\sum_{j=1}^{|V|} e^{\hat{y}_j}$ normalizes the vector ($\sum_{j=1}^{|V|} \hat{y}_j = 1$) to give probability

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)



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

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

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.

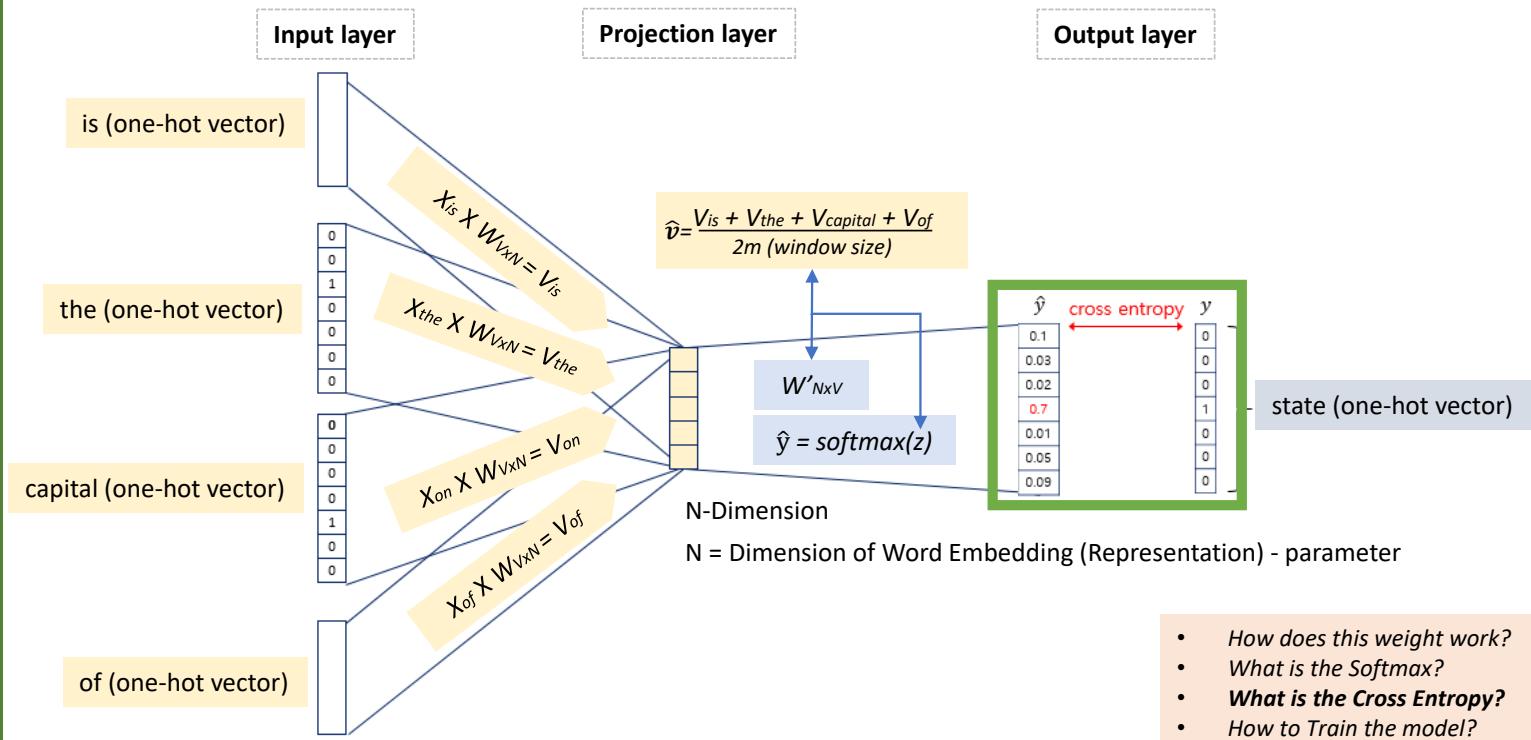
$$H(\hat{y}, y) = -y_j \log(\hat{y}_j)$$

