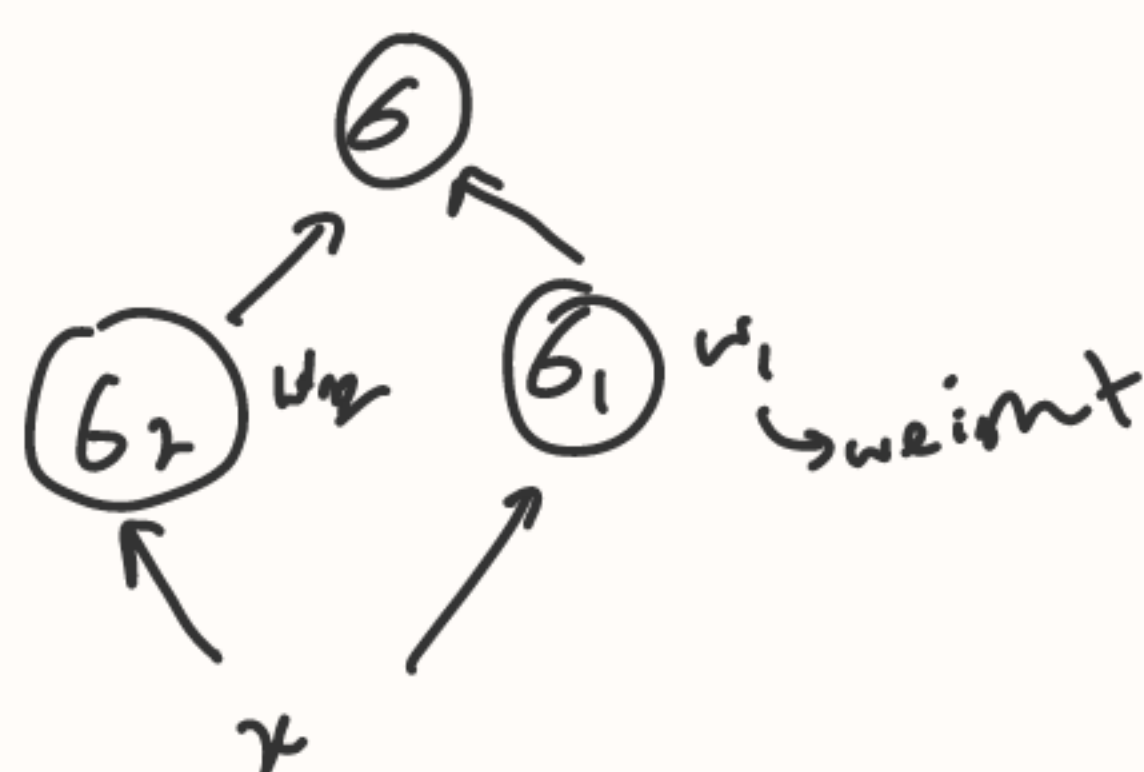
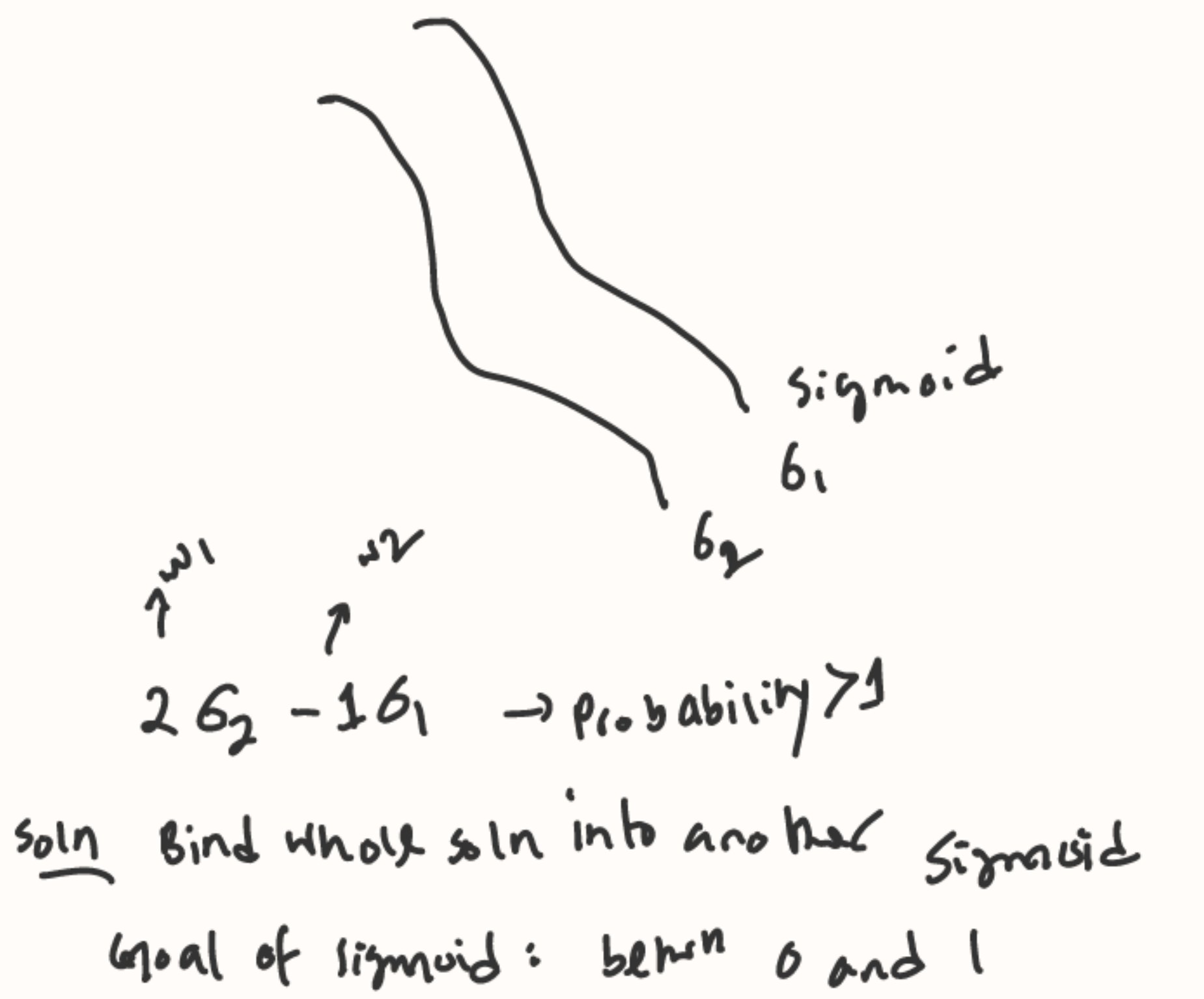


