# FOLT Lecture 3, NNs for NLP

## Intro

- A language model (LM) measures the probability of a text sequence using the chain rule

$$P(w_1 w_2 \dots w_T) = \prod_{i}^{T} P(w_i | w_{<i})$$

After (3-gram)

$$P(\text{I saw a cat on a mat}) =$$

P(I) $\longrightarrow$ P(I)
· P(saw | I) $\longrightarrow$ · P(saw | I)
· P(a | I saw) $\longrightarrow$ · P(a | I saw)
· P(cat | ~~I~~ saw a) $\longrightarrow$ · P(cat | saw a)
· P(on | ~~I saw~~ a cat) $\longrightarrow$ · P(on | a cat)
· P(a | ~~I saw a~~ cat on) $\longrightarrow$ · P(a | cat on)
· P(mat | ~~I saw a cat~~ on a) $\longrightarrow$ · P(mat | on a)

ignore      use

$$P(w_T | w_1 w_2 \dots w_{T-1}) \approx P(w_T | \underbrace{w_{T-n+1} \dots w_{T-1}})$$

all previous words          $n-1$ previous words

$$P(w_T | w_{<T}) \approx \frac{N(w_{T-n+1} \dots w_{T-1} w_T)}{N(w_{T-n+1} \dots w_{T-1})} = \frac{N(n\text{-gram})}{N((n-1)\text{-gram})}$$

We usually measure the performance of LMs with the intrinsic evaluation method: Perplexity
exp: How hard is recognizing (30,000) names at random? Perplexity = 30,000
We previously talked about N-gram language models some of their Pros

- Easy to understand
- Cheap (with modern hardware)
- Fine in some applications and when training data is scarce
  And some of their cons :
- Assume the vocabulary is known
- Sparsity problems
- Zero count
- Cannot use long context ( large n(-gram) leads to higher sparsity)
  this last point is the challenge to when we have to model Semantic similarity

## Semantic word Similarity

Two words with $N(w_1) >> N(w_2)$ $(N(cat) >> N(kitten))$
Can $P(cat|\text{saw a})$ give us information about $P(kitten|\text{saw a})$ ?
Not in N-gram models
But in Neural LMs Yes

## PART1: Neural Network Overview

State-of-the-Art NLP Methods are :

- Deep learning approaches(• (often) free from linguistic features • very large neural models • pre-training over large raw text • works well with unstructured data • deliver high-quality results)
- Neural network architectures(• LSTM (Long Short-Term Memory) • GRU (Gated Recurrent Unit) • CNN (Convolutional Neural Networks) • Transformers)
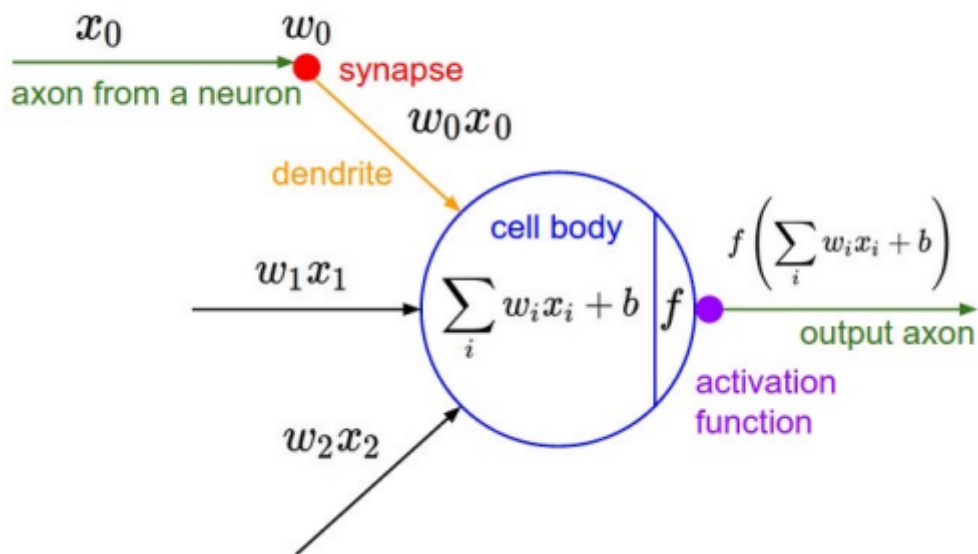  Neural networks: originally inspired by modelling biological neural systems

## Neural networks and Deep Learning Pros and Cons

### Pros

● State-of-the-art performance (current best performance) ● Architecture easily adaptable to multiple problems, e.g., Transformers ● Reduce need for feature engineering, e.g., counting n-grams

● Require a lot of cleaned training data and computational resources (E.g., Reinforcement learning from human feedback (RLHF) à ChatGPT, GPT-4)

● Long training time with a lot of Hardware and multiple rounds of human feedback

● Interpretability is a challenge (DL is a Blackbox)
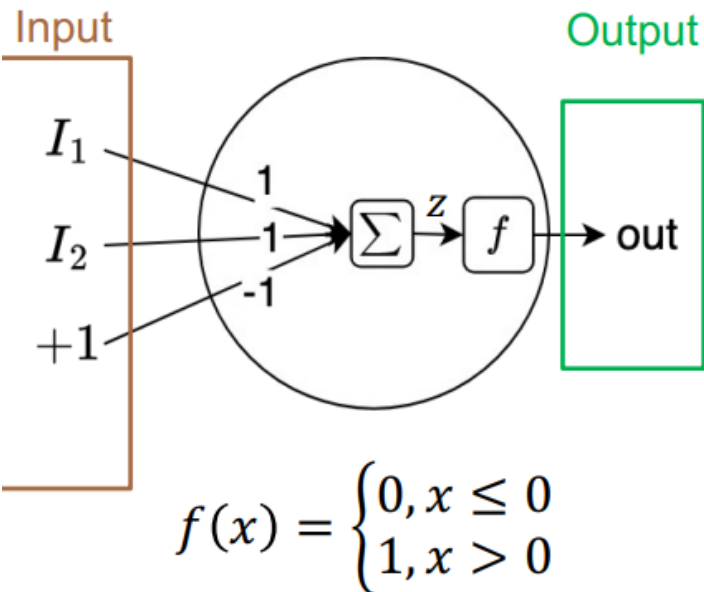
# PART2: Perceptrons



