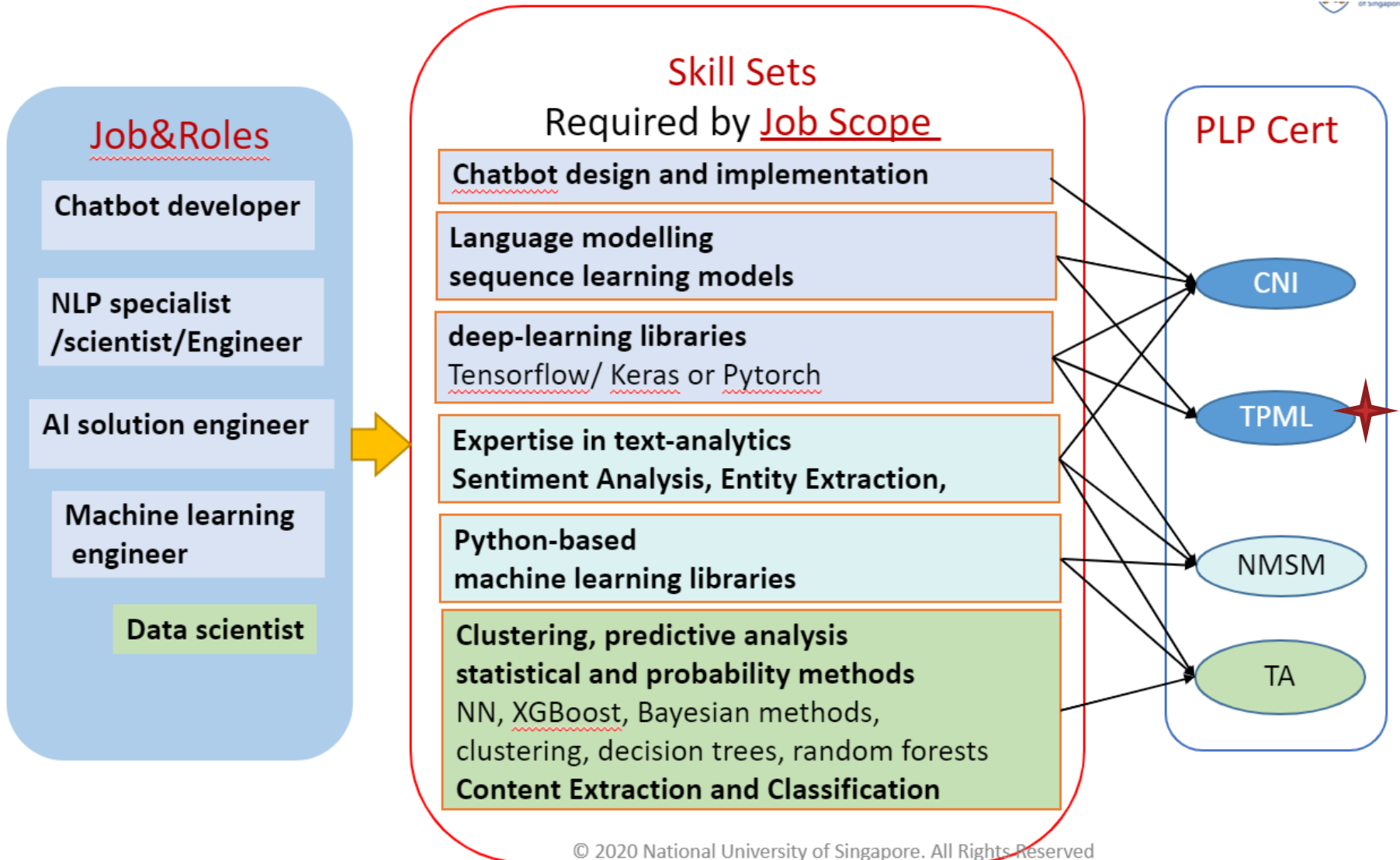


# Text Processing Using Machine Learning

## Neural Nets for Text Processing

Dr Wang Aobo  
aobo.wang@nus.edu.sg

# PLP Cert



# Learning Outcomes

- **Recent** techniques for NLP
  - Deep learning
- Getting a basic grasp of deeplearning library for NLP
  - Knowledge can be applied to any tensor libraries
- Understanding the **underlying** techniques and knowing how to implement them.

# Agenda

- **Overview of NLP**
  - Applications
  - Statistical Modelling vs. Deep Neural Nets
- Deep Learning Basics for NLP
  - MLP
  - Deep Learning Training Routine (quiz)
    - Workshop: Basic NN on Colab
- Word2Vec & DL Specific
  - CBOW and SkipGram
  - ActivationFunction/LossFunction
  - Optimiser/Learning Rate
    - Workshop: Word2Vec from Scratch

# What can NLP do?

## 2 basic Use Cases:

### 1. Automatically put text into categories- **Classification**

- Sentiment detection
- Spam email detection
- Emotion detection

### 2. Extract specific information from the text- **Extraction**

- Named Entity extraction from sentence



# What can NLP do?

Other *Fancier* Use Cases:

1. **Object Classification/Clustering & Recommendation**
2. **Search Engine**
3. **Question Answering System**
4. **Voice Assistant**
5. **Machine Translation**
6. **Grammar Error Correction & Language Learning**
7. **Chatting Robots**
8. ***“Fake Articles”* Generation & Detection**
9. ....

# How does NLP Work

- The whole task here is...

## Documents

Lost glamor  
R  
o  
High tea at Raffles!  
R  
o  
2 Not what it was, but still  
a  
w  
R  
o  
ir  
th  
Amazing service  
Rated 5 by travel-gini  
on Feb 26, 2013  
Great location with a  
little bit of history, the  
staff make this hotel  
though

*smart*



**Doc/Feature Matrix**

	amazing	service	lost	glamour	...
Doc1	1.5	2.1	0	0	
Doc2	0	0	3.1	1.5	
Doc3	0	0	0	1.9	
Doc4	0	0	0	0	
...					



ML + CPU  
Supervised  
Unsupervised

# Classic NLP vs. Deep learning

- **Frequency TF-IDF and PPMI vectors are**
  - long ( $|V| > 100,000$ )
  - sparse (lots of zero)
  - efficient for simple tasks with reasonably large of dataset
  - interpretable as designed by human intelligence
  - **difficult to capture contextual dependency**

	he	drink	hold	...	<i>tfidf</i> drink apple	<i>tfidf</i> hold apple	<i>tfidf</i> apple juice	...	<i>PPMI</i> drink apple	<i>PPMI</i> hold apple	<i>PPMI</i> apple juice
<b>sent0</b>	0.01	0.38	0.00	...	0.87	0.00	0.92	...	4.23	0.00	8.90
<b>sent1</b>	0.01	0.00	0.28	...	0.00	0.87	0.00	...	0.00	2.45	0.00



# Statistical modelling

- **All models are wrong, but some are useful**
  - for Document/Sentence Classification
    - SVM/KNN/Decision Tree/RandomForest/Naïve Bayes/MaxEnt
  - for Sequence Labeling
    - HMM/CRF
  - for Language modeling
    - N-gram/RandomForest/MaxEnt
  - for Machine Translation
    - Phrase/Tree-based Model + beam search

# How does NLP Work

- The whole task here is...

## Documents

Lost glamor  
R  
o  
2  
v  
ir  
th

High tea at Raffles!  
R  
o  
2  
v  
ir  
th

Not what it was, but still  
a  
v  
ir  
th

Amazing service  
Rated 5 by travel-gini  
on Feb 26, 2013

Great location with a  
little bit of history, the  
staff make this hotel  
though

*smart*



Doc/Feature Matrix

	amazing	service	lost	glamour	...
Doc1	1.5	2.1	0	0	
Doc2	0	0	3.1	1.5	
Doc3	0	0	0	1.9	
Doc4	0	0	0	0	
...					



ML + CPU  
Supervised  
Unsupervised

- Able to generate this Matrix **automatically** ?
- Able to generate this Matrix **task-independently** ?
- Yes! Learn to generate from data

# How does NLP Work

- The whole task here is...

## Documents

Lost glamor  
 R High tea at Raffles!  
 o R 5 by travel-gini  
 2 Not what it was, but still  
 a  
 v R  
 o  
 ir  
 th  
 Amazing service  
 Rated 5 by travel-gini  
 on Feb 26, 2013  
 Great location with a  
 little bit of history, the  
 staff make this hotel  
 though

*smart*



## Doc/Feature Matrix

	amazing	service	lost	glamour	...
Doc1	1.5	2.1	0	0	
Doc2	0	0	3.1	1.5	
Doc3	0	0	0	1.9	
Doc4	0	0	0	0	
...					



ML + CPU  
 Supervised  
 Unsupervised

## Word Indexing Matrix

	amazing	service	lost	glamour	...
amazing	1	0	0	0	
service	0	1	0	0	
lost	0	0	1	0	
glamour	0	0	0	1	
...					

*naive*



DNN + GPU<sup>s</sup>  
 Supervised

=

	?	?	?	?	?
Doc1	11.5	2.1	5.70	-30.2	
Doc2	-3.40	0.34	3.1	1.5	
Doc3	5.8	0.560	5.9	1.9	
...					

# Classic NLP vs. Deep learning

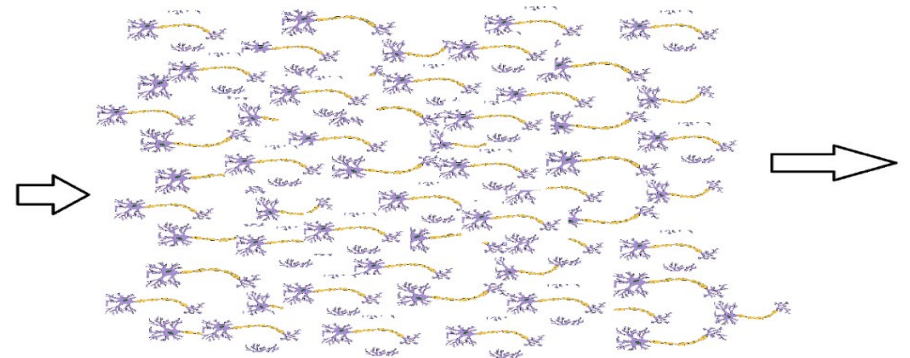
- Deep learning can **create** feature vectors that are
  - short (often fixed-sized <2000, decided empirically)
  - dense (most are non-zeros)
  - **non-Interpretable** as decided empirically without human intelligence
  - able to capture **contextual dependency**
  - beneficial to **all** tasks (classification/sequence labeling/Translation/QA)

# Agenda

- Overview of NLP
- **Deep Learning Basics for NLP**
  - MLP
  - Deep Learning Training Routine (quiz)
- Workshop: Deep Learning from Scratch
- Word2Vec & DL Specific
  - CBOW and SkipGram
  - ActivationFunction/LossFunction
  - Optimiser/Learning Rate
  - Workshop: Word2Vec from Scratch

# Deep Learning Basics

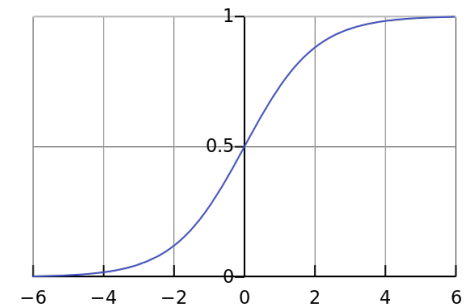
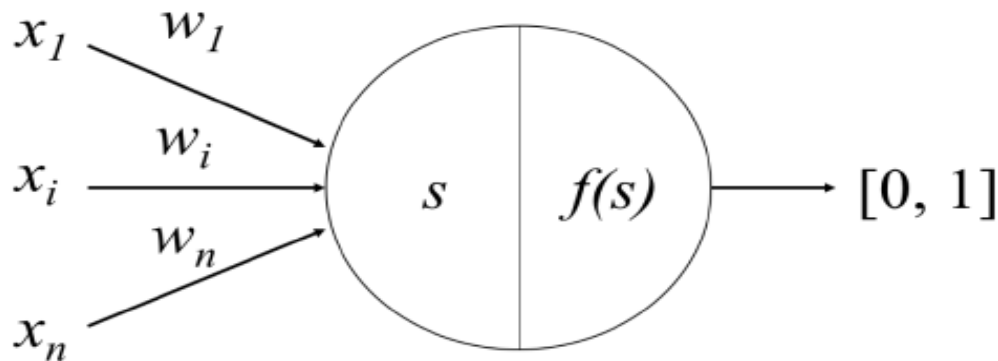
- Perceptron
  - "system that depends on **probabilistic** rather than deterministic principles for its operation, gains its reliability from the properties of statistical measurements obtain from a large population of elements"
    - Frank Rosenblatt (1957)
  - 100 billion **perceptron** in our brain
  - BERT ~ 110 million to 17 billion **para**
  - GPT ~ 1.5 billion to 175 billion **para**



(\*Image from [Akshay Chandra Lagandula's blog](#))

# Recap Perceptron

- Given a **set of inputs**  $\mathbf{x}$ , perceptron
  - learns  **$\mathbf{w}$  vector** to map the inputs to a real-value output between  $[0,1]$
  - through the summation of the dot product of the  $\mathbf{w} \cdot \mathbf{x}$
  - with a transformation function (aka. **activation function**)



**Summation**

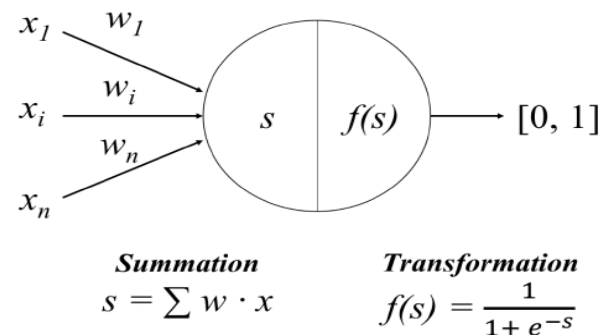
$$s = \sum w \cdot x$$

**Transformation**

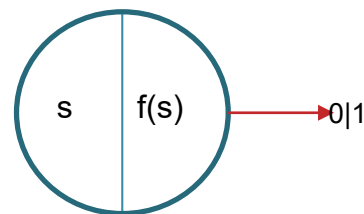
$$f(s) = \frac{1}{1 + e^{-s}}$$

# Perceptron

- Word level classification
  - Positive=1;Negative=0
  - let  $n=5$



	Like x1	Hate x2	Good x3	Enjoy x4	Bad x5
Like	1	0	0	0	0
Hate	0	1	0	0	0
Good	0	0	1	0	0
Enjoy	0	0	0	1	0
Bad	0	0	0	0	1