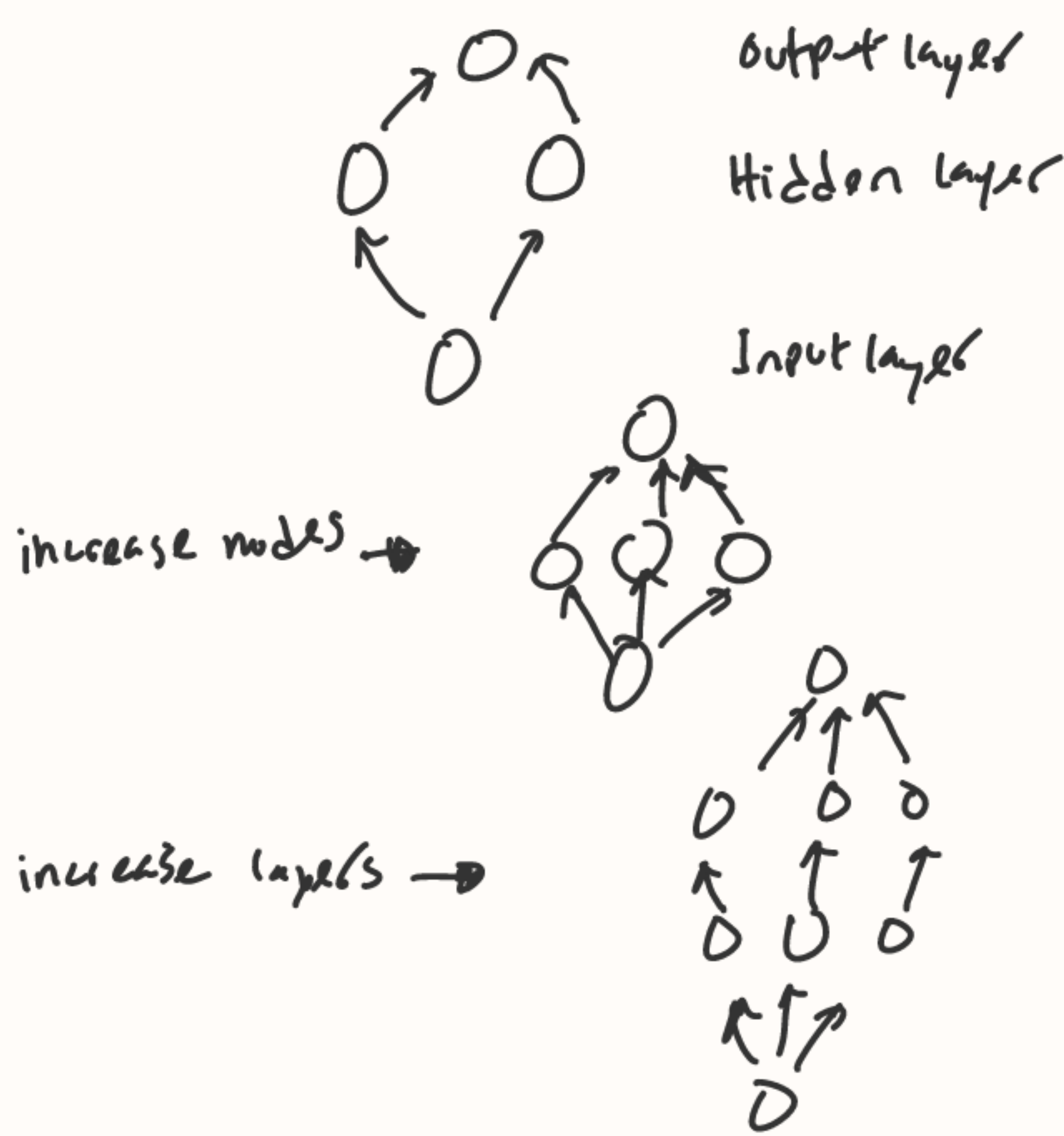
Neural Networks in NLP

Classify tweets based on no. of chars



Logistic Regr: $y = \sigma(wx+b)$

Neural network: $y =$



Fully connected - neural network needs not be

feed forward: uniform direction of calculations
NN need not be
prev weights → learn next weights

Recurrent neural networks

LR-linear prob - curve descending - global minima
Unintended to find a soln.

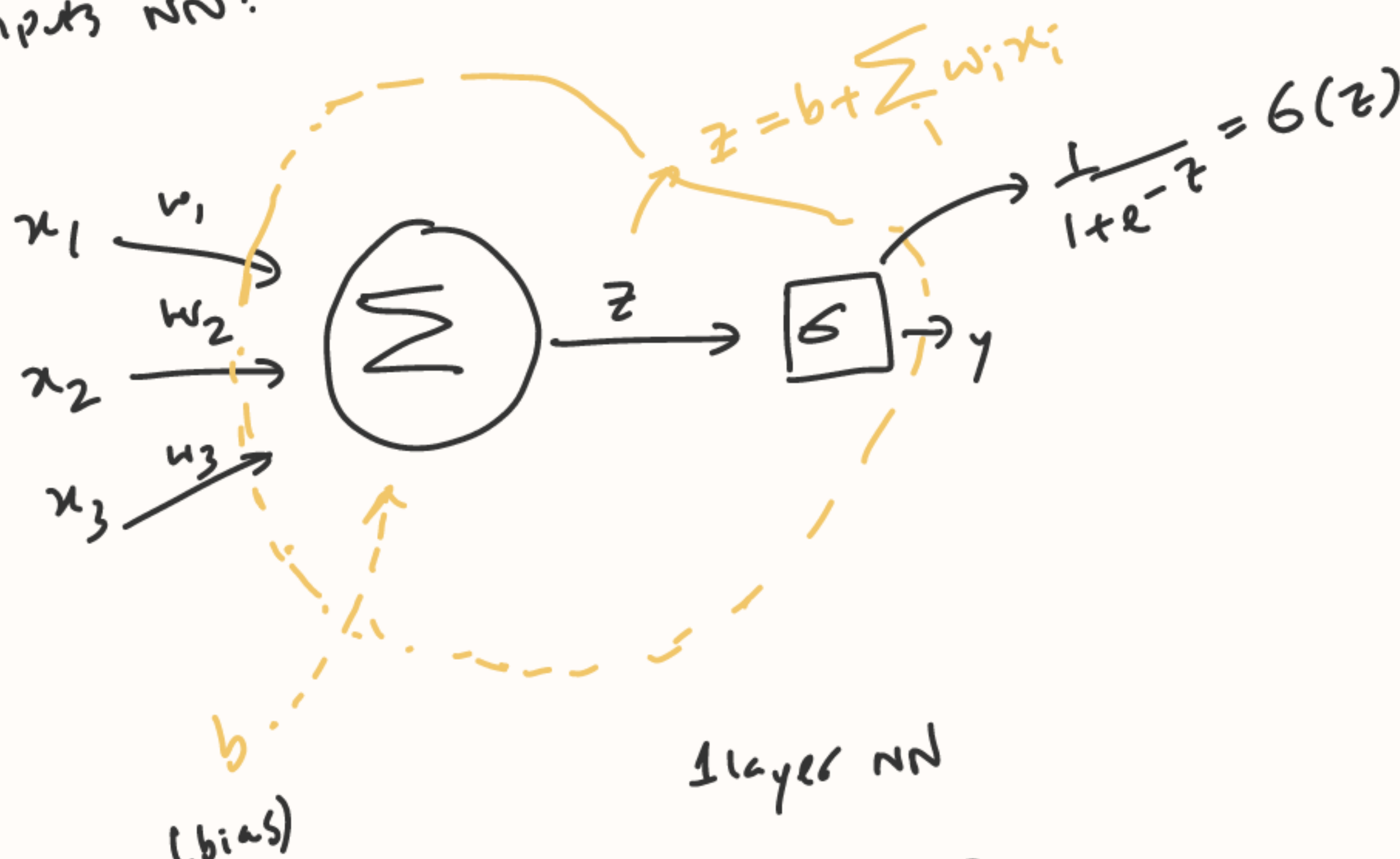
Here, instead of global minima, we have a local minima

