

# Linear regression Formula

Model equation

$$y = \beta_0 + \beta_1 x + \epsilon$$

y is Predicted value/output

$\beta_0$  is the Y intercept (constant term)

$\beta_1$  coefficient for the input feature  
(slope)

$\epsilon$  is the error term

(difference between Predicted value  
& actual value)

• Cost function (Mean Squared Error)

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

MSE is the Mean Squared error to minimize

n is number of data points

$y_i$  is actual output,  $\hat{y}_i$  is a predicted output

# Logistic regression

- Hypothesis function: (sigmoid)

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$h_{\theta}(x)$  is the probability that the output is 1 given  $x$ .

$\theta$  is parameter vector (weights)

$x$  is feature vector (input data)

$T$  is the linear combination of input features.

- Cost function (Binary-Cross Entropy):

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

→ Flip to next Page

# Explanation of Binary Cross-Entropy

$J(\theta)$  is Cost function to minimize

$m$  is number of training examples

$y^{(i)}$  is Actual Label (0 or 1)

$h_{\theta}(x^{(i)})$ : Predicted Probability of the  $i$ th data-point.

## Gradient Descent

update rule:

$$\theta = \theta - \alpha \frac{\partial J(\theta)}{\partial \theta}$$

$\theta$  is the model parameters (weights)

$\alpha$  is the learning rate (step size)

$\frac{\partial J(\theta)}{\partial \theta}$  is the Gradient of the Cost function  $J(\theta)$ .

## L2 Regularization (Ridge):

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \theta_j^2$$

$\lambda$  is the regularization strength it controls overfitting

$\sum_{j=1}^p \theta_j^2$  is the sum of squared model parameters

\* The rest of the equation is MSE (Mean Squared Error)

## L1 Regularization (LASSO):

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\theta_j|$$

$|\theta_j|$  is the absolute value of the parameters (includes Sparsity).

- Distance Formula (Euclidean distance)

$$d(x, x') = \sqrt{\sum_{i=1}^n (x_i - x'_i)^2}$$

$d(x, x')$  is the distance between points  $x$  and  $x'$ .

$x_i, x'_i$  is the  $i$ th feature of points  $x$  and  $x'$ .

## Support vector machines (SVM)

- Decision boundary:

$$f(x) = w^T x + b = 0$$

$f(x)$  is the decision function

$w$  is the weight vector

$x$  is the input feature vector (defines the hyperplane)

# Support vector machine

• Hinge loss:

$$L = \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

$L$  is hinge loss

$y_i$  is the true label (+1 or -1)

$\mathbf{x}_i$  is the input feature vector  
for the  $i$ th data point.

$\mathbf{w}^T$  is the weight vector.

## Principal Component Analysis: PCA

Variance Explained

$$\text{Var}(z) = \frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})^2$$

$\text{Var}(z)$  is the variance of  
the transformed Principal  
Components

$n$  is the number of data points

$z_i$  is the value of  $i$ th Principal Component

→ Continued for Principal Component Analysis

$\bar{z}$  is the Mean of the Principal Components

## Cross - Validation

K-Fold Cross - Validation:

$$CV_{\text{error}} = \frac{1}{K} \sum_{i=1}^K MSE_i$$

$CV_{\text{error}}$  is the Average validation error across K folds

K is the number of folds

$MSE_i$  is the MSE (Mean Squared Error) on the ith fold  
(MSE is explained on Pg. 1)

# Evaluation Metrics

Accuracy:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

TP is True Positives

TN is True negatives

FP is False Positives

FN is False negatives

Precision:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

• Precision measures the proportion of true positives to sum of true positives & false positives

• Recall:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Recall measures the proportion of true positive predictions

out of all actual positives.

F1 Score:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

\* x means multiplication

F1 is the harmonic mean  
of Precision & Recall,  
used when there is  
an uneven class distribution

### Confusion Matrix

True Positive: Correctly classified  
Positive

True Negative: Correctly classified  
negative

False Positive: Incorrectly  
classified as  
positive

False negative: Incorrectly  
classified as  
negative

# Decision Trees

Gini Impurity: (for classification)

$$G = 1 - \sum_{k=1}^K p_k^2$$

$p_k$  is the proportion of classes in class K

K is number of classes

Entropy: for classification

$$H = - \sum_{k=1}^K p_k \log_2(p_k)$$

$p_k$  is the proportion of data points in class K.

Information gain :

$$IG = H_{\text{parent}} - \sum_i \frac{|D_i|}{|D|} H(D_i)$$

IG is the information gain  
from splitting the  
dataset

$H_{\text{parent}}$  is the Entropy of the  
parent node.

$H(D_i)$  : Entropy of the child  
node (i)

$|D_i|$  is the number of samples  
in child node (i)

$|D|$  is the total number of  
samples in the Parent  
node

# Random Forests

Final Prediction: (for classification)

$\hat{y}$  = majority vote of individual trees

Final Prediction: (for regression)

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T \hat{y}_t$$

T is the number of trees in the forest

$\hat{y}_t$  is the prediction from the t-th tree

Bayes theorem:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

$P(y|x)$  is the posterior probability of class Y → given feature x

$P(x|y)$  is the Posterior Probability of observing feature  $x$  given class  $y$ .

$P(y)$  is the Prior Probability of class  $y$

$P(x)$  is the Probability of Observing feature  $x$

Gaussian Naive Bayes (for continuous features);

$$P(x_i|y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

$x_i$  is the feature value

$\mu_y$  is the mean of the feature  $x_i$  for class  $y$ .

$\sigma_y^2$  is the Variance of the feature  $x_i$  for class  $y$ .

# Computer Vision Formulas

## Convolution Operation (in CNN's)

Formula (2D convolution):

$$(I * K)(x, y) = \sum_{i=-K}^K \sum_{j=-K}^K I(x+i, y+j) \cdot K(i, j)$$

$I(x, y)$  is the input image pixel at position  $(x, y)$

$K(i, j)$  is the Kernel (filter) value at position  $(i, j)$

\* is the convolution operator

$(x, y)$  is the position in the output feature map

The kernel (filter) slides over the image, applying the dot product between the image patch and the kernel to produce the output feature map

# Pooling Operation (Max Pooling)

Pooling is used to downsample the image while still retaining important information.

Max Pooling formula:

$$P(x, y) = \max(I(x+i, y+j)), \text{ for } i, j \in [0, f-1]$$

$P(x, y)$  is the pooled output at position  $(x, y)$

$I(x+i, y+j)$  is the image pixel values in the  $f \times F$  neighborhood

$f$  is the size of the pooling window (usually  $2 \times 2$  or  $3 \times 3$ )

The maximum value within each pooling window is taken to reduce the image size.

Average Pooling

$$y_{i,j} = \frac{1}{K} \sum_{m,n} x_{(2i+m), (2j+n)}$$

check  
other  
page  
for  
variables.

## ReLU (Rectified Linear Unit) Activation Function

ReLU function:

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

$x$  is the input value

$f(x)$  is the output value after applying ReLU

ReLU replaces negative values with 0, retaining only positive values which helps in learning complex patterns.

### Softmax activation function (for image classification)

Softmax formula:

$$\sigma(z_i) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

$z_i$  is the logit or score for class  $i$

$\sigma$  is the Probability of class  $i$   
( $z_i$ )

$n$  is the total number of classes.

The softmax function normalizes the output so that the Probabilities for all classes sum up to 1.

Cross Entropy Loss (for classification tasks)

Cross-Entropy  
Loss formula

$$L = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

$L$  is the loss value

$y_i$  is the Actual Label (0 or 1 for binary classification)

$n$  is the total number of classes

- the loss increases as the predicted probability for the true class decreases.

Bounding Box regression  
(for object detection):

Bounding  
box

coordinates:

$$(x_{\min}, y_{\min}, x_{\max}, y_{\max})$$

$x_{\min}, y_{\min}$ : coordinates of the top left corner for the bounding box

$x_{\max}, y_{\max}$ : coordinates for the bottom right corner of the bounding box.

IoU (Intersection over Union),

$$IoU = \frac{\text{Area of intersection}}{\text{Area of Union}}$$

IoU is the overlap between the predicted and ground truth bounding boxes.

This metric is used to evaluate the accuracy of object localization.

