

Date: \_\_\_\_\_

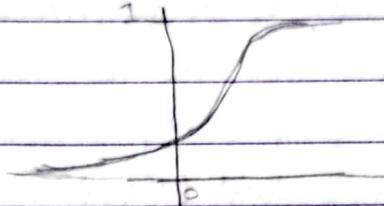
Sun Mon Tue Wed Thu Fri Sat

## Activation

functions.

- 1- Sigmoid: Squashes input in range of 0 to 1  
used for binary classification

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

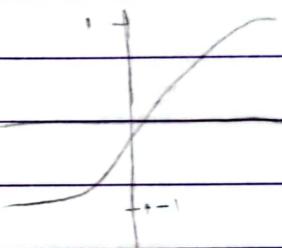


- vanishing gradient

- 2- Hyperbolic Tangent (Tanh): Squashes input into range between -1 & 1, often used in RNN

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

- vanishing gradient



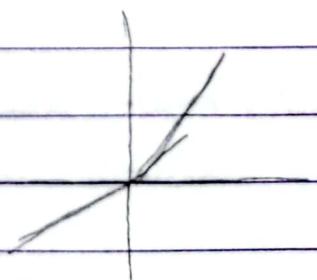
- 3- Rectified Linear Unit (ReLU): outputs the input if positive  
negative input converted to 0, commonly used in  
deep learning models

$$\text{ReLU}(x) = \max(0, x)$$

Problem: neurons dying. Best for HL

- 4- Leaky ReLU: Similar to leaky ReLU but allows small  
gradient when negative input

$$\text{Leaky ReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases}$$



fixed problem of dying neurons of ReLU.

less flexible

Best for hidden layer

Date: \_\_\_\_\_

Sun Mon Tue Wed Thu Fri Sat

5- ~~PRR~~ Parametric ReLU (PReLU): extends Leaky ReLU by learning slope of negative part.

when times perform

Leaky ReLU:  $\alpha$  is fixed.

PReLU:  $\alpha$  is learned to find optimal  $\alpha$

- Increased complexity longer training time

6- Exponential Linear Unit (ELU): Similar to ReLU but smoother outputs - helps in training stability - Also faster convergence

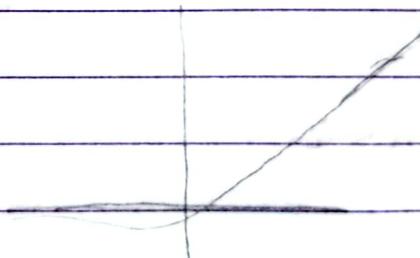
$$\text{ELU} = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{otherwise.} \end{cases}$$

- Slightly more complex than ReLU or Leaky ReLU

1- Swish: proposed by google applies Sigmoid to the input and multiply itself

$$\text{Swish}(x) = x \cdot \text{Sigmoid}(x)$$

good alternative to ReLU



3- Softmax: used in multiclass classification, converts scores into probabilities summing to 1.

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

- Output layer

- cannot be used

for regression

Date: \_\_\_\_\_

Sun Mon Tue Wed Thu Fri Sat

- 9- Softplus: Smooth approximation of ReLU, helpful in optimization due to its differentiability

$$\text{Softplus}(x) = \log(1+e^x)$$

- 10- Gaussian Error Linear Unit (GELU): A smooth approximation of ReLU with gaussian noise claimed to improve performance.

$$\text{GELU}(x) = x \cdot \text{CDF}(x)$$

$\text{CDF}(x)$  = the cumulative distribution function of the standard normal distribution.

- good in hidden layers - faster by perform better than  
ReLU

## Loss functions

1. Mean Squared Error: Average squared difference between actual and predicted values.

- used in regression
- sensitive to outliers.
- Not ideal for skewed distribution.
- Best where outliers are minimal & normal dist

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

2. Mean Absolute Error: Average difference between actual & predicted values.

- less sensitive to outliers.
- & Not differentiable, cannot be used in optimization algorithms.
- Not capture overall error as well as MSE
- Best for where outliers are a concern.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

3. Huber Loss: Combines MSE & MAE.

- uses MAE for small error & MSE for large
- remains differentiable for optimization
- hyperparameter delta as threshold for transition
- Regression with outliers

$$\begin{aligned}
 P(A) &\rightarrow [P(A) \ 0 \ 0] \\
 &\rightarrow [0 \ P(B) \ 0] \\
 &\rightarrow [0 \ 0 \ P(C)]
 \end{aligned}$$

Date:

Sun Mon Tue Wed Thu Fri Sat

$$\text{HL} = \begin{cases} \frac{1}{2}(y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta|y - \hat{y}| - \frac{1}{2}\delta & \text{otherwise.} \end{cases}$$

4- Binary Cross entropy Loss: Measure average difference between predicted probability of correct class

$y_i \ 1$

- Performance of binary class
- Not suitable for multiclass.

$$\text{Binary CE} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

5-Categorical cross entropy: Generalization of binary CE for multiclass classification.