Not guaranteed to have a soln
soln: defensive technique (to avoid local minima)

- random initialization

weight values in NN
(prev. weight value b)

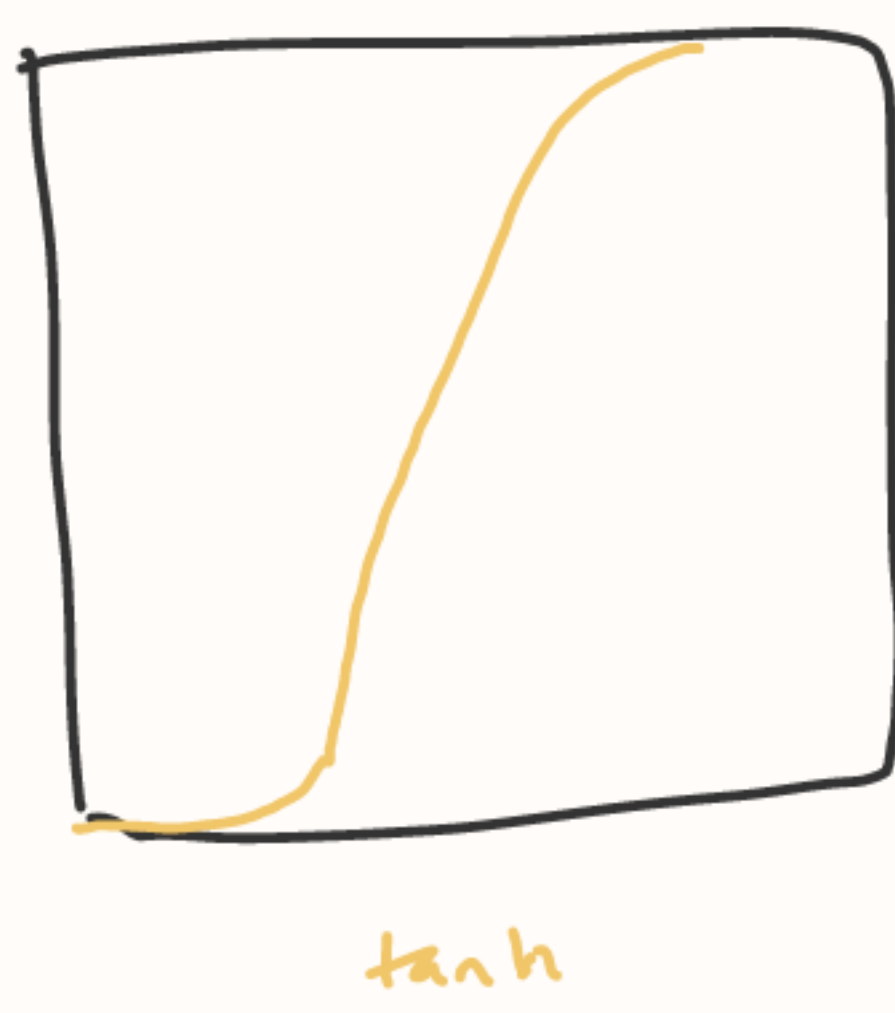
3 inputs NN:



single unit → 1 node
perceptron: outputs a binary outcome (0 or 1)
linear activation function
(sometimes non linear)

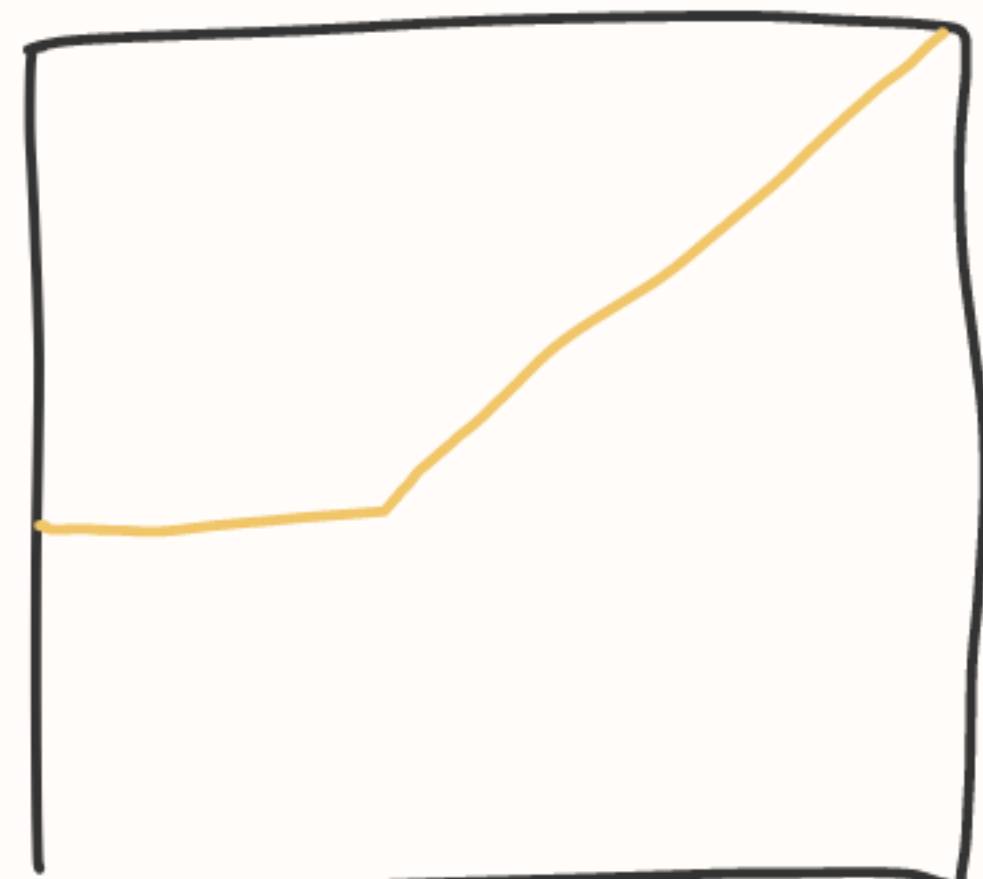
for last layer | Sigmoid activation func. → binary layer → non linear A.F.
Softmax for multinomial classification

for hidden layer | tanh
ReLU



$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

more suitable to handle outliers
(projects outliers to mean value)



$$y = \text{relu}(z) = \max(z, 0)$$

-ve value → will be 0
> 0 → that value kept
→ handle outliers
Apply sigma in tanh for high values of z → 1
vanishing gradient (gradient becomes 0) → differentiated = 0
ReLU keeps same values of z
gradient doesn't become 0

output layer \rightarrow softmax

Neural Network (NN)