Bounding box Loss  $L_2$  loss:

$$L_{bbox} = (x_{pred} - x_{true})^2 + (y_{pred} - y_{true})^2$$

$x_{pred}, y_{pred}$  is the predicted center coordinates of the bounding box

$x_{true}, y_{true}$  is the ground truth center of the bounding box.

The loss penalizes the difference between predicted or

# Affine transformation

for image augmentation

these are transformations applied to images.

Affine transformation Matrix;

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} e \\ f \end{bmatrix}$$

$(x, y)$  are the original coordinates

$(x', y')$  are the transformed coordinates

$a, b, c, d$  are the Transformation Matrix Elements (e.g. Rotation, Scaling)

$e, f$  are the translation values

Multi-scale  
Processing

Image Pyramid Scaling formula:

$$I'(x,y) = I(sx,sy)$$

$I(x,y)$  Original Image

$I'(x,y)$  Rescaled Image

$s$  is the scale factor (e.g. 0.5  
*for downscaling*)

This allows detection  
of objects in  
various sizes  
images.

# Optical Flow (for motion detection)

Optical flow tracks motion in videos by estimating movement of pixels between frames.

Optical Flow  
Equation

$$I_x u + I_y v + I_t = 0$$

$I_x, I_y$  are the image intensity gradients in the  $x, y$  directions

$I_t$  is Temporal Intensity Gradient  
(change in Pixel intensity over time)

$u, v$  are the optical flow velocities  
(displacement in  $x$  and  $y$ )

# Histogram of oriented gradients (HOG)

HOG is used for feature extraction by capturing the distribution of gradient directions

Gradient Magnitude and direction:

$$M(x, y) = \sqrt{I_x^2 + I_y^2}, \Theta(x, y) = \tan^{-1}\left(\frac{I_y}{I_x}\right)$$

$M(x, y)$  is the Gradient magnitude at Position  $x, y$

$\Theta(x, y)$  is the gradient direction.

$I_x, I_y$  are the image gradients in the  $x$  and  $y$  directions

HOG captures these gradients and groups them into orientations bins to describe

SSD (single shot multibox detector)

SSD multibox Loss function:

$$L = \frac{1}{N} (L_{\text{conf}} + \alpha L_{\text{loc}})$$

$L_{\text{conf}}$  is how well the predicted class probabilities match the ground truth.

$L_{\text{loc}}$  is the difference between the predicted and ground truth bounding box coordinates.

$\alpha$  is the balancing factor between confidence loss and localization loss.

$N$  is the number of default bounding boxes

(is the confidence score for the object)

Predicted box =  $(x_c, y_c, w, h, c)$

$x_c, y_c$  center coordinates,  $w, h$  width & height

Similar are: YOLO and R-CNN

# Term frequency (TF)

The term frequency measures how frequently a term appears in a document

$$TF(t,d) = \frac{f_{t,d}}{n_d}$$

$t$  is the term being evaluated

$d$  is the specific document

$f_{t,d}$  is the count of term  $t$  in document  $d$

$n_d$  is the total number of terms in document  $d$

# Inverse Document Frequency (IDF)

Measures the importance of a term across multiple documents

$$IDF(t, D) = \log\left(\frac{N}{f_{t,D}} + 1\right)$$

$t$  is the term

$D$  is the set of documents

$N$  is the total number of documents

$f_{t,D}$  is the number of documents containing the term  $t$

TF-IDF (Term frequency-Inverse Document frequency)

(Combines TF & IDF to evaluate the importance of a term in a document relative to a corpus)

$$TF-IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$

Measures the similarity between the two vectors  
(e.g document representations)

$$\text{Cosine Similarity}(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$$

$A \cdot B$  is the dot product of vectors A and B

$\|A\|$  is the length of vector A

$\|B\|$  is the length of vector B

Dot Product

$$A \cdot B = \sum_{i=1}^n A_i B_i$$

This calculates the sum of

## Magnitude of a vector

$$\|A\| = \sqrt{\sum_{i=1}^n A_i^2}$$

this computes the euclidean length of the vector

## Range of values

the value of cosine similarity ranges from -1 and 1

1 shows that the two vectors are similar

0 shows that both vectors are orthogonal (no similarity)

-1 shows that both vectors are opposed

## Word embeddings (e.g word2vec)

This is used to convert

words into vectors so

numbers, so computers

can process them

## Skip-Gram model formula

$$J = - \sum_{t=1}^T \sum_{j=-c}^c \log P(w_{t+j} | w_t)$$

$J$  is the object function to minimize

$T$  is the total number of words in the corpus

$w_{t+j}$  is the context of words within a window size  $c$ .

$P(w_{t+j} | w_t)$  is the probability of context word given the target word, typically modeled with softmax activation function.

Skip Gram's are used to generate word embeddings in NLP.

## Sequence-to-Sequence models (seq2seq)

this converts one sequence of data to another.

Loss function (cross entropy)

$$L = - \sum_{t=1}^T \sum_{k=1}^V y_{t,k} \log(p_{t,k})$$

L is the loss function to minimize

T is the length of the sequence

V is the size of the vocabulary

$y_{t,k}$  is the ground truth output

(1 if the t-th output is k, else 0)

$p_{t,k}$  is the predicted probability of word k being the t-th output.

Hidden State update:

$$h_t = \sigma(W_h h_{t-1} + W_x x_t + b)$$

$h_t$  is the hidden state  
at time  $t$

$h_{t-1}$  is the hidden state  
at the previous time  
step.

$W_h$  is the weight matrix for  
the hidden state

$W_x$  is the weight matrix for  
the input  $x_t$ .

$b$  is the bias vector

$\sigma$  is the activation function  
(e.g. tanh, ReLU)

RNN's are used to process  
word input as it can  
update current <sup>keep track of past</sup> state words and

# BERT (Bidirectional Encoder representation from transformers)

Self-attention mechanism:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

$Q$  is the Query matrix

$K$  is the Key matrix

$V$  is the Value matrix

$d_K$  is the dimensionality of the key vectors for scaling.

The softmax function ensures that the attention weights sum to one.

It is helpful for operations with text input and is a more effective model of RNN's.

Cross entropy loss

Formula:

$$L = -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

\* refer to other page  
for variable meaning

# Reinforcement Learning Algorithms

## Markov Decision Process (MDP)

States ( $S$ ) Possible situations the agent can be in

Actions ( $A$ ) choices available to the agent

Transition function ( $T$ ) Probability of moving from one state to another given an action

Reward function ( $R$ ) Immediate reward received after transitioning from one state to another

Discount Factor ( $\gamma$ ) A value between 0 & 1 that determines the importance of future rewards.

MDP's provide a formal framework for modeling decision making problems where outcomes are

## Value function

Formula

$$V(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid s_t = s \right]$$

use: This estimates the expected return (total future rewards) starting from state  $s$ . helps us evaluate long-term value of states in a policy.

Action Value function (Q-Function)

Formula:

$$Q(s, a) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid s_t = s, a_t = a \right]$$

use: Estimates the expected return of taking action  $a$  in state  $s$  and following afterward.

Bellman Equation & signals gradients

Value function Bellman equation

$$V(s) = \max_a [R(s, a) + \gamma \sum_{s'} P(s'|s, a)V(s')]$$

Action-value function Bellman equation

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \underbrace{\max_{a'} Q(s', a')}$$

Use: This provides a recursive relationship to compute value functions enabling iterative methods for policy evaluation and improvement

Policy

What: It is balancing the choice between exploring new actions (to discover their effects) and exploiting known actions that yield high rewards

## Use of Policies:

This is essential in RL  
to ensure the agent learns  
effectively over time and  
avoids local optima.

## Temporal Difference Learning (TD Learning)

Formula:

$$V(s) \leftarrow V(s) + \alpha(R + \gamma V(s') - V(s))$$

Use: Combines ideas from  
monte carlo methods & dynamic  
programming to update value  
estimates based on other  
learned estimates enabling  
learning from incomplete  
episodes.

## Q-Learning update rule

Formula:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(R + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

use of Q-Learning:

This allows an agent to learn the value of actions directly from interactions with the environment, enabling off-Policy Learning.