- Measure the average difference between the predicted probability distribution by the one hot encoded distribution.
- Sensitive to imbalanced classes.

$$\text{C-CE} = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m y_{ij} \log(\hat{y}_{ij})$$

6-Sparse Categorical cross Entropy: Variant of CCE with sparse targets - only considers loss for predicted class.

- More efficient than CCE
- Not good for all multiclass classification tasks.
- Best for large number of classes

$$n = 1 \sum^h \log(p_i)$$

7. Kullback-Leibler Divergence: Measure difference between 2 probability distribution

- used in IR, Style transfer where comparing distribution is crucial
- Not symmetric  $KL(A||B) \neq KL(B||A)$
- can be undefined in some cases.

$$KL(P||Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right)$$

8. Hinge loss: Used in SVM, Penalizes model for misclassification

- May not be as smooth as other loss function
- Slower convergence
- Best for SVM where maximizing

$$\text{Hinge} = \max(0, 1 - y \cdot \hat{y})$$

9. Contrastive loss: Used in similarity learning

- Enforces smaller distance in similar datapoint (+ve pair)
- Enforces larger distance in dissimilar datapoint (-ve pair)

- Careful selection of the +ve pairs
- computationally expensive for large datasets
- Best for Siamese networks, facial recognition, IR

$$L(y, \hat{y}) = (1-y) \cdot \frac{1}{2} \hat{y}^2 + y \cdot \frac{1}{2} \max(0, m - \hat{y})^2$$

10. **Triplet loss:** used in triplet networks, minimizes the distance between anchor by positive example while maximizing distance between anchor by negative example.

- Computationally expensive
- Best for Image retrieval by facial recognition
- Doesn't require class label - only need to know if 2 points are similar.

$$L(x_a, x_p, x_n) = \max(0, d(x_a, x_p) - d(x_a, x_n) + \alpha)$$

11. **Adversarial loss:** core concept in GANs.

- 1. **Discriminator loss:** This measure discriminator's ability to distinguish real & fake data.  
often: binary-crossentropy.
- 2. **Generator loss:** This is based on how well the generator has fooled the discriminator  
normally negative of discriminator's loss
- Training instability: finding right balance between the generator by discriminator. Sometimes one overpower other.
- Some cases generator stuck in loop ↗
- Best for GAN: Image, text, for generation

## Optimizers

### 1- Stochastic Gradient descent: (SGD)

- Updates parameters by taking step in direction of negative gradient of loss function.
- Slow convergence, noisy updates.
- Sensitive to learning rate.
- =  $\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$
- Simple to implement
- computationally inexpensive
- Best for baseline for ~~other~~<sup>comparison</sup> with others

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$

### 2- Mini Batch Gradient Descent:

- Similar to SGD but updates parameters in mini batches instead individually.
- More efficient, less noisy
- Problem in choosing optimal mini batch size.

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$

where  $J(\theta_t)$  = loss of mini batch

## 3- Adam:

- Adapts learning rate for each parameter based on past gradients
- faster convergence, better performance
- can be aggressive leads to large parameter updates
- good general purpose optimizer

$$m_t = \beta_1 \cdot m_{t-1} (1 - \beta_1) \cdot \nabla J(\theta_t)$$

$$v_t = \beta_2 \cdot v_{t-1} (1 - \beta_2) \cdot \nabla J(\theta_t)^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_{t+1}} + \epsilon} \cdot m_{t+1}$$

- $\beta_1$  &  $\beta_2$  : hyper parameter for momentum & forgetting
- $m_t$  &  $v_t$  : mean & variance of past gradients

## 4- RMSprop:

- Improves on SGD by adapting Learning Rate of each parameter by using history of squared gradients
- Solves the vanishing gradient problem.
- Can be slow to converge
- Might not be as good as adam

Date: .....