tanh and ReLU

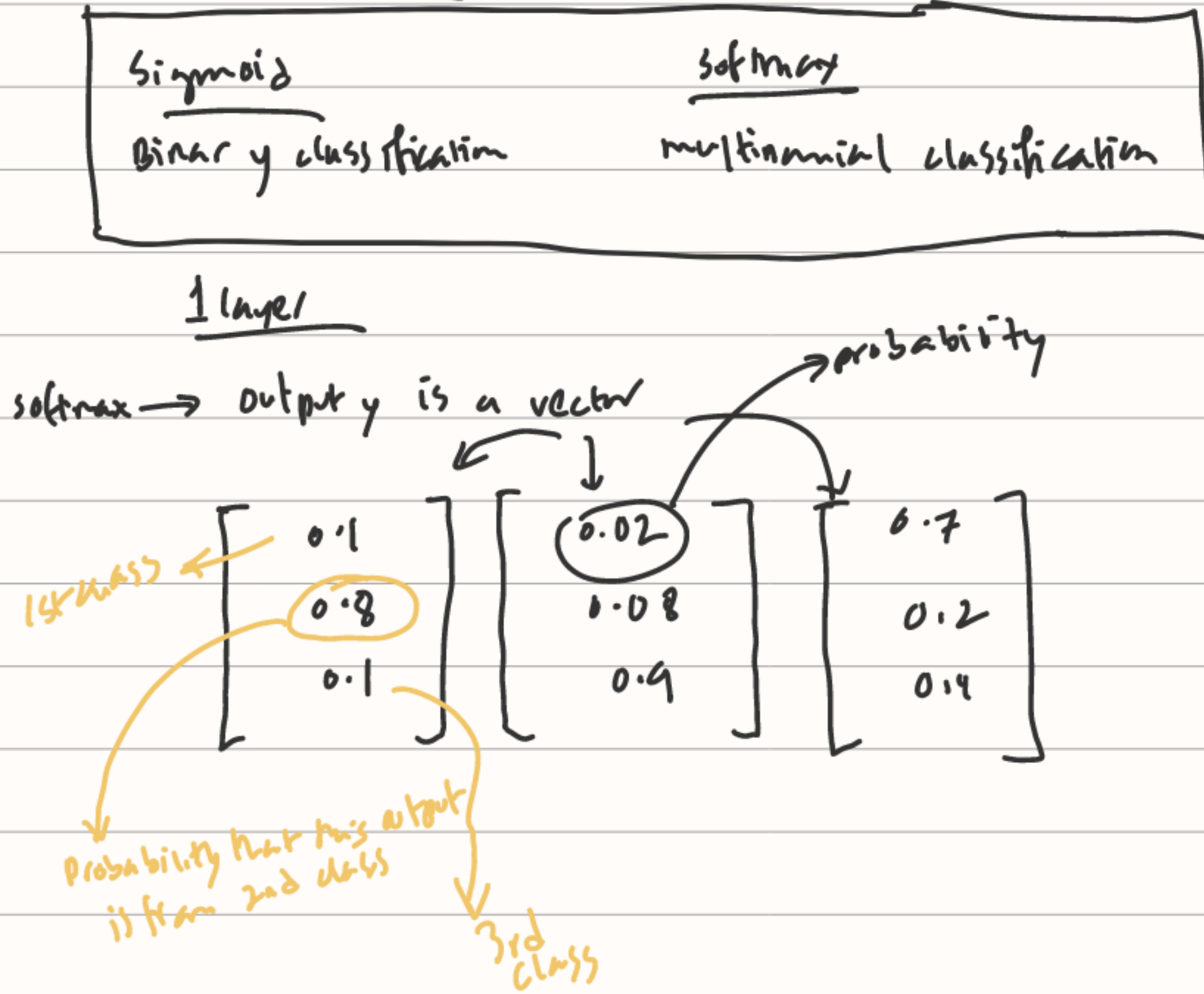
Perceptrons

XOR problem

Soln: use intermediate layers

Feed forward NN — output of one layer as input of next layer

Multinomial Logistic Regression — softmax (MLR)



2 layers

output, hidden layer

multilayer NN \approx MLR

- add as many hidden layers as possible
- not possible to define as many
- tanh or ReLU in intermediate layers
- softmax/sigmoid but use tanh or ReLU



Why do you use hidden layers?

Trying many diff combinations of features to classify

Hidden layer used for feature representation

which feature contribute more to classification

multilayer notation

$z^{[1]}$ output after 1st layer

$a^{[0]}$ input
g activation function
b bias

≡ Necessity of non linear Activation function
tanh, ReLU, softmax...

If g_1, g_2 both linear:

$$\begin{aligned}
 z^{[1]} &= w^1 x + b^1 \\
 z^{[2]} &= w^2 z^{[1]} + b^2 \\
 &= w^2 (w^1 x + b^1) + b^2 \\
 &= (w^2 w^1) x + ((w^2 b^1) + b^2) \\
 &= w^1 x + b^1
 \end{aligned}$$

basically 1st layer (no feature learning happening)

Non linear A/c activates feature representation

w/o them, no trying out combo of features

hidden layers don't learn anything

Repeating the bias unit

x features
 b bias
 $x_0 = 1$ ~~not~~ ± 1 each input last node

Text classification

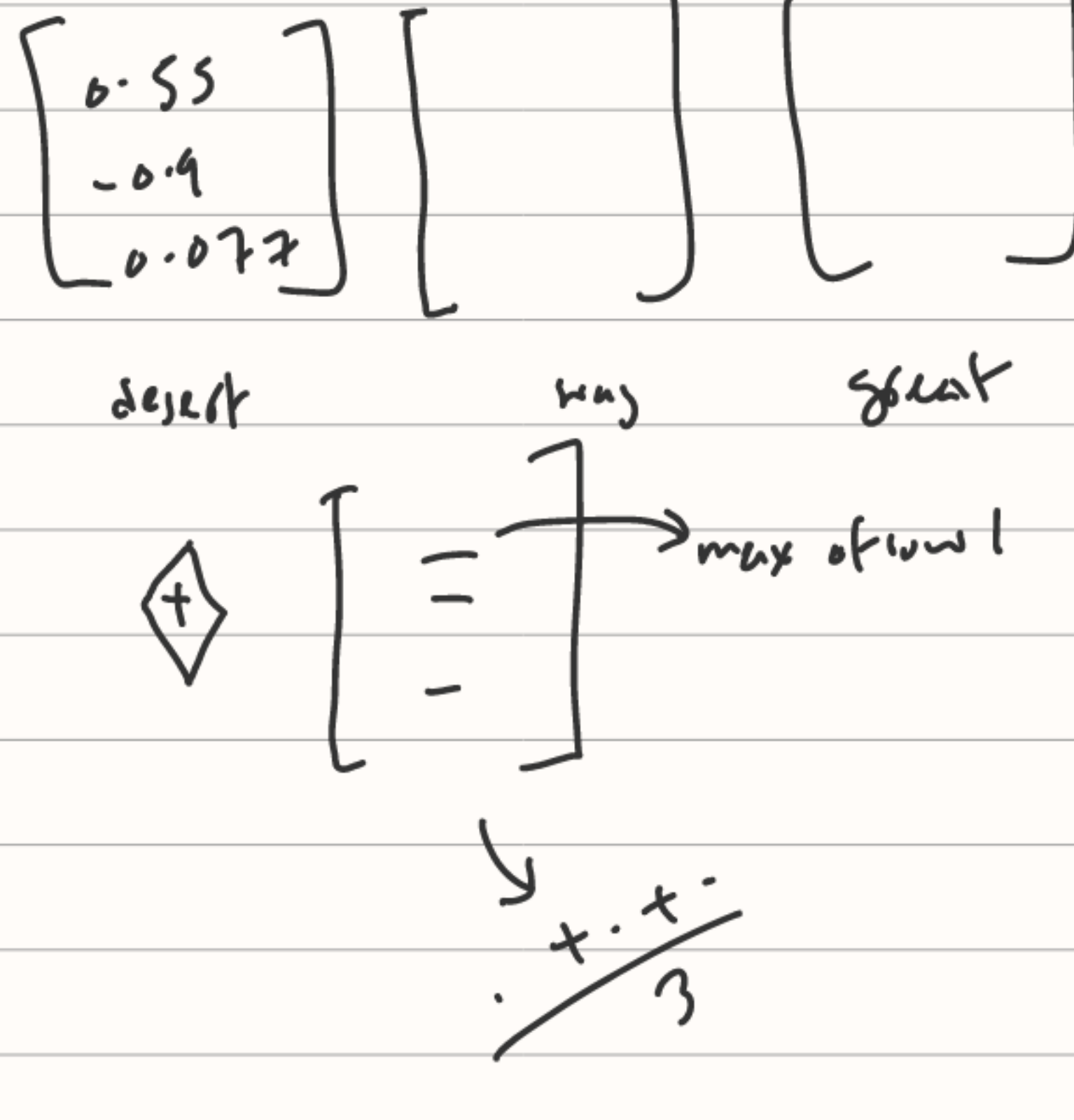
Lang modeling

≡ Why feeding handcrafted features?

