

* Programming :- We give data and program as input to the computer to get desired output

Machine learning :- We give data & desired output (or) expected output to computer to get program as an output

* Supervised - Training data with desired output values that are provided before training

Unsupervised - Training data without desired outputs
Semi-supervised - Training data with small desired outputs and iteratively train the ~~the~~ model using the most confident predictions

→ Train initial ~~the~~ model on a small set of labeled data.

→ Use this model to generate predictions on a larger set of unlabeled data.

→ Select the most confident predictions and add them to labeled data set

→ Retrain the model on expanded dataset

→ Repeat step 2 and 4 until desired output is achieved

* types of learning :- i) Active learning ii) Federated learning
iii) Reinforcement Learning

i) Active learning - Train data with and without labels; interactive labeling of the training data by asking a human expert to label the least confident data samples

Ex:- Pool-based active learning ; Diversity sampling ;

Stream-based active learning

ii) Federated learning - Training data is localized ; learn a central model using distributed local data without sharing data

* Reinforcement learning: an agent learns to interpret its environments by interacting with it through actions are rewarded (or) penalized.

Examples of Supervised learning

* Predict whether a patient, hospitalized due to heart attack, will have a second heart attack.

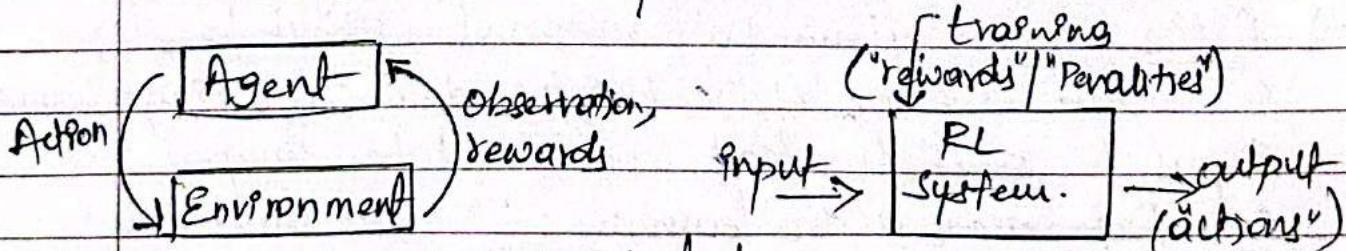
Identify strongest predictors of heart attack

Predictions based on many data modalities :- demography, diet, clinical measurements, clinical history

* Predict ~~stocks~~ the price of stocks in 6 months from now, based on i) Performance of the company (ii) other economic metrics (iii) other exogenous variables.

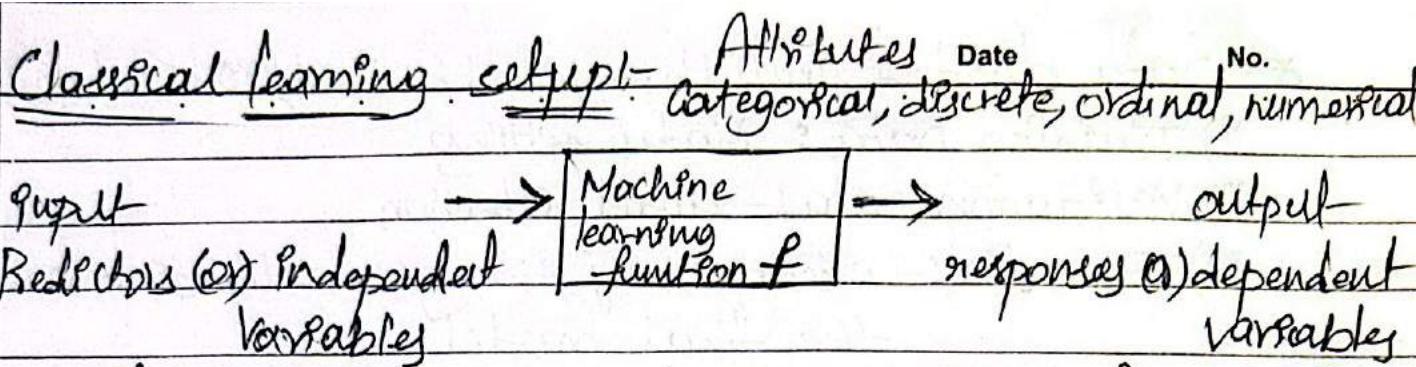
* Identify handwritten numbers in a digitalized image & identify animals in image

* Examples of Reinforcement learning: - Define the measures to be taken by the public health authorities under a pandemic situation.



Examples of Federated learning:-

Given data from distributed hospitals learns a central model that can propose the optimal patient treatment without sharing any data.



Classification:- It's task of learning a target function f , that maps an input (attribute set) X to one of the predefined class labels Y .

Regression:- It's task of learning a target function f that maps an input (attribute set) X to ~~one of the~~ a value in range of real values Y .

- * The Bias variance decomposition- The metric to optimize is Mean Squared Error (MSE)
- * How good is f ? - If training error drops to zero as model complexity increases lead over-fitting (not good)
 - i) High bias \rightarrow under-fits
 - ii) High variance \rightarrow over-fits
 - iii) Low bias & low variance \rightarrow good ; This can be achieved by selecting appropriate model complexity, using regularization techniques & selecting it's right features. In order to achieve a goodness of machine learning function, can be evaluated based on it's ability to model relation b/w the input features & targeted variables. & bias, variance properties

* Model selection & assessment-

Model selection- estimating the performance of different models in order to choose the best one.

~~Model validation~~ Model assessment- having chosen a final model, estimating its prediction error on new data.

* Data divided into three parts :- Test, Validation, Train

i) Training data :- Model fitting

ii) Validation data :- Model selection

iii) Test data :- assessment of the generalization error of the final model

* Hold out and subsampling - Two common techniques to create training and test sets.

Holdout - Reserves 50% for training 25% each for test & validation.

Random subsampling - Repeated holdout

Stratified subsampling - Reserves class balances in the samples.

* K-fold cross validation :-

1) Divide data into K equal parts (or) folds.

2) Iterate K times :

i) Take the k^{th} fold as validation set, remaining as training set

ii) Train the model using training data

iii) Test the model using validation set

3) Average the accuracy of all K iterations & use it as the performance metric.

* Comparing learning curves :-

High Bias :- As training examples increase,

training error increases too; model is not fitting the data with an appropriate curve. Bias leads to comparable errors in cross validation & training. Adding more data does not help.

MATRIXAS

Increase no. of features and more complex features.

(high variance) Test error grows slowly but well within the desired range as no. of training examples increases. High variance leads to overfitting & hence high test error the model is forced to learn more generalized properties. ~~But get more data and reduce no. of features~~

- * K-fold cross validation pitfall: Data leakage - If information from test set is inadvertently include in the training set during K-fold Cross Val. This can be happened if the data processing (or) feature selection is done on entire dataset, rather than just the training set.
- * Nested Cross validations - It involves an outer loop of K-fold cross val for evaluating the generalization performance of the model & inner loop for tuning the hyperparameters of the model. The inner loop is performed on each outer loop and is used to select the best hyperparameters for the model.

We deploy the feature set obtained by applying the winning procedure to the whole dataset.

- * Bootstrapping: Create B replications of same size of training set of size N by selecting examples uniformly with replacement. ~~From the each replacement B~~
 - Train on each replication b
 - Test on whole training set

Leave One Out Bootstrapping: Train on each replication b test each sample not in b.