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”

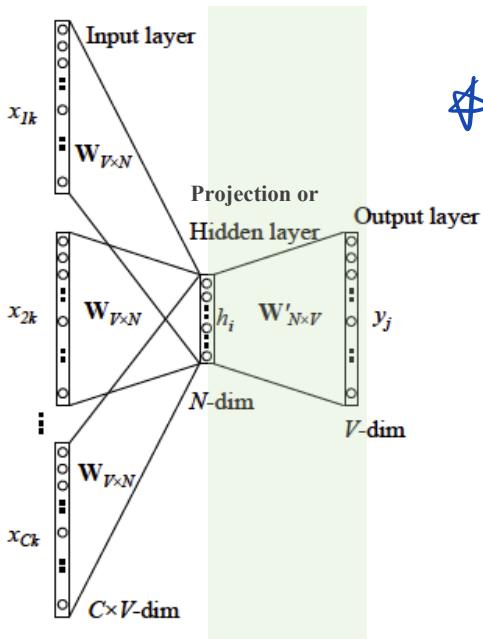


3 Deep Learning for NLP

CBOW – Neural Network Architecture (ReCAP with Optimizer)

Predict center word from (bag of) context words.

Summary of CBOW Training (Review your understanding with equations)



8-1. Optimization Objective Function can be presented:

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

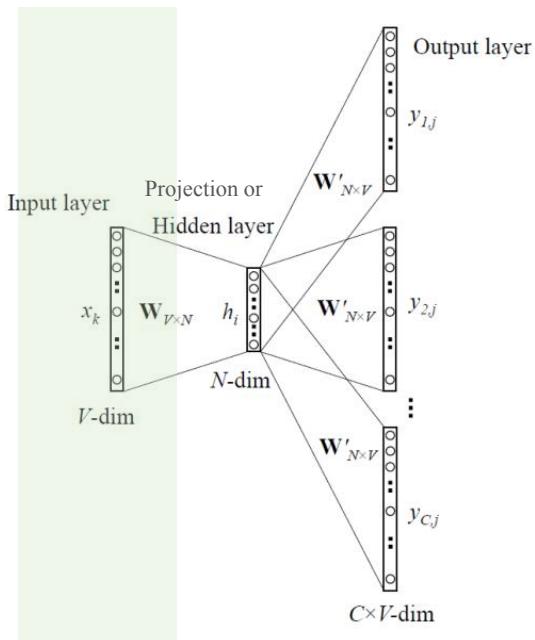
u_i =the output vector representation of word w_i

3 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)



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

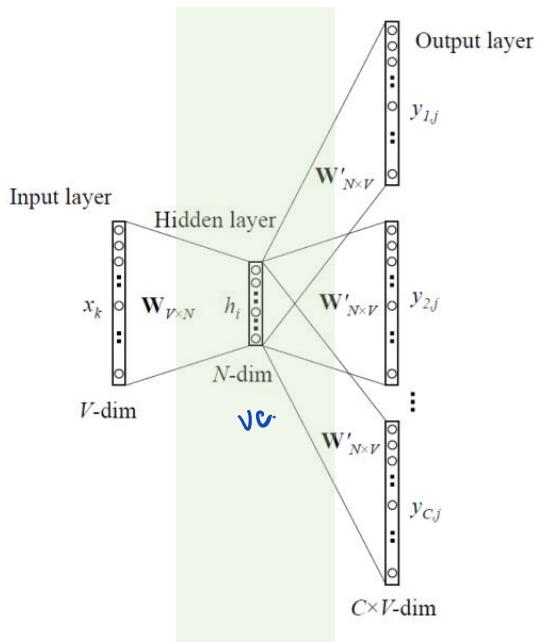
2) Parameter Matrix (to Output Layer)
 $\mathbf{W}' \in \mathbb{R}^{N \times V}$

3 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)



3. 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.