we learnt embedding of each token

bert, word2vec, glove — trained embeddings — outside source

pooling: combining all embeddings, maybe learn a new embedding



which pooling method to use? — based on empirical evidence

Pre trained: Learning from an already learned embedding
transfer learning — using learning from another source

units in NN

Issue: texts come in diff sizes

Neural LM

LM Predicting next tokens based on prev tokens
↓
Lang. modeling

Transformer based LM / better than n-gram based LM

Embedding for each token from outside source

unknown sample in N-gram LM → cannot predict a word that didn't appear
(word not in training corpus) in the training sample

In neural LM, not mandatory for a word to be in training sample.
similar tokens - cat and dog.
what goes for cat, goes for dog.
can predict

training

I ate English breakfast

I had good lunch at a Chinese restaurant.

testing

Lunches are good in _____ restaurant. ✓ can predict

Lunches are good in _____ cafe. ✗ cannot predict

Neural LM knows restaurant and cafe similar. → can predict.

output layer : $z_2 = Ux a_1$

hidden layer (ReLU A.F) $a_1 = \sigma(z_1)$

Embedding layer $w \cdot e + b$

$$\begin{array}{l} \text{Eqns} \\ \left\{ \begin{array}{l} e = [Ex_{t-3}; Ex_{t-2}; Ex_{t-1}] \\ h = \sigma(w e + b) \\ z = U h \\ \hat{y} = \text{softmax}(z) \end{array} \right. \end{array}$$

noise / word not in pretrain corpus → return noise / random embedding

2 layer network

Every layer has weights

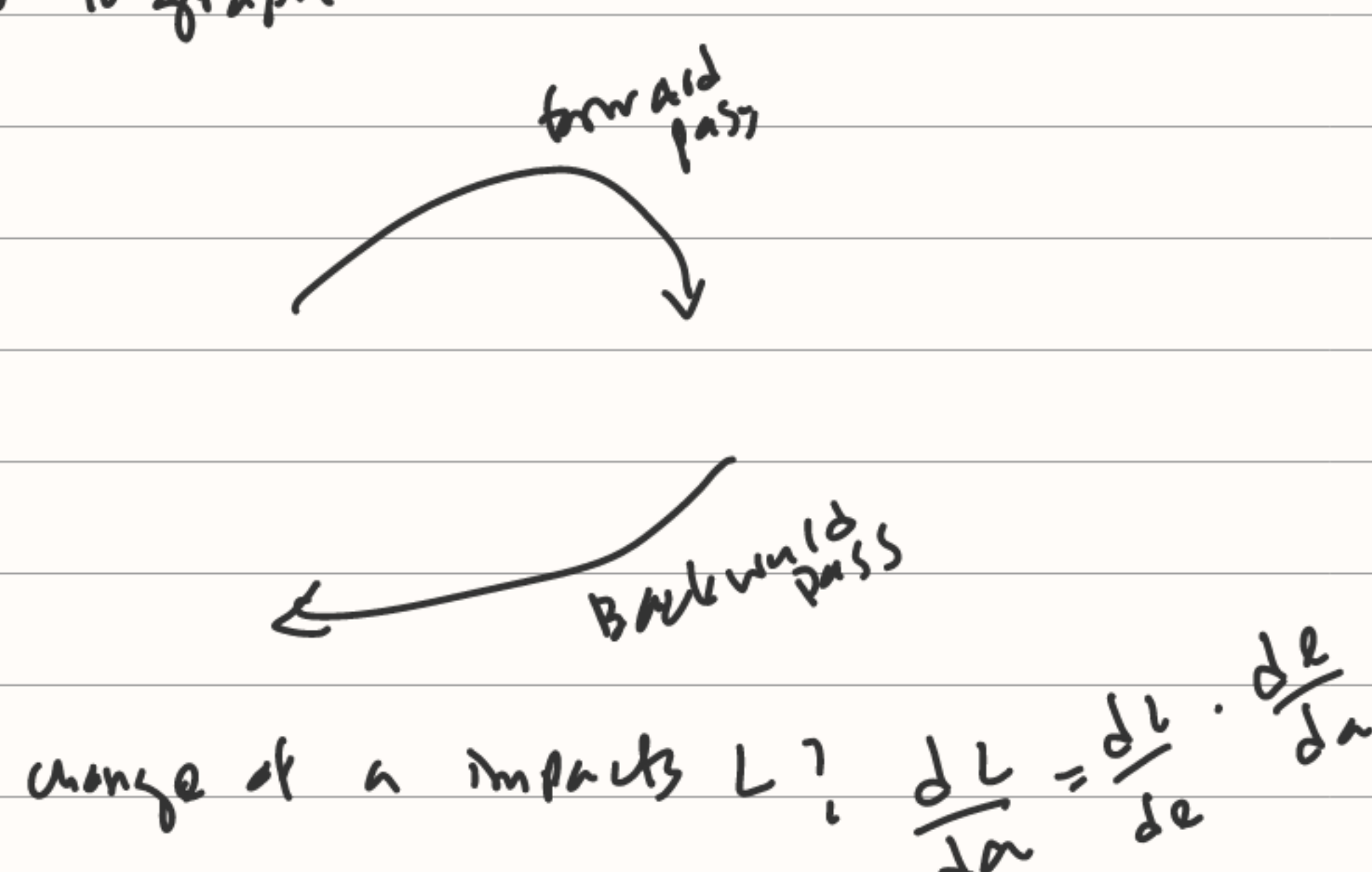
update weights based on loss function

(depends on output layer) | loss function only on last layer

Backpropagation: calculate weight of first layer

computation graphs

Eqn to graph



Training LM

large corpus, gold label? ✗

↓
unsupervised learning approach
(not annotated)

sliding window

self supervised learning

Train your model to predict next word

BERT Paper

(NSP (Next Sentence Prediction) - minimize loss function
MLM (Masked Language Modeling)

Lunches are <tok> in restaurants.