Problems: #1 The avg no. of distinct observations in each bootstrap sample is about $0.632N$, it's bias slightly below 1.0 as two-fold cross val; if the learning curve has considerable slope at sample $N/2$, leave one out bootstrap will be overestimate of the true error.

Q. 6.3.2 Bootstrap! - extension of leave-one-out Bootstrap.

Pulls the leave-one-out Bootstrap estimate down towards the training error rate, and hence reduces its upward bias.

Lec-2

Example of supervised learning:- i) Quantitative, ii) Categorical,
iii) Quantitative :- stock price ; Categorical :- heart attack
↳ outcome measurement :> Predict based on a set of features.

* Regression when predict quantitative outputs.

Classification - ii. v. Categorical outputs (qualitative).

* Regression- Statistical measure that attempts to

determine the relation b/w dependent variables Y & independent variables X ; used to predict next value of dependent variable Y from the value of independent variable X .

→ Dependent variable Y :- Is the variable whose value changes as the consequence of change in other values in system.

→ Independent variable :- x_1, x_2, \dots, x_n are input to the system & may take diff values freely. Also called as predictors.

* Simple Linear Regression- Mathematical relation b/w the variables X & Y is a linear relationship.

X - cause ; Y - effect in cause-and-effect relation. Model relation b/w two variables by fitting a linear equation to the observed data.

* How do we find linear model?

We attempt "best fit line" by minimizing the actual and predicted values. But positive differences could offset negative ones so we take the squared difference.

- * Ordinary Least Squares (OLS) minimizes the Residual Sum of Squares (RSS) of observed values to the straight line.
- * When there is single input variable it is referred to simple linear regression, when there are multiple input variable then it is Multiple Linear Regression.

* Cognitive function - The abilities of AI model to perform tasks that are typically associated with human cognition, such as perception, reasoning, learning, decision-making, language processing.

Gauss-Markov theorem - When all assumptions hold, OLS will produce better coefficient estimates compared to all other linear model estimation methods.

* Gradient Descent - It's an optimization algorithm that finds the linear regression coefficients iteratively. This is used by model to reduce the loss function & achieving the best fit line. To check if it's working is to check the error rate decreases as for each iteration.
Working process of gradient descent :-

- i) It assigns an initial random value to the coefficients
- ii) It iteratively updates the coefficient values proportional to the negative of the gradient of the function.
- iii) Finally it finds a locally (or globally) minimum for the loss function.

How good the fit is? - The coefficient determination can be used to determine how well the model fits the training data. Denoted by R^2 ; $R^2 = \frac{\text{Variance explained by model}}{\text{Total variance}}$

Problems with fitting the data - Overfitting, Underfitting

Overfitting: - the model models the training data too well;

If we have large ~~data~~ no. of features & the test score is worse than the training score.

Underfitting: - Model that can't model the training data nor generalize to new data. If we have fewer features and the test score & training score is poor for both.

Regularization: - Not satisfied with Ordinary Least Squares (OLS): - (i) Low performance (ii) Interpretation. Regularizing the coefficient estimates \rightarrow i) Ridge Regression (ii) Lasso Regression

Ridge Regression: - Uses a penalty term, called a regularization parameter to limit the complexity of the model, which prevents overfitting. It works by adding the sum of the squares of the coefficient to the loss function which is then minimized during training.

The penalty term adds small bias to the model.

Lasso Regression: - adds absolute values of the coefficients. This has the effect of shrinking some of the coefficients to zero, effectively performing variable selection and reducing no. of features in the model.

Conclusion: - Ridge adds a penalty term to limit the size of the coefficients while lasso adds a penalty term that eliminates some of the coefficients all together.

Example of Ridge Regression:- Average ^{data} credit card debt balance is measured for each individual. We have information on several predictors; age, no. of cards, income, credit rating.

*Evaluation Metrics:- 1) Mean Squared Error (2) Root Mean Squared Error (3) Mean Absolute Error (4) Mean absolute Percentage error

1) MSE- is calculated as the mean (or) average of the squared differences between predicted & expected target values in a dataset. The higher the MSE the larger the error in the model.

2) RMSE- The square root of the MSE, which means that the units of the error are the same as the units of the target value that is being predicted. The higher the RMSE the larger the error in the model.

3) MAE- Average of the absolute error values and like RMSE the units of the error does not match the units of the target value that is being predicted.

4) MAPE- The percentage equivalent of MAE. Each residual is scaled against the actual value.

Name	Residual operation?	Robust to outliers?
MSE	Square	No
RMSE	Square	No
MAE	Absolute value	Yes
MAPE	Absolute value	Yes

*Logistic Regression:- Statistical model that uses a logistic function to model binary dependent variable. It models the probability that Y belongs to a particular Category / class. The output of Logistic regression problem can be only 0 and 1.
How to fit?:- We use a method called "Maximum Likelihood"

Types of Logistic Regression

Date _____

No. _____

- i) Binary logistic Regression: - The categorical response has only two possible outcomes.
Example:- Spam or Not
- ii) Multinomial logistic Regression: - Three or more categories without ordering.
Ex:- Predicted which food is preferred more (Veg./Nonveg./Vegan).
- iii) Ordinal logistic Regression: - Three (or) more categories with ordering.
Ex:- Movie rating from 1 to 5.

Rating	Students
1	10
2	15
3	20
4	25
5	30

Rating	Students
1	10
2	15
3	20
4	25
5	30

Rating	Students
1	10
2	15
3	20
4	25
5	30

Rating - 1st: Ishan (10), 2nd: S. A. (15), 3rd: A. (20), 4th: B. (25), 5th: C. (30)

Rating - 1st: D. (10), 2nd: E. (15), 3rd: F. (20), 4th: G. (25), 5th: H. (30)

Rating - 1st: I. (10), 2nd: J. (15), 3rd: K. (20), 4th: L. (25), 5th: M. (30)

Rating - 1st: N. (10), 2nd: O. (15), 3rd: P. (20), 4th: Q. (25), 5th: R. (30)

Rating - 1st: S. (10), 2nd: T. (15), 3rd: U. (20), 4th: V. (25), 5th: W. (30)

Rating - 1st: X. (10), 2nd: Y. (15), 3rd: Z. (20), 4th: A. (25), 5th: B. (30)

Type of Logistic Regression-

Date

No.

i) Binary logistic Regression :- The categorical response may have only two possible outcomes.

Example :- Spam or Not

ii) Multinomial logistic Regression :- Three or more categories without ordering.

Ex :- ~~Prediction~~ predicted which food is preferred more (Veg./Nonveg./Vegan).

iii) Ordinal logistic Regression :- Three (or) more categories with ordering.

Ex :- Movie rating from 1 to 5.

Deep Learning :- Deep learning algorithms attempt to learn feature representations by using a hierarchy of multiple layers.

ANN :- (i) Single layer Perception (ii) Single layer of artificial neurons, with each one connected to every input feature.

(i) Summation function. E.g. dot product.

(ii) Bias used to shift activation function ($f(x+b) > 0$), to either right (or) left.

MLP - More complex problems may require intermediate layers.

Activation function :- Non-linear, network with one hidden layer can learn perfectly any classification problem. Problem is to find the proper no. of neurons & learn their weights.

$$(1) H = N$$

Linear activation function :- Not flexible enough to capture complex relationships b/w data variables.

Non-linear activation :- Makes it easier for the model to adapt as data varies.

Sigmoid activation :- Convenient for models that predict range (0) probabilities. Ex:- Range: (0,1). Softmax function is a more generalized logistic activation function which is used for multiclass classification.