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

4. Calculate the score value for the output layer by multiplying the parameter matrix W'

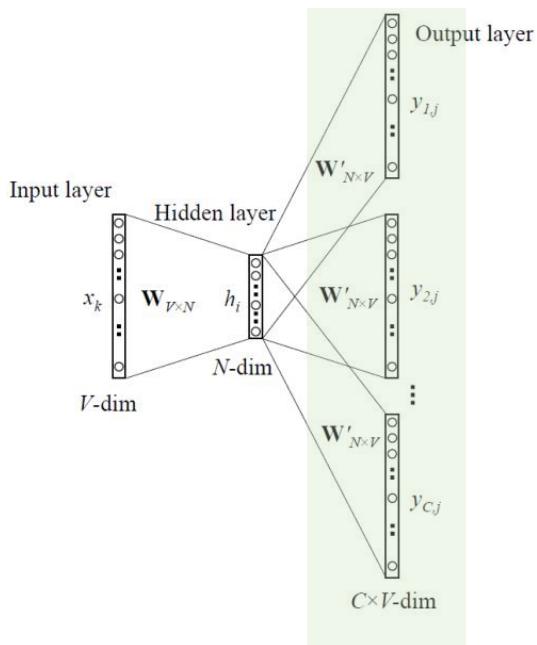
$$\mathbf{z} = \mathbf{W}'_{\mathbf{v}_c} \quad \mathbf{w}' \mathbf{v}_c$$

3 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.

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

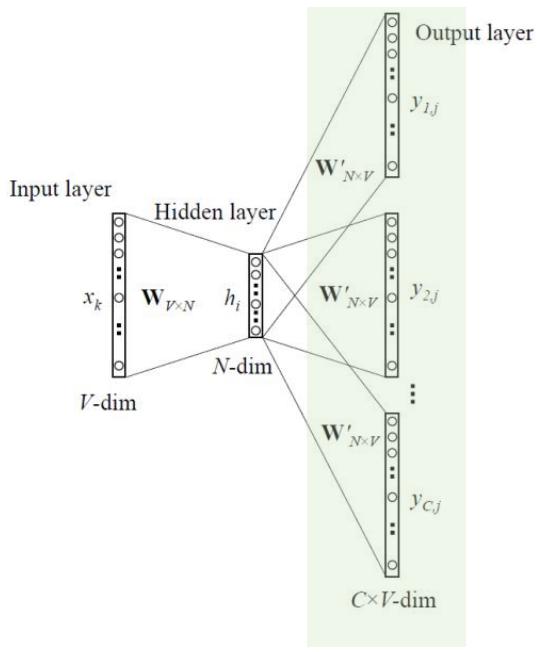
and compare with the ground truth (one-hot vector)
 $y^{(c-m)}, \dots, y^{(c-1)}, y^{(c+1)}, \dots, y^{(c+m)}$

3 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.

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

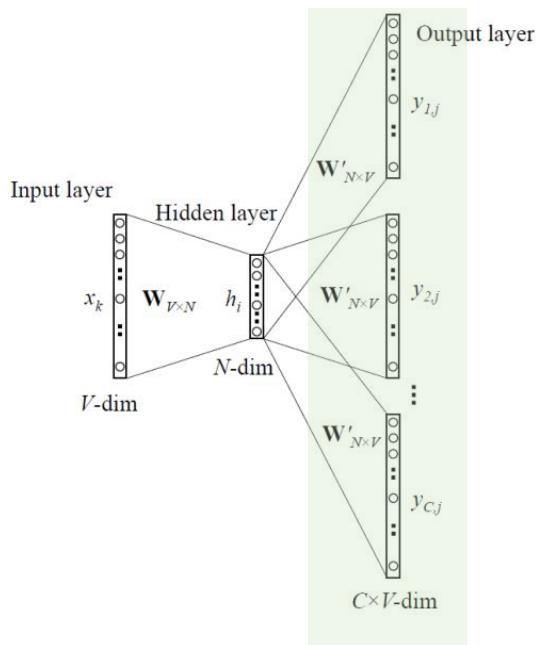
$u_i = \text{the output vector representation of word } w_i$