## 1) Let's try to compute AND

| AND | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

Input    Output



$$f(x) = \begin{cases} 0, x \leq 0 \\ 1, x > 0 \end{cases}$$

| AND | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

Predict output using a perceptron

o  We feed input into the perceptron

o  Try to predict output

$$f(x) = \begin{cases} 0, x \le 0 \\ 1, x > 0 \end{cases}$$

**AND**

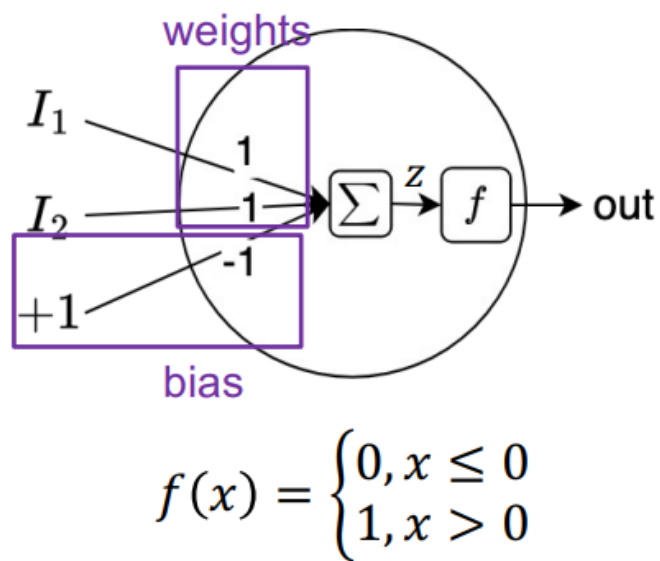| I₁ | I₂ | out |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

Predict output using a perceptron:

Parameters $\theta$ of a perceptron:

weights and biases

○ We multiply input by weights

○ Take sum and add the bias

$$z = I_1 \times 1 + I_2 \times 1 + 1 \times -1$$
$$= 0 \times 1 + 0 \times 1 + 1 \times -1$$
$$= -1$$



$$f(z) = \begin{cases} 0, z \le 0 \\ 1, z > 0 \end{cases}$$

activation function

**AND**

| I₁ | I₂ | out |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

Predict output using a perceptron:

Parameters $\theta$ of a perceptron:

weights and biases

○ We get $z = -1$

○ Then we feed $z$ into

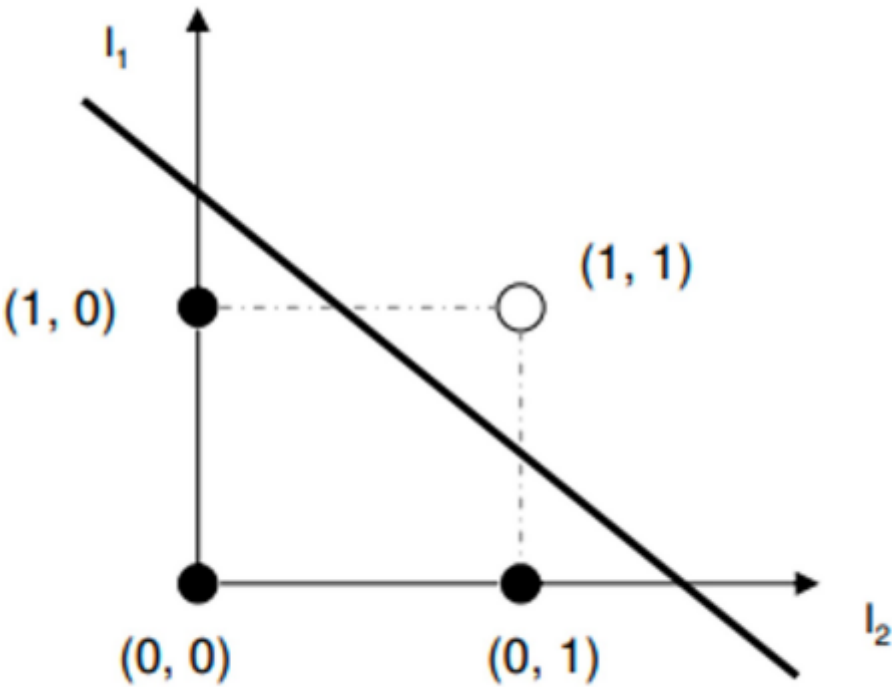the activation function $f(z)$

○ Because $z = -1 < 0, f(z) = 0 =$ out

**AND**

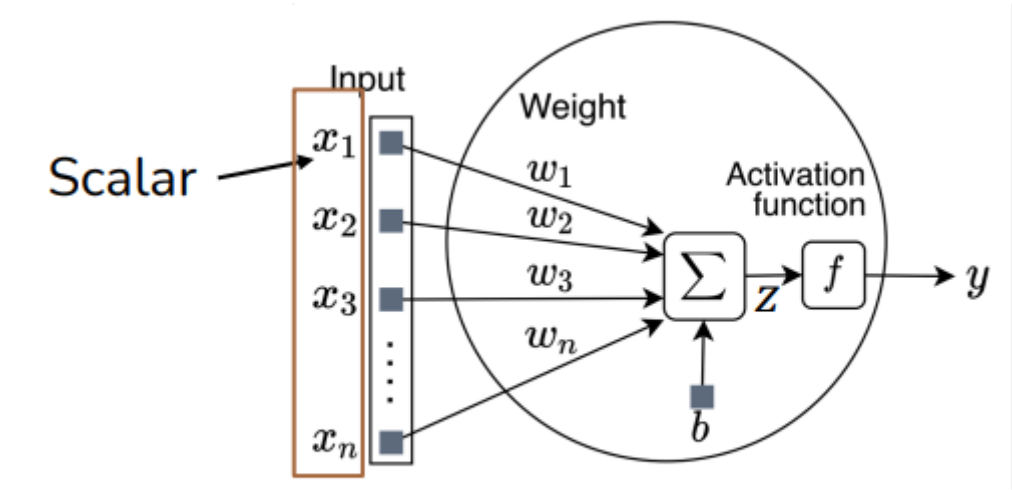| I₁ | I₂ | out |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

$$f(z) = f(I_1 \times 1 + I_2 \times 1 + 1 \times -1)$$
$$= f(0 \times 1 + 1 \times 1 + 1 \times -1)$$
$$= f(0) = 0 = \text{out}$$