## Deep Q-Networks (DQN)

use: Combines Q-learning with Deep neural networks to approximate the Q-Function, enabling the handling of high-dimensional state-spaces like images.

## Policy Gradient Methods

Formula:

$$\nabla J(\theta) = \mathbb{E}_{T \sim \pi_\theta} \left[ \sum_{t=0}^T \nabla \log \pi_\theta(s_t, a_t) R \right]$$

use: Directly optimize the policy by adjusting its parameters ( $\theta$ ) to maximize expected returns, useful for complex action spaces and continuous actions.

# Deep Learning formulas

## Gradient Descent

This is a optimization algorithm used to minimize the loss function by updating model parameters iteratively.

Formula:

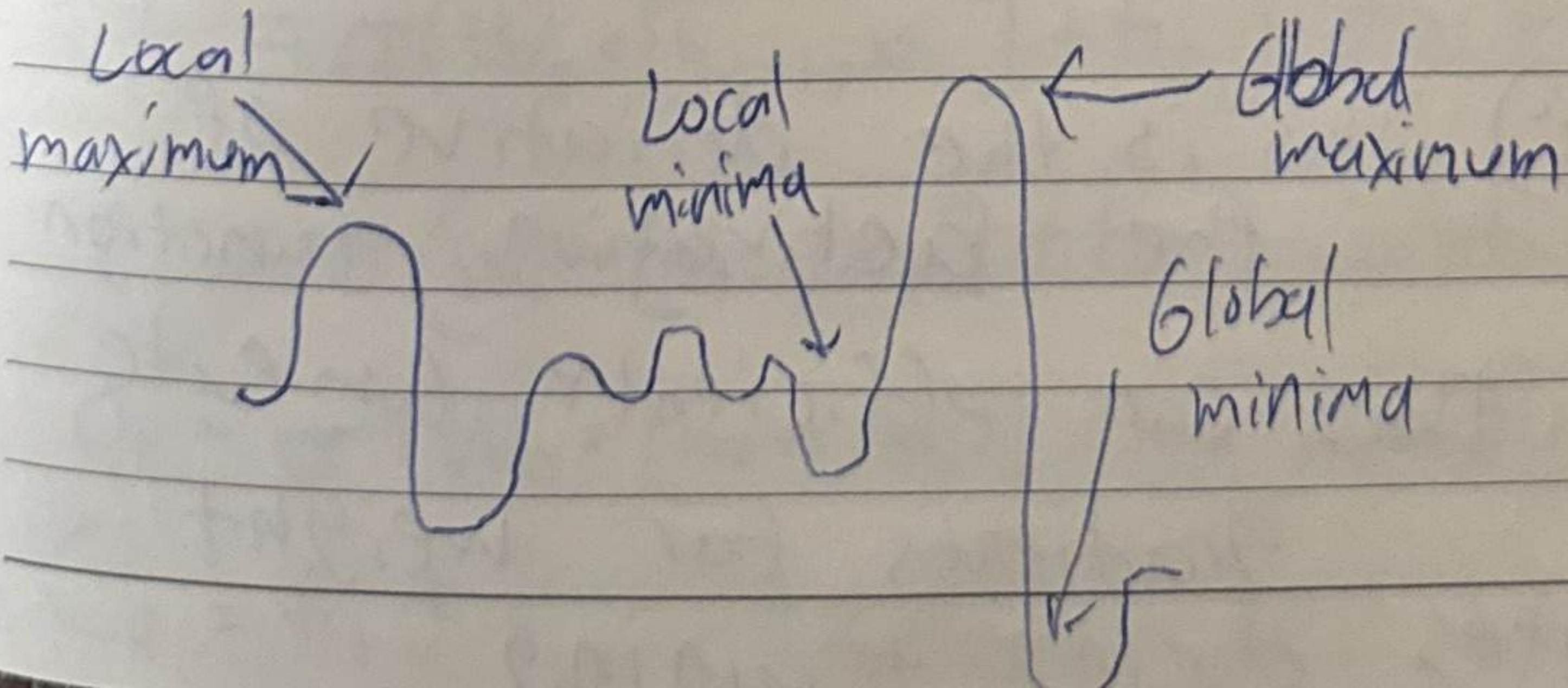
$$\theta = \theta - \alpha \nabla J(\theta)$$

$\theta$  is the model Parameters

$\alpha$  is the learning rate

$J(\theta)$  is the loss function

this algorithm guides the model by minimizing error.



# Back Propagation through Neural Network

An algorithm for calculating gradients of the loss function with respects to weights using chain rule.

Formula

For a layer  $l$ ,

$$\delta^l = \nabla_a L \circ \sigma'(z^l)$$

$\delta^l$  is the error for layer  $l$ .

$\nabla_a L$  is the Gradient of the loss with respect to activations

$\sigma'(z^l)$  This is the derivative of the activation function

This can efficiently compute gradients for weight

## Activation function

\* check other pages for  
ReLU & sigmoid

TANH activation function:

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

These introduce non linearity  
enabling the model to  
learn complex relationships

Long short-term memory (LSTM) for  
MLP

LSTM cell equations meaning is

on next page →

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \text{ Forget gate}$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \text{ input gate}$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \text{ output gate}$$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t \text{ cell state}$$

$$h_+ = o_+ * \tanh(c_+) \text{ hidden state}$$

These LSTM manage long term dependencies in sequential duty effectively.

## Auto encoders

Structure

Encoder:

$$z = f(x) = \sigma(W_e x + b_e)$$

Decoder:

$$\hat{x} = g(z) = \sigma(W_d z + b_d)$$

Useful for dimensionality and anomaly detection

## Optimizers

Adam optimizer:

Importance: Different

help navigate

effectively, improving training

Adam optimizer formula:

$$\hat{m}_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$\hat{v}_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta = \theta - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

LSTM cell equation meanings

$f_t$  is the forget gate activation at time  $t$  (a value between 0 and 1)

$N_f$  is the weight matrix for the forget gate

$h_{t-1}$  is the hidden state from the previous time step ( $t-1$ )

$i_t$  is the input at time  $t$ .

~~Explanation of~~  
 $b_f$  is the bias term for forget gate.

$\sigma$  is the sigmoid activation function

$O_f, I_f, \tilde{I}_f, \tilde{C}_f$  are same meaning for different gates

$W_o, W_f, W_i, W_c$  are also same meaning

$b_o, b_f, b_i, b_c$  are also same meaning

tanh is the hyperbolic tangent function, which outputs values between -1 and 1.

$C_t$  is the cell state at time t.

$C_{t-1}$  is the cell state from previous

M Gates:

\* is the symbol for element wise multiplication

$h_t$  is the hidden state at output at time  $t$ ,

$c_t$  is the current cell state at time  $t$

## Auto encoders Explanation

$W_e$  is the weight matrix for encoder.

$b_e$  is the bias vector for encoder.

$h$  is the encoded representation  
(latent space of input)

$h = f(x) = \sigma(W_e x + b_e)$  is the equation for the encoder where  $f$  is the activation function.

\* These mean the same thing for the Decoder's as well.

## Adam Optimizer explained

### Parameters

$\theta$  are the Parameter (weights) of the model that are being optimized.

### Learning rate

$\alpha$  is the learning rate, which controls how much to update the parameters during training (it is set to a small value e.g 0.001)

### Exponential Decay rate:

$\beta_1$ -decay rate for first moment (e.g 0.9)

$g_t$  is the gradient of the loss function with respect to parameters at time step  $t$ .

Moment estimates

$m_t$  the first moment estimates  
(mean of the gradients)

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$v_t$  is the second moment estimate  
(the uncentered variance of the gradients)

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Bias corrected moment estimates

$\hat{m}_t$  the bias corrected first moment estimate

$\hat{v}_t$  is the bias corrected second moment estimate

Parameter update:

The parameters are updated using:

$$\theta = \theta - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

$\epsilon$ , is the small constant added for numerical stability usually around ( $10^{-8}$ ).