Tanh activation :- Negative input is mapped to strongly negative values and zeros \rightarrow inputs \neq zero are mapped near to zero. Takes negative values. Range: (-1,1).

ReLU activation :- Both the function and its derivative are monotonic. Problem is negative inputs turn to zero immediately so can not map negative values appropriately.

Leaky ReLU Function :- leak helps to increase the range of the ReLU function. usually, a is 0.01 (leaky) or it may have any other value. Nonzero in nature both ReLU functions & derivatives

Forward Propagation :- In here, the input features are fed into the first layer of neural network which process the inputs and passes the result to next layer. Each layer ~~is~~ of the network consists of multiple neurons which calculates the weighted sum of their input and apply an activation function and gives an output. The output of each neuron is passed as input to the next layer until the final layer of network is reached. The output of final layer represents the predicted output of the neural network for given input features.

let x be the input features; $w \& b \rightarrow$ weights & bias of neuron in each layer of the network; f activation function.

y output of final layer of neural network.

$$\text{Output of first layer} \rightarrow z_1 = w_1 * x + b_1$$

$$a_1 = f(z_1)$$

$$\text{Output of second layer} \rightarrow z_2 = w_2 * a_1 + b_2$$

$$a_2 = f(z_2)$$

Date: No.

Back propal- The idea is to propagate the error back to all neurons. After forward propagation, the predicted output of neural network is compared with true output to compute errors. Then errors are propagated back through the network to update the weights and bias of the neurons in the output layer.

Regression loss functions - Huber loss: If the absolute difference b/w the actual & predicted value is less than or equal to a threshold value, δ , then MSE is applied. If errors are sufficiently large, then MAE is applied.

0-1 loss - Counts how many class labels the predictions are correct. It is not differentiable hence cannot optimize efficiently. Consider alternative functions that can be optimized also referred to as Surrogate loss function.

Binary Cross entropy - Classification neural networks outputs a vector of probabilities \vec{g}_i .

Categorical Crossentropy - Classification neural networks output a vector of probabilities \vec{g}_i :

$$\vec{g} \rightarrow \text{softmax} \rightarrow \text{Cross-entropy loss}$$

Regularization (L_2 or) L_1 :- L_1 not preferred as it will end up compressing the network L_2 to make the derivatives computation easier.

Regularization term adds penalties big weights, added to the objective weight decay value determines how dominant regularization is during gradient computation. Big weight decay coefficient \rightarrow big penalties for big weight.

Data augmentation! - Increase the size of the training data whenever possible. May result in improved predicted performance. (concrete cost of collecting and labelling data. Enables rare event prediction. Prevents data privacy issues)

Early stopping! - Use validation error to decide when to stop training. Stop when the monitored quantity has not improved after n subsequent epochs n is patience.

Validation error too high! - Training error zero, val error remains high: overfitting. model too complex; reduce layers and hidden nodes apply early stopping & dropout. Set threshold on training performance: increase training examples use regularization to penalize weights if weights get too large.

Stochastic Gradient Descent! - Randomly sample a small batch of data from entire dataset at each iteration and compute gradient of the loss function with respect to the parameters of the batch. The parameters are then updated using the computed gradient, and the process is repeated with another randomly sampled batch.

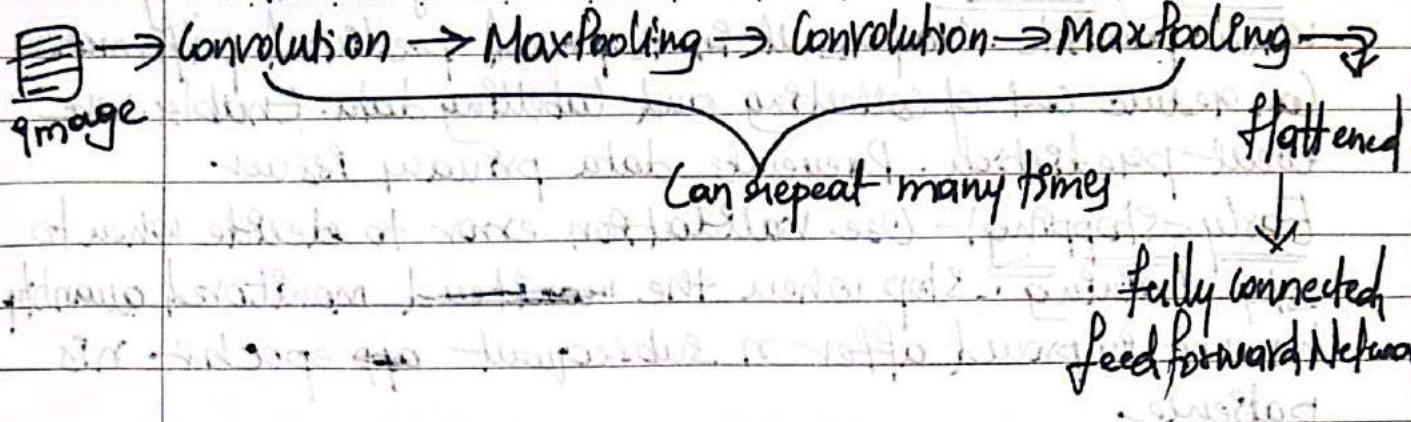
CCS-CNN: We need to detect the small patterns

Multilayer perception! - fully connected networks, each neuron in one layer is connected to all neurons of next layer. fully connectivity makes them prone to overfitting.

Hierarchical patterns! - Assemble patterns of increasing complexity. Using smaller & simpler patterns captured by filters

Date _____ No. _____

Whole CNN - :



- * Max pooling: takes the max value in a 2×2 (non overlapping)
- * Avg. Pooling: takes the avg in a 2×2 (non overlapping)

- # Dark Background: - maxpooling selects the brightest pixels
- # Light Background: - minpooling selects the darkest features
- # Image Smoothing: - average pooling.

- * Why pooling? → Subsampling pixels will not change the object. We subsample the pixels to make image smaller, fewer parameters to characterize image.
Reduced no. of connections, shared weights on the edges, Maxpooling further reduces the complexity.

- * Flattening: image refers to convert image from 2D to 1D by stacking its pixel values row by row (or) column by column. It is often done as preprocessing step before feeding into a ML model.

- * ReLU is applied after each convolution: faster & moreover effective training by mapping the negative values to zero and maintaining positive values.

- * final layer: softmax on the previous layers, yielding the probability of output being of a particular class.

Training a CNN: Two steps - forward phase:- The input is passed through the network; each layer will cache any input and intermediate values needed for the backward phase.

Backward Phase! Gradients are backpropagated and weights are updated.

(Crossentropy^{loss}) is used as the loss function.

Training a CNN: Maxpooling:- A maxpooling layer cannot be trained as it does not have any weights. We can propagate gradients by double the width & height of the gradient by assigning each value to where the original max value was in the corresponding pooled block.

Training a CNN: Convolution:- Changing the filter weights affect the entire output image for that filter. Every output pixel uses every pixel weights during convolution. Increasing any of the filter weights by 1 would increase the output by the value of the corresponding image pixel. The derivative of a specific output pixel (i, j) with respect to a specific filter weight is just the corresponding image pixel value.