↓
token

unfreeze embedding → modern LM
(update embedding layer)

LM 7 done

RNNs and LSTMs

Recurrent Neural Networks

--- dash line on h_t means hidden layer dependent on sth.

like feedforward neural network

Diff betwn RNN and feedforward NN:

non-linear
(cycles in network)
dependency

linear representation
of network

Training on RNN \xrightarrow{LM} Autoregressive generation

Self-supervising

Automatically generating sth. (Don't specify gold label)

Teacher forcing

Sampling technique: give prev work, ask LM to predict next word

Text generation / LM

If it can generate a word token, loss is minimal

word2vec, fasttext \rightarrow outside source

use predicted output to calculate loss

RNN for sequence labeling

Each token of sequence has a gold label

POS-tagging (parts of speech)

RNN for sequence classification

I love to play cricket — positive (overall)
no need for each token

n — last token

map output to 3 classes — FFN
(+, -, neutral)

prob of RNN: loss of info by the time it comes to the last token

Autoregressive Generation with RNNs

Teacher forcing before AG with RNNs (Task) \rightarrow text gen

(training mechanism)

fine tuning neural

~~Soln~~ Average pooling / max pooling

Bidirectional RNN \rightarrow RNN1, RNN2

stacked RNN \rightarrow RNN1, RNN2, RNN3

more layers \rightarrow learn more features
abstract / sophisticated

Limitations RNN

- Info passed is local
(Does not perform well for distant info)

Provide info in current context
Carrying info from forward

- vanishing gradient problem
back propagation: past info depending on future info
(backward pass)

Remember important part, forget rest

↓
LSTM
removes info not needed
add info which might be used to take future decisions
use Gates

Focus on important parts of a sentence → Attention

- feed forward layer
- sigmoid activation function → output 0/1
erase remember
-

Hadamard Product: Element wise multiplication
→ Accumulation of dot products

forget gate: forgets info not needed
binary masking (sigmoid function)

higher weights to info for keeping

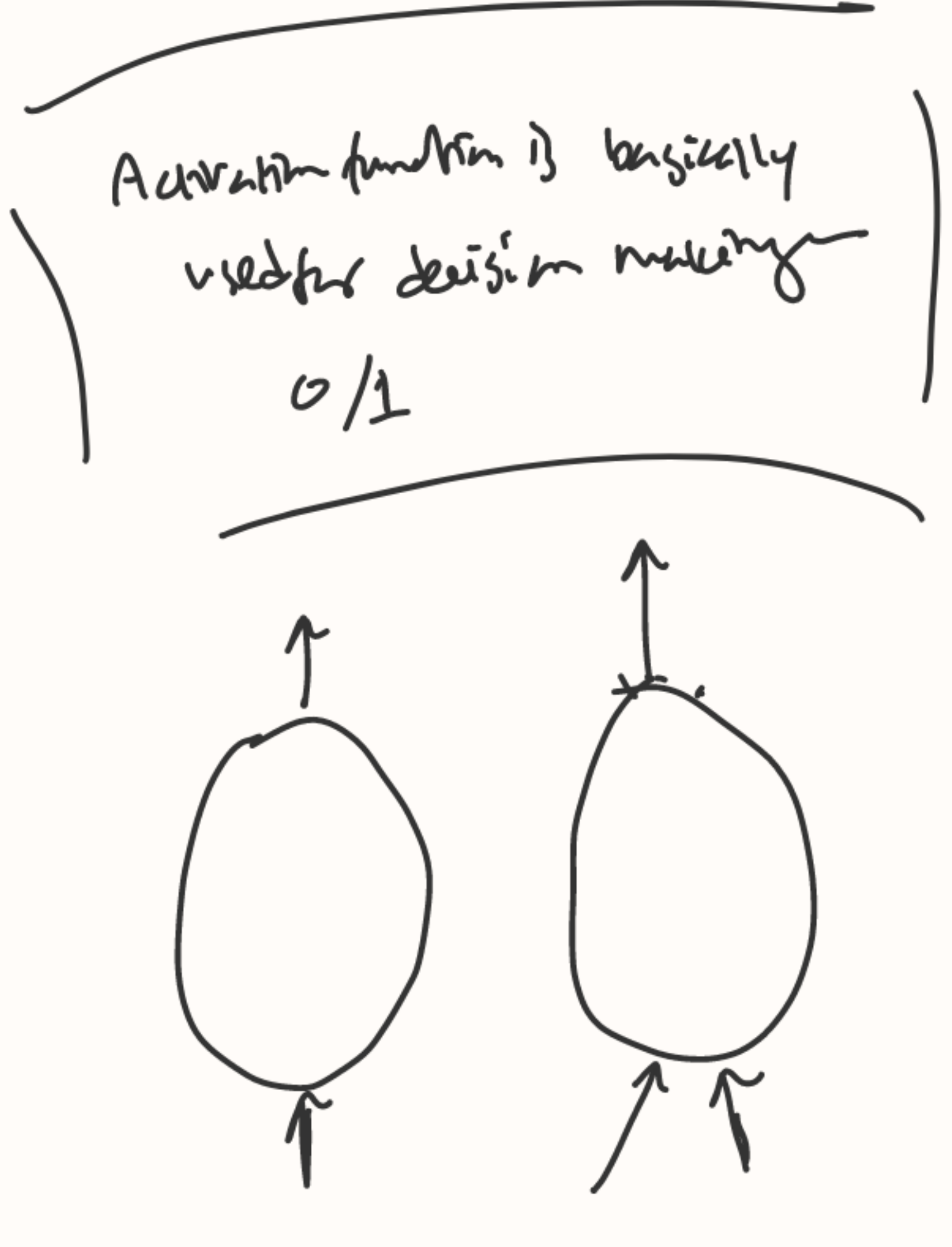
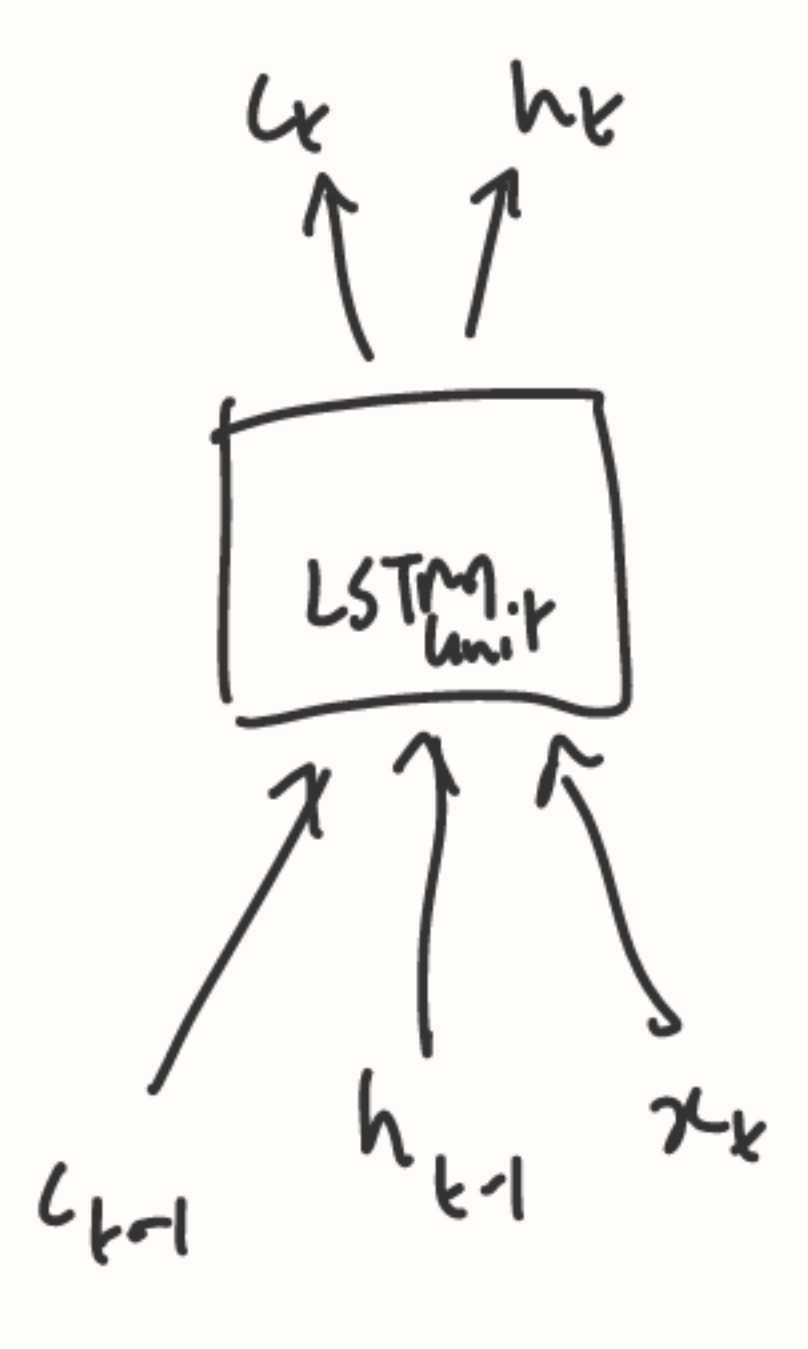
Info calculated for current context → use tanh activation func.
Keeps weights as it is
1 → as it is info
0 → 0
(use less info)