## Diagrams of Nets in ML

### Feed forward Neural Networks (FNN)

- single-layer Perceptron
- multi layer Perceptron

### Convolutional neural networks (CNN)

- Standard CNN
- Residual Networks (ResNet)
- Inception Networks
- U-Net (used for image segmentation)
- SqueezeNet

# Recurrent Neural Networks (RNN)

- Standard RNN
- Long short-term memory (LSTM)
- Gated Recurrent Unit (GRU)

# Generative Models

- Generative adversarial Networks (GAN)
- Variational Autoencoders (VAE)

# Transformer Models

- Vanilla transformer
- BERT (Bidirectional Encoder Representation from transformers)
- GPT (Generative Pre-trained transformer)
- T5 (Text-to-text transformer)
- Vision transformer (ViT)

# Autoencoders

- Denoising Autoencoder
- Sparse Autoencoder
- Variational Autoencoder (VAE)

## Graph Neural Networks

- Graph convolutional network (GCN)
- Graph Attention Networks (GAT)

## Reinforcement Learning networks

- Deep Q-Networks (DQN)
- Policy Gradient Methods
- Actor-Critic methods

## Specialized networks

- Capsule networks
- Attention networks
- Neural Turing machines
- Self-organizing Maps (SOM)

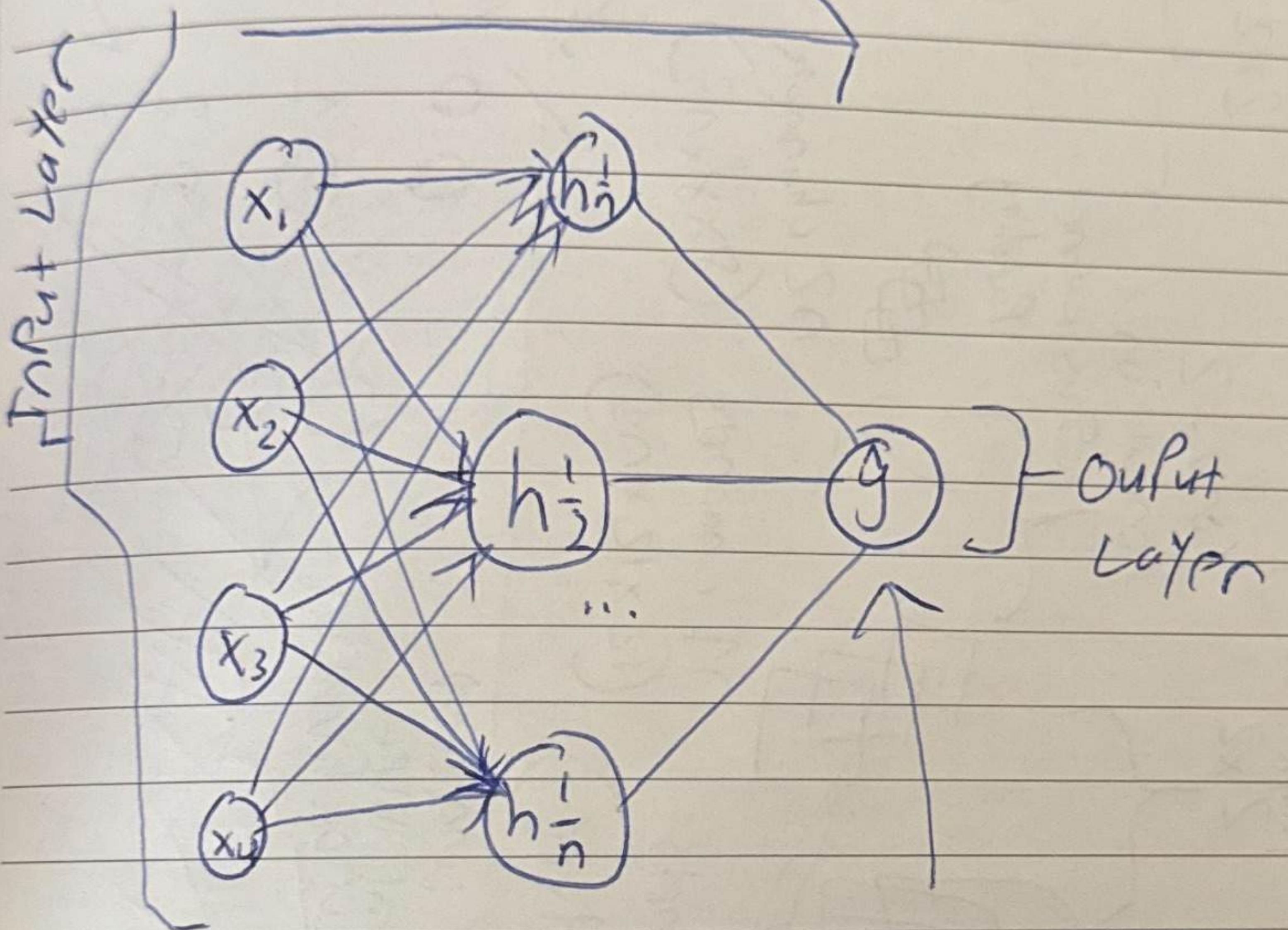
## Hybrid models

- Convolutional Recurrent Neural Networks (CRNN)
- Neural Architecture Search (NAS) Networks

# Models / visual representations

## Feed forward Neural Network

### Multi-LAYER Perceptron



SINGLE  
LAYER

HIDDEN  
LAYER

$$f(x) = \begin{cases} 1 & \text{if } \sum w_i x_i + b \geq 0 \\ 0 & \text{if } \sum w_i x_i + b < 0 \end{cases}$$

$x$   $w^1$   $\rightarrow$

$x_2$   $w^2$   $\rightarrow$   $\oplus$

$x_3$   $w^3$   $\rightarrow$

Summation & bias

$$\sum_{i=1}^m (x_i w_i) + \text{bias}$$

Conv - 1  
5x5 kernel  
padding

Max-Pooling  
 $2 \times 2$

Conv - 2  
convolution

5x5 kernel  
padding

Max Pooling  
 $2 \times 2$



Input

(28x28x1)

n1 channels  
(24x24x1)

fc - 3

Fully Connected  
layer  
activation

ReLU activation

0 0 1

0 2 2

0 1 2

1 1 1

1 1 1

1 1 1

1 1 1

1 1 1

$n^2$  channels

(8x8x1)

n2 channels

(12x12x1)

n2 channels

(8x8x1)

$n^2$  channels

(4x4x1)

$n^2$  channels

(2x2x1)

$n^2$  channels

(1x1x1)