Sun Mon Tue Wed Thu Fri Sat

$$v_t = \beta \cdot v_{t-1} + (1 - \beta) \cdot (\nabla J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \cdot \nabla J(\theta_t)$$

### 5- AdaGrad:

- Adapts LR for each parameter
- Based on square root of the sum of past squared gradients
- lead to very low LR with each update
- Not suitable for sparse model.
- Rarely used

$$g_t = g_{t-1} + (\nabla J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{g_t + \epsilon}} \cdot \nabla J(\theta_t)$$

### 6- Adadelta:

- Extension of AdaGrad
- Reduces the problem of vanishing LR
- Only look at a window of past gradients
- Requires tuning

$$\Delta \theta_t = \frac{\sqrt{E[\Delta \theta^2]_{t-1} + \epsilon} \cdot \nabla J(\theta_t)}{\sqrt{E[g^2]_{t+\epsilon}}}$$

## 7- Adamax:

- Similar to Adam but uses infinity norm of past gradients instead of the mean
- In certain cases large gradients can effect Adam's performance
- Need to find optimal settings.

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot \nabla J(\theta_t)$$

$$u_t = \max(\beta_2 \cdot u_{t-1}, \|\nabla J(\theta_t)\|_\infty)$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{u_t + \epsilon} \cdot m_t$$

## 8- NAdam:

- Combines Nesterov momentum technique, with Adam,
- Might not provide significant benefit over adam
- Can be used for faster convergence

$$\theta_{t+1} = \theta_t + \hat{g}_t + \beta_1 \cdot \frac{(1 - \beta_1) \cdot \nabla J(\theta_t)}{(1 - \beta_1^2) \cdot (\nabla J(\theta_t))^2}$$

Date: .....

Sun Mon Tue Wed Thu Fri Sat

9- RAdam: Solve the issue in Adam where the gradient of mean become 0

- Rectifies the variance of estimation
- leads to better generalization

10- Ranger: (RAdam + lookahead)

## AutoEncoders

- Special type of NNs where input is same as output
- Auto encoder consists of
  - 1- Encoder: compresses the input by produces the code.
  - 2- Code: Efficient representation of input data latent vector.
  - 3- Decoder: Reconstructs the input only using this code.
- Unsupervised learning , unlabelled data.

### Why use Autoencoder:

- Useful for dimension reduction.
- powerful feature detector (feature selection)
- can be used for unsupervised pre-training of DNNs
- Capable of generating new data similar to training data, can be used as generative model.

### PCA vs Autoencoder:

- Non linear activation by multiple layer ,
- Can learn from convolution layer.
- More efficient
- gives multiple representation as output of each layer

Date:

Sun Mon Tue Wed Thu Fri Sat

- There are 2 configuration of AutoEncoders:

### 1- Under complete:

- Dimension of latent space is less than the dimension of input space
- Model forced to produce a compressed representation of input
- used for dimension reduction, feature extraction
- less prone to overfitting
- May struggle to capture all variation of input

### 2- Over complete:

- Dimension of latent space is more than dimension of input space.
- More parameter than undercomplete.
- produces a latent <sup>representation</sup> space of more complex features than input
- Prone to overfit.
- Regularization, Dropout by sparsity used for generalization

## Types of AutoEncoders:

### 1- Vanilla AE: - Basic AE ,

- fully connected dense layers

### 2- Sparse AE: - Sparse <sup>constraints</sup> ~~constraint~~

- Good for feature extraction

Date:

Sun Mon Tue Wed Thu Fri Sat

3- Denoising AE: - Trained to reconstruct clean data from noisy input

4- Variational AE: - Probabilistic AE

- learns prob distribution over latent space
- Widely used for generative modeling tasks

5- Convolutional AE:

- Conv layers instead of fully connected layers
- used for image related tasks
- capture hierarchies of features

6 Recurrent AE:

- Incorporates RNN layer in either Encoder or decoder or both.
- Good for sequential data, time series, text, audio

7- Adversarial AE:

- concept of AE with adversarial training
- Encoder, decoder by A discriminator.
- GAN

8- Spiky lab: - sparse coding  
- feature are binary

## GAN

= GAN consists of 2 NNs trained simultaneously in competitive manner

1- Generator: Takes random noise as input and tries to generate data sample similar to real data

2- discriminator: Receives real sample by fake sample from generator and distinguishes between real & generated