0	7	0
0	80	31
33	14	0

→

0	0	0
0	1	0
0	0	0

→ 80

Deploying a Pre-trained CNN:-

- #) Pre-train a neural network model on a source dataset.
- #) Create a new neural model by:
 - Copying the pre-trained model
 - Removing the output layer
 - adding an output layer
 - initializing randomly the model parameters of the new output layer
- #) Train new model on the target dataset
- #) The output layer will be trained from scratch while the

Parameters of all other layers are fine-tuned based on the parameters of the source model.

Lec-6 - RNN; Basic RNN task -

- * Predict the future from the past
- * Map the past x into a fixed-length vector $h^{(t)}$
 $(x^{(t)}, x^{(t-1)}, x^{(t-2)}, \dots, x^{(1)}) \rightarrow h^{(t)}$
- * How much of the past should one store/medal in $h^{(t)}$

Three common vanilla RNNs:

- * Version 1: *) output at each time step
- *) Recurrent Connection b/w hidden units
- * Version 2: *) output at each time step
- *) Recurrent Connection only b/w from the output at one time step to the hidden units at the next
- * Version 3: *) Recurrent Connection b/w hidden units
- *) Read an entire sequence \mathcal{S} ; produce a single output

Version 2) - no need to compute the output for previous time step first why? we can give the actual output.

* Teacher Forcing - Applicable to RNNs that have connections from their outputs to their hidden states at the next time step. At training time, we feed the correct output $y^{(t)}$ drawn from the training set as input to $h^{(t+1)}$.

Paper said

During training time, we will feed in the intended output to the next hidden state and teach it how to produce the output time $t+1$ (compare it with intended output of time $t+1$). We can add noise to it cause it's not always be the exact output.

Input Output Scenarios:

Date

No.

- * Single-Single \rightarrow Feed-forward Network
- * Single-Multiple \rightarrow Image captioning (Image only)
- * Multiple-Single \rightarrow Sentiment classification
- * Multiple-Multiple \rightarrow Translation ; Image captioning ;
(Previous words too).

Injecting the Image to RNN:

- * Int-inject - Image features are used to initialize the hidden state of the RNN. The image features are passed through a fully connected layer to obtain an initial hidden state, which is used to initialize the hidden state of RNN.
- * Pre-inject - The image features are concatenated with the input sequence at each time step before being fed into the RNN. This allows RNN to take the image features into account when processing the sequence.
- * Par-inject - The image features are passed through a parallel network that runs alongside of RNN. The output of the parallel network is then concatenated with hidden state of RNN at each time step.
- * Merge inject - The image features and the hidden state of the RNN are combined using a merge layer, such as concatenation before being passed through a fully connected layer to obtain the new hidden state.

Back propagation through time (BPTT)

Date _____ No. _____

- * Treat unfolded network as one big feed-forward network
- * This network takes \vec{x}_n , the whole sequence as input.
- * We compute the gradients through the usual Backpropagation
- * We update weights, shared weights

Truncated BPTT

- * Run forward & backward propagation through segments of the sequence instead of the whole sequence
- * Back propagation through a segment & make a gradient step on the weights
- * Next Batch: - still have the hidden state from the previous batch & carry the hidden state forward.

Moving Gradients: The effect of change in L to the weight in some layers becomes so small due to increase model complexity with more hidden units, that it becomes zero after a certain point.

Exploding Gradients: a large increase in the norm of the gradient during training due to an explosion of long-term components, can result in the gradients growing exponentially.

Gradient scaling: normalizing the error gradient vector such that the vector norm equals a defined value, such as 0.1.

Gradient clipping: forcing the gradient value to a specific minimum (or) maximum value if the gradient exceeds an expected range.

- * In some cases large range of error gradients is permitted in the output layer compared to hidden layers.

Gated RNNs:- Creates paths through them that have derivatives that neither vanish or explode. Most effective sequence models used in practical applications: i) Long short-term memory, ii) Gated Recurrent Unit.

- * Allow the network to accumulate information over a long duration
- * Once this information is used the neural network can forget old state
- * Instead of when to clear the state manually, neural network learns to decide where to do it.

LSTM:- Uses this idea of constant error flow to create a constant error carousel which ensures that gradients do not decay. Memory cell that acts as an accumulator over time. Instead of computing a new state as a matrix product with the old state, it rather computes the difference b/w them. Gradients are well behaved.

Structure:- Each cell has the input and output as an ordinary recurrent network but also more parameters and a system of gating units that controls the flow of information.

- * The sigmoid layer outputs numbers between 0-1 which determines how much each component should be let through
- * \times gate : point wise multiplication
- * $+$ gate : point wise addition

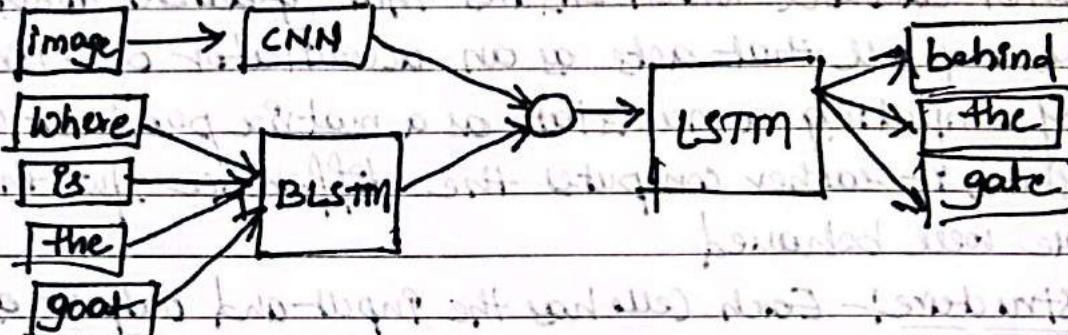
Gated Recurrent Unit:- * A single gating unit that combines the forget and input into a single update gate.

- * It also merges the cell state and hidden state
- * This is simpler than LSTM
- * There are many other variants too
- * Has fewer parameters than an LSTM & has been shown to outperform LSTM on some tasks

LSTM vs GRU!

- *¹) GRU uses lesser framing parameters and hence needs less memory.
- *²) GRU executes faster than LSTM where as LSTM is more accurate on larger datasets and longer sequences.
- *³) One can choose LSTM when dealing with long sequences and accuracy is concerned.
- *⁴) GRU is used when we have less memory availability.

Multimodal learning-



Lec 7!

Autoencoders:- Neural networks that are designed to learn encodings/features spaces from data. How? By reproducing the input from a learned encoding.

Encoder! - compresses the input into a latent space of usually smaller dimensionality.

Decoder! - reconstructs the input from the latent space with as close as possible.

- *¹) Unlike PCA we can achieve the non-linearity by using activation functions. It has been shown ^{NO} without activation function achieve PCA capacity:
- $$x \rightarrow f(x) \rightarrow z \rightarrow g(z) \rightarrow x'$$

Given data x and we would like to learn function f (encoder) and g (decoder) where $f(x) = s(wx+b) = z$ & $g(z) = s(w^Tz+b) = x'$

z_p is a latent space representation, S_p is non-linear identity function.

Date _____
No. _____

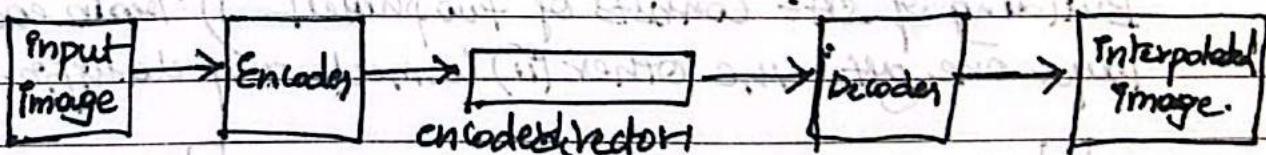
* Train the autoencoders using gradient descent as any other ANN using traditional squared error loss function. If input p is interpreted as bit vectors then cross entropy

* Two types of autoencoders :- Undercomplete, Overcomplete.