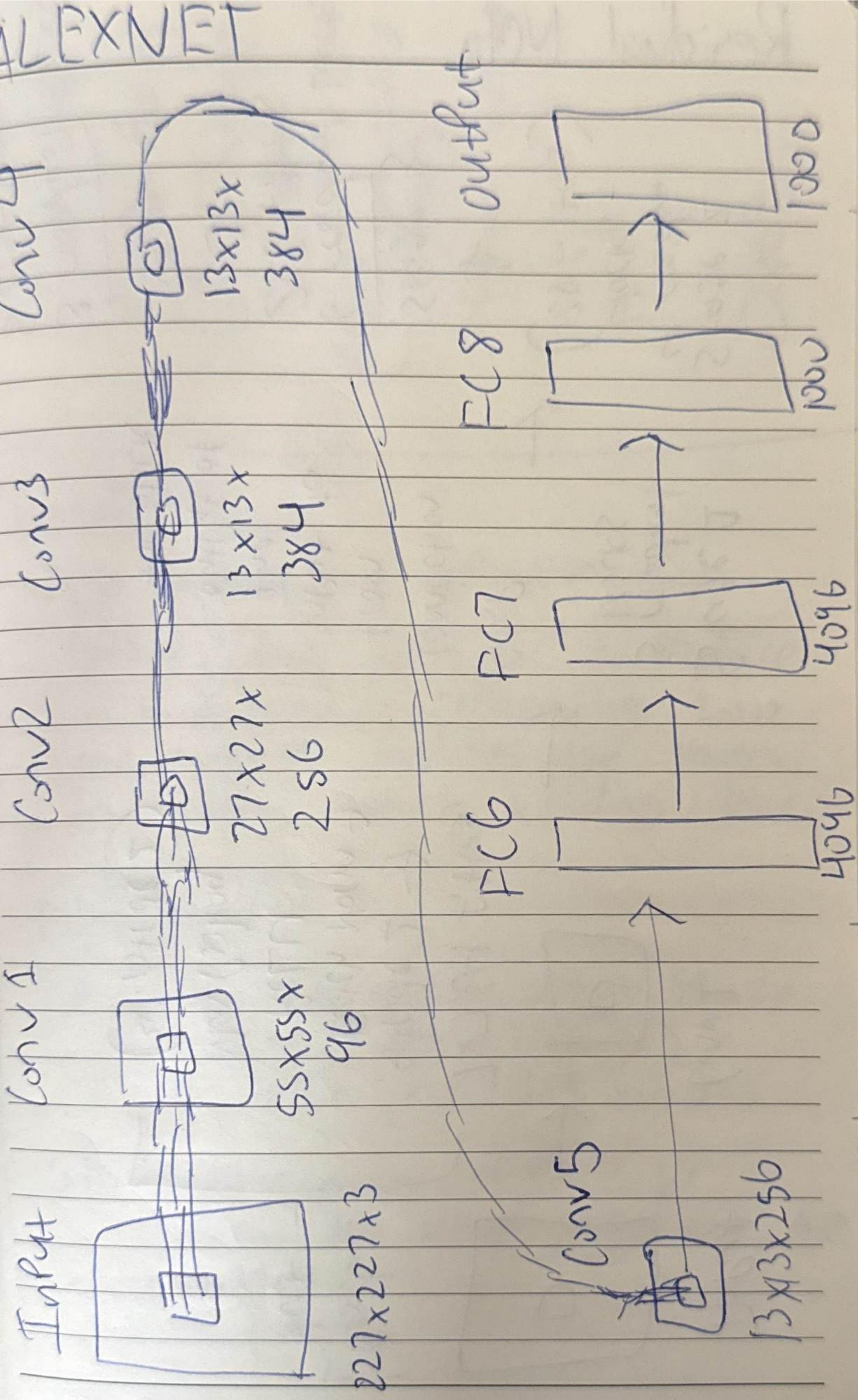
fully connected

with neural net

Dropout

out net

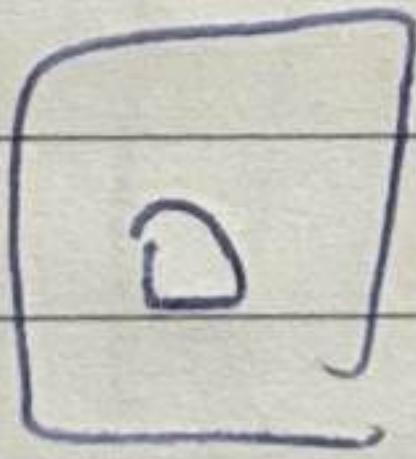
out net



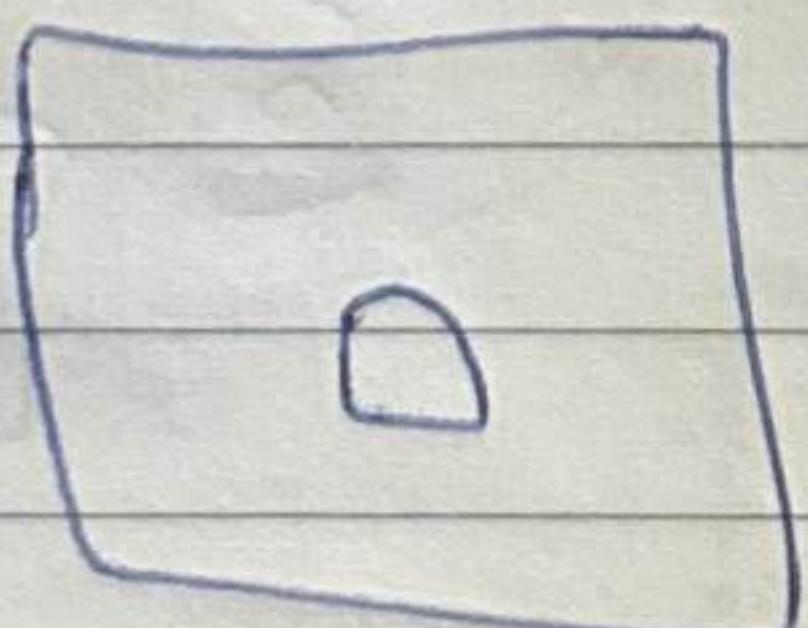
Residual  $\approx$

Stage 1  
3 residual  
blocks

Conv1



Input



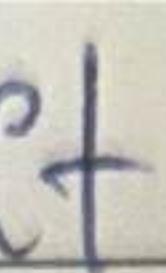
$7 \times 7$  64 filters,  
Stride 2  $\rightarrow$   
Batch Norm  $\rightarrow$   
ReLU  
Max Pooling  
 $(3 \times 3, \text{ stride } 2)$

1 Residual  
block -

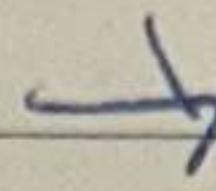
Skip  
connection



from  
input to  
the  
output of  
the block

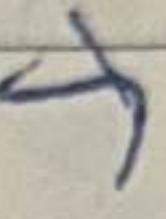


Stage 2  
4 residual  
blocks  
(28 filters)



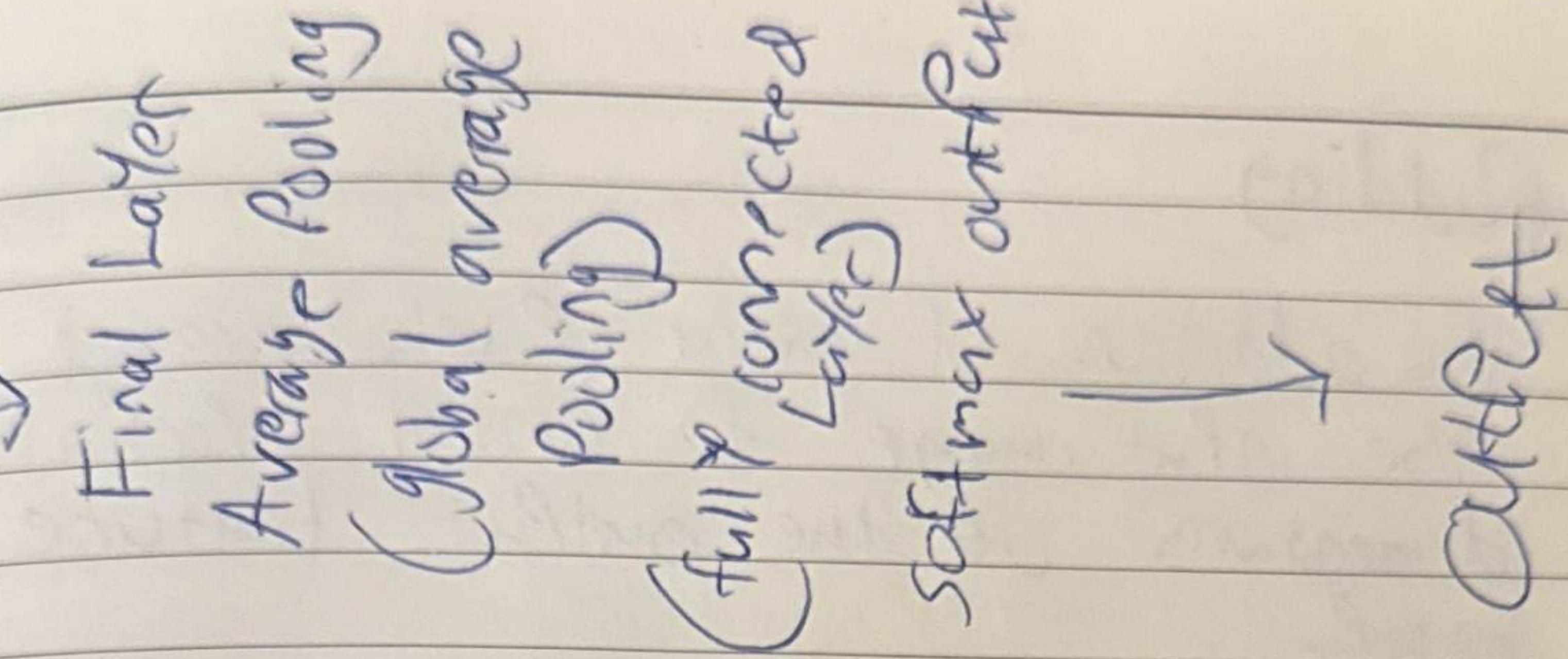
Stage 3

6 residual blocks  
256 filters



Stage 4

3 residual blocks  
512 filters



## Convolution

• Is a mathematical operation that combines two functions to produce a third function, representing how one function modifies the other.  
In CNN it involves a sliding filter.