- During training gen by dis are trained iteratively
- trained on minimax game framework. generator minimizes the likelihood of discriminator correctly classifying samples while discriminator maximizes its accuracy
- Application : Image generation, style transfer, translation

Problems:- Training instability: one network overpower other

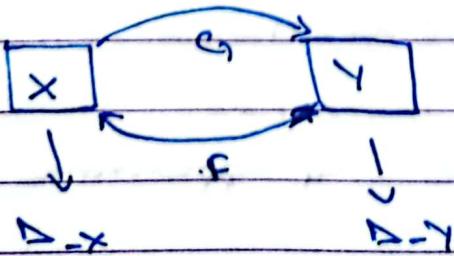
- Mode collapse: generator stuck in loop, generates similar data instead diverse realistic samples
- Hard to find balance in both Network's loss

## Conditional GAN:

- Both Network receive conditional information.
- Conditions: guide for generation process.
- Generator produce realistic sample while following conditions.
- Discrimination learns to distinguish between real & generated data while considering conditions.
- minimizes 2 losses
  - 1- adversarial loss
  - 2- Conditional loss
- Used when specific characteristics desired in generated sample.

## Cycle GAN:

- learns to translate image from one domain to another.
- 2 generators, 2 discriminators.
- generator



- Cycle consistency loss: ensures translation process is consistent in both directions.
  - helps to preserve semantic content during translation
- used in style transfer, object transfiguration.

### Types of GANs:

- 1- Vanilla GAN: Basic GAN.
- 2- CGAN: - Conditional information
  - class labels, text description other info
  - Allows for more controlled generation of samples
- 3- WGAN: - Wasserstein loss as metric
  - Addresses some instability by mode collapse
- 4- Cycle GAN: - unpaired image translation
  - cycle consistency loss
  - 4 Networks.
- 5- DCGAN: - Deep convolution NN as disc by gen
  - More stable by effective
- 5- Style GAN: - control style by appearance of generated image
  - fine grain control like facial expression

- 7- StarGAN: - Multi domain image translation
- 8- BigGAN: - Scaled up version of GAN
  - Used to improve image quality / resolution.

## RNNs

- Short-term dependency: immediate or near immediate relationships between adjacent elements.
- long-term dependency: relationship between elements separated by several steps. or time intervals.
- long-term dep are particularly because traditional ANNs often struggle to capture them due to vanishing gradients.

RNNs: Cases of ANN designed to handle sequential data by

- Maintains hidden state that capture info about previous state inputs.
- Connection that form a cycle allowing them to exhibit temporal dynamic behaviour
- used for sequences, time series, NLP,

working:

input: each element of the sequence is represented as a vector

Recurrent connections:

- At step  $t$  input vector  $x_t$  and previous hidden state  $h_{t-1}$  are combined to produce new hidden state  $h_t$
- $$h_t = f(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$W_{hh}$  &  $W_{xh}$  are weight matrices.

$b_h$  is bias       $f$  = activation function

Output: hidden state  $h_t$  can be used to generate output  $y_t$  which can be used to predict

BPTT:  
- Back propagates through time  
- unfolds the network over time

## Strengths:

- Capable of modeling sequential dependencies
- can handle variable length of inputs by can produce variable length of outputs.
- Maintains a hidden state that retains info about previous inputs.

## Weakness:

- vanishing grad / Exploding gradients
- difficulty in retaining info over long sequences.
- Training is challenging due to nature of sequential data.

## LSTM: - Type of RNN

- Solve vanishing gradient problem
- Capture long term dep in data
- Gated mechanism

- forget Gate: Takes input for current batch  $x_t$  and previous hidden state  $h_{t-1}$

- ~~backpropagates~~ produces a forget gate vector  $f_t$
- determines how much percentage of long term memory will be remembered