# Perceptron

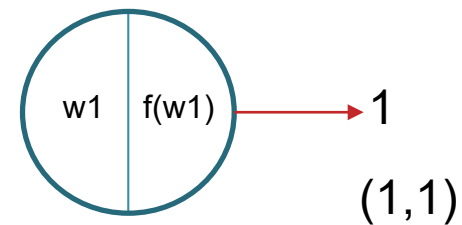
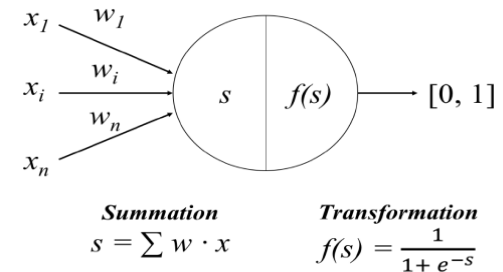
- Word level classification
  - Positive=1;Negative=0
  - let  $n=5 = \text{vocabulary\_size}$

	Like x1	Hate x2	Good x3	Enjoy x4	Bad x5
Like	1	0	0	0	0
Hate	0	1	0	0	0
Good	0	0	1	0	0
Enjoy	0	0	0	1	0
Bad	0	0	0	0	1

(1,5)

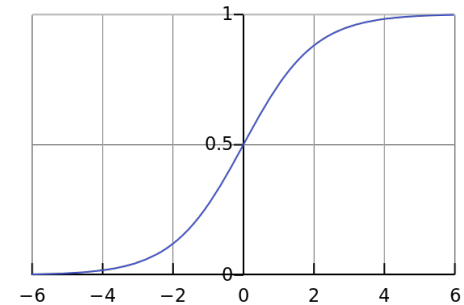
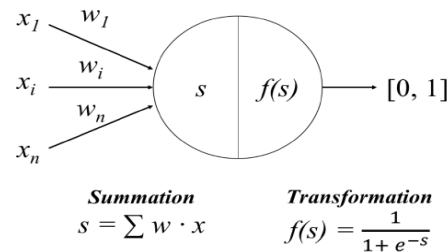
weight
w1
w2
w3
w4
w5

(5,1)



# Perceptron

- Word level classification
  - Positive=1;Negative=0
  - let  $n=5 = \text{vocabulary\_size}$
  - $\text{Batch\_size} = 1$

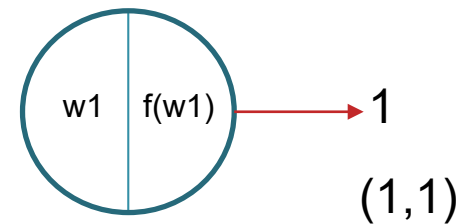


	Like x1	Hate x2	Good x3	Enjoy x4	Bad x5
Like	1	0	0	0	0
Hate	0	1	0	0	0
Good	0	0	1	0	0
Enjoy	0	0	0	1	0
Bad	0	0	0	0	1

(1,5)

weight
6
-6
6
6
-6

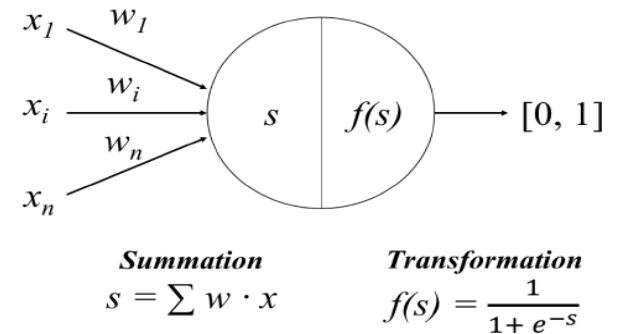
(5,1)



- How about unknown words?

# Perceptron

- Word level classification
  - Positive=1;Negative=0
  - let  $n=26 \ll \text{vocabulary\_size}$

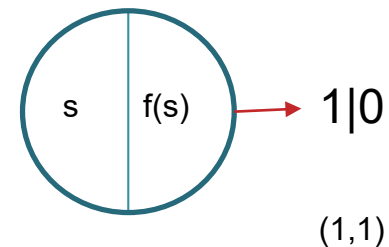


	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	1	1	0	...	0

(1,26)

weight
w1
w2
w3
...
w26

(26,1)

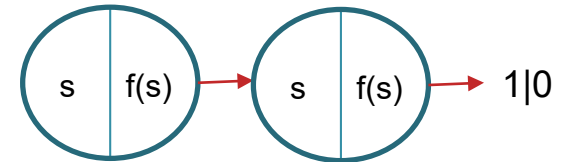


# Multi-Layer Perceptron

- Word level classification
  - let  $n=26 \ll \text{vocabulary\_size}$
  - More parameters
  - More layers

	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	0	1	0	...	0

weight
w1
w2
w3
...
w26

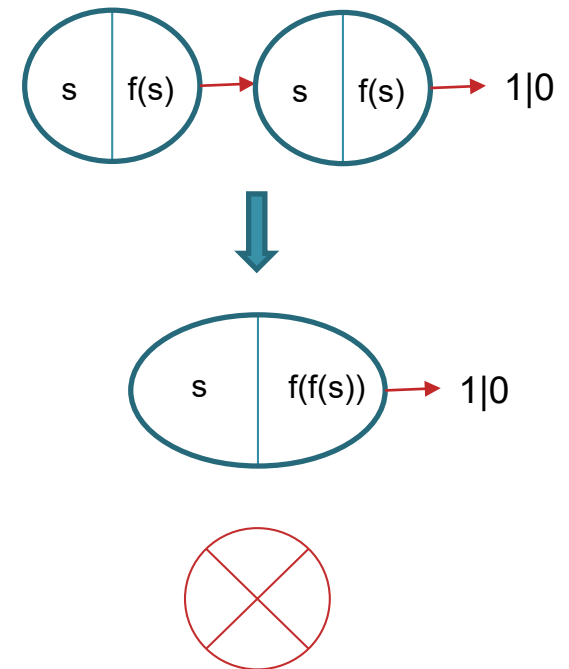


# Multi-Layer Perceptron

- Word level classification
  - let  $n=26 \ll \text{vocabulary\_size}$
  - More parameters
  - More layers

	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	0	1	0	...	0

weight
w1
w2
w3
...
w26



# Multi-Layer Perceptron

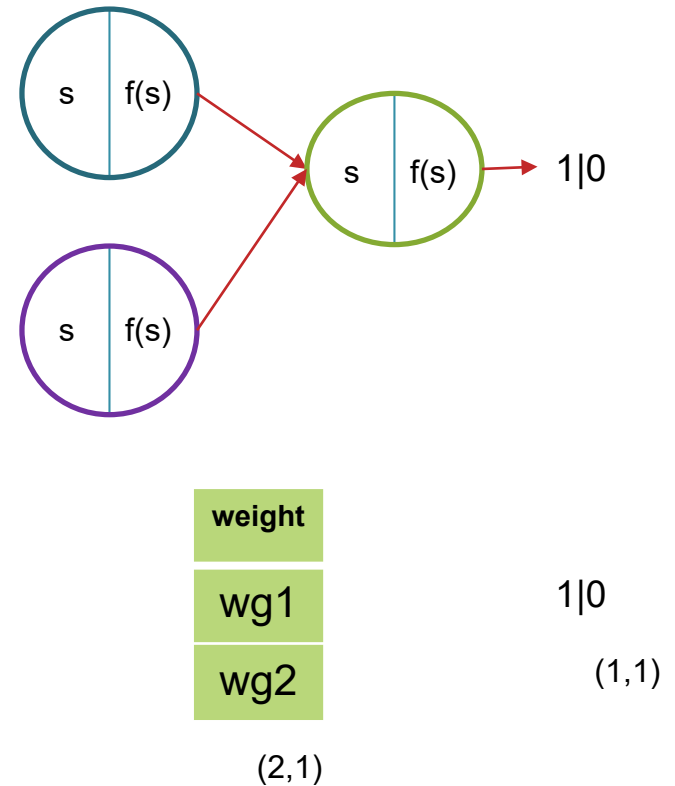
- Word level classification
  - let  $n=26 \ll \text{vocabulary\_size}$
  - More layers
  - More parameters

	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	0	1	0	...	0

(1,26)

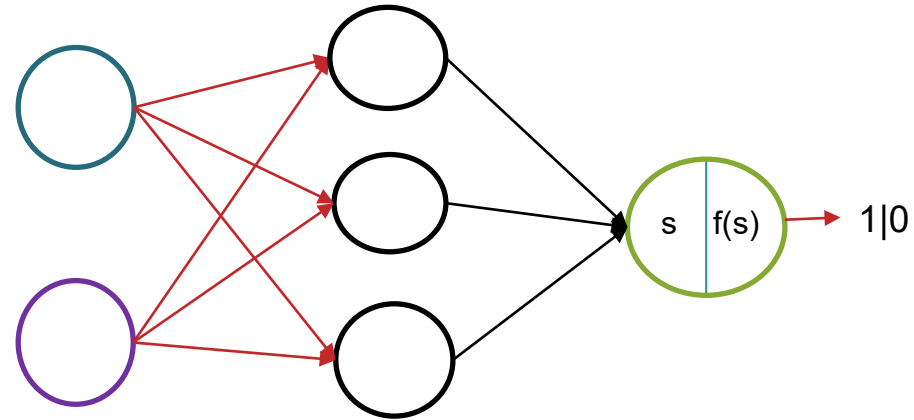
weight	weight
wb1	wp1
wb2	wp2
wb3	wp3
...	...
wb26	wp26

(26,2)



# Multi-Layer Perceptron

- Word level classification
  - let  $n=26 \ll \text{vocabulary\_size}$
  - More parameters
  - Even More Layers



	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	0	1	0	...	0

(1,26)

weight	weight
wb1	wp1
wb2	wp2
wb3	wp3
...	...
wb26	wp26

(26,2)

weight	weight	weight
w	w	w
w	w	w

(2,3)

weight
wg1
wg2
wg3

(3,1)

Y\_pred  
(1,1)

# Deep leaning Training Routine

**Repeat the following until desired**

- Initialize weights vector
  - Random
  - One-hot encoding
- Forward Propagation
- Compute and log the loss
- Back Propagation
- **Optimizer**



# Deep leaning Training Routine

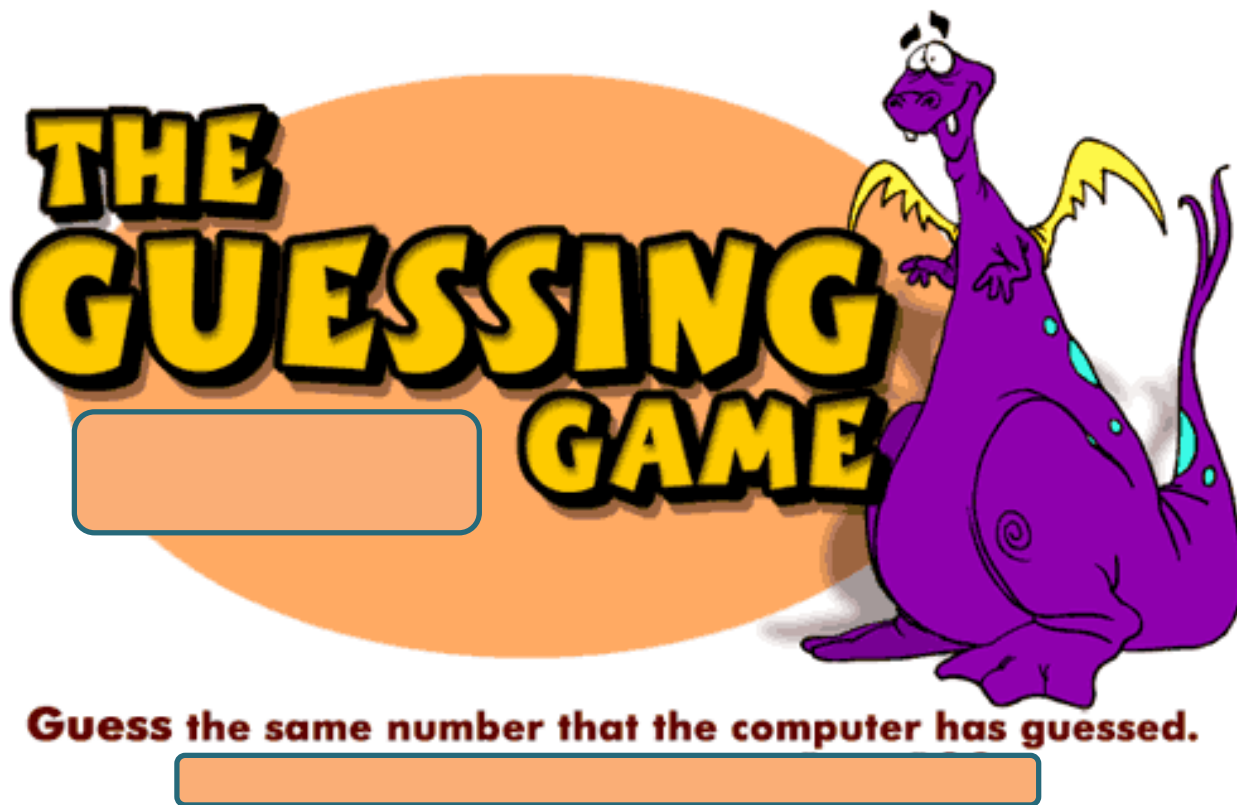
- **Optimizer** (searching action with a **Strategy** )



**Guess the same number that the computer has guessed.  
The number will range from 1 to 100.**

# Deep leaning Training Routine

- **Optimizer** (“brute force” searching action)



# Deep leaning Training Routine

Repeat the following until desired

- ...