## Filter (kernel)

A small matrix used to detect specific features (edges & textures) in the input (e.g.  $3 \times 3, 5 \times 5$ )

Stride: The number of pixels by which the filter slides over the image.

## Padding

The addition of extra pixels around the input image to control spatial dimensions of the output feature map.

Valid Padding: No padding, output is smaller than input

Same Padding: Padding is added to keep the output size the same as input size.

Feature map: the output produced by applying a filter to the input.

Pooling: A downsampling operation is used to reduce the spatial dimensions of the image.

Max Pooling: Takes the maximum value from a region.

Average Pooling: Takes an average value

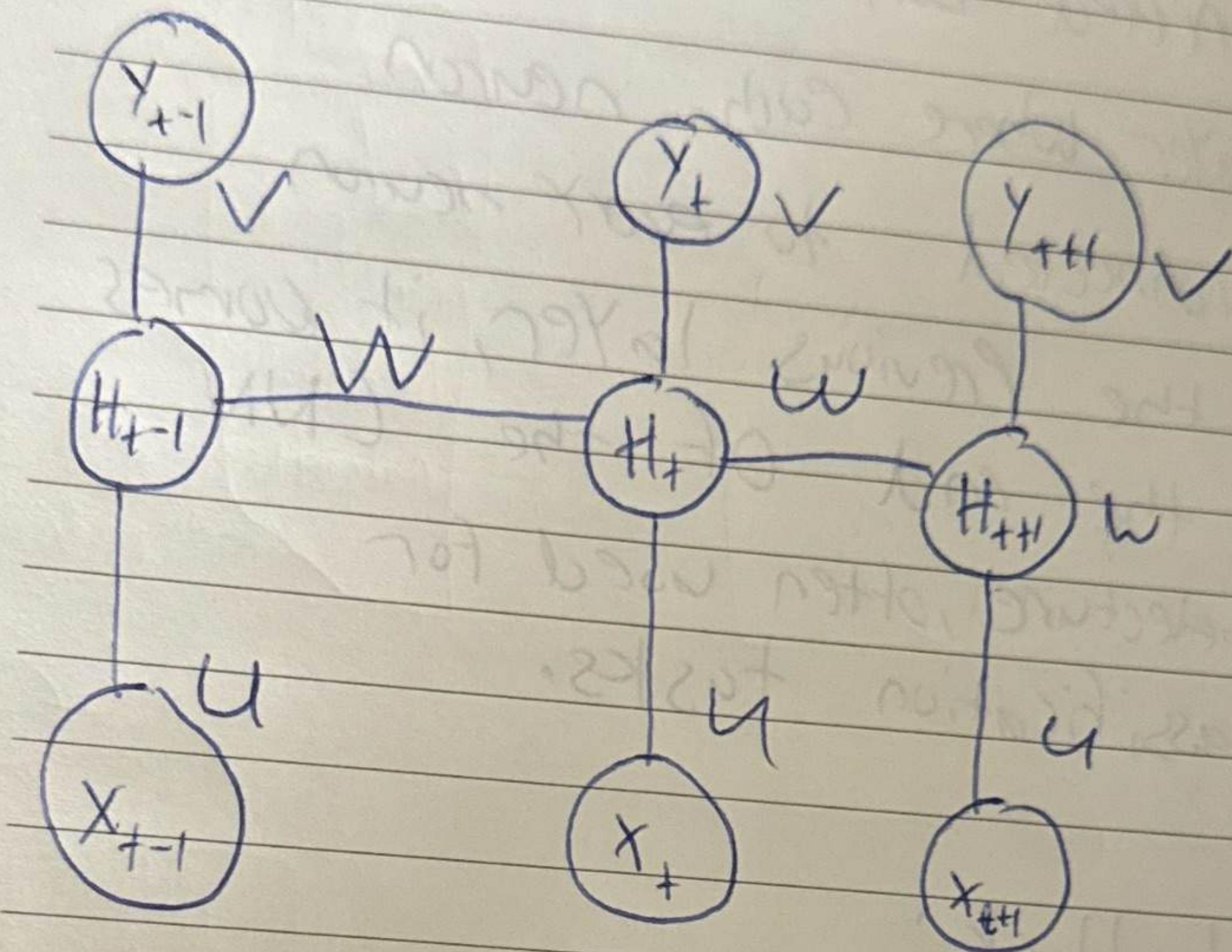
## Fully Connected Layer:

A layer where each neuron is connected to every neuron in the previous layer, it comes at the end of the CNN architecture, often used for classification tasks.

## Residual block:

A building block in architectures like ResNet that include skip connections, allowing the gradient to flow through the network more effectively

# Recurrent Neural Network



$u$  = weight vector for hidden layer

$v$  is weight vector for output layer

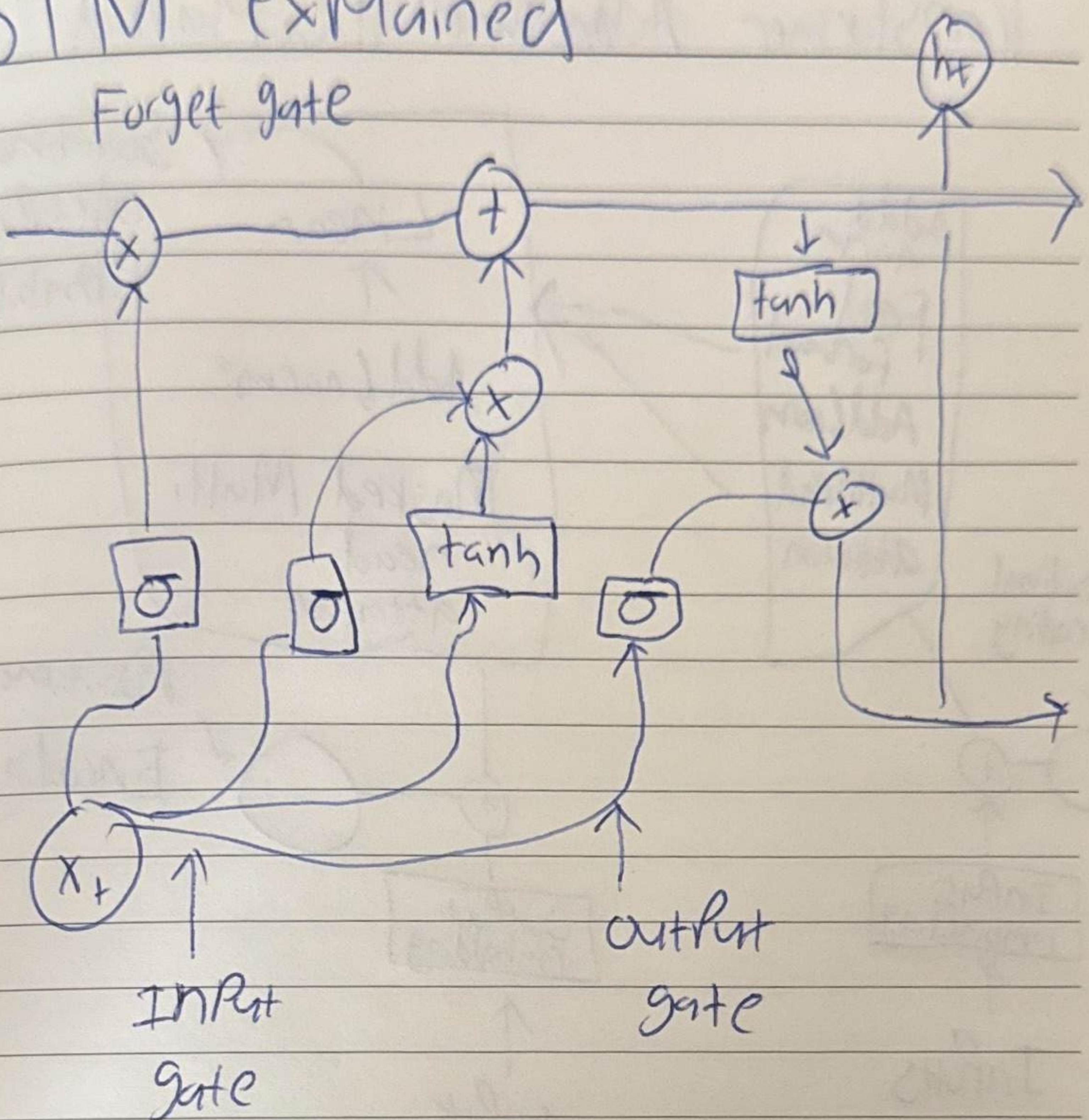
$w$  is the same weight for many timesteps