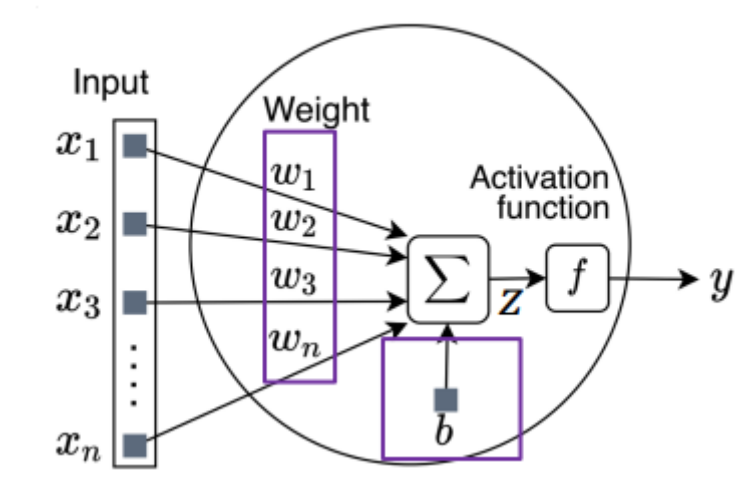rest of input output combinations computed similarly and we obtain:



## 2) Perceptrons

o Classic form



- Input $x$: usually a real value vector
Each element is a scalar = real number • E.g., input is a vector with $n$ dimensions
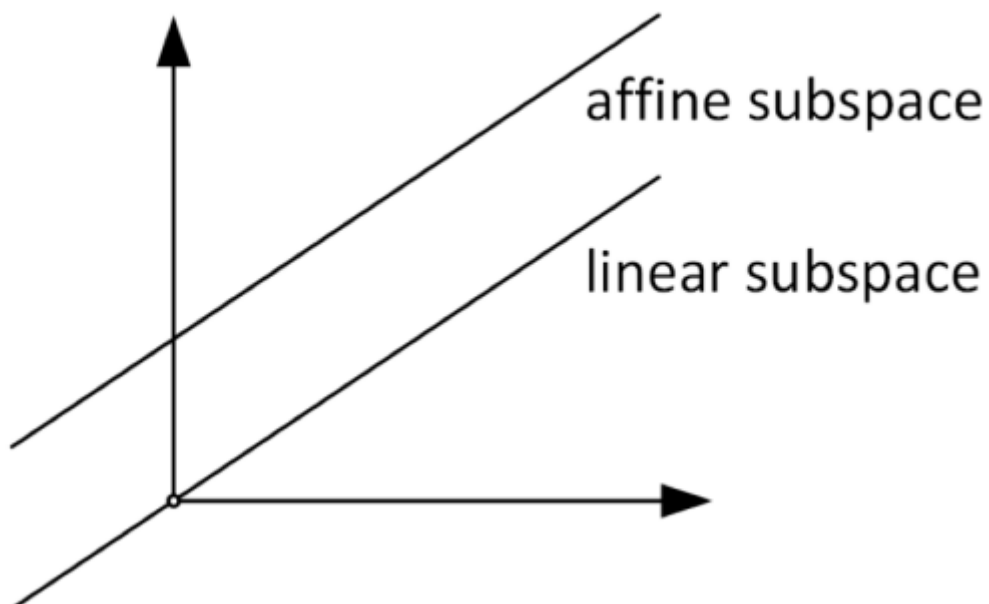


- Parameters $\theta$ :
-Weights $w = w_1, w_2 \ldots, w_n$ (please don't confuse with words input x is our "words") Each weight associated with each input indicate each input's "importance" • E.g., $n$ scalar inputs, $n$ weight values or weight vector of $n$
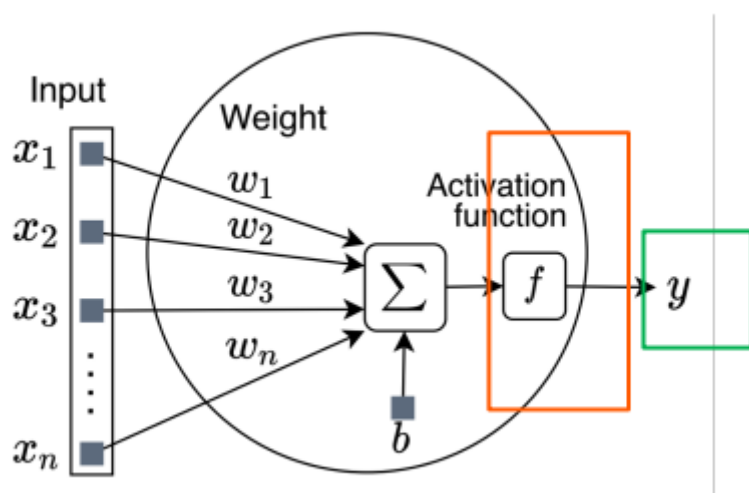-Bias $b$ : a scalar

- Intermediate output $z = (\sum_i^n w_i x_i) + b$:
  - **Weighted** summation of the inputs
    - E.g., $z = I_1 \times 1 + I_2 \times 1$
  - With an additional **bias**
    - E.g., $z = I_1 \times 1 + I_2 \times 1 + 1 \times -1$

**Why Using Bias?**

o A neuron with weights only must fix the origin to zero.