- **Optimizer** (“brute force” **searching** action + **Strategy**)

- Gradient Descent and Delta rule

$$\text{New weight} = \text{Old weight} - \text{Derivative Rate} * \text{Learning rate}$$

- learning rate a constant (usually very small)
- to avoid big steps

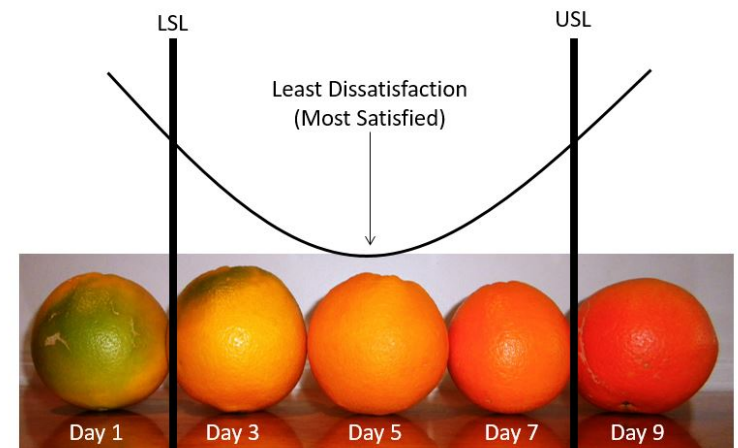


**Guess the same number that the computer has guessed.**

# Deep leaning Training Routine

“brute force” searching action with **strategy** until desired

- **Compute and keep the Cost/Loss**
  - Define and compute the  $loss(Y_{pred}, Y_{correct})$
- **Back Propagation**
  - Compute the **partial derivatives rate** of *Loss function* wrt *W* for all layers
  - Chain rule applies



- **Optimizer**

$$\text{New weight} = \text{Old weight} - \text{Derivative Rate} * \text{Learning rate}$$

Business Performance Improvement

BIZ-PI.com



# Optimization

- Example
  - $y = wx$  with  $n=5$  training examples (forward)
  - Loss function:  $L = \frac{\sum(y_p - y_t)}{n} = \frac{\sum(wx - y_t)}{n}$
  - when  $w$  randomly initialized as 3 (guessing)

$$\frac{dL}{dw} = \sum_{i=0}^n x/n = 2$$

$$w^{new} = w^{old} - \eta \frac{dL}{dw}$$

$$\eta = 0.5 \quad w_{init} = 3$$

1st iter

$$w^{new} = 3 - 0.5 * 2 = 2$$

$$y = 2x \quad \text{Loss function} = 0$$

$x$	$y_{pred}$ ( $w=3$ ) $y=3x$	$y_{true}$	$y_{pred} - y_{true}$
0	0	0	0
1	3	2	1
2	6	4	2
3	9	6	3
4	12	8	4
Loss	-	-	2

# Chain Rule

$$\text{loss}(Y_{\text{pred}}, Y_{\text{correct}}) = Y_{\text{pred}} - Y_{\text{correct}} = \underbrace{x * W_{bp}}_{f_1(x, W_{bp}) = y_1} * \underbrace{W}_{f_2(y_1, W) = y_2} * \underbrace{W_g}_{f_3(y_2, W_g) = y_3} - 1$$

$$\frac{\partial \text{Loss}}{\partial w_{bp}} = \frac{\partial(y_3)}{\partial y_2} \frac{\partial(y_2)}{\partial y_1} \frac{\partial(y_1)}{\partial w_{bp}}$$

	a x1	b x2	c x3	...	z x26
like	0	0	0	...	0
hate	1	0	0	...	0
good	0	0	0	...	0
enjoy	0	0	0	...	0
bad	0	1	0	...	0

(1,26)

weight	weight
wb1	wp1
wb2	wp2
wb3	wp3
...	...
wb26	wp26

(26,2)

weight	weight	weight
w	w	w
w	w	w

(2,3)

weight
wg1
wg2
wg3

(3,1)

y  
(1,1)

$$W_{bp}^{\text{new}} = W_{bp}^{\text{old}} - \eta \frac{\partial \text{Loss}}{\partial w_{bp}}$$

# Deep learning Training Routine

Repeat the following until desired

- Initialize weights vector  $W$  for all layers
- **Forward Propagation**
  - *reaching the final layer to get  $Y_{pred}$*
- **Compute and keep the cost/loss**
  - Define and compute the ***loss***( $Y_{pred}, Y_{correct}$ )
- **Back Propagation**
  - **partial derivatives rate of *Loss function wrt  $W$  for all layers***
- **Optimizer**
  - ***New weight = Old weight – Derivative Rate \* Learning rate***



# Agenda

- Overview of NLP
- Deep Learning Basics for NLP
- **Workshop: NN Basics on Colab**
- Word2Vec & DL Specifics
  - CBOW and SkipGram
  - ActivationFunction/LossFunction
  - Optimiser/Learning Rate
- Workshop: Word2Vec from Scratch



# Text Processing using Machine Learning

## Word2Vec & DL Specifics

Dr Wang Aobo  
aobo.wang@nus.edu.sg

# Features from Word Vectors

- **Count from Data**
  - Word Co-occurrence + SVD
  - Count-based model
- **Learn from Data**
  - CBOW and SKIPGRAM
  - NN Methods
  - Predictive Model
- **Count and Learn from Data**
  - GLOVE: Global Vectors for Word Representation
  - Count + SGD

NMSM  
Day 2

TPML  
Day 1

# One-Hot Encoding (Sparse Representation)

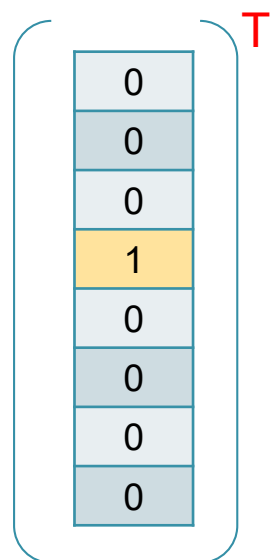
Vocabulary of the corpus (big enough)

	give	she	at	talk	have	ramen	a	drink
give	1	0	0	0	0	0	0	0
talk	0	0	0	1	0	0	0	0
have	0	0	0	0	1	0	0	0
drink	0	0	0	0	0	0	0	1
a	0	0	0	0	0	0	1	0
at	0	0	1	0	0	0	0	0

$v('talk')$

# Lookup Function

$v('talk')$



One-Hot Encoding

$1 \times |V|$



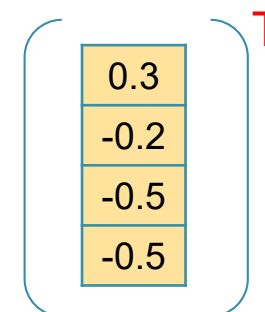
0.1	0.2	-0.4	0.9
0.2	0.1	-0.3	0.9
0.2	-1.4	0.3	-0.1
0.3	-2.0	0.5	-0.5
0.2	-1.1	0.3	-0.7
0.9	-1.3	0.4	-0.9
0.3	-3.0	0.5	-0.2
0.5	-0.1	0.2	0.1

Word Embeddings

$|V| \times d$



$v('talk')$   
Embedded



Input

$1 \times d$

Learn the Matrix through Making Prediction

38

# Word2Vec (CBOW)

Learn the Matrix through “**classification**” task

**Sentence:** the bulk of linguistic questions concern the distinction between a and m. a linguistic account of phenomenon ...

of  
linguistic  
questions  
concern  
the  
dis-  
tinction  
between  
a  
and  
m.  
a  
linguistic  
account  
of  
a  
phenomenon

the bulk \_\_\_\_\_ linguistic questions  
bulk of \_\_\_\_\_ questions concern  
of linguistic \_\_\_\_\_ concern the  
linguistic questions \_\_\_\_\_ the dis-  
questions concern \_\_\_\_\_ dis- tinction  
concern the \_\_\_\_\_ tinction between  
the dis- \_\_\_\_\_ between a  
dis- tinction \_\_\_\_\_ a and  
tinction between \_\_\_\_\_ and m.  
between a \_\_\_\_\_ m. a  
a and \_\_\_\_\_ a linguistic  
and m. \_\_\_\_\_ linguistic account  
m. a \_\_\_\_\_ account of  
a linguistic \_\_\_\_\_ of a  
linguistic account \_\_\_\_\_ a phenomenon  
account of \_\_\_\_\_ phenomenon gen-  
of a \_\_\_\_\_ gen- erally

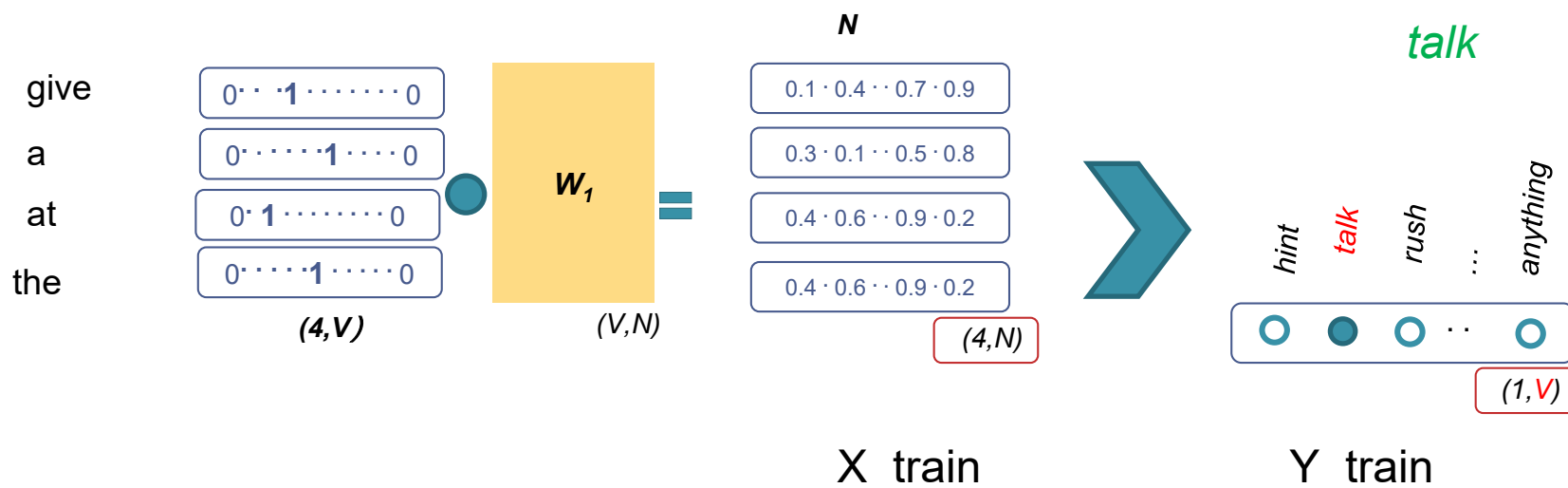


**window\_size = 2**

## Word2Vec (CBOW)

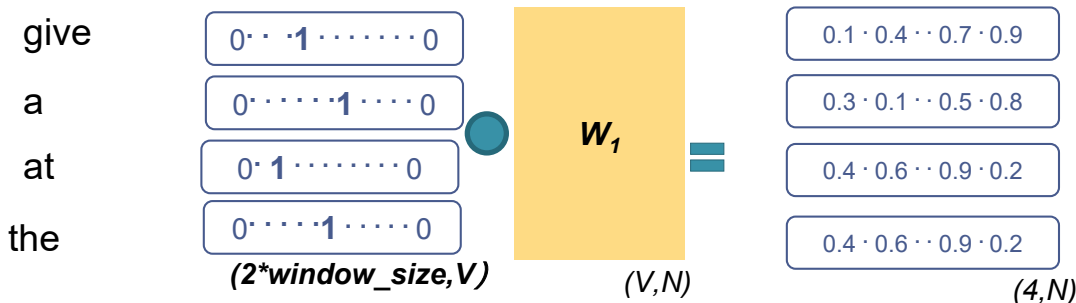
**Task:** Iterate through every word with a given window; learn  $W$  such the models can predict what's the word given only the context words as inputs.

give a \_\_\_\_\_ at the



# Word2Vec (CBOW)

give a talk at the



Re-shaping  
the rows

**AVG () + Relu ()**

Re-shaping  
the values

0.1 · 0.1 ··· 7 · 0.4

$(1, N)$

Re-shaping  
the columns

$W_2$

$(N, V)$

Re-shaping  
the Probs

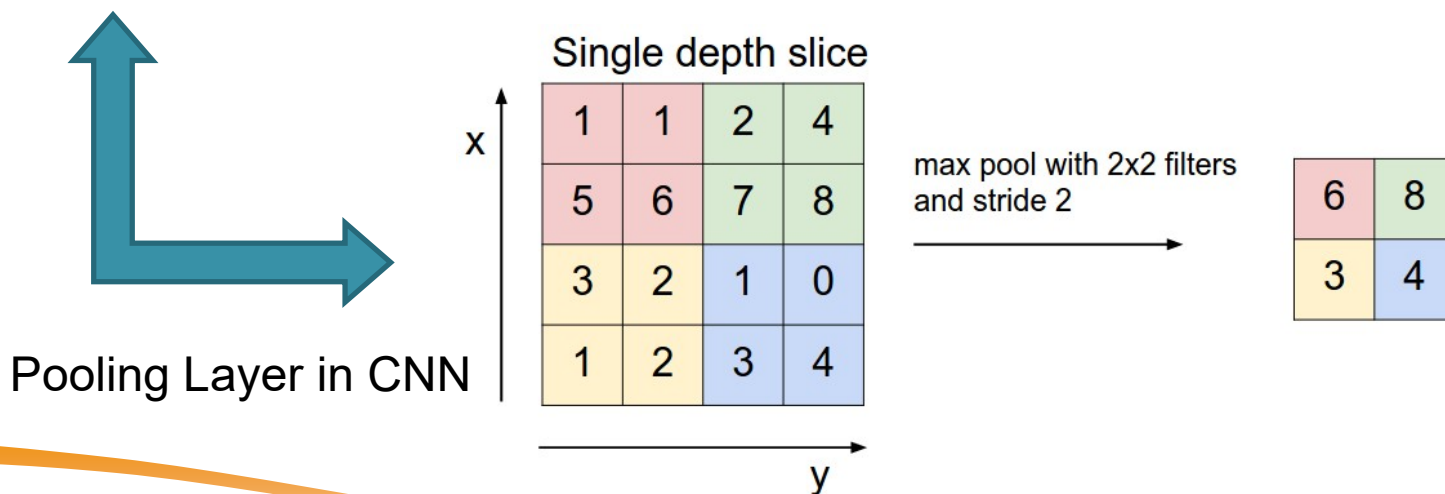
**Softmax ()**

hint talk rush ... anything

$(1, V)$

# Naïve Sentence Embedding

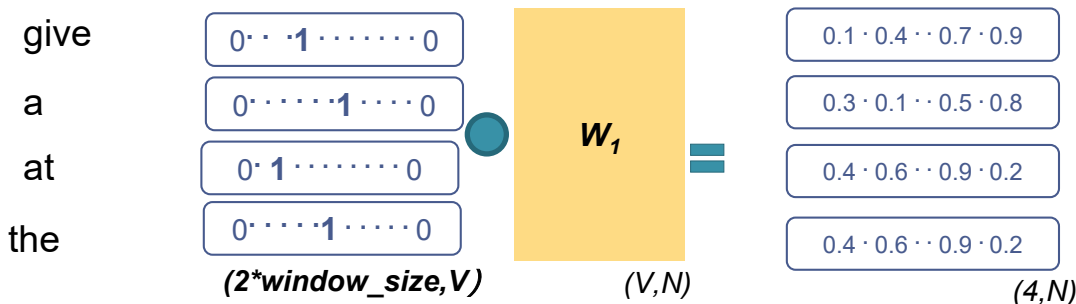
<b>v('give')</b>	0.5	-1.3	0.6	1.1
<b>v('a')</b>	0.3	-0.2	0.7	0.5
<b>v('at')</b>	0.3	2.3	-0.8	1.0
<b>v('the')</b>	1.7	-0.2	-0.1	1.0
	<b>AVG()/MAX()/MIN()/Concat()</b>			
<b>Sent0 ("give a at the")</b>	0.7	0.2	0.1	0.9