Undercomplete- Hidden layer p is smaller than input layer. Hidden nodes capture good features of the training set distribution.

Overcomplete- Hidden layer p larger than input, no compression in hidden layer, no guarantee that model will capture good features, a latent space of higher dimensionality models a more complex distribution.

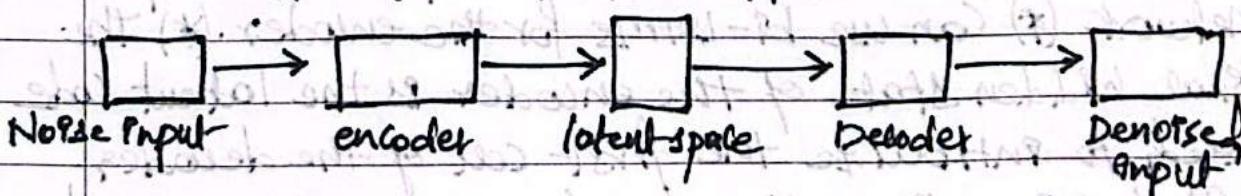
Image fading- Take two random images and generate their encoding by passing them through encoders. The encodings are flat vectors of length 128 each. Generate interpolations b/w the encodings. Generate images for each of the interpolation by passing them through decoders.



Denoising autoencoder- A more robust model.

Encode the input & do not mimic the identity function

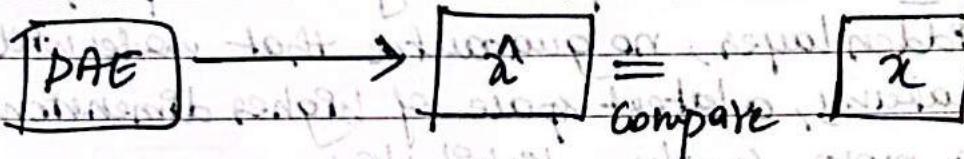
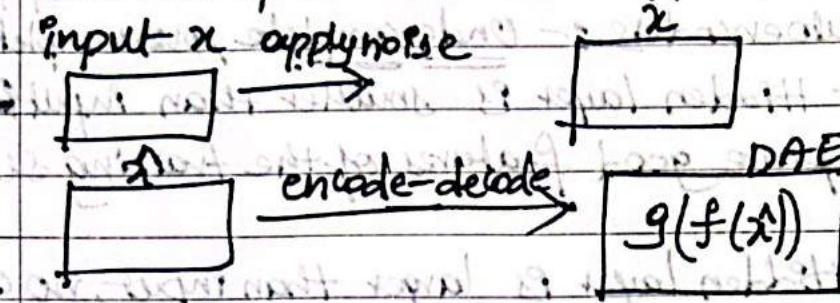
• Undo the corruption process applied to the input.



* Instead of trying to mimic the identity function by minimizing, A DAE instead minimizes

MATRIXAS

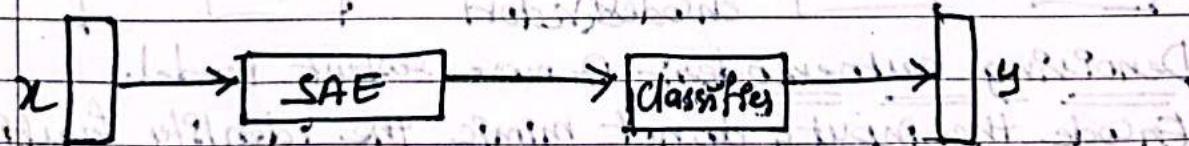
Denoising autoencoders - process: - Using the corrupted version of the input data as input to the network, it is trained to produce the original uncorrupted input by the output.



Stacked autoencoders - we want harness the feature extraction quality of an AE; we can build a deep classifier, whose input is the output of a SAE.

Benefit - Our deep models are not randomly selected but rather smartly selected.

Building of SAE consists of two phases - i) Train each AE layer one after the other ii) Connect any classifier



Autoencoders with RNNs - *). LSTM for encoders & decoders network (*). Can we bi-LSTM for the encoder. (*) The final hidden state of the encoder is the latent code used to initialize the first cell of the decoder.

- * At training time: the decoder receives the ground truth previous word as input at each time step
- * At inference time: - The decoder receives the predicted word at the previous time step.

MATRIXAS

Date _____ No. _____

Imitation of autoencoders:-

- *) language translation: encoder + decoder: learn the translation tasks together.
- *) Encoder: reads & encodes a source sentence into a fixed-length vector; Decoder: outputs a translation from the encoded vector.
- *) A Neural network needs to be able to compress all the necessary information of a source sentence into a fixed length vector. This may make it difficult to cope up with long sentences → Solution is Attention Mechanism.

Attention Mechanism [Bahdanau 2015]: In autoencoders, attention mechanism can be used to help autoencoders learn to focus on the most important parts and ignoring the irrelevant parts. This also helps autoencoder to identify patterns in the input data, which can be used to generate more meaningful representations of the data. The autoencoder does not encode a whole input sentence into a single fixed length vector. It encodes the input sentence into a sequence of vectors and choose a subset of these vectors adaptively, while decoding the translation.

Decoder:- implements an attention mechanism by soft searching through annotations generated by the encoder.

Encoder:- uses a BiLSTM to capture richer context from the input sentence by producing annotations that are concatenations of the forward & backward ~~propagation~~ hidden states at each step.

How does decoder pay attention to annotations?

It uses an alignment model with a trainable feed forward network to score how relevant each annotation is for the decoders hidden state. It calculates weights of each annotation using the alignment score. The ^{MATRIXES} decoder receives

as input at time point i a context vector that is a weighted sum of the annotations.

Date : No. . . .

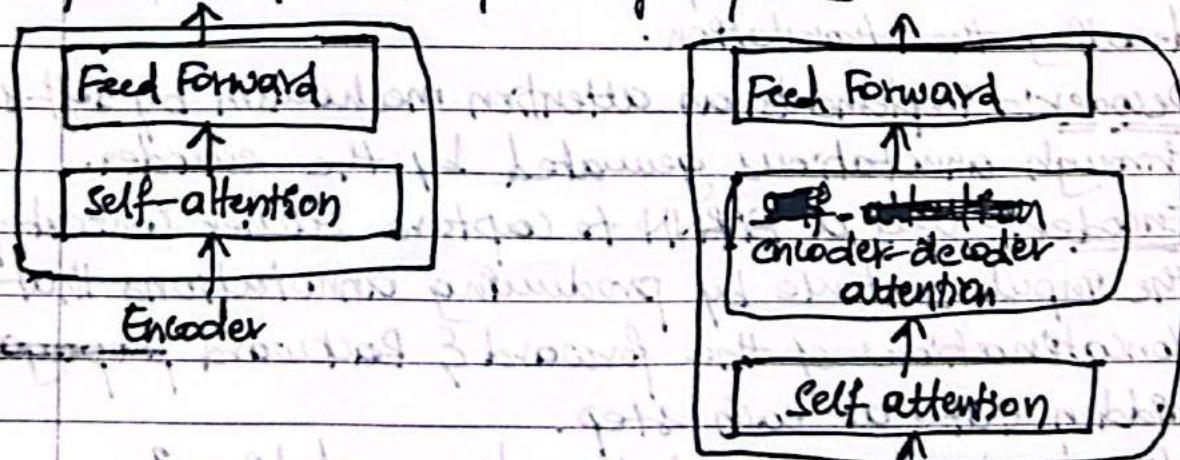
Decoder at step i : takes the previous state, current state to generate the current hidden state. Which is why we use previous state to calculate the alignment scores.