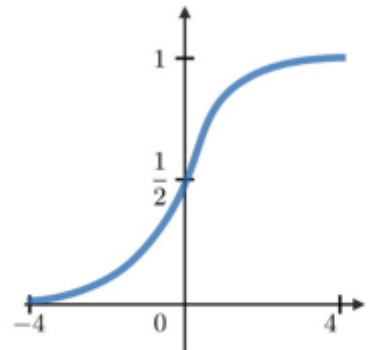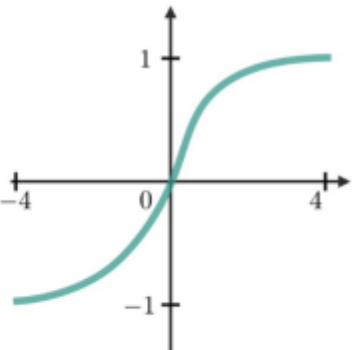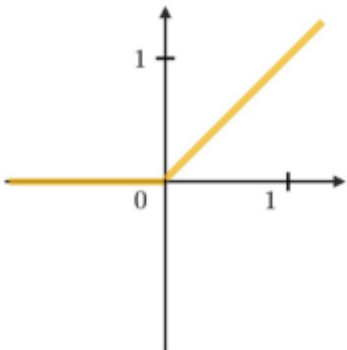o The bias allows for modelling the set of **affine** functions, which is a superset of **linear** functions.
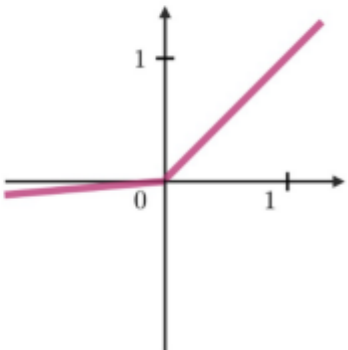


- Then, pass $z$ through an activation function $f$
  - E.g.,
    $$f(z) = \begin{cases} 0, z \leq 0 \\ 1, z > 0 \end{cases}$$
  - Controls the value range of the output y

**3)Activation function $f(.)$**

o Identity (linear function): $f(z) = z$

o Activation functions are used introduce non-linear complexities to a model

| Sigmoid | Tanh | ReLU | Leaky ReLU |
|---|---|---|---|
| $g(z) = \dfrac{1}{1+e^{-z}}$ | $g(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | $g(z) = \max(0, z)$ | $g(z) = \max(\epsilon z, z)$ with $\epsilon \ll 1$ |



## Example

Compute "AND"

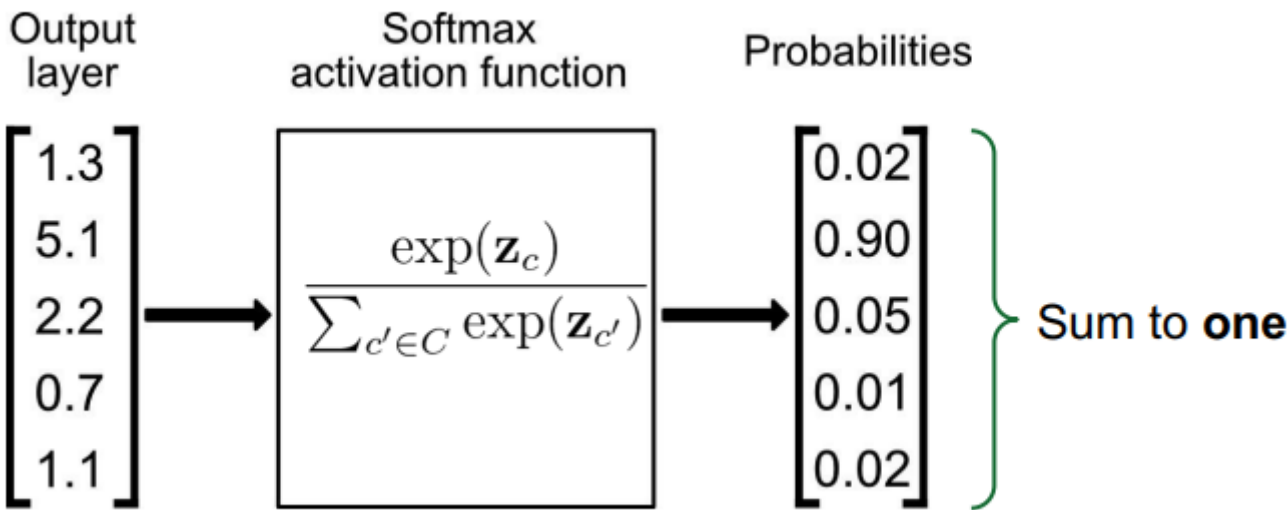$$f(z) = \begin{cases} 0, z \leq 0 \\ 1, z > 0 \end{cases}$$

$$f(z) = \begin{cases} 0, z \leq 0 \\ z, z > 0 \end{cases}$$
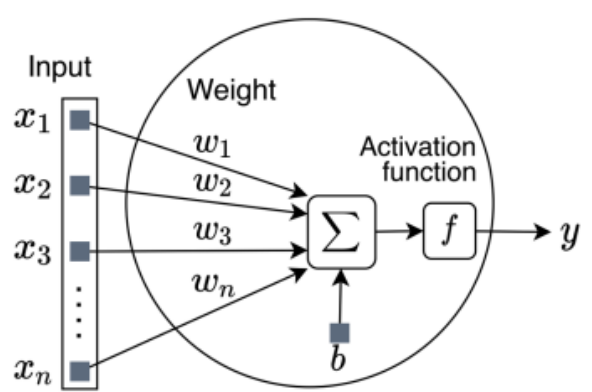$$= \max(0, z)$$
$$= \text{ReLU}$$

## 1a)Softmax

**Take a vector and compute a probability distribution out of the vector**
○ Each score resides in [0, **1**]
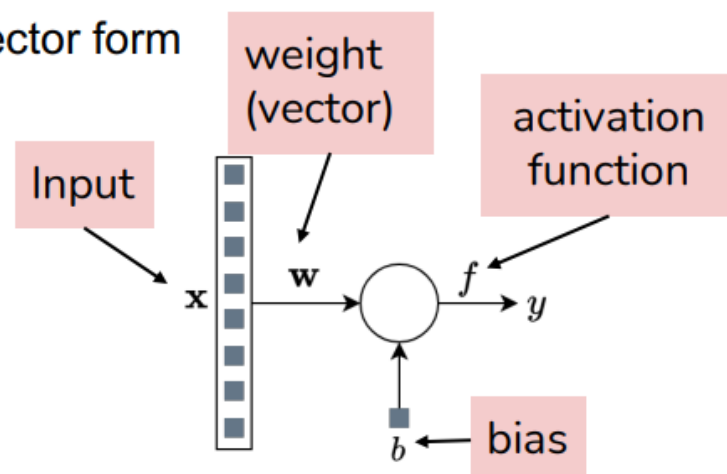● Act on a layer, not individual node