# Word2Vec (CBOW)

give a talk at the



Re-shaping  
the rows

**AVG () + Relu ()**

Re-shaping  
the values

$0.1 \cdot 0.1 \cdots 7 \cdot 0.4$

(1, N)

Re-shaping  
the columns

$W_2$

(N, V)

Re-shaping  
the Probs

**Softmax ()**

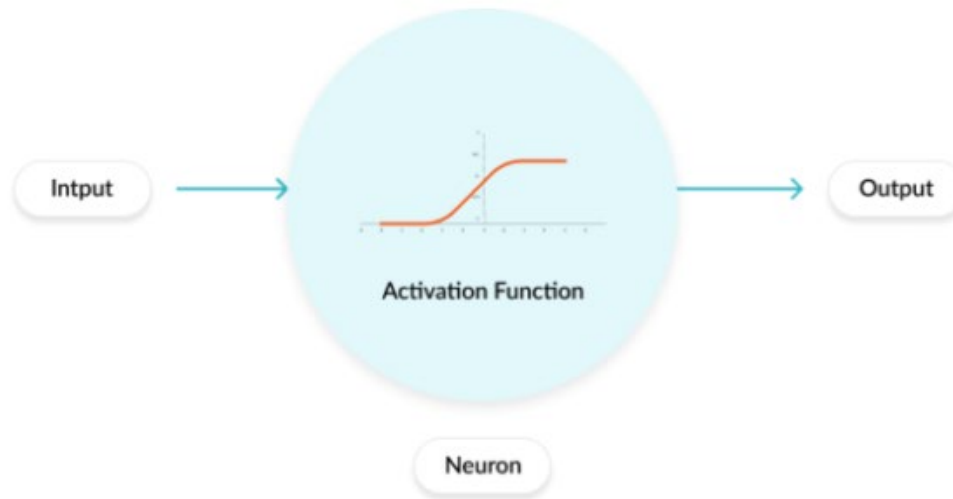
hint  
talk  
rush  
...  
anything

(1, V)

# Activation Function

- **Activation Function**

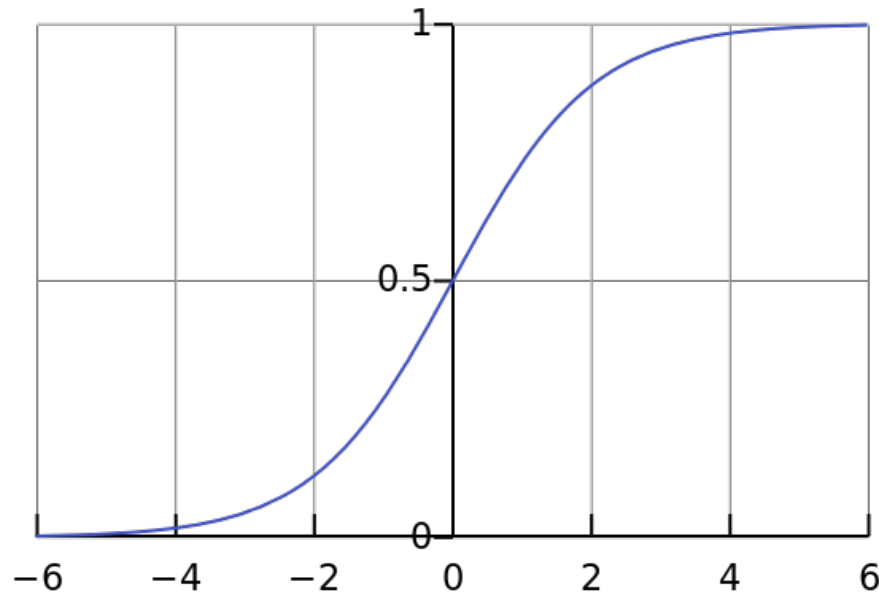
- Smoother decision function is expected
- Activations bound in  $(0,1)$
- Support backpropagation
- Need to be non-linear



# Activation Function

- **Sigmoid** function
  - smooth output between 0 and 1
  - interpreted as a probability of “Yes”

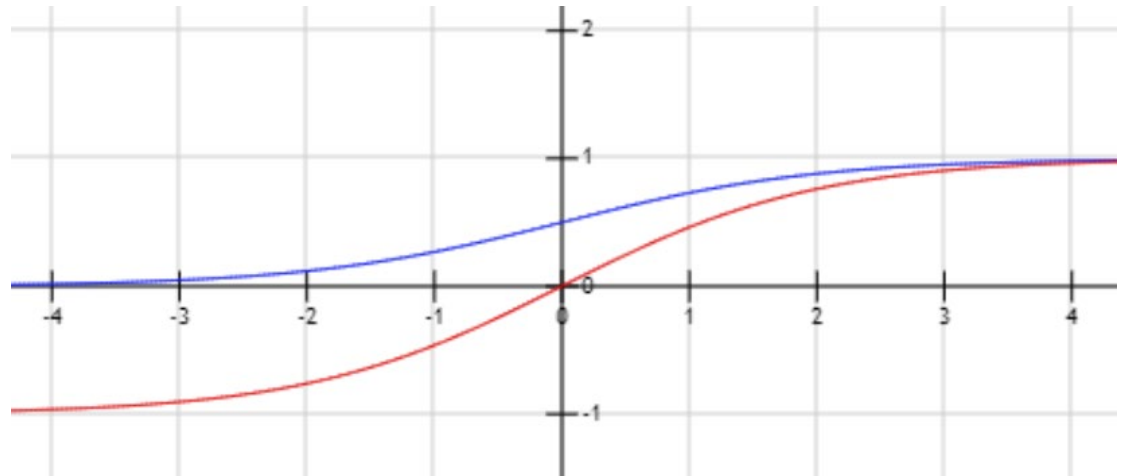
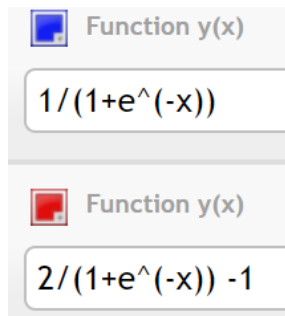
$$S(x) = \frac{1}{1 + e^{-x}}$$



(Image from Wikipedia)

# Activation Function

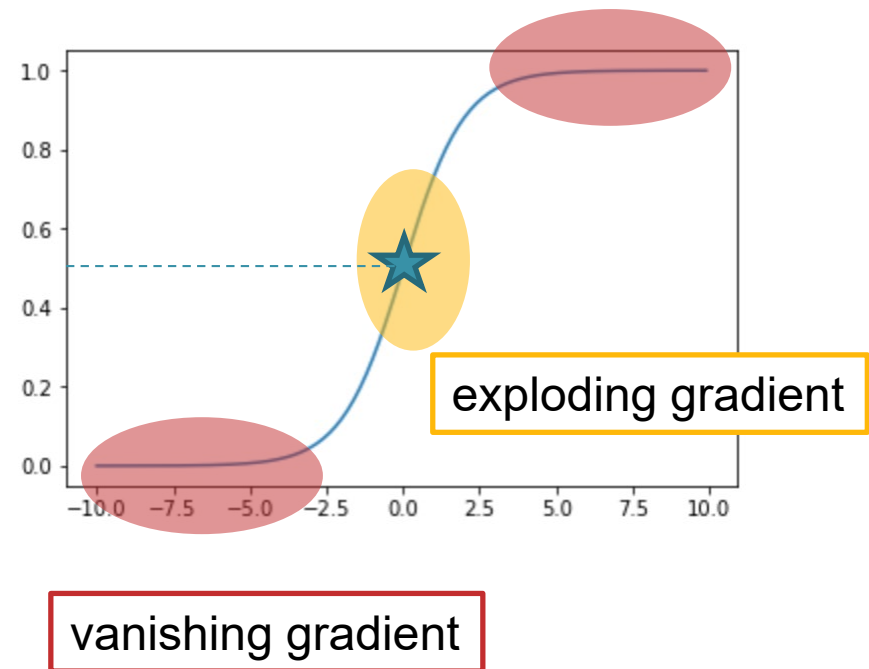
- Tanh Function
  - smooth output between -1 and 1
  - 0 centroid



$$\tanh(x) = 2 \operatorname{sigmoid}(2x) - 1$$

# Activation Function (Sigmoid)

```
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 def sigmoid(x):
7     return 1/(1+np.exp(-x))
8
9 # Generate points from -10 to +10,
10 # in steps of 0.1
11 x = np.arange(-10, 10, 0.1)
12 y = sigmoid(x)
13
14 # Plot the graph.
15 plt.plot(x, y)
16 plt.show()
```



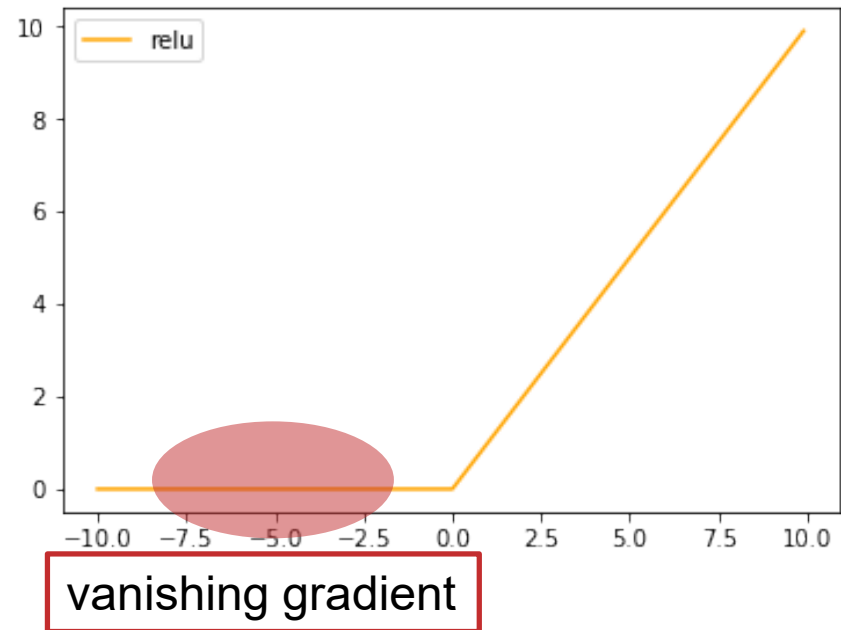
# Activation Function (ReLU)

- **Intuitions can be useless for NN**
  - Smoother decision function is expected
  - Activations bound in a range
  - “0” centroid



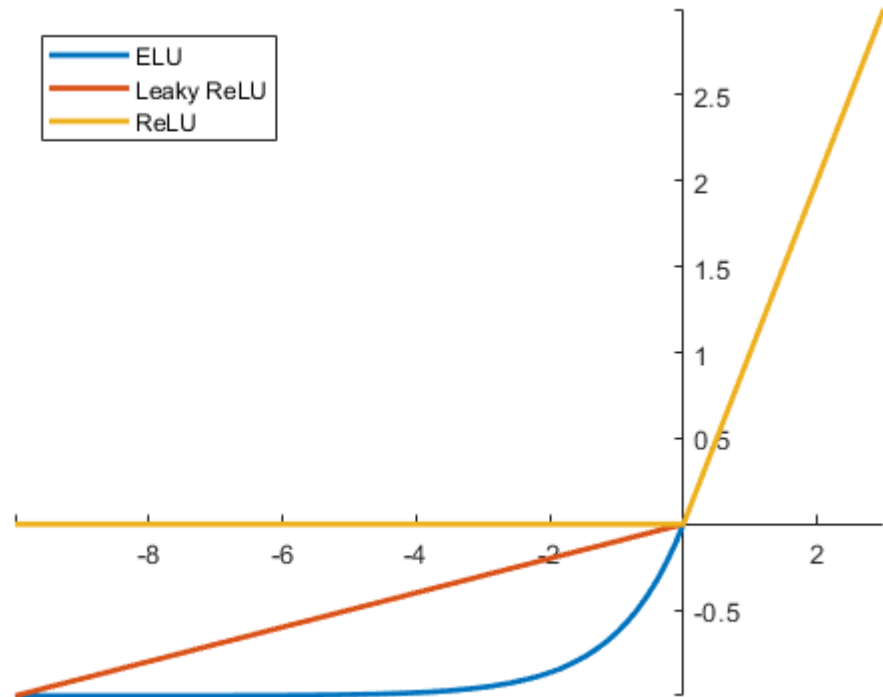
- **All you need is just**
  - Simple
  - Fast