Transformers: first to relying entirely on self attention to compute the representations of its input g output, without using any RNN & CNN units.

Encoder: All are identical. R_n structure, but don't share similar weights.

- 1) self attention layer helps the encoder look at other words in the input sentence as it encodes a specific word.
- 2) Output of self attention layer is directly fed into a feed forward neural network.
- 3) The same feed forward network is independently applied to each position.

Decoders: has both the same layers but in b/w them it has another layer (attention layer) that helps decoder focus on relevant parts of input sentence.



Simultaneous with a previous page, with
however, but different, a new, but a different, in the same
sort of problem, due to the evolution of the
need to change, understand, if there is a, but, a, but, a
subset which will, for example, be, with Matrikas

Encoding: *1) Each encoder receives an input vector, each of 512 bits.
*2) Each embedding flows through each of the two layers of encoder.
*3) These paths have dependencies when in self-attention layers but not through feed-forward NN, which is the same for all output vectors.

* Hence, the latter can be executed in parallel.

Self attention- Consider a statement: "the animal didn't cross the street because it is too tired". The self-attention layer solves this by associating "it" with "animal".

Step 1- Define three vectors for each input embedding by multiplying with three matrices. (i) Query vector (ii) Key vector (iii) Value vector: the default dimensionality of these three vectors is 64.

Step 2- Self attention score is the dot product of the query vector with key vectors, with respective word.

Self attention score for the word r_1 at the position #1 is $q_1 \cdot k_1$ & second, $q_2 \cdot k_2$, etc. The score defines how much focus to place on other parts of input sentence as we encode a word at a certain position.

Step 3-4- Normalization: Divide the scores by 8, this leads to have more stable gradients. Then pass the result through a softmax operation.

Step 5-6- Scaling & summarization:- *1) Multiply each value vector by softmax score. Maintain the values of the word(s) we want to focus on & drown out irrelevant words. Output of self-attention layer at each position is sum of the weighted value vectors.

Multiple heads: It focus on different positions Date _____ No. _____

It offers attention layer multiple representation subspaces
It tries to catch different definitions of relevance
If we add all attention head to the picture, it is hard to interpret.

- (*) How do we account for the order of the words?
- the transformer adds a vector to each input embedding
 - these vector follows a specific pattern that the model learns, which help it to determine the position of each word in the sequence
 - adding these values to the embedding provides meaningful distance b/w the vectors once they are projected into Q/K/V & during dot product attention

- (*) Layer Normalization: so that each feature has the same average and standard deviation. Normalizes the input across the features.

Decoder: (1) The encoder starts by processing input sequence

- (2) the output of the top encoder is then transformed into a set of attention vectors K & V.
- (3) These are then used by each decoder in "encoder-decoder attention" layer helping the decoder to focus on appropriate places in the input sequence.
- (4) Encoder-decoder layer: works just like multihead self attention, except it creates query matrix from the layer below it, & takes the key & value matrix from output of the encoder. stack Self attention layer: only allowed to attend to earlier positions in the output sequence by masking future positions before the softmax step in self attention calculation

The final layers - (1) The decoder stack outputs a vector of floats

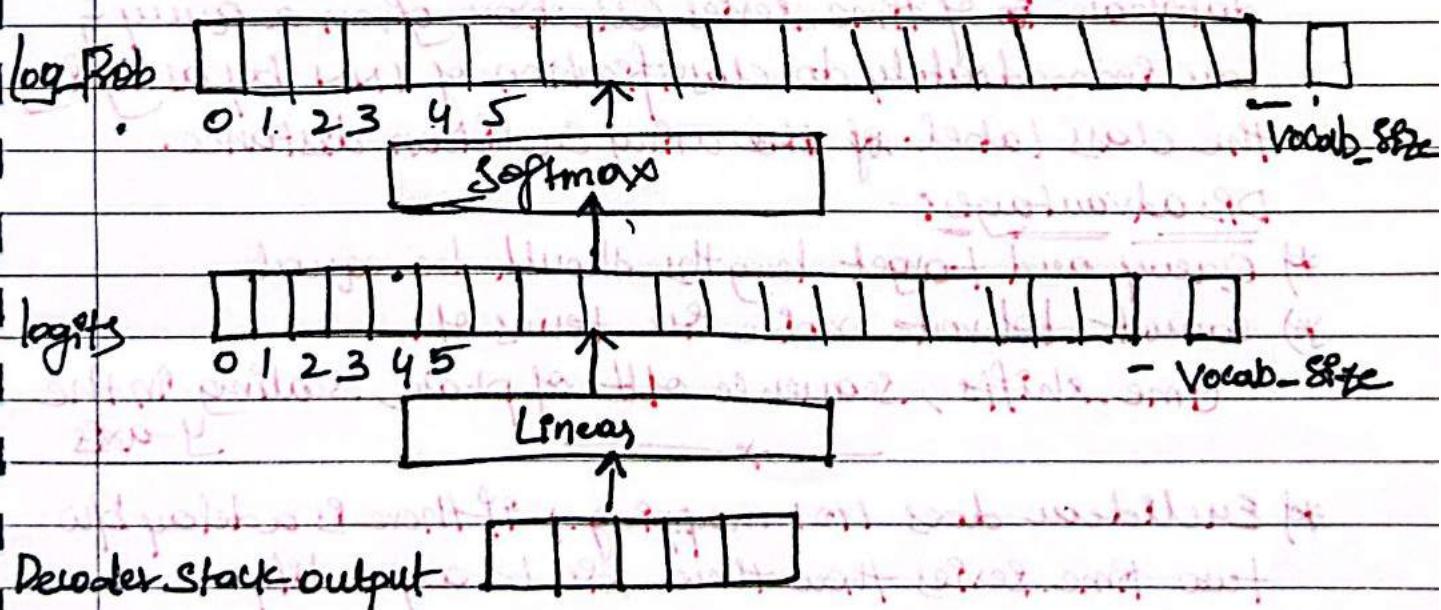
(2) How do we turn that into a word?

- final layer followed by a softmax layer
Linear layer - a simple fully connected neural network that projects the vector of the stack of decoders onto a large vector called logits vector.

(3) Assuming the model knows 10,000 unique words learned from the training set

(4) Logits vector: ~~is~~ has 10,000 cells

(5) The softmax layer then turns the scores into probabilities
(6) the cell with highest probability P_3 chosen, and the word associated with P_3 produced as output for this time step.



Loss function - translate "merci" to "thanks".

* use one-hot encoding to represent the output vocabulary
* Cross-entropy loss (or) Kullback-Leibler divergence

- *) How do we use transforms for classification?
- *) Add an FNN and a softmax layer on top.

lec-8 Time Series

Date _____

No. _____

Data Series:- Sequence of points ordered along some dimension.

- *) How do we compute the distance b/w two time series?
How do we define a similarity?

The L_p norm/Euclidean distance and after apply k-means or clustering :- View each time series as a point in the n -dimensional space ($n=\text{length}$); Define the distance b/w time series X and Y .