Cell State = LTM

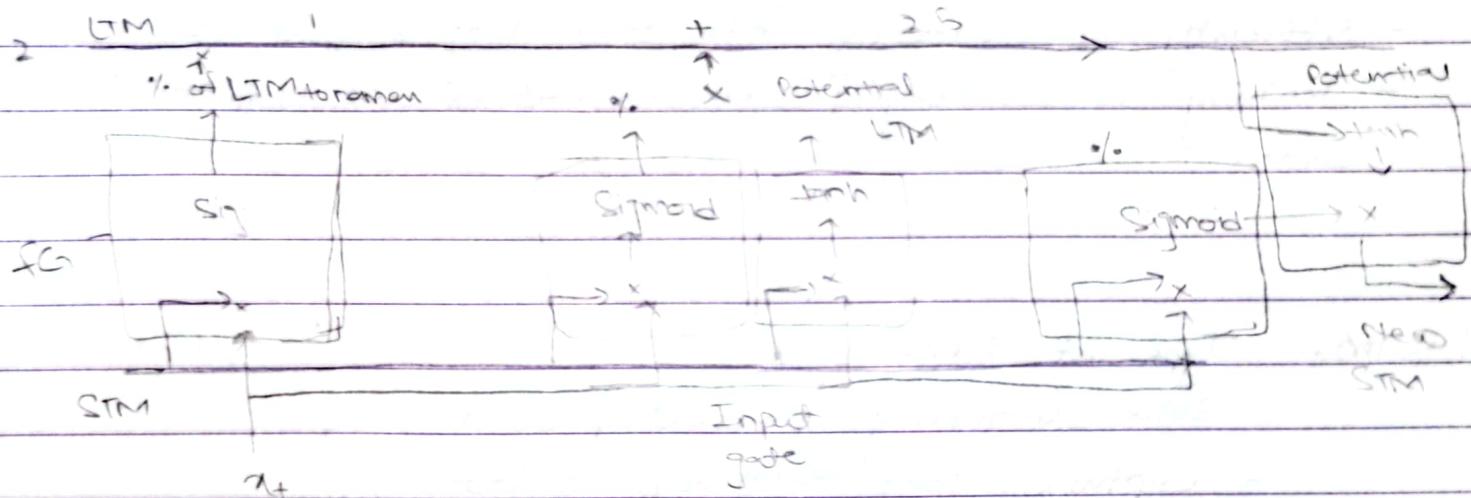
Hidden State = STM

Date:

Sun Mon Tue Wed Thu Fri Sat

2- Input gate: combines the short term memory with the input to create a potential long term memory  
And how much percentage of that potential LTM should be add to LTM

3 Output gate: updates the short term memory  
- uses new LTM to create a potential STM  
- uses input by STM to find how much percentage of potential STM should be passed on.



**Strength:**

- effective in capturing long range dep
- Mitigating vanishing gradient
- Maintains memory cell to store long sequences

**Weakness:**

- more complex than RNNs
- Require more computational power
- Sensitive to noisy input

Date:

Sun Mon Tue Wed Thu Fri Sat

GRU: Type of RNN, similar to  $\rightarrow$  LSTM

- Solve vanishing gradient by capturing long term dep.
- Simplify architecture of LSTM by combining the forget gate & input gate into single update gate
- And merging cell  $q$ , hidden state,

update gate: uses input  $x_t$  and  $h_{t-1}$  with sigmoid to determine how much of previous hidden state to retain and how much to update.

- Determine trade-off between retaining vs updating hidden state.

Reset gate (optional): Some variations include this gate decides how much of past info to forget

Date:

Sun Mon Tue Wed Thu Fri Sat

### Strength:

- Simpler architecture to LSTM • easier to train
- fewer parameters. more efficient
- Effective longterm dep

### Weakness:

- May not capture long term dep as good as LSTM
- less control over memory

Bi-Directional RNN: Type of RNN that processes in both forward by backward direction

- Allows them to capture information from both past by future context.
- Input sequence is processed by 2 separate RNNs.
  - 1. forward 2. backward.
- each independently update their hidden state
- The hidden state of each RNN at each time step captures info from both past by future
- Both info is then concatenated to produce final representation.

Date:

Sun Mon Tue Wed Thu Fri Sat

### Strength:

- Better <sup>under</sup> ~~better~~ stand context
- Improved performance
- Robust to noise

### Weakness:

- More complex, high comp cost
- access to future may lead to overfitting.

## GNN

- Graph neural Network one class of DNN used to inference on data in Graph
- Used to predict node level, edge level, graph level task
- GNN do what CNN failed at

### Why CNN Failed in graphs:

- Arbitrary size
- complex topology
- unfixed node ordering

GNN operate on graph structure : node  $\rightarrow$  entity  
relation : edge  $\rightarrow$  relationship

- Message passing: at each layer nodes exchange information
  - achieved through propagation.
  - Node aggregate its neighbour info to update its own representation
- Node imbedding: Representation of graph in low dimension
- Aggregate function: used to combine info from neighbouring node