

1 Model