The Euclidean distance:- Euclidean distance b/w two time series & ranking them according to the distance low to high.

And the disadvantages are query & target-lengths should be equal.

→ Distance function, then we can perform k-nn search on database of time series collection given a query & can immediately do classification of INN by assigning the class label of the using Euclidean distance.

DPs advantages:-

- *) Query and target lengths should be equal
- *) Cannot tolerate noise in terms of:
 - Time shifts, sequence off. of phase, scaling in the y -axis

- *) Euclidean does 1 to 1 mapping. If there is a delay b/w two time series then there will be a penalty.

Dynamic Time wrapping:- DTw allows sequences to be stretched along the time axis.

- Insert 'stutters' into sequence
- Then compute Euclidean distance
- It does it in minimum way does that cost of the matching is minimized using Dynamic programming

MATRIXAS

Computing DTW:-

Date

No.

- *) DTW computed by dynamic Programming (DP)
- *) Given two sequences. $X = \{x_1, x_2, \dots, x_m\}$; $Y = \{y_1, y_2, \dots, y_n\}$
- *) Use an $n \times m$ DP matrix f to store the scores
- *) Computing the difference of the two points that are corresponding to the cell (Pairwise difference of two time series points) and add the minimum value looking at three adjacent cells & the final distance.
Equations $f(i, j) = \|x_i - y_j\| + \min \begin{cases} f(i-1, j-1) \\ f(i-1, j) \\ f(i, j-1) \end{cases}$ i.e., the value of the top-right cell.
- *) Each cell (i, j) is a pair of indices
- *) cell values are computed using some norm, $P=2$
- *) first cell $= (x_1 - y_1)^2$; we then moved diagonal then the second cell $= (x_2 - y_2)^2 + (x_1 - y_1)^2 \rightarrow$ Pairwise (adding the previous cell value to the second cell when it's computed). So this is identical to Euclidean but without square root. That means, we are doing just one-to-one mapping. When $P=2$ diagonal then the equation will be $f(i-1, j-1)$.
- *) Actually DTW allows any path not the diagonal path.
- *) Pairwise distance of the particular cell and then look on to the left side of the cell and downwards to cell & check the value in the cell and finds the smallest value and propagate it to add to the cell value. If we go from the left we are applying a y-stutter if we go from the bottom we apply x-stutter that means we are repeating the same y value and x value respectively if we go diagonal we are doing one-to-one mapping, not repeating anything. The top-right cell contains final output & it's the DTW score. If for each cell we store also their path (or) the transition (the adjacent cells having smaller value) we can propagate back to the initial cell.

MATRIX AS

Properties of DTW:- 9) Boundary Conditions (ii) continuity

Date _____

No. _____

(ii) Monotonicity.

Boundary condition:- The first two time series and the last two time series should match no matter what.

Continuity:- We should not skip any time point. Each point should atleast match other point otherwise series & viceversa.

Monotonicity:- The path should always go to the right and upwards; it cannot go back (or) down.

* What happens if it go backward (or) down?

— There will be a crossover in the matching of points

* How can we speed DTW?

Global constraints:- Applying global constraints (or) apply bands around the diagonal.

* We can apply how the working path can move. It can only move inside band & it cannot go outside the band and we have band of fixed length of δ_1 and we can also have band that δ_1 rigid at the beginning and as it moves it becomes more liberal and again it becomes rigid. "Itakura" fits better and mostly used.

Pros & Cons:- Advantages:- 1) Target & Query lengths may not be equal length

* Can tolerate: time shifts, sequence out of phase, scaling on the y-axis (how?) (z-normalization)

Disadvantages:- 1) Computational complexity; may not be handled: masking of noisy values;

* If we have noise, noise will go through the matching and it propagate add. it to the sum of the Euclidean distance

Z-normalization - Needed when interested in detected trends and not absolute values. If we want to detect trends. If we care only about trend not fact that they are different in years, what we can do is Z-normalize the whole time series or each one of them. Take one time series, compute the mean & standard deviation and for each and every value in time series we subtract the mean and divide by standard deviation.

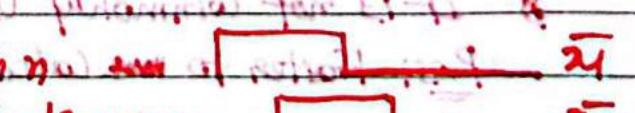
2) we observe patterns are similar but we are using two different scales, for two patterns for two time series. If we immediately run Euclidean (or DTW) we will get high values. If we apply Z-normalization it comes to the same scale.

If care about patterns then Z-normalization is must.

Longest common subsequence - Using dynamic programming but what it does, it checks the pairs of points. If the diff b/w the pairs of points is less than a small value then counts 1 to the match. So we add 1 to the score. If it is higher than that then we don't add anything to the match but we propagate the value from either the left or bottom. In this we can avoid which have high difference in value.

Time Series Summarization techniques -

Preface Aggregate Approximation (PAAP) - Represent the original time series of length n , as a square of boxes. Each box begins of the same length resulting in a new time series of length $N \leq n$.

X1: original time series of length n \rightarrow  \rightarrow 21

It can be represented in the N -dimensional space. \rightarrow  \rightarrow 22

\rightarrow  \rightarrow 23

Example: Let $X = [13 \ 14 \ 44 \ 45 \ 37]$; X can be mapped Date No.

from its original dimension $n=9$ to a lower dimension, e.g.

$N=3$, as follows: $[13 \ 14 \ 44 \ 45 \ 37]$

$\downarrow \downarrow \downarrow$

$1 \ 4 \ 5$

What does it do (Paaes):- It takes original time series &

breaks it down into segments non overlapping and

it computes average of within segment.

Pros-
1) Very fast to compute it. As efficient as other

approaches
2) Supports non-Euclidean measures.

Cons- It assumes as fixed-length non-overlapping summary windows over the time series.

* The solution to PAA vs APCA:- It does the opposite to

PAA, if we have a grow data & do simple PAA we

will have a fixed weight for each segment, that we have to define then we compute average/mean of the within segment

The problem:- We need two represents, length & value.

We need to come up exactly with one algorithm that we

want to use for segmentation and there are different time series segmentation algorithms.

Where should you do this segmentation where should you

stop one segment? (where time segment is begin?)

So we have to define how much fine grained we want

it to be, but this is hard to implement.

* It is not commonly used because harder to implement.

Pros- Faster to calculate, More efficient than other.

approaches, supports non-Euclidean measures.

Cons- Performance depends upon segmentation technique used.

Symbolic Approximation (SATX):— Similar principle to PAA,
start with original time series & use segments to represent it. But SATX represents segments with symbols (rather than real numbers)
Small memory footprint.

Definition of how it works— While we define those green lines (cates), we decide how many cates we want and then for each band how cates we assign a letter (discrete value).

Procedural— \rightarrow First, we have to normalize & apply PAA to it by defining no. of segments. Next, discrete into symbols

* How do we define the bands?

We are making an assumption that the time series we are having is a Gaussian distribution. Then we can divide the Gaussian distribution into areas of equal size.

* How do we decide the alphabet? We will decide it, we will try different alphabet size. To the first, we are solving and see which alphabet size solves the problem & we want to map the time series to SATX before we apply classification then we try diff. size of alphabets which performs well.
Problem— The symbiotic of those values.

Time Series classification:

If we use INN we can use any distance measures;

Pros— accurate, Robust, Simple; Cons— time and space complexity; results are not interpretable.