3 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

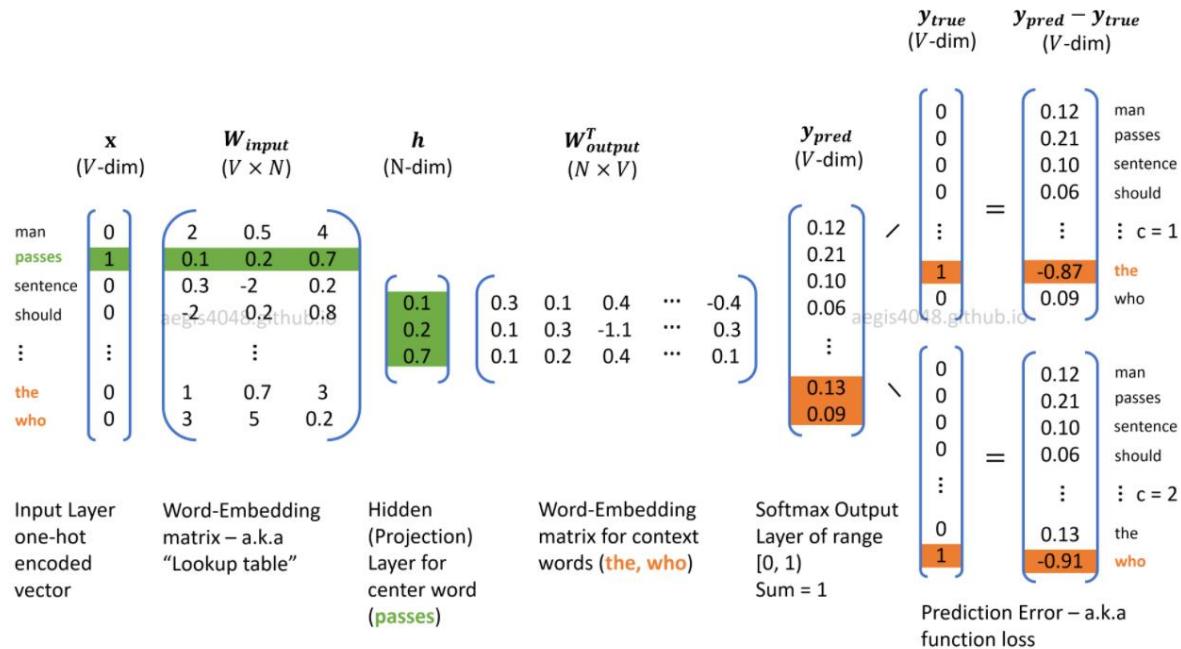


$$\begin{aligned}
 J &= - \sum_{j=0, j \neq m}^{2m} \log P(u_{c-m+j} | v_c) \\
 &= \sum_{j=0, j \neq m}^{2m} H(\hat{y}, y_{c-m+j})
 \end{aligned}$$

4 Deep Learning for NLP

Word2Vec-SkipGram Overview

With a simple diagram



4 Deep Learning for NLP

Key Parameter (2) for Training methods: Negative Samples

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

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

*Each che
absolute negative one*

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 |
| eat | building | 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