Add gate

info to be retained
(+) info to be erased

updated info (next context)
Output gate: context with updated info
hidden layer calculated using
prev hidden layer + current input + current context

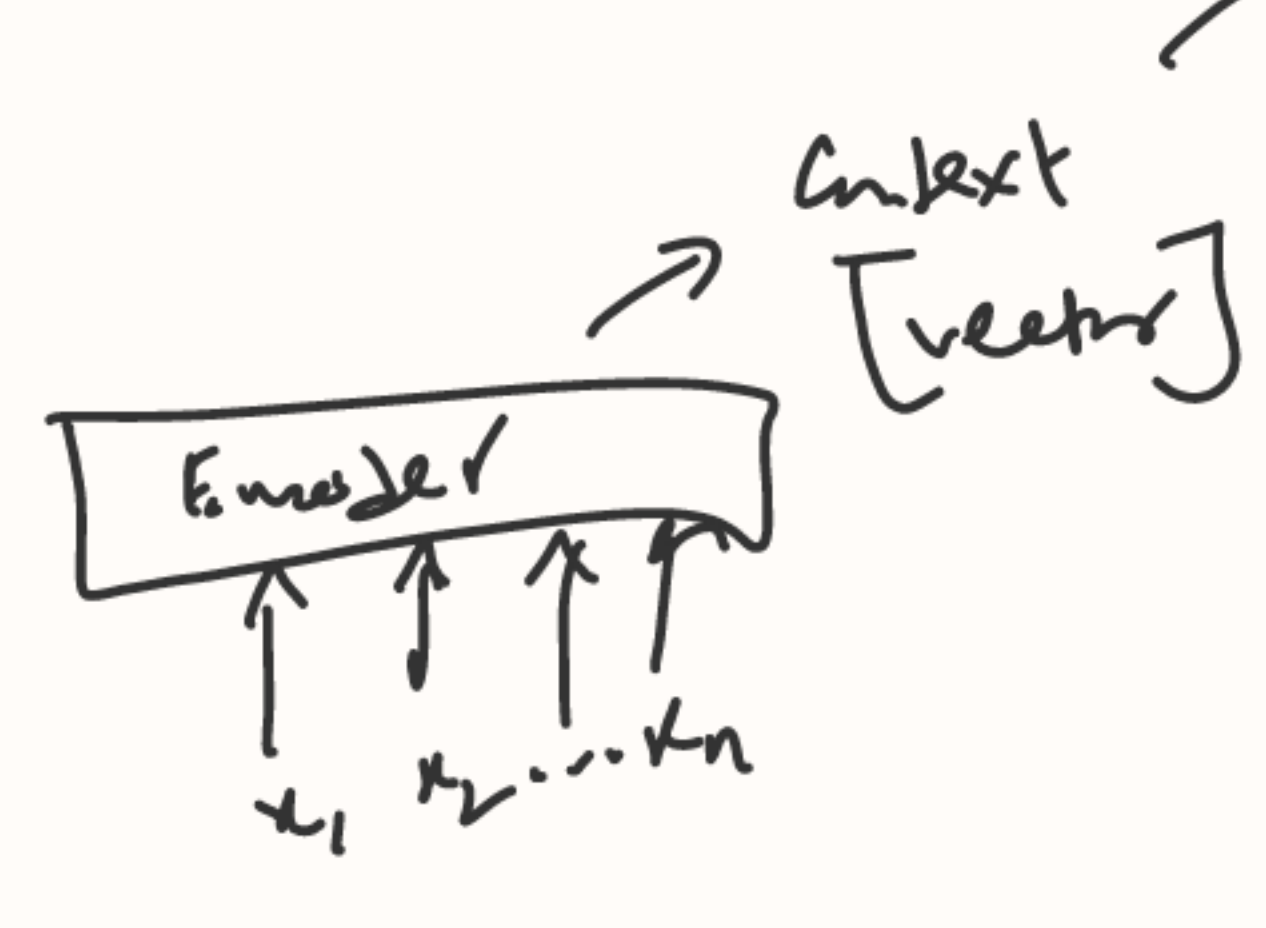
Feed forward neural network
RNN → knows most recent info
LSTM → most relevant info/context



RNN NLP Architectures

1. Sequence Labeling Task
2. Sequence Classification
3. Language Modelling
4. Encoder-Decoder

seq-seq modeling
translation length! = input length
 y_1, y_2, \dots, y_m



Translation task

Inference == prediction

Training the encoder/decoder model

Sample wise loss (not elementwise loss)

Training - loss - weight submit
validation set - loss calculate acc for seq?

Decoder
Hidden layer:

$$h_t^d = g(y_{t-1}^d, y_{t-1}^e, c)$$

$$c = h_n^e = h_0^d$$

$$h_t^d = g(y_{t-1}^d, h_{t-1}^d, c)$$

$$y_t = \text{softmax}(u_t)$$

$$y_t = \text{softmax}(z_x)$$

LSTM performs uniformly better than RNNs.

CH10: Transformers and LMs

30K-11ac vocab : Age 30 person

740 words learn everyday

reading enriches vocab

Distributional hypothesis - words that appear together have a similar meaning

Let LM read books

RNN, LSTM, Encoder-Decoder

These architectures cannot learn complex association of words

Soln: Transformers

LSTM does sm like attention

Building block: Attention - to focus on the relevant part of the sentence

Attention mechanism -

self attention - Building block of transformers

What part of RNN does LSTM improve?

Remember dist into by forgetting adding more info

LSTM also not accurate (probabilistic model)

could forget

uncertainties in LM

The notion of letting LM read and learn association of words: Pre-training
(fitting LM with info)

Downstream task: Fine-tuning

To make contextual representation of each word based on its surroundings - word embedding
(don't use frequency based methods like tf-idf)

Parallelization? Parallelization



How self attention is different from attention?

hidden layers not connected to each other
compare each of words (not the hidden layers)
no recurrent connection
compare hidden unit

Backward looking self att: (only look at preceding words)

undo

Autoregressive generation

not dependent on each other

Bidirectional looking self att: (look all words)

dot products

Create a representation of context for each word based on surrounding words - self att.