How about feature-based classification? We shapelets or "attributes" or "features" for splitting a node in the decision tree.

Shapelets— Time series Subsequence, maximally representative of a class, discriminative from other classes.

III. Chaplets are: the short time series

* If we have a shapelet, how do you quantify it?

So, any subsequence has a feature and what is the feature value? It can use many diff ways to represent.

One way to take the shapelet is: use a sliding window and sum D_E, Euclidean dist, less for each slide take the minimum one will be the best match.

* If we have a set of shapelets we can build a decision tree or a shapelet tree. How does it work?

It works same as the regular Decision tree. How can

apply? So, we have the dictionary and then each node

is decision node could be binary. If the best match is higher than some value, if yes you go to one side of the tree if not we go other side of the tree. And keep on testing every shapelet until we find the decision node.

How do we define shapelet?

* We apply any transformation (SPK, PAA) & do KNN.

* As soon as we find shapelet, select top k most informative shapes as features and learn any suitable classifier (svm, random forest) using the transformed dataset.

* Take random no. of shapelets & cluster them, take means of those clusters by using K-means & start putting those k-shapelets to improve them & check their utility in terms of classification.

ROCKET: - It keeps similar or better classification results on huge time series datasets and its much faster than any other.

LocNet: - It generates a set of convolution kernels that are random. The convolution of each kernel with an

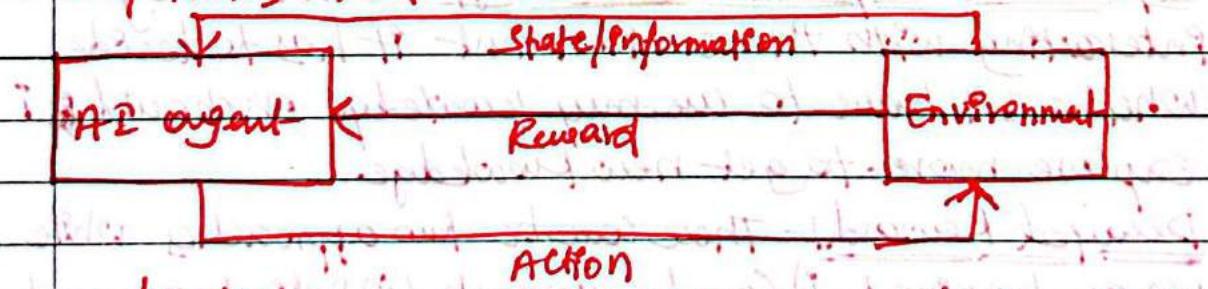
Input time series produces a feature vector. They ~~don't~~ ^{Date} ~~don't~~ ^{No.} define kernels, they don't train them. Next, they look at only two features they look at is the proportional positive value after convolution is applied & also maxpooling. And then they concatenate these two values. total no. of features if they feed it to the regressor.

- *) Dilated Convolution - If we have 3×3 size Kernel & then take image scan Kernel across the image we will get a number by doing dot product the normal convolution. So if we apply dilation what it does is, it moves with dilation of 2 that make our convolution area wider, actually add row (or) column after each and every point of the convolution Kernel. It very useful to manage with noise & time series of different scales that's why dilation is needed for time series

L-9 Interpretability & Explainability:

L-9 RL - where we learn through experience of trial & error method.

- *) RL - The art of learning, so there is an intelligent agent which interacts with environment. E.g. it learns by trial & error. & gradually learns how to take decisions.



Goal - Make optimal sequence of decisions in an unknown environment by Trial & Error.

- *) What do the AI agent learn? Our goal is to maximize the total accumulated reward over a period of time by taking the sequence of the rewards we get.

- *) Three Main Paradigms:-
(i) Multiarmed Bandits
(ii) Contextual Bandits
(iii) Markov Decision process

Aimed

MCTS Baudits - No change in environment, the environment remains static, the stat doesn't change we just take some actions & based on it we get some rewards & we learn which actions are good.

Contextual Baudits - Environment has some context by state but there is no co-relation from one state to another state.

Markov Decision Process:- Has both states as well as dynamics b/w one state to another state. The next state is dependant on current state.

Main characteristics of RL :- (i) Optimization (ii) Exploration
(iii) Exploitation (iv) Decay Reward (v) Generalization.

Optimization - How do we design rewards? We have to design in such a way that what the behaviour the agent is learning is what we want it to learn.

Exploration - Exploitation Trade off - When an agent is interacting with the environment it has to decide whether it has to use my knowledge or should i explore more to get new knowledge.

Decayed Reward - There can be two approaches while we are learning :- (i) Greedy approach (ii) patient approach.

Greedy - Just always take the action which is giving you the best reward immediately.

Patient - Taking an action that will lead to better accumulated reward on the long run, not immediately.

Generalization - In which each pixel combination is a new state for the agent. We took a few states and we can generalize for a large set of states, this is achieved using deep Reinforcement learning.

Where we train a weight vector of the neural network \vec{w} :
the weight vector \vec{w} is different to give result for any state

Multi-Armed Bandits - Setup - what we learn is how the exploration and exploitation trade off is done. We need to explore more to get more experience about the environment.

Action-value - How to explore? and how to exploit?

~~We need some value which can help~~ The action value tells us, this action has this much value. If you take this action this is the amount reward you get.

If you already know that for this action this is the reward from going to get ~~then~~ that is optimal ~~value~~ action.

* When we want to exploit or explore? When we want to exploit, whatever Q value which we have till now we will see which action gives us the best value and we will take that action. When we want to explore we will take a random action & then we will see what's the reward & update the Q value.

Baye Algorithms which helps in exploration exploitation :-

* (i) ~~E-Greedy~~, (ii) Upper Confidence Bound algorithm
E-Greedy:- If we want to explore we do it with the probability ϵ (epsilon). We decide some probability at the beginning. If we want to explore we will take a random action and update the Q value based on the reward we get. When we want to exploit we will choose the max value what we received till now. Probability of exploitation $\pi = 1 - \epsilon$.

Problems - Inefficient exploration - It explores with the same probability from starting to end

Inefficient exploitation - Spend too much / little time on clearly bad / good actions

Upper Confidence Bound) :- we want to take an action only if the reward for the action is large enough (i) when we are not confident about that action.

Principle's for UCB :-

i) optimism in face of uncertainty - Act as if the environment is as plausible as possible. If unsure, then overestimate. $\hat{Q}(a)$.

ii) statistical upper confidence bound - we can use some statistical bounds. The average value of the random variable is bigger than the particular value there. Inequality followed by that. ~~These upper confidence bounds help to capture the~~

* How it's applied? In each iteration we choose the action which maximizes the UCB value & update Q value based on the reward. We get:

* How it works? If you are sure the UCB value is high if you are not sure the UCB value is high.

* How good they are performing? This can be done by using regret analysis.

* Regret Analysis :- UCB is more efficient than E-greedy as time increases in E-greedy and lower time in UCB.

Content of Markov Decision Process :- We have a state and we have relationship b/w one state and other state. The environment is in one state and the agent takes action and get reward. As a result of this action the agent moves to next state, next action, next reward. This trajectory.

* Given current states and the history of actions taken till now, what is the probability moving to next state. So some dynamics involved in that.

* Markovian assumptions - The current state is enough to give the ~~assumption~~ of whole history of the past. Information

Deep RL: Generalization The no of features increase the state space increases exponentially that's called curse of dimensionality
→ We will learn the approx Q values and using the approx based RL algorithms they are proven to be more effective.

Date _____

No. _____