## Vector form of a perceptron

**

○ Classic form



$$y = f\left(\left(\sum_{i}^{n} w_i x_i\right) + b\right)$$

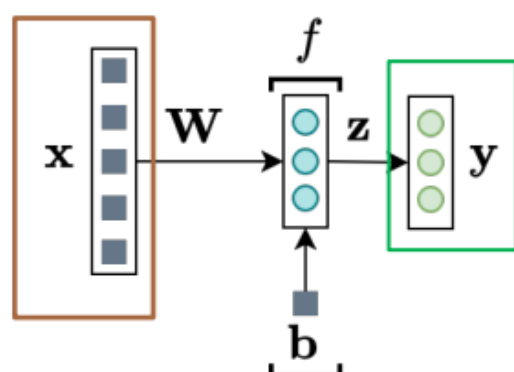○ Vector form



weight (vector)
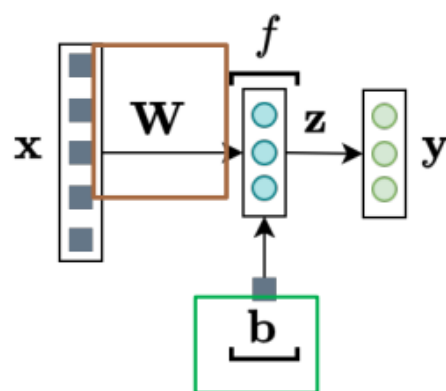
activation function

Input

bias

$$y = f(\mathbf{w}^\top \mathbf{x} + b)$$

Input: $\mathbf{x} \in \mathbb{R}^n$, Output: $y \in \mathbb{R}$
Parameters $\theta$: Weights: $\mathbf{w} \in \mathbb{R}^n$, bias: $b \in \mathbb{R}$

## 2) Perceptron : multiple Outputs



One-layer perceptron
or one-layer neural network

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

○ Input: $\mathbf{x} \in \mathbb{R}^n$ → vector
○ Output: $\mathbf{y} \in \mathbb{R}^m$ → vector



One-layer perceptron
or one-layer neural network

$$\mathbf{y} = \qquad \mathbf{xW} + \mathbf{b}$$

○ Input: $\mathbf{x} \in \mathbb{R}^n$ → vector
○ Output: $\mathbf{y} \in \mathbb{R}^m$ → vector

Parameters $\theta$:
○ Weights: $\mathbf{W} \in \mathbb{R}^{n \times m}$ → matrix
○ bias: $\mathbf{b} \in \mathbb{R}^m$ → vector

$$\mathbf{y} = f(\mathbf{z}) = f(\mathbf{xW} + \mathbf{b})$$

One-layer perceptron

or one-layer neural network

- o Input: $\mathbf{x} \in \mathbb{R}^n$ → vector
- o Output: $\mathbf{y} \in \mathbb{R}^m$ → vector

Parameters $\theta$:

- o Weights: $\mathbf{W} \in \mathbb{R}^{n \times m}$ → matrix
- o bias: $\mathbf{b} \in \mathbb{R}^m$ → vector

Activation function $f(\cdot)$: element-wise

- o Apply to each element of its input
- o E.g., each element in the intermediate output $\mathbf{z} \in \mathbb{R}^m$

$$\mathbf{y} = f(\mathbf{z}) = f(\mathbf{xW} + \mathbf{b})$$

One-layer perceptron

or one-layer neural network

Example:**

$$n = 3 \qquad m = 2$$

Example: $y = xW + b$, $d_{in} = 3$, $d_{out} = 2$

$$\begin{pmatrix} x_{[1]} & x_{[2]} & x_{[3]} \end{pmatrix} \begin{pmatrix} W_{[1,1]} & W_{[1,2]} \\ W_{[2,1]} & W_{[2,2]} \\ W_{[3,1]} & W_{[3,2]} \end{pmatrix} + \begin{pmatrix} b_{[1]} & b_{[2]} \end{pmatrix} = \begin{pmatrix} y_{[1]} & y_{[2]} \end{pmatrix}$$

*Ivan H., Deep learning for NLP, TUDarmstadt*

$$\mathbf{y} = f(\mathbf{z}) = f(\mathbf{xW} + \mathbf{b})$$

- o Input: $\mathbf{x}=[1, 0.5, -2] \in \mathbb{R}^3$ → vector with 3 dimensions
- o Output: $\mathbf{y} \in \mathbb{R}^2$ → vector

Parameters $\theta$:

- o Weights: $\mathbf{W} = \begin{bmatrix} 1 & 0 \\ -1 & -1 \\ 10 & -9 \end{bmatrix} \in \mathbb{R}^{3 \times 2}$ → matrix
- o bias: $\mathbf{b} = [0,0] \in \mathbb{R}^2$ → vector zeros → ignore

Activation function $f(\cdot)$: element-wise

- o Assume $f(\cdot)$ is identity → ignore

$$y = xW$$

$$= [1, 0.5, -2] \begin{bmatrix} 1 & 0 \\ -1 & -1 \\ 10 & -9 \end{bmatrix}$$

$$= [\ 1 \times 1 + 0.5 \times -1 + -2 \times 10, \ 1 \times 0 + 0.5 \times -1 + -2 \times -9]$$

$$= [-19.5, 17.5]$$

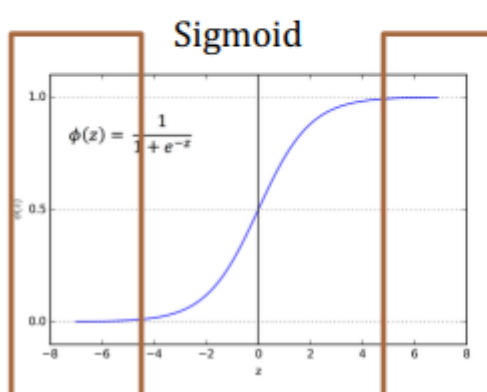If activation function is Sigmoid $= \dfrac{1}{1+e^{-x}}$

$$y = f(xW)$$

$$= f([1, 0.5, -2] \begin{bmatrix} 1 & 0 \\ -1 & -1 \\ 10 & -9 \end{bmatrix})$$

$$= f([\ 1 \times 1 + 0.5 \times -1 + -2 \times 10, \ 1 \times 0 + 0.5 \times -1 + -2 \times -9])$$

$$= f([-19.5, 17.5]) = \text{Sigmoid}([-19.5, 17.5]) = [3.39e\text{-}09, 0.99] \approx [0, 1]$$