$x$  is the word vector for input word

$y$  is the word vector for output

# LSTM explained

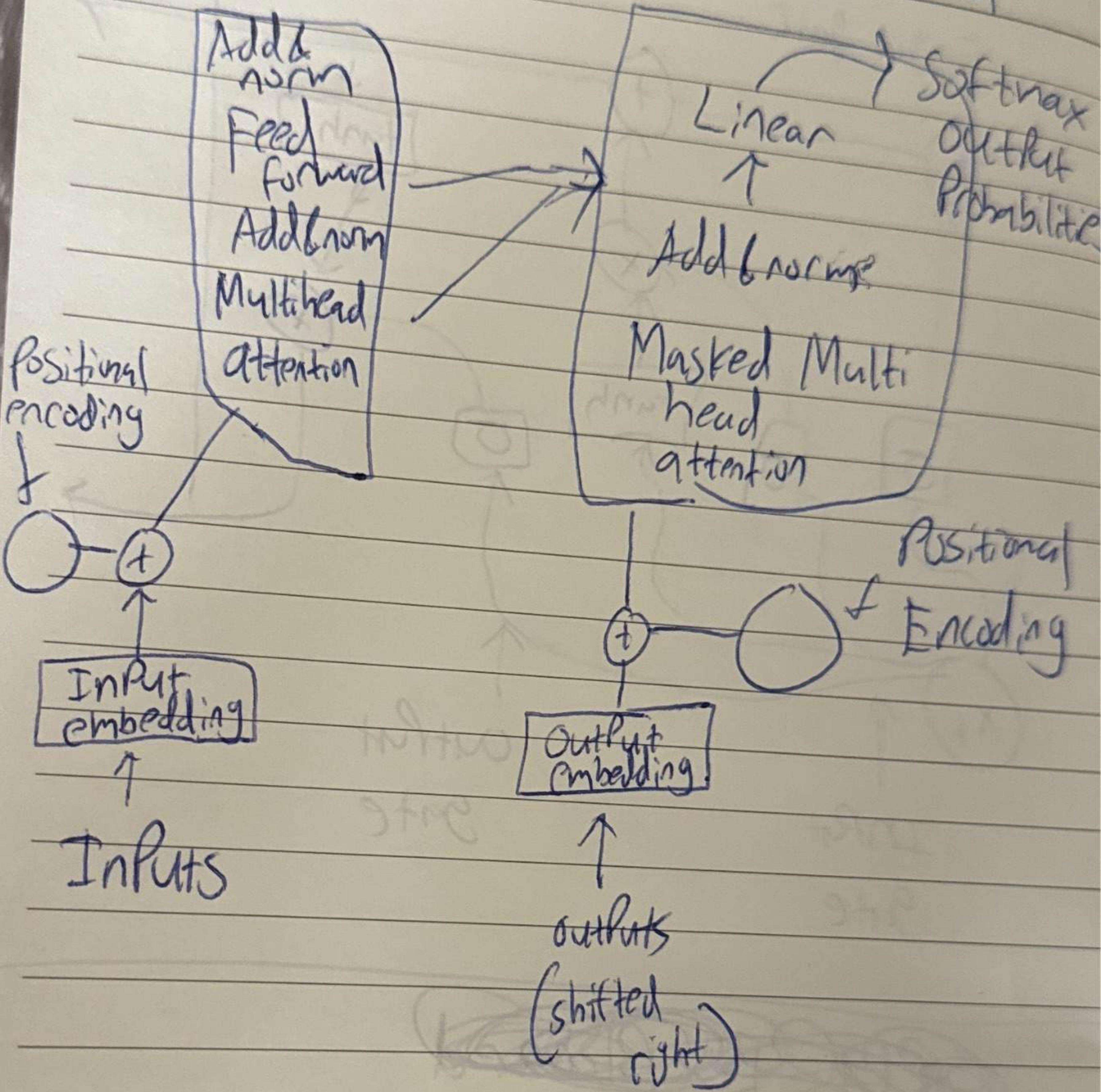
## Forget gate



~~GRU~~ LSTM

!

# Transformer Architecture Explained



# DATA SCIENCE Fundamentals

## for ML

### 1. DATA COLLECTION

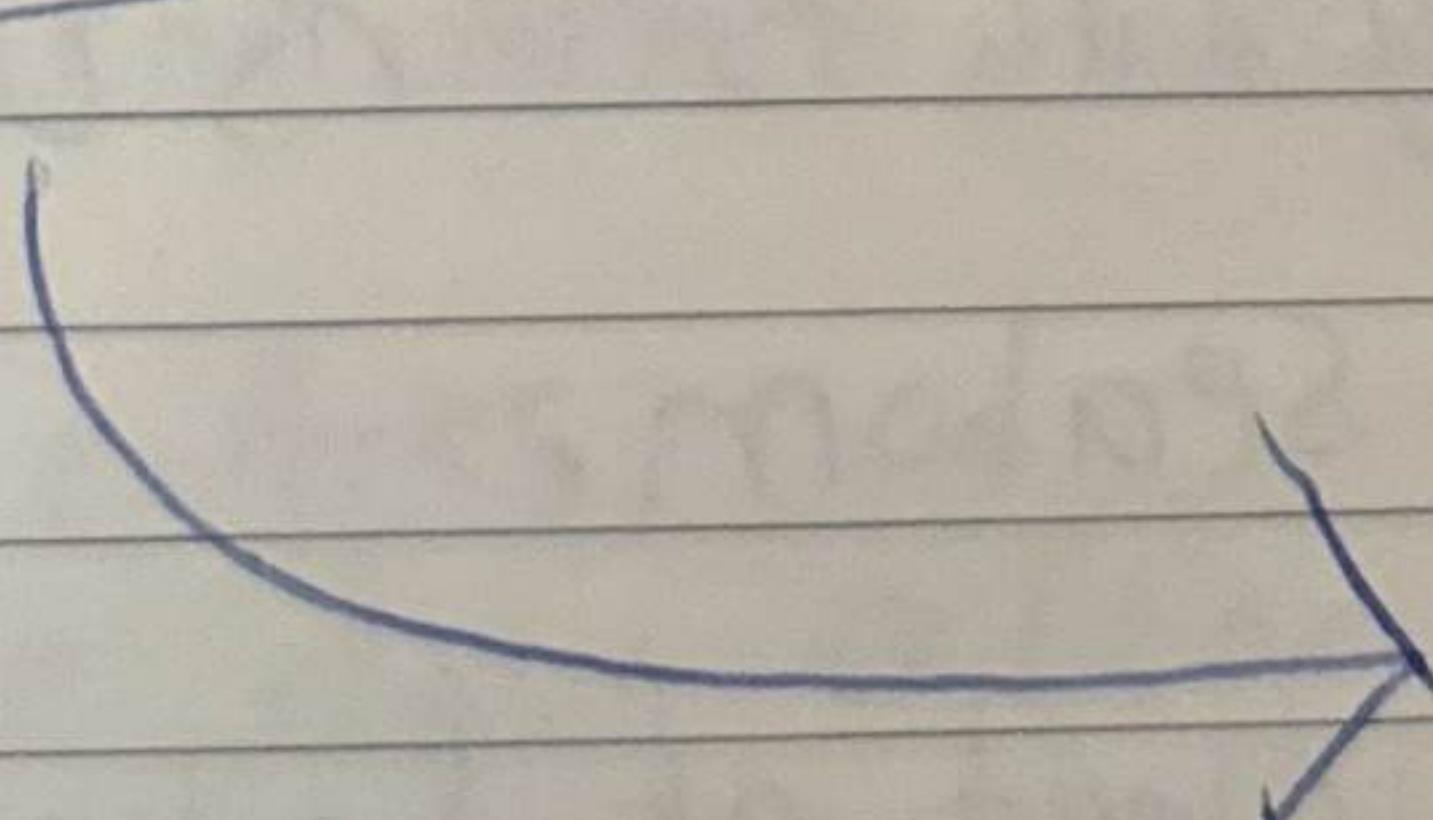
Gather relevant data from various sources:

Databases: SQL or MySQL, where it is stored in structured formats

API's: Interfaces to retrieve data from external websites

Web scraping: extracting data from web directly

CSV/excel: import data from flat files



## 2. DATA (Learning)

Ensure Data quality & integrity

Handling missing values

- Removing records:

Delete rows or columns with excessive missing values

- Imputation techniques:

- Mean, Median, Mode & range

- K-Nearst Neighbors

- Interpolation

Use heat maps for visualizing

missing data patterns using libraries

like seaborn.