$$A(x) = \max(0, x)$$



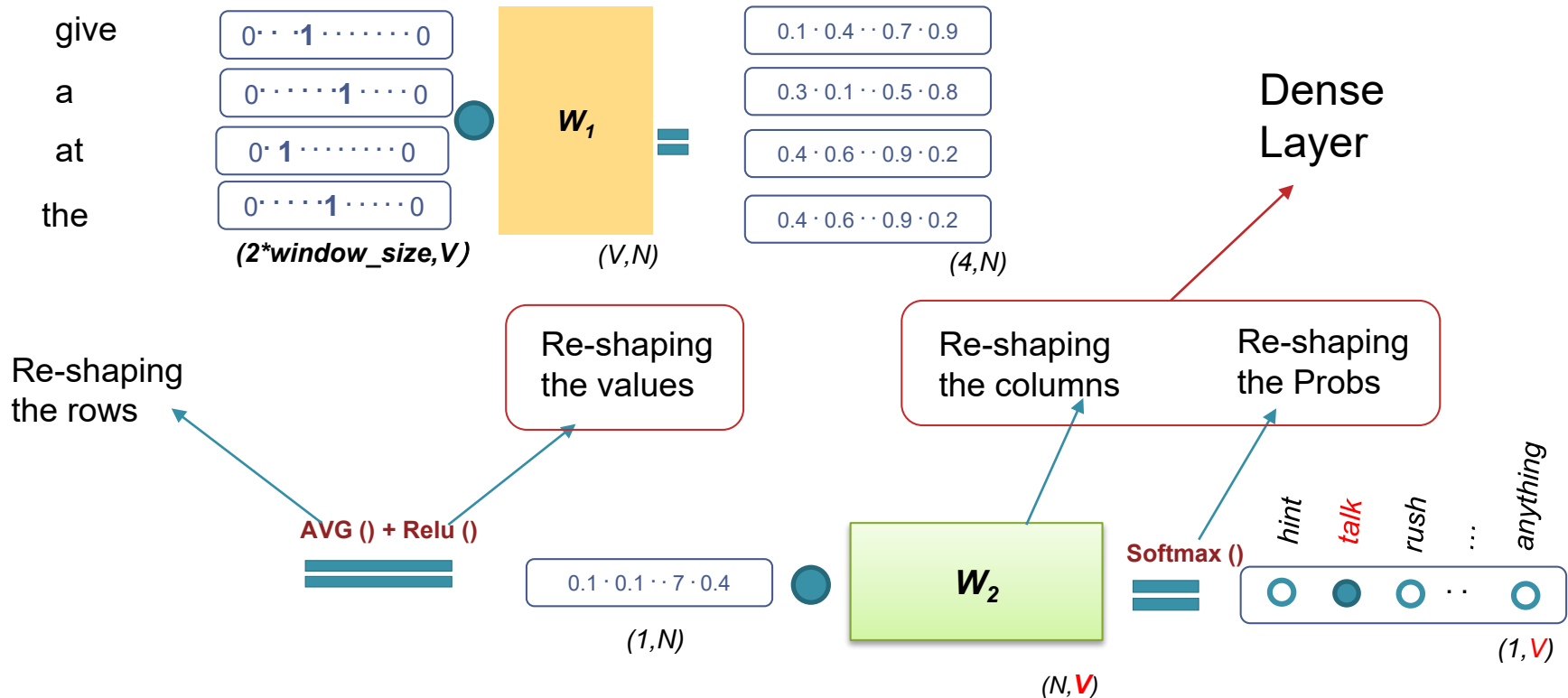
# Activation Function (ReLU)

- ReLU's Family



# Word2Vec (CBOW)

give a talk at the

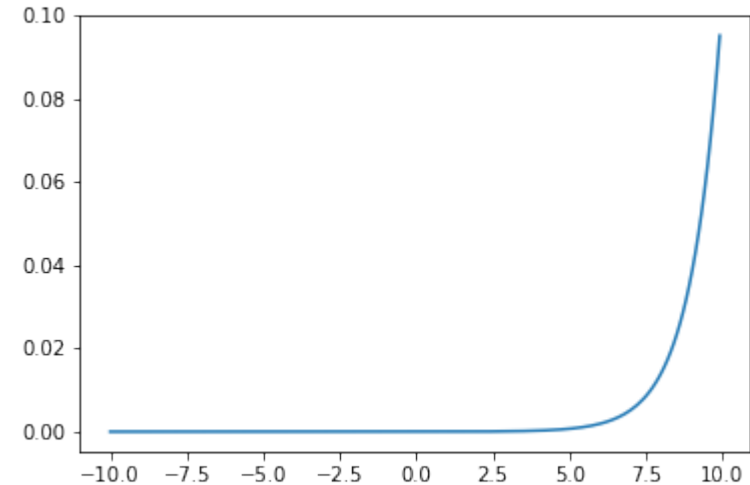




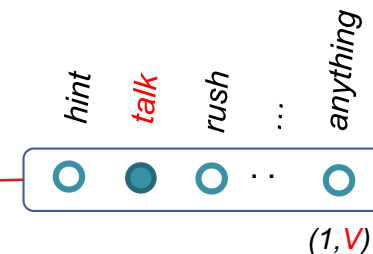
# Last Layer Activation Function

- SoftMax

```
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 def softmax(x):
7     return np.exp(x) / np.sum(np.exp(x), axis=0)
8
9 x = np.arange(-10, 10, 0.1)
10 y = softmax(x)
11
12 plt.plot(x,y)
13 plt.show()
```

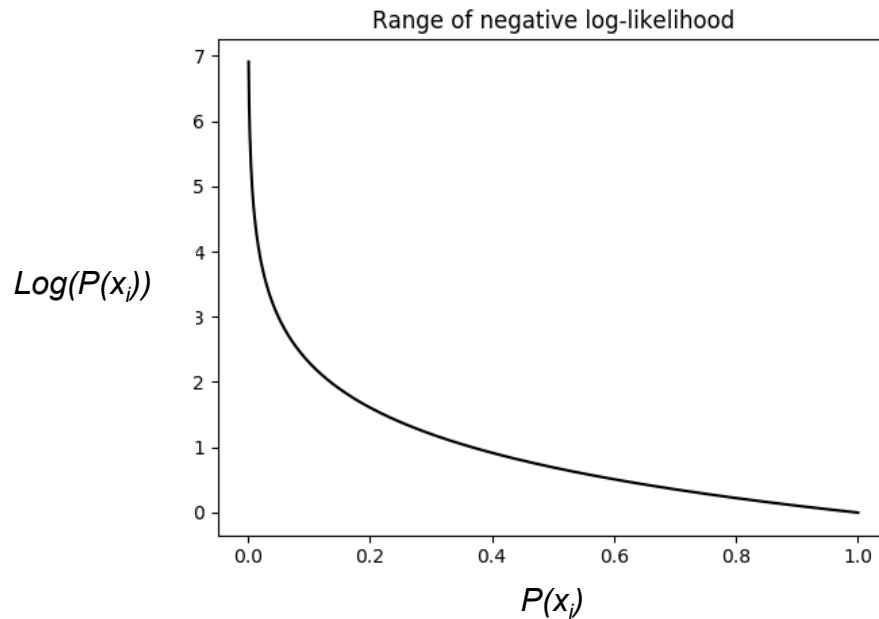


Vector of  $x_i$  as Input  
Never be a single  $x$

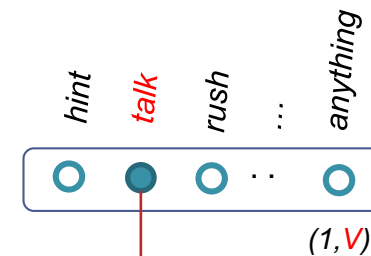


# Last Layer Activation Function

- NegLogSoftMax



Max  $\{ p(x_i) \}$  as input  
Never be a vector of  $p(x_i)$

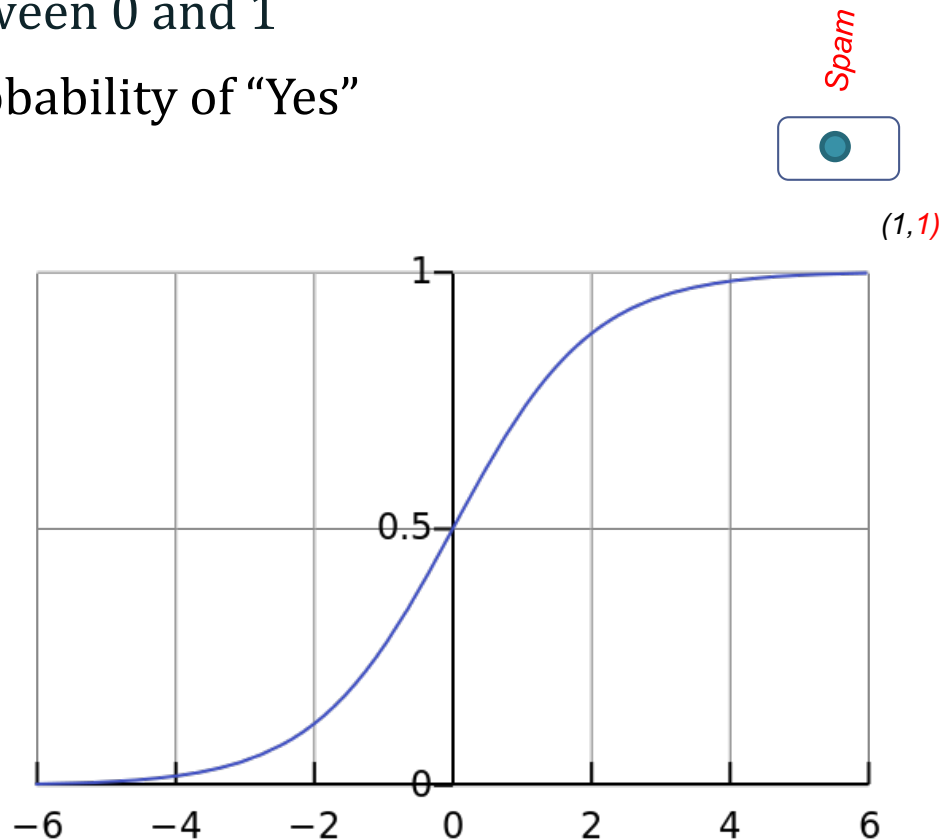


# Last Layer Activation Function

- **Sigmoid function**
  - smooth output between 0 and 1
  - interpreted as a probability of “Yes”

A single  $x$  as input  
Never be a vector of  $x_i$

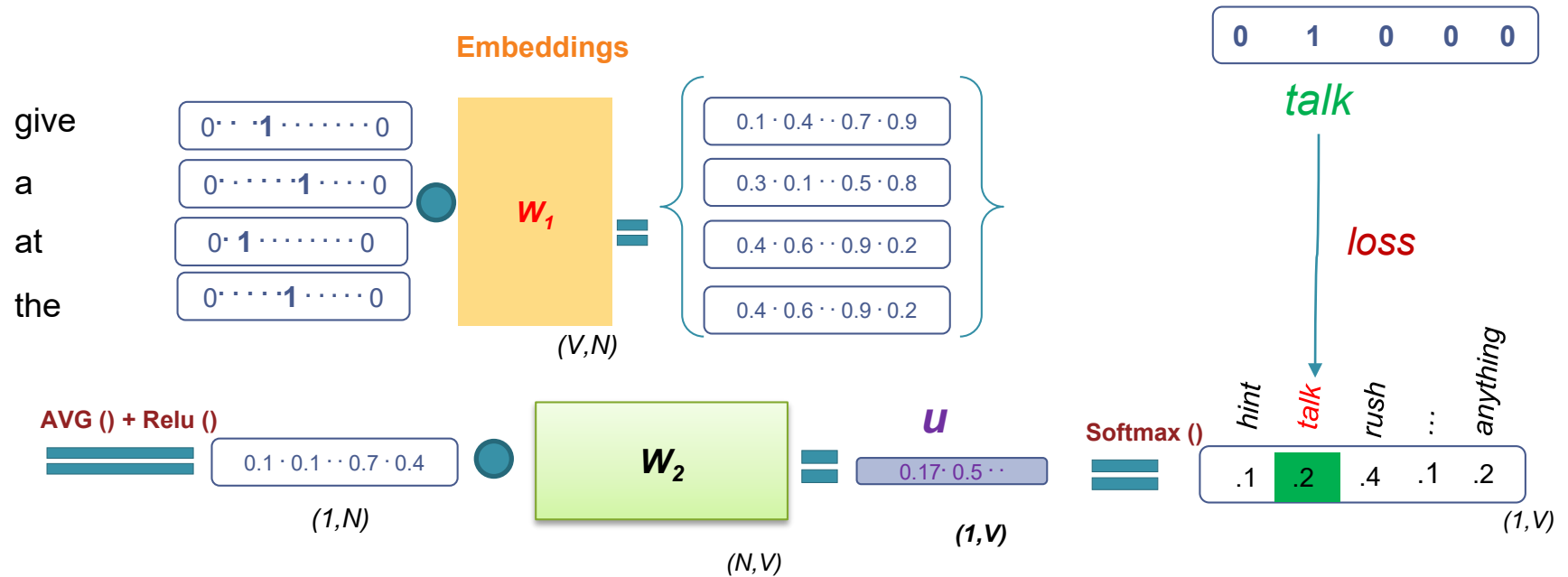
$$S(x) = \frac{1}{1 + e^{-x}}$$



(Image from Wikipedia)