Sigmoid

$\phi(z) = \dfrac{1}{1+e^{-z}}$

## 4) Multi-layer Perceptron

### Non-linearly Separable Problems

● Single-layer perceptron cannot solve non-linearly separable problems
○ Such as the most famous simple example is the boolean XOR operator



| XOR | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

## 1)XOR from OR and NOT AND

OR

NOT AND

XOR

| XOR | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

Compute OR



(1, 0)   (1, 1)

?

(0, 0)   (0, 1)

---

OR

NOT AND

XOR

| XOR | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

Compute OR



(1, 0)   (1, 1)

?

(0, 0)   (0, 1)

---

OR

NOT AND

XOR



| XOR | | |
|---|---|---|
| $I_1$ | $I_2$ | out |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$x_1 = I_1,$
$x_2 = I_2$

$v_1 = f(-1.5 + 2x_1 - x_2)$
$v_2 = f(-1.5 + 2x_2 - x_1)$

$f(z) = \begin{cases} 0, z \le 0 \\ 1, z > 0 \end{cases}$

$y = -0.5 + v_1 + v_2$

(1, 0)   (1, 1)

?

(0, 0)   (0, 1)

## 2)Multi-layer Perceptron

Two-layer neural network
or
One hidden layer neural network

$$\mathbf{z}_0 = \mathbf{x}$$

hidden layer



Two-layer neural network
or
One hidden layer neural network

$$\mathbf{z}_0 = \mathbf{x}$$
$$\mathbf{z}_1 = f_1(\mathbf{z}_0 \mathbf{W}_1 + \mathbf{b}_1)$$

Activation function

output layer



Two-layer neural network
or
One hidden layer neural network

$$\mathbf{z}_0 = \mathbf{x}$$
$$\mathbf{z}_1 = f_1(\mathbf{z}_0 \mathbf{W}_1 + \mathbf{b}_1)$$
$$\mathbf{z}_2 = f_2(\mathbf{z}_1 \mathbf{W}_2 + \mathbf{b}_2)$$

Two-layer neural network
or
One hidden layer neural network

$$\mathbf{z}_0 = \mathbf{x}$$
$$\mathbf{z}_1 = f_1(\mathbf{z}_0 \mathbf{W}_1 + \mathbf{b}_1)$$
$$\mathbf{z}_2 = f_2(\mathbf{z}_1 \mathbf{W}_2 + \mathbf{b}_2)$$

output of each layer



Two-layer neural network
or
One hidden layer neural network

$$\mathbf{z}_0 = \mathbf{x}$$
$$\mathbf{z}_1 = f_1(\mathbf{z}_0 \mathbf{W}_1 + \mathbf{b}_1)$$
$$\mathbf{z}_2 = f_2(\mathbf{z}_1 \mathbf{W}_2 + \mathbf{b}_2)$$

input of each layer
=
output of previous layer

$$\mathbf{z}_0 = \mathbf{x}$$
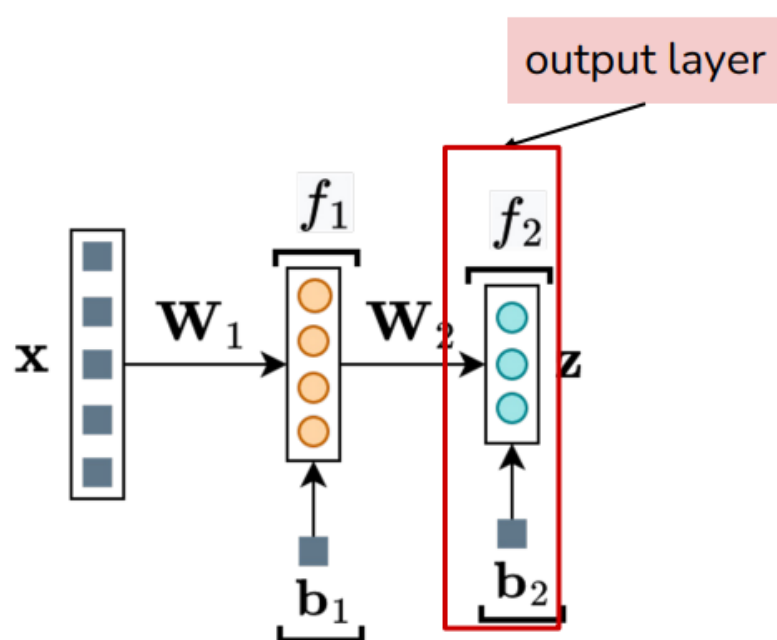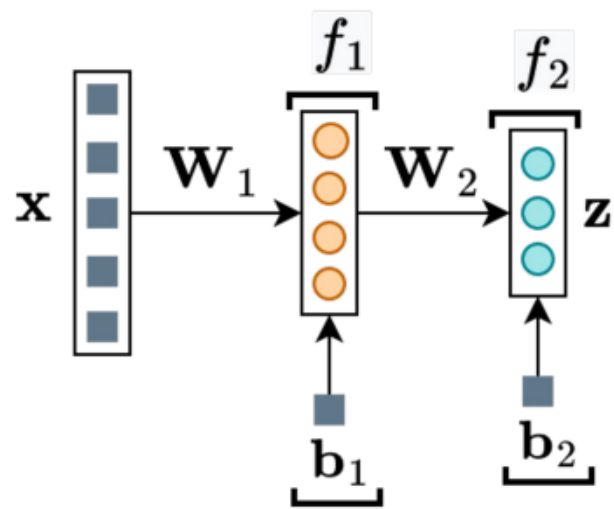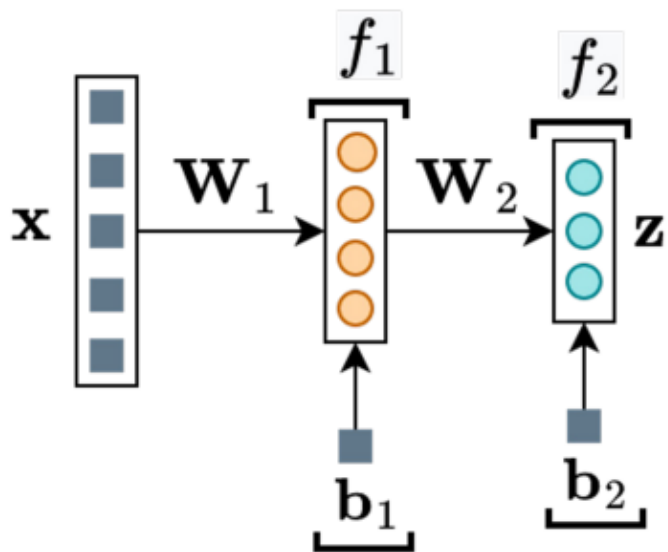
$$\mathbf{z}_1 = f_1(\mathbf{z}_0 \mathbf{W}_1 + \mathbf{b}_1)$$