Keep format of data consistent.

# Handling outliers in DATA

use IQR Method to identify outliers outside of the range:

$$[Q1 - 1.5 \times IQR, Q3 + 1.5 \times IQR]$$

Capping: Limit the value of outliers by capping them at certain thresholds.

Transformations: APPLY transformations to reduce the effect of outliers  
eg: Log transformation.

Cross-Reference Data to make sure it is accurate.

Documentation of cleaning processes in data.  
Automate repetitive data cleaning tasks.

Involve domain experts.

Always test the data after cleaning the data.

# ~~DATA~~ Transformation

## Normalization

This rescales the data to a specific range typically  $[0,1]$ . This is important for various tasks.

## Methods:

### Min-Max scaling:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Used, when the data does not follow a Gaussian distribution & when it is important to have features on same scale.

Ex: Converting range of ages from a range of  $(0,1)$ .

## Standardization of DATA

this transforms data to have a average of 0 and a standard deviation of 1. This is important for algorithms that assume normally distributed data, such as

LINEAR REGRESSION,

LOGISTIC REGRESSION,

Principal component analysis.

Methods:

Z-score Normalization:

$$x' = \frac{x - \mu}{\sigma}$$

$\mu$  is the mean of the data and  $\sigma$  is the standard deviation of the data

This is effective when features have different units & scales.

## Categorical Variables

Categorical variables need to be converted into numerical format.

### Methods:

- One-hot encoding:

Creates binary columns for each value, useful for data without intrinsic ordering (e.g color, type)

Ex. `['Red', 'blue', 'green']` create three columns is-red, is-blue, is-green.

- Label encoding:

This is good for data which categories have a clear order (e.g 'low', 'medium', 'high').

Example: `["low"=0, "Medium"=1, "high"=2]`

Other important tips:

### Feature Extraction:

- Create new features from data to enhance the model's ability to capture patterns.

use TF-IDF for text data.

### Feature engineering

use knowledge to create new features that can improve a model's performance.

## DIMENSIONALITY REDUCTION

- Reduce the number of features while preserving essential information, which can help prevent overfitting.

Methods:

PCA, t-distributed Stochastic neighbour embedding

### DATA Integration

- Combining data from various sources to create one coherent dataset.

# DATA SPLITTING IN ML

Generalization: This can help the model make predictions on new unseen data.

Model evaluation: This can help model evaluation on new unseen data.

Avoiding overfitting: This can help avoid overfitting as model's need to train on other data than the train set.

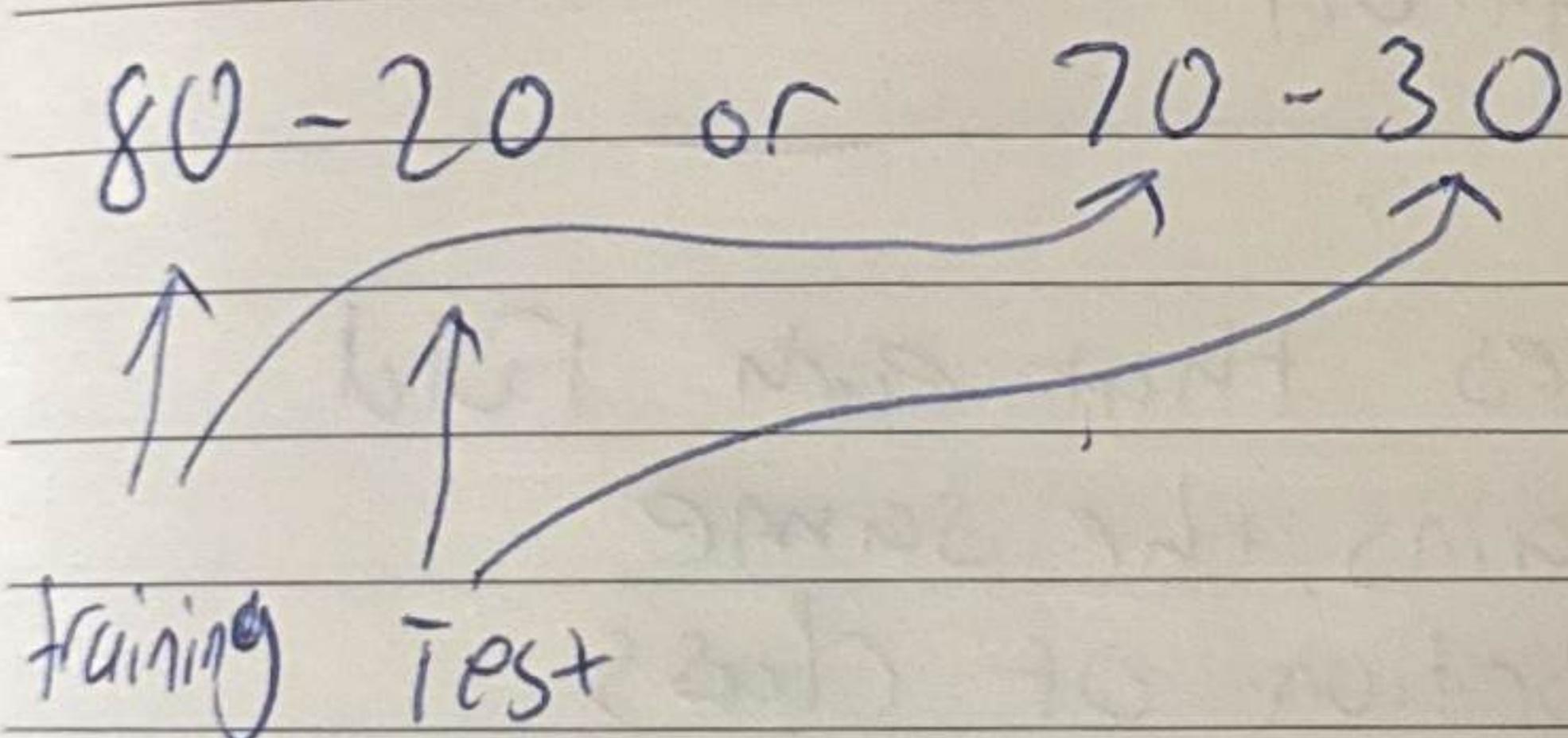
Training set: This set of the data is used to train the ML Model.

Validation set: This is used to tune hyperparameters and make choices about the model architecture.

Test set: this is used after the model has been trained & tuned, this provides a final evaluation on how well the model is doing in real-life.

### SIMPLE Train-Test Split:

Divide the Set like this:



\* Do not use this if the dataset is small.

### K-Fold cross validation

Dataset is divided into K-Folds, then it is trained K times using K-1 Folds for training and Rest for validation.

The issue is that this is computationally expensive.

Reduces Variance in Performance estimates.

- All data points are used for training & validation.

### Stratified K-Fold cross validation

This ensures that each fold maintains the same proportion of class labels as the original dataset.

Essential for imbalanced datasets.

Pros: Preserves class distribution

Cons: slightly more complex to implement than standard K-fold.