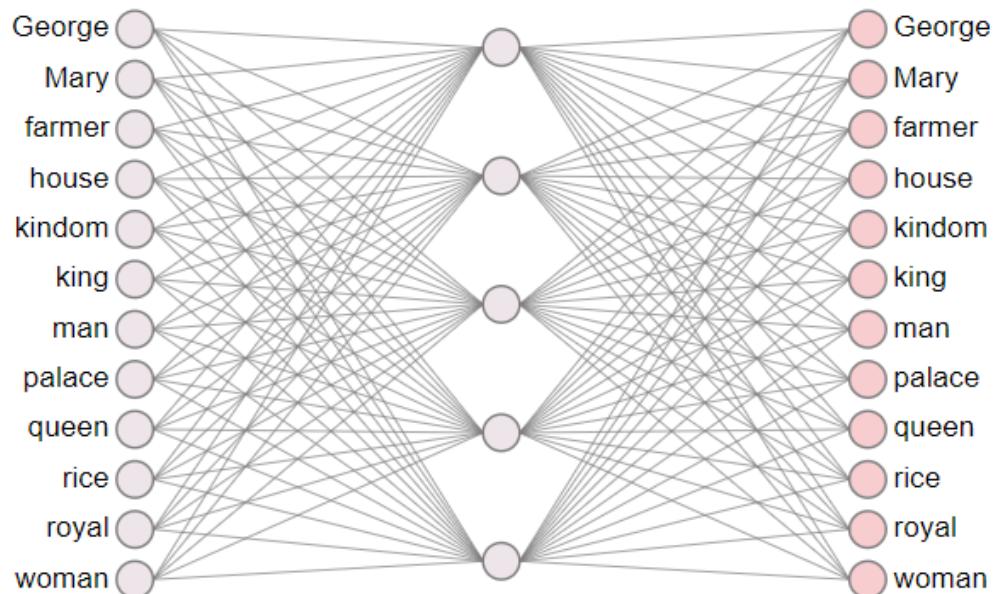
With a simple diagram



4 Deep Learning for NLP

Application

Application #1: Embedding Pretraining



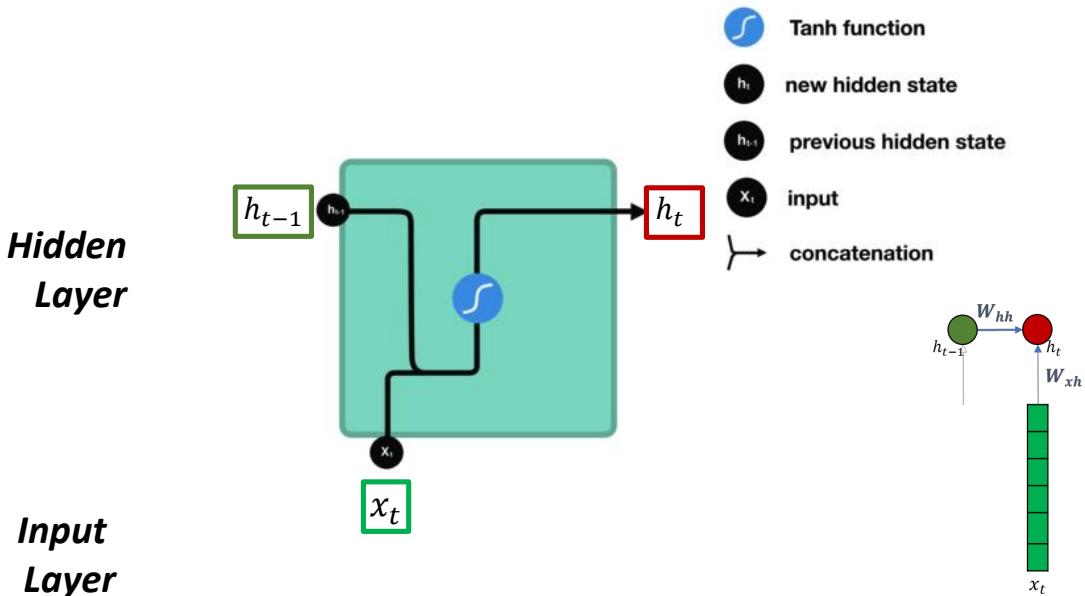
0 LECTURE PLAN

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.



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

New hidden state Previous state input
 A function with parameters W

/ 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