$$\mathbf{z}_2 = f_2(\mathbf{z}_1 \mathbf{W}_2 + \mathbf{b}_2)$$

Two-layer neural network
or
One hidden layer neural network

input of each layer
=
output of previous layer

$$\mathbf{z}_2 = f_2(f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2)$$

hidden layer



Three-layer neural network
or
Two hidden layer neural network

## Hidden Layers

o Decides the expressiveness of the whole neuron network

o Hidden layer can project the data onto another vector space in which they are now linearly separable.



hidden layer

## Do We Need More Hidden Layers?

● In theory, two-layer neural networks (one hidden layer) can approximate any functions if we increase the number of hidden units and cover all examples .

● In practice, it is impossible to include all examples, the network can fit any training data but not generalise to new data points by reasoning from the observed data. So, we usually add more layers to increase the complexity

● To prevent linearity → approximate complex functions

$$\mathbf{z} = f_2(f_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2)$$
$$= ((\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2)$$
$$= (\mathbf{W}_1\mathbf{x})\mathbf{W}_2 + (\mathbf{b}_1\mathbf{W}_2 + \mathbf{b}_2)$$
$$= \mathbf{W}'\mathbf{x} + \mathbf{b}'$$



## Artificial Neural Networks

▪ Common neural networks



# PART3:Training a Neural Network

o How does the model get the parameters?
Training on the training data.
o Use stochastic gradient descent algorithms (more details in DL courses) to optimize
-Maximum likelihood estimate (MLE) ( Compute AND, maximize the likelihood of output 1 given input (1, 1) , Compute XOR, maximize the likelihood of output 1 given input (0, 1) or (1, 0) )
- I.e., minimize the error between the prediction of the model and the "correct" output

## 1) Optimize Parameters

○ Step 0. Initialize model
○ Step 1. Take a batch of training data
○ Step 2. Perform forward pass (inference) to obtain the prediction
o Forward pass (inference)
o Feed the input into neural networks
o Compute from input to output
o Get the prediction / output



$$z_2 = f_2(f_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2)$$

○ Step 3. Backward pass (training) from MLE
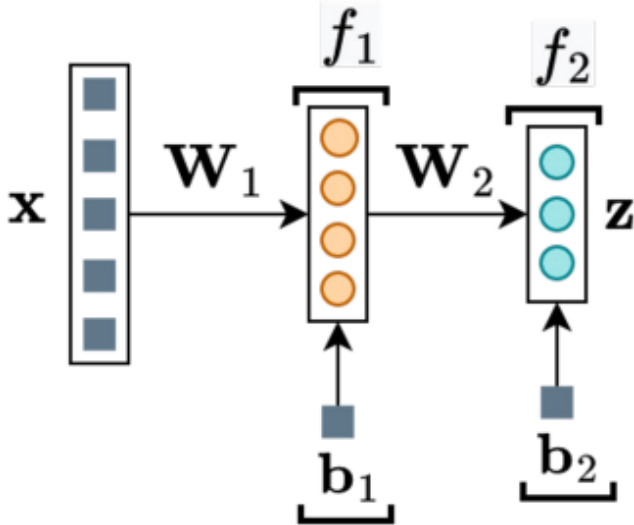o Backward pass (training)
o Compute error of the predictions

o Compute changes to parameters in each layer, backward from output to input

o Add the changes to parameters → update

○ Step 4. Update parameters

o **Epoch** or **parameter updating steps**. In the context of training, one epoch refers to one time when the model sees the entire training examples

   o Set a maximum number of epoch

o **Validation performance**. Performance evaluation on the validation set.

   o Early stopping: If the validation performance does not increase after k (parameter-)updating steps

      o k is called patience

## 2)Hyperparameters

o There are "micro-decisions" to be made for each ML (NLP) methods

   o N-gram LMs

      o N in n-gram, which can be 1, 2, 3, 4, …

      o $\alpha$ in Add-$\alpha$ smoothing, $\alpha < 1$ or $\alpha = 1$ (Laplace smoothing)

   o Neural networks:

      o Types of neural layers: most studies recently use Transformers

      o Number of neural layers, e.g., Llama-2 7B uses 32 layers

      o Activation functions, e.g. ReLU is one of the most used

      o Learning rate

# Learning Rate



| Too low | Just right | Too high |
|---|---|---|
| A small learning rate requires many updates before reaching the minimum point | The optimal learning rate swiftly reaches the minimum point | Too large of a learning rate causes drastic updates which lead to divergent behaviors |

o How to choose hyperparameters?