How to calculate a_3 ?

1st step: $score(x_i, x_j) = x_i \cdot x_j$

2nd step: $\alpha_{ij} = \text{softmax}(score(x_i, x_j))$

att. weight

3rd step: $a_i = \sum_{j \leq i} \alpha_{ij} x_j$

actual input to self att. calculate
(in att. hidden input to self att. calculate)

which word is more relevant → give more weights

prev: outsource embeddings.

Transformer have their own embeddings.

Query, Key, value - 3 roles of self att

Each token multiplied with softmax

A cow will come.
key: will, value: come, focus/query: will

Each input has these 3 representations
(x_1, x_2, x_3) (q, k, v)

steps of self att:

1) Comparison: compare query of x_3 with key of x_1 - dot product
key of x_2 - ~
key of x_3 - ~

2) Normalize: send this dot product to softmax activation function

3) Calculate value: Multiply each vector with their corresponding softmax value.

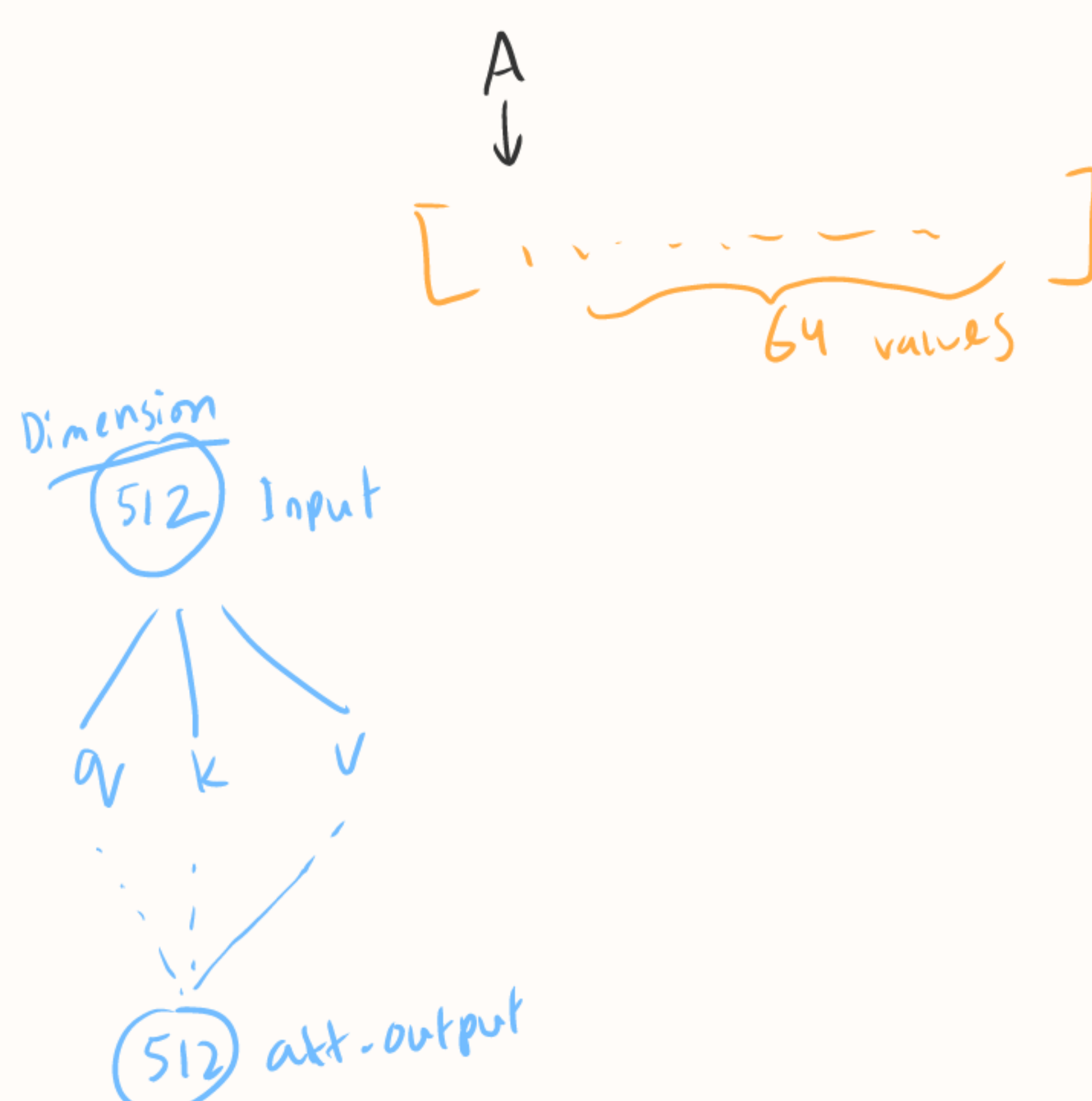
Add \rightarrow Get a single vector

Output of self-att.

Dimension

$$d_k = d_v = 64$$

Each word represented by 64 values



Transformer remember long context - very good autoregressive generator

RNN and LSTM cannot look at very distant info - fail after a few tokens

Magic out the future $\rightarrow -\infty$
(lower A matrix)

1 self att unit

// A cat is chasing a rat.

Syntax relation

Semantic relation

Multihread att.

Transformer Blocks

Residual connection: original info pass to the next layer

layer norm: to keep values in a calculable range

γ - game } learnable curve
 β - offset }

" 16 transformer blocks inside GPT3 "

cat sat m mat

$$q_{cat} = [0.2 \ 0.3 \ 0.45 \ 0.35]$$

$$k_{cat} = [\quad \quad \quad n \quad \quad]$$

$$v_{cat} = [\quad \quad \quad n \quad \quad]$$

$$q_{sat} \cdot k_{cat} = 3.45$$

$$q_{sat} \cdot k_{sat} = 120$$

$$\alpha_{21} = \frac{3.45}{3.45 + 120} \approx 0.45$$

$$\alpha_{22} = \frac{120}{3.45 + 120} \approx 0.55$$

$$a = \sum \alpha_{ij} v_j$$

attention \uparrow
weight \uparrow
value \uparrow

$$a_2 = (0.45 \times v_{cat}) + (0.55 \times v_{sat})$$

Post-norm vs pre-norm transformer

\downarrow
traditional

layer normalization w/ raw input (GPT1 GPT2 GPT3)

word embedding for each token
and position embedding

\rightarrow To prevent exchange of position which can change meaning of the sentence

David beat John w/ a stick
 \curvearrowright
token exchange

John beat David w/ a stick.

positional embedding used in negative sentiment

NOT - good \rightarrow negative sentiment

Language modeling head

\downarrow
Addition top of transformer block
to match dimension

vocab size = no. of tokens in dataset = #cols of embeddings

LM = next token/word prediction task
(Lang. modeling)

Transformer has a long context window.

Works w/ Transformer

↓
LM trained (generalized)

Sentiment Analysis
Ques-Answering
Summarization } Downstream task solver

RAG - use external sources w/ LM
(to prevent hallucination)
↓

Retrieval Augmented Generation

- Knowledge graphs
- vector database
- prompt engineering (chain of command)
↓
Divide query into a chain

Toxic Lang Generation - data filtering

Privacy Issues - remove private datasets from training data

Datasheets/model cards

Finals

After mid

2 Ques from each 3 chapters

RNN, LSTM,
Transformers

maps

Diagrams of underlying architecture

Self supervised training

Dimension