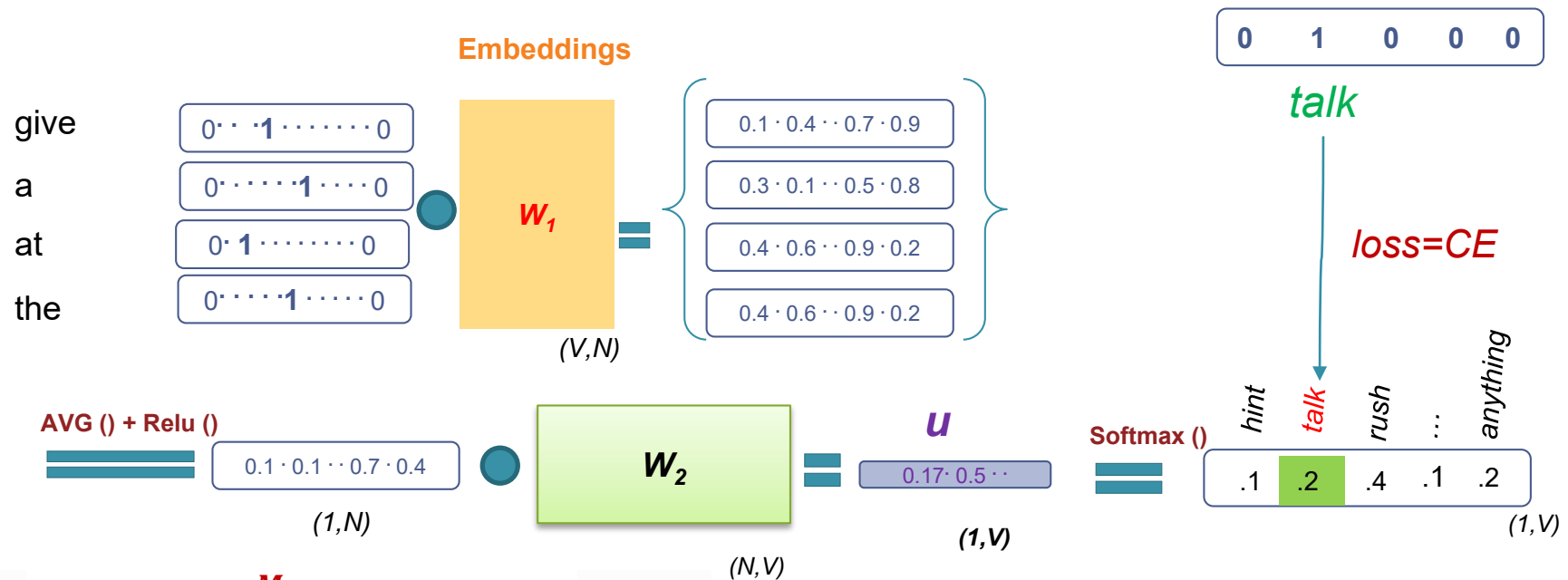
# Word2Vec (CBOW)

## Loss Function



# Loss Function

## (Categorical) Cross Entropy Loss Function



$$Loss = - \sum_{c=1}^M y_{o,c} \log(p_{o,c})$$

### Note

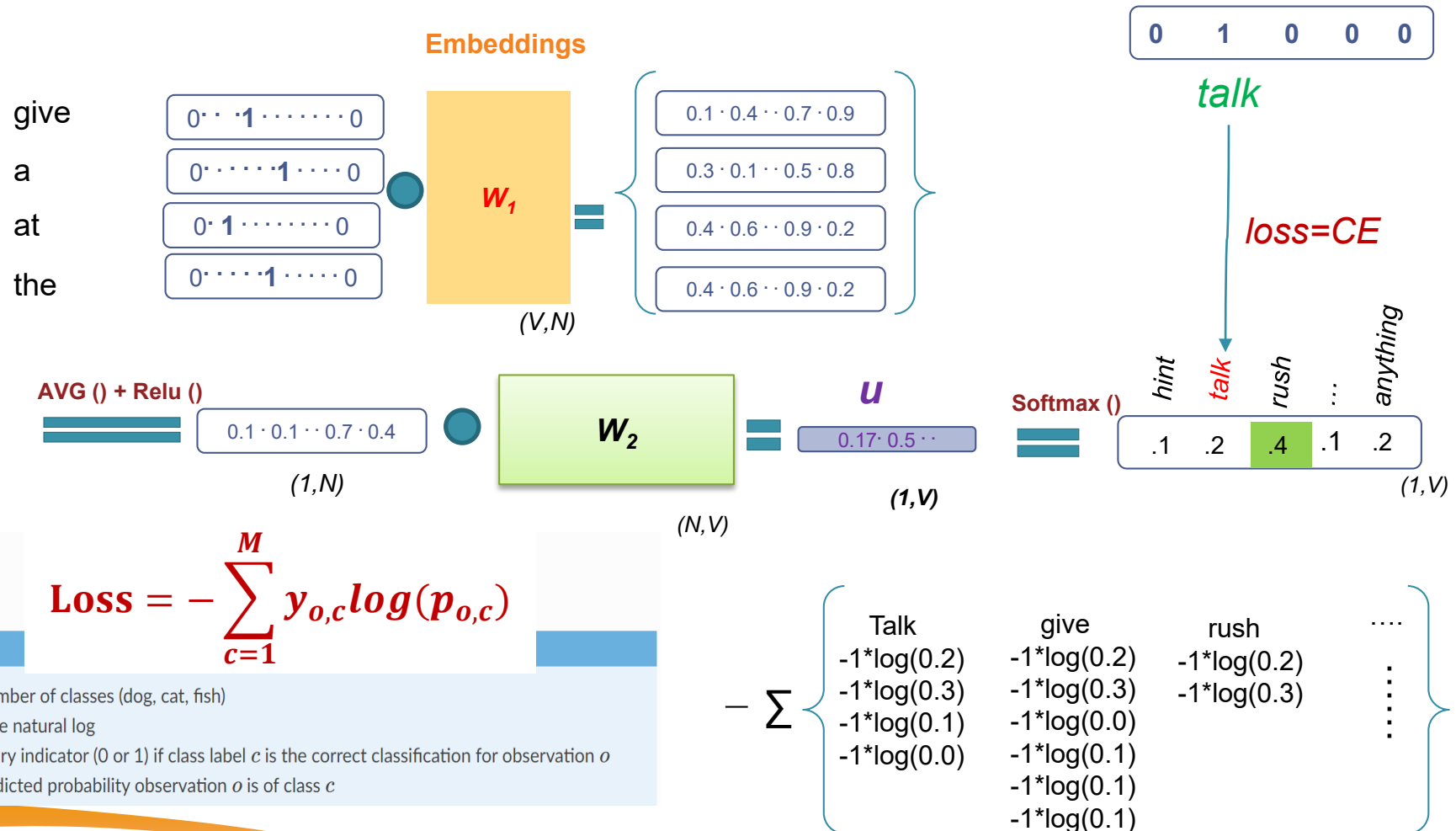
- M - number of classes (dog, cat, fish)
- log - the natural log
- y - binary indicator (0 or 1) if class label c is the correct classification for observation o
- p - predicted probability observation o is of class c

Talk

$$- 0 \cdot \log(0.1) - 1 \cdot \log(0.2) - 0 \cdot \log(0.4) - 0 \cdot \log(0.1) - 0 \cdot \log(0.2)$$

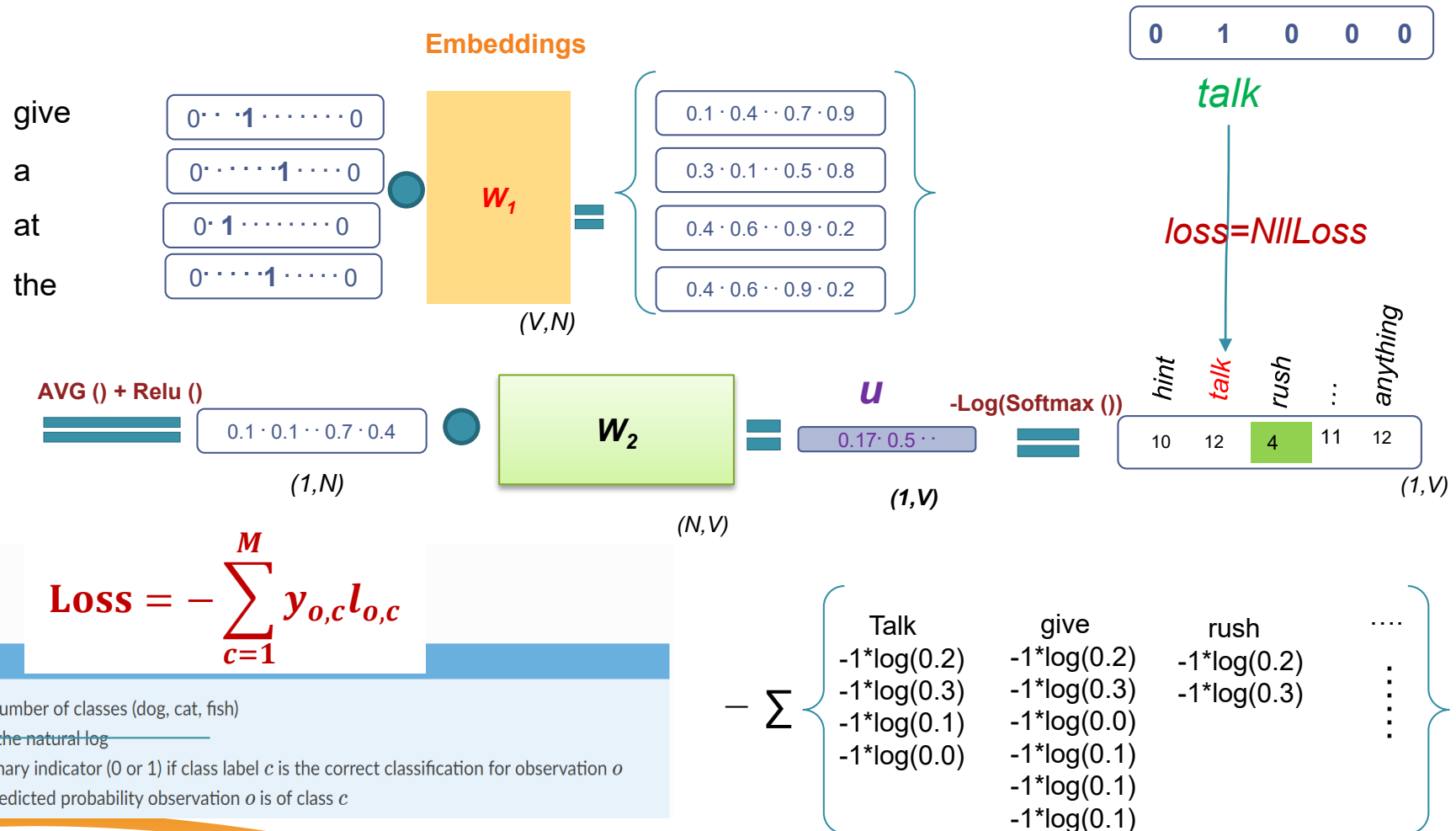
# Loss Function

## (Categorical) Cross Entropy Loss Function



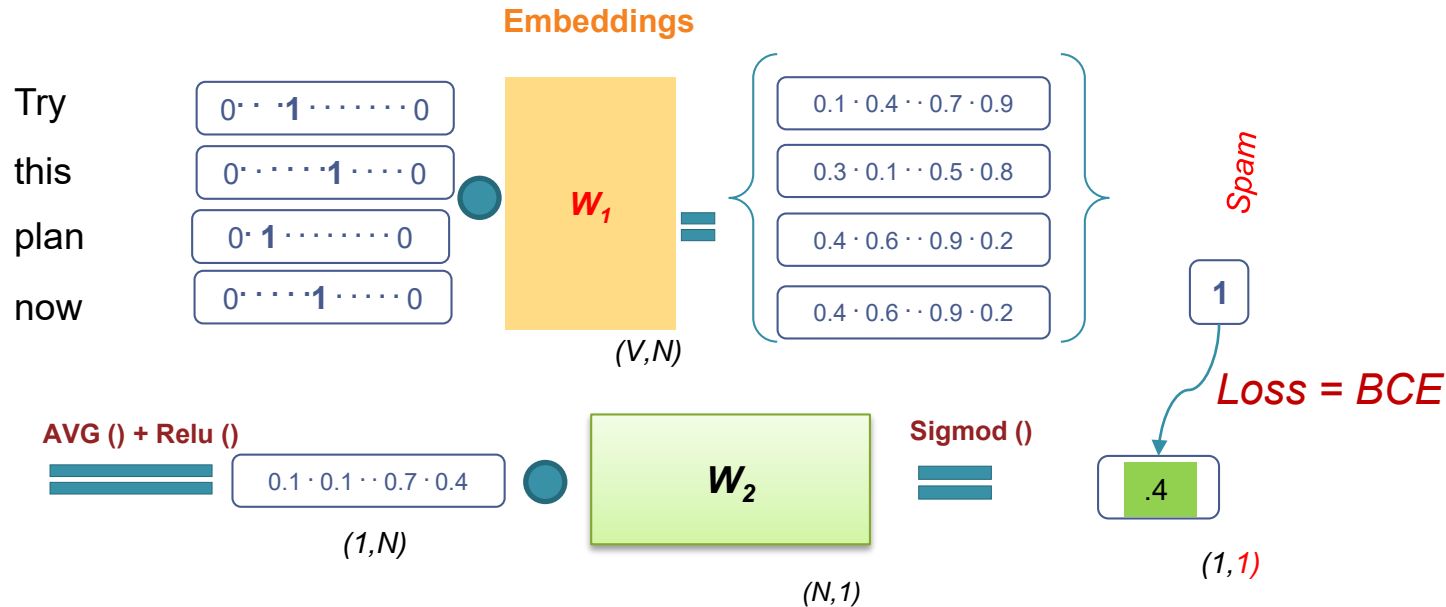
# Loss Function

## (Categorical) Negative Log Loss Function



# Not For CBOW

## Binary Cross Entropy Loss Function



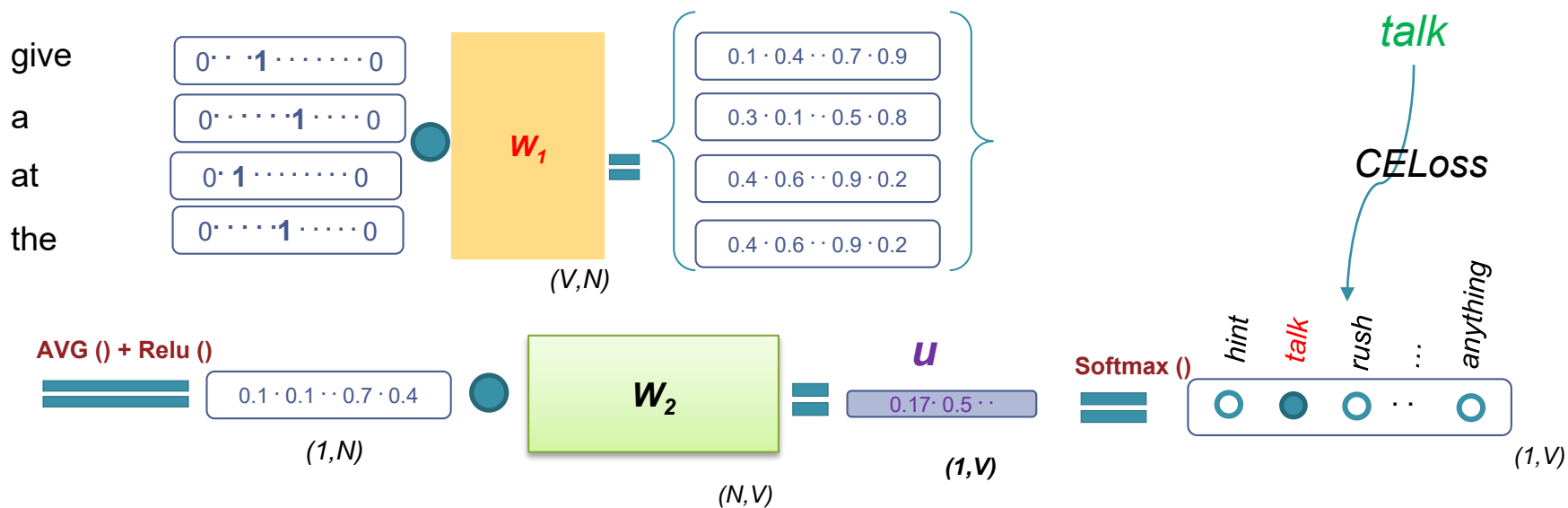
$$BCE = - \sum_{i=1}^{C'=2} t_i \log(s_i) = -t_1 \log(s_1) - (1 - t_1) \log(1 - s_1)$$

$$- \sum \left\{ \begin{array}{ll} \text{spam} & \text{not spam} \\ -1 * \log(0.4) & - (1-0.4) * \log(1-0.4) \\ -1 * \log(0.3) & - (1-0.3) * \log(1-0.3) \\ -1 * \log(0.1) & - (1-0.1) * \log(1-0.1) \\ -1 * \log(0.0) & - (1-0.0) * \log(1-0.0) \end{array} \right\}$$



# Word2Vec (CBOW)

## Backpropagation SGD



$$W_1(new) = W_1(old) - \frac{\partial loss}{\partial w_1} * lr$$

$$W_2(new) = W_2(old) - \frac{\partial loss}{\partial w_2} * lr$$

# Backpropagation

## SGD vs Adam

Faster  
but sometimes  
not converging

*# Vanilla SGD*

```
x += - learning_rate * dx
```

**x** is a vector of parameters and  
**dx** is the gradient

Well generalised  
but slower  
Used together with  
*momentum*

*# Adam*

```
m = beta1*m + (1-beta1)*dx  
v = beta2*v + (1-beta2)*(dx**2)  
x += - learning_rate * m /  
      (np.sqrt(v) + eps)
```

**x** is a vector of parameters and  
**dx** is the gradient  
**m** is the smoothen gradient  
**v** is the 'cache' used to normalize **x**  
**eps** is smoothing term (1e-4 to 1e-8)  
**beta1, beta2** are hypers (0.9, 0.999)

# Backpropagation

- **Learning Rate**

- Fixed Learning Rate based on Experience
  - Vanilla model: MLP/RNN/CNN  $10^{-2} \sim 10^{-3}$
  - Typical Models: LSTM/CNN  $10^{-3}$
  - Complex Models: BERT  $10^{-5}$

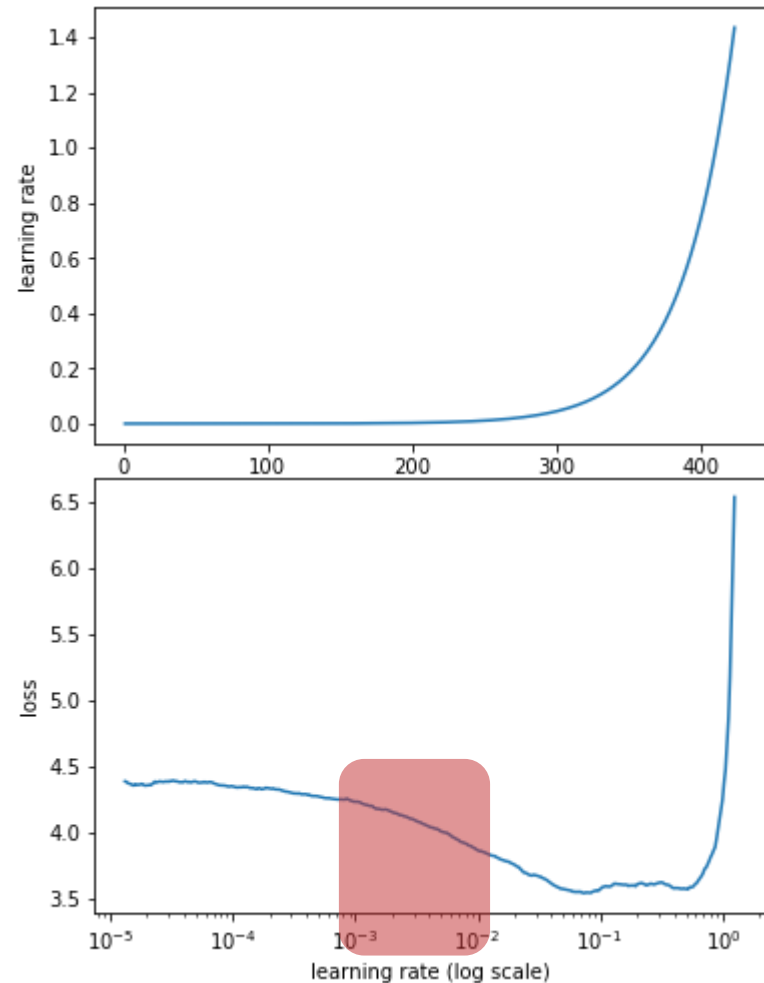
# Backpropagation

- **Fixed Learning Rate**
  - Estimate Fixed Learning Rate

Increase *lr* as iterations going up

save and plot the *loss* as per *lr*

get the *lr*  
fastest decrease in the loss

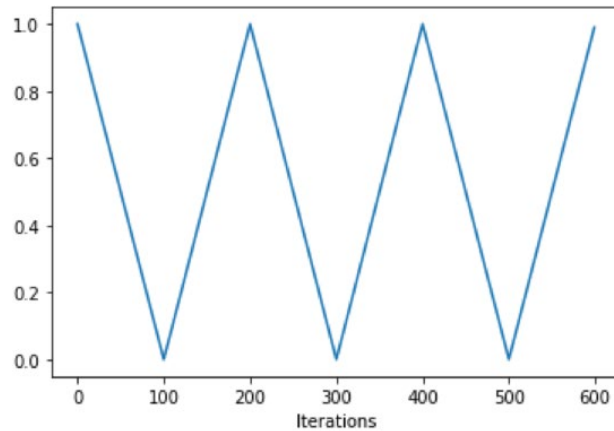


# Backpropagation

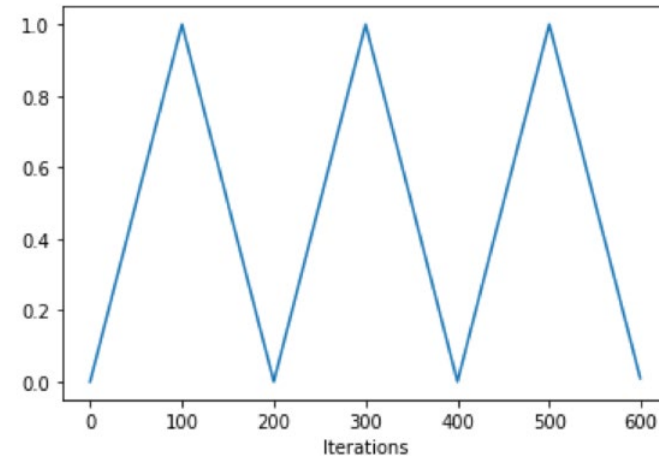
- Dynamic Learning Rate

- Estimate Initial Learning Rate
- Program Dynamic Learning Rate

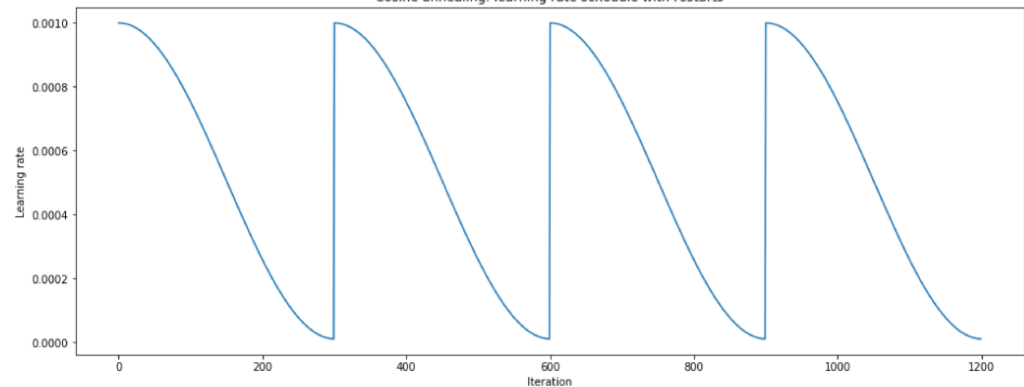
Function:  $\frac{\text{iterations}}{\text{step\_size}} - 2(\text{cycle}) + 1$



Function:  $1 - \frac{\text{iterations}}{\text{step\_size}} - 2(\text{cycle}) + 1$



Cosine annealing: learning rate schedule with restarts

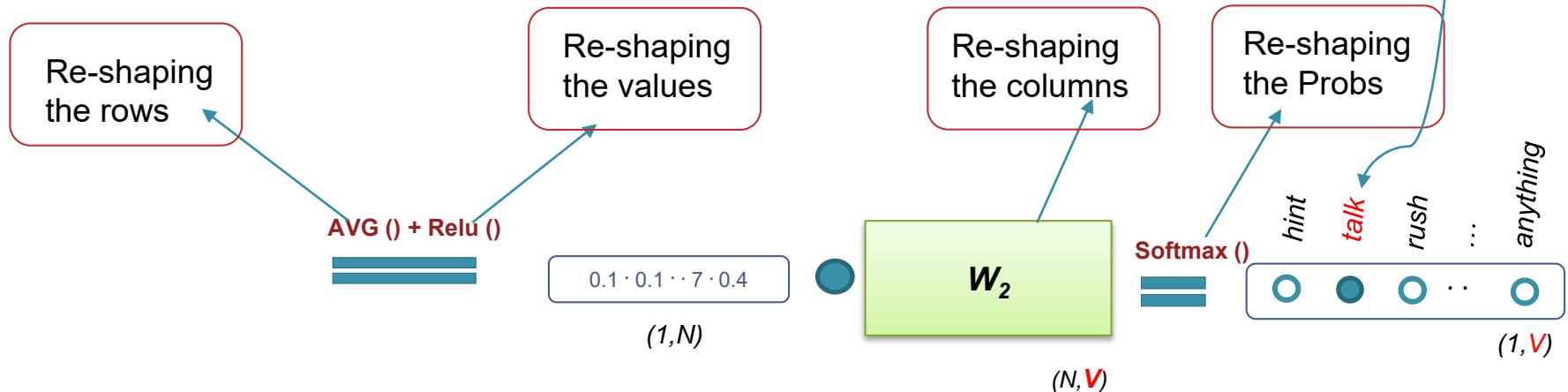
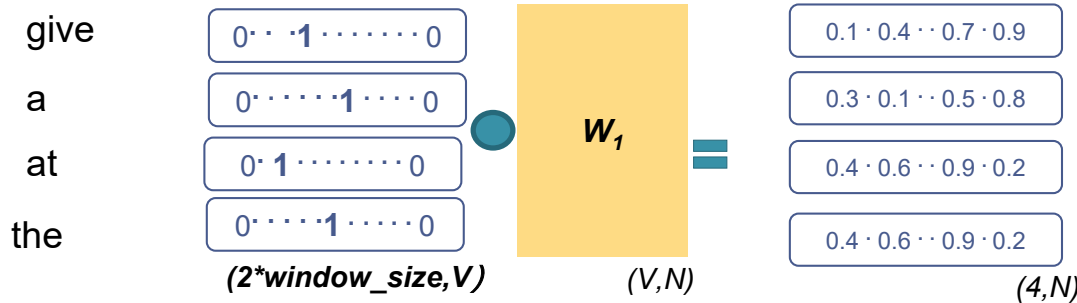


$$\eta_t = \eta_{\min}^i + \frac{1}{2}(\eta_{\max}^i - \eta_{\min}^i) \left( 1 + \cos\left(\frac{T_{\text{current}}}{T_i} \pi\right) \right)$$

**Annealing,**  
i.e. taking a  
partial step

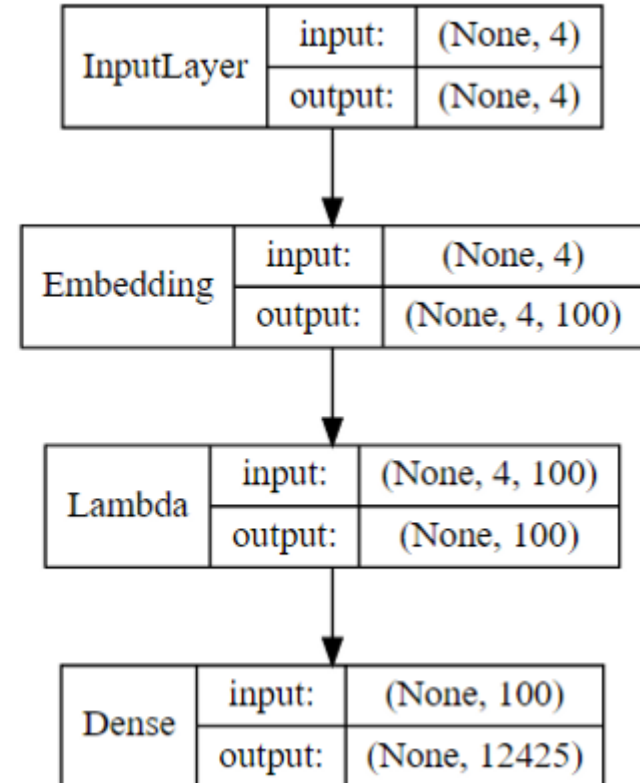
# Word2Vec (CBOW)

give a talk at the



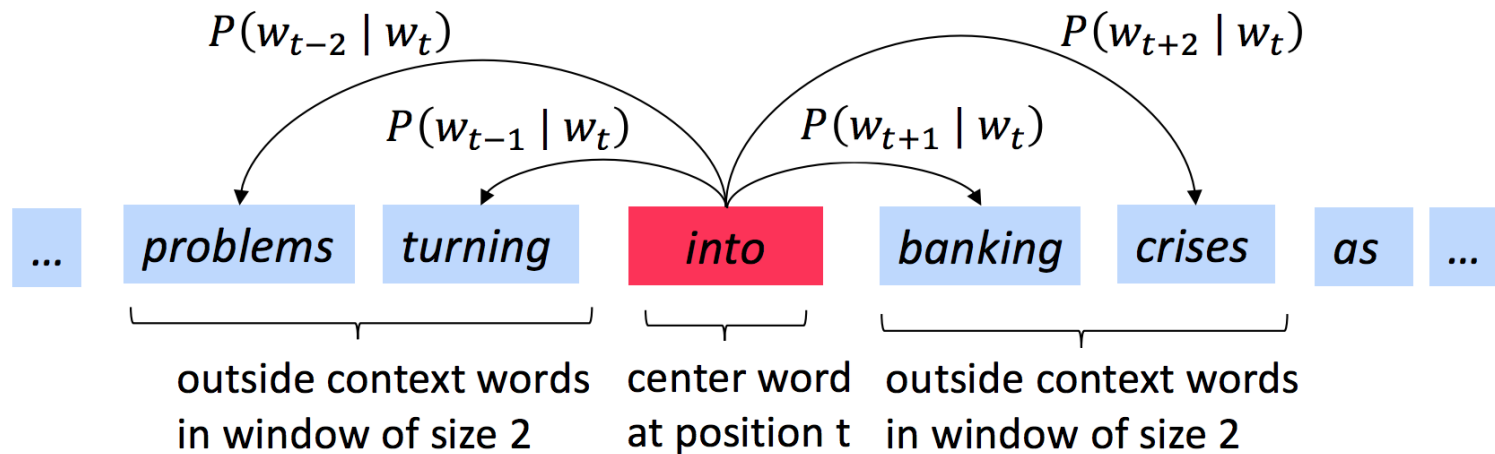
# CBOW Model Summary

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 4, 100)	1242500
lambda_1 (Lambda)	(None, 100)	0
dense_1 (Dense)	(None, 12425)	1254925
Total params: 2,497,425		
Trainable params: 2,497,425		
Non-trainable params: 0		



# Word2Vec (Skipgram)

**Task:** Iterate through each word with a given window; for each word predict the context words within the window



(E.g. from Manning (2018) Stanford cs224n course)



# Word2Vec (Skipgram)

**Sentence:** language users never choose words randomly , and  
language is essentially non-random .

## In-/Outputs:

```
[  
  ('never', ['language', 'users', 'choose', 'words']),  
  ('choose', ['users', 'never', 'words', 'randomly']),  
  ('words', ['choose', 'words', ',', 'and']),  
  ('randomly', ['words', 'randomly', 'and', 'language'])  
]
```

# Word2Vec (Skipgram)

**Sentence:** language users never choose words randomly , and  
language is essentially non-random .

**Windows:**

`['language', 'users', 'never', 'choose', 'words']`

```
('never', 'language', 1),  
( 'never', 'users', 1),  
( 'never', 'choose', 1),  
( 'never', 'words', 1),  
( 'never', ' ', 0),  
( 'never', 'non-random', 0),  
( 'never', 'is', 0),  
( 'never', 'is', 0)
```

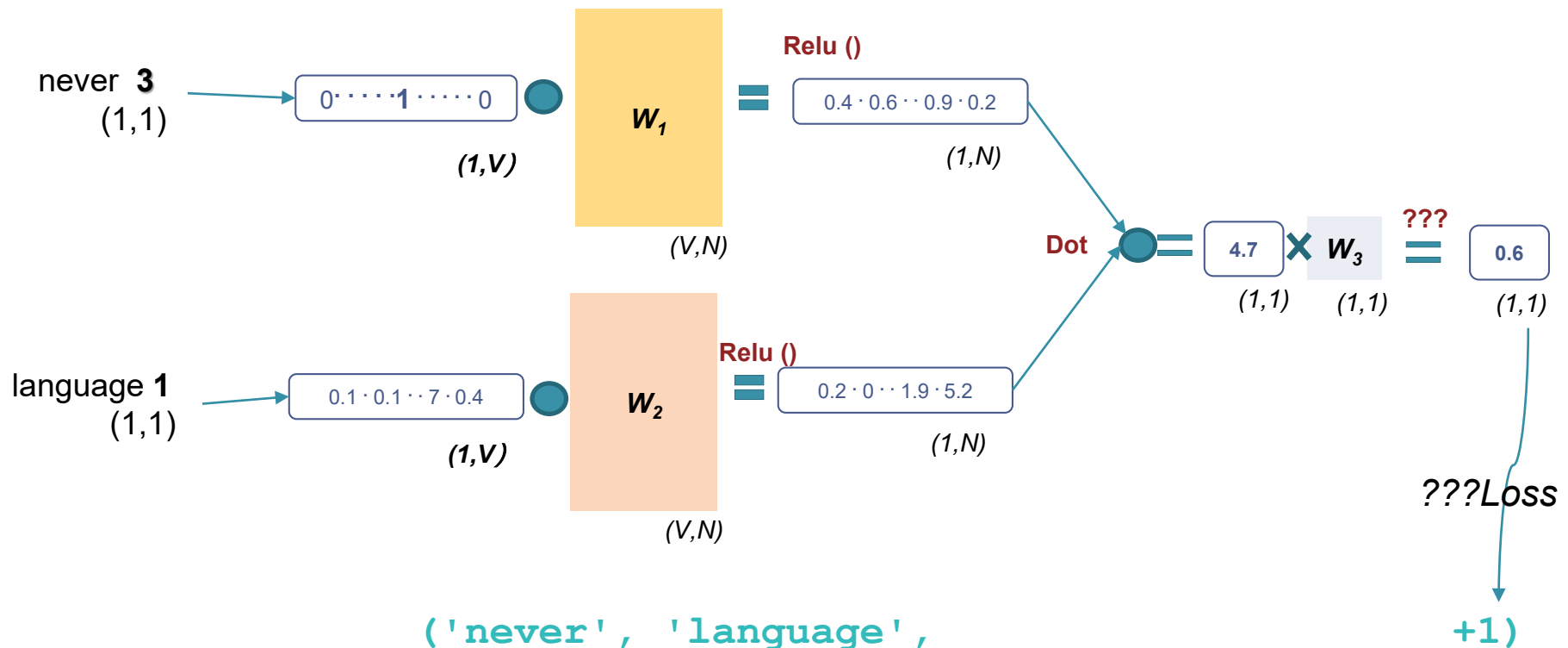
aka.  
negative  
sampling



# Word2Vec (Skipgram)

**Sentence:** language (1) users (2) never (3) choose (4) words  
(5) randomly (6), (7) and (8) language (1) is (10)  
essentially (11) non-random (12) . (13)

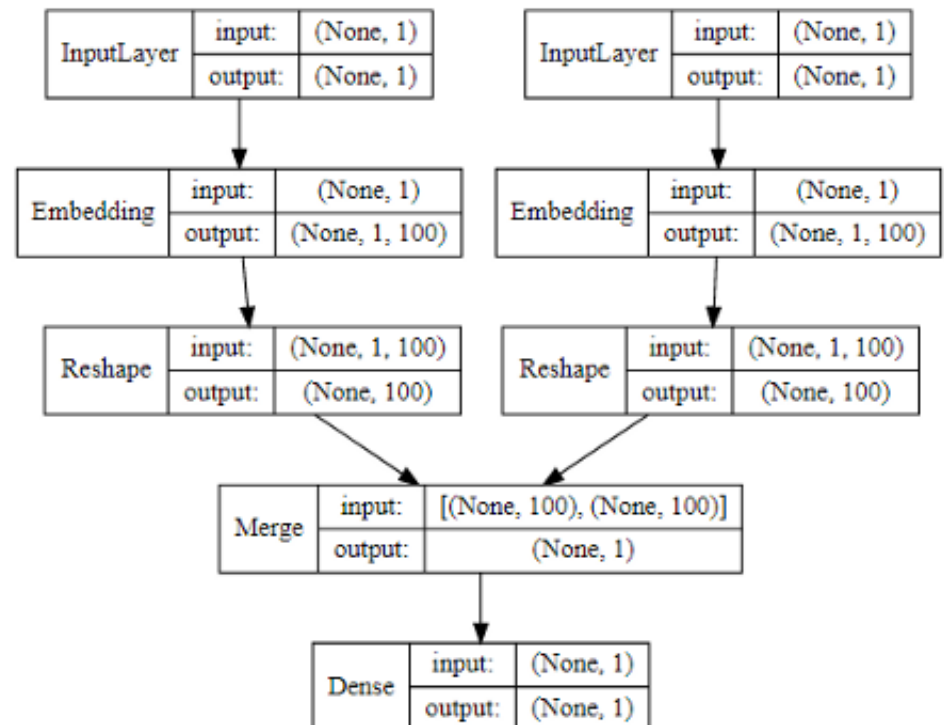
['language', 'users', 'never', 'choose', 'words']



# Word2Vec (Skipgram)

- Skipgram Model

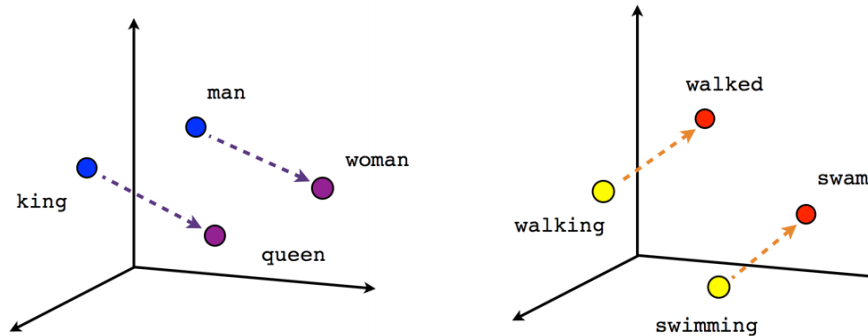
Layer (type)	Output Shape	Param #
merge_2 (Merge)	(None, 1)	0
dense_3 (Dense)	(None, 1)	2
Total params: 2,485,002		
Trainable params: 2,485,002		
Non-trainable params: 0		



# How to Choose Context?

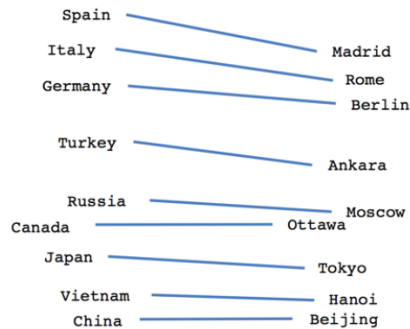
- **Different contexts lead to different embeddings**
- **Small context window:** more syntax related
  - *I like...*
  - *She likes ...*
- **Large context window:** more semantics related
  - ***stackoverflow** great **website** for **programmers***

# Properties of Word Embeddings



Male-Female

Verb tense



Country-Capital

## Ingredients

<b>Corpus of text</b>	As large as possible
<b>Annotations</b>	0
<b>Initialize weights (aka Embeddings)</b>	1x per word
<b>Deep Learning Model</b>	1x
<b>Cost Function</b>	Appropriately
<b>GPU</b>	Lotsa of it

# When to use pre-trained embeddings?

- Generally, when you don't have much training/annotated data
- **Useful:** Use as inputs to model for classification task, e.g. tagging, parsing, ranking (based on similarity)
- **Less Useful:** Machine Translation / Sequence generating tasks
- **Not Useful:** Generic Language Modeling, for those, we have sentence embeddings...

# Limitations

- **Sensitive to “tokens”** (cat vs cats)
- **Inconsistent across space**, embeddings for the same words trained with different data are different
- **Can encode bias** (stereotypical gender roles, racial bias)
- **Not interpretable**



# Embedding Bias

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{king}} - \vec{\text{queen}}$$

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}$$

	“He” Occupations	“She” Occupations
Cosine Similarity	[“retired”, “doctor”, “teacher”, “student”, “miller”, “assistant”, “lawyer”, “baker”, “judge”, “governor”, “butler”]	[“doctor”, “ <b>teacher</b> ”, “ <b>nurse</b> ”, “actress”, “student”, “miller”, “reporter”, “retired”, “lawyer”, “actor”, “artist”]
Inner Product Similarity	[“cleric”, “photographer”, “ <b>skipper</b> ”, “chaplain”, “accountant”, “inspector”, “rector”, “investigator”, “psychologist”, “treasurer”, “supervisor”]	[“ <b>librarian</b> ”, “ <b>housekeeper</b> ”, “nanny”, “accountant”, “sheriff”, “envoy”, “tutor”, “salesman”, “butler”, “footballer”, “solicitor”]

# Summary Word2Vec

## Steps

1. Define task that we want to predict
2. Go through each sentence and create the task's in-/outputs
3. Iterate through task's I/O, put the inputs through the embeddings and models to create predictions
4. Measure cost of the predicted and expected output
5. Update embedding weights accordingly (\*backprop)
6. Repeat Step 3-5 until desired.

Neither **GloVe** or **Word2Vec** has been shown to provide definitively better results rather they should both be evaluated for a given dataset.

# Train Word2Vec from “scratch”

**CBOW AND SKIPGRAM**

# References

1. Role, Franois, and Mohamed Nadif. "Handling the impact of low frequency events on co-occurrence based measures of word similarity." Proceedings of the International Conference on Knowledge Discovery and Information Retrieval (KDIR-2011). Scitepress. 2011.
2. Yogan, Jaya Kumar, et al. "A review on automatic text summarization approaches." Journal of Computer Science 12.4 (2016): 178-190
3. Peters, Matthew E., et al. "Deep contextualized word representations." arXiv preprint arXiv:1802.05365 (2018)
4. Radford, Alec, et al. "Improving language understanding by generative pre-training." URL [https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language understanding paper. pdf](https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language%20understanding%20paper.pdf) (2018).
5. Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.
6. Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." arXiv preprint arXiv:1801.06146 (2018).
7. Pennington, Jeffrey, Richard Socher, and Christopher Manning. "Glove: Global vectors for word representation." Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014.