Try different values, and choose one using a validation dataset

Try a systematic and replicable procedure for finding the best values; report this procedure

## PART4:Neural Language Models

# Semantic Word Similarity

- Two words with $N(w_1) \gg N(w_2)$
  - where $N(w_i)$ is count of word $i$, $\gg$ denotes *much larger than*
  - E.g., $N(\text{cat}) \gg N(\text{kitten})$
- Can $P(\text{cat}|\text{saw a})$ give us information about $P(\text{kitten}|\text{saw a})$?
- N-gram language models: no
- Neural language models: yes
  - Can be extremely time-consuming (hours-days)
  - More complex / less transparent
  - Robust over longer contexts
  - Usually lower perplexity than n-gram

# Recap: Language Models

- A language model (LM) measures the probability of a text sequence. Using the chain rule, we get:

$$P(w_1 w_2 \dots w_T) = P(w_1)P(w_2|w_1) \dots P(w_T|w_1 w_2 \dots w_{T-1}) = \prod_i^T P(w_i|w_1 \dots w_{i-1}) = \prod_i^T P(w_i|w_{<i})$$

P(I saw a cat on a mat) =

P(I)
· P(saw | I)
· P(a | I saw)
· P(cat | I saw a)
· P(on | I saw a cat)
· P(a | I saw a cat on)
· P(mat | I saw a cat on a)

$$P(w_T|w_{<T}) = P(w_T|w_1 w_2 \dots w_{T-1}) = \frac{N(w_1 \dots w_{T-1} w_T)}{N(w_1 \dots w_{T-1})}$$

Instead of looking up the count of a word or n-gram, neural language models do:

o encoding the previous words

model-specific

Idea: get a vector representation for the previous words (computers work with numbers)

We can use different model architecture, e.g., RNN, CNN, Transformers

o transforming it into a distribution over the next word

Idea: get the probability distribution for the next token

get probability
distribution for
the next token

process context
(previous history)

o Formally, we compute $P(w_T|w_{<T})$ using neural network

o Input $w_{<T}$: previous words

o Expected output $P(w_T|w_{<T})$: a probability distribution for the next word $w_T$

$$P(w_T|w_{<T}) = \boxed{\text{NN}(\text{enc}(w_{<T}); \; \theta)}$$

o where $\theta$ (trained parameters or weights of a model) do

o encoding the previous words – enc($\cdot$)

o transforming it into a distribution for the next word – NN($\cdot$)



**1) Get Word Embeddings(later in course )**

**|V| tokens**

**d-sized vector**

Transform **h** linearly from size **d** to **|V|** - the vocabulary size

Linear layer

softmax

P( * | **I saw a cat on a**)

get probability distribution for the next token

Neural network

**h**: vector representation of context **I saw a cat on a**

process context (previous history)

Input word embeddings

**I saw a cat on a**

## 2)Get vector representation of context



**|V| tokens**

**d-sized vector**

Transform **h** linearly from size **d** to **|V|** - the vocabulary size

Linear layer

softmax

P( * | **I saw a cat on a**)

get probability distribution for the next token

Neural network

**h**: vector representation of context **I saw a cat on a**

process context (previous history)

Input word embeddings

**I saw a cat on a**

## 3)Predict Next Word (Token)

|V| tokens

Transform **h** linearly from size **d** to |V| - the vocabulary size

**Linear layer** → softmax → $P( * | \text{I saw a cat on a})$

get probability distribution for the next token

Neural network

**h**: vector representation of context **I saw a cat on a**

← Input word embeddings

process context (previous history)

I saw a cat on a

---

We have:
- **h** - vector of size **d**

We need:
- vector of size |V| – probabilities for all tokens in the vocabulary

Transform linearly from size **d** to size |V|

**d** features (vector dimensionality)

|V| tokens

$\mathbf{h}^\mathsf{T}$ × Linear layer → softmax

transform vector representation of context into a probability distribution

$$y = \text{softmax}(\mathbf{h}\mathbf{W} + \mathbf{b})$$

$$\mathbf{h} \in \mathbb{R}^d, \ \mathbf{W} \in \mathbb{R}^{d \times |V|}, \ \mathbf{b} \in \mathbb{R}^{|V|}$$

## Training

$$P(w_T | w_{<T}) = \text{NN}(\text{enc}(w_{<T}); \theta)$$

- o  Training: updating $\theta$ (trained parameters or weights of a model)

  - o  Using *stochastic gradient descent* algorithms to maximize likelihood

  - o  **Maximize likelihood estimate** (MLE): maximize the probability of the "correct" next word predicted by the model, i.e., **minimize error** between the model's predicted probabilities and the "correct" next word
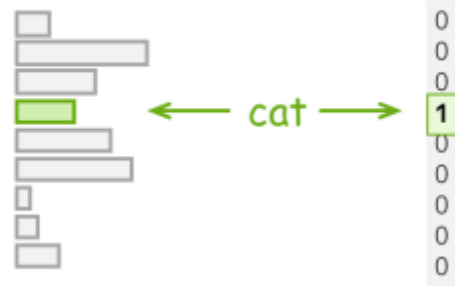
we want the model
to predict this

vord predicted by

correct" next word

Training example: **I saw a cat** on a mat <eos>

Model prediction: p( * | **I saw a**)    Target

$$\text{Loss} = -\log(p(cat)) \rightarrow \min$$

0
0
0
1
0
0
0
0

— cat —

decrease
increase

decrease

negative log-likelihood
or cross-entropy

## Inference

o  Inference: use the trained model (trained parameters/weights $\theta$) to compute $P(w_T | w_{<T})$

I was _____

P( * | I was)    get probability
distribution

| there | 0.10 |
| in | 0.09 |
| glad | 0.07 |
| happy | 0.05 |
| confused | 0.03 |
| ... | ... |
| hungry | 0.01 |
| ... | ... |

I was _____

P( * | I was)    sample from the
distribution

| there | 0.10 |
| in | 0.09 |
| glad | 0.07 |
| happy | 0.05 |
| confused | 0.03 |
| ... | ... |
| hungry | 0.01 |

I was _____

generate

| there | 0.10 |
| in | 0.09 |
| glad | 0.07 |
| happy | 0.05 |
| confused | 0.03 |
| ... | ... |
| hungry | 0.01 |

I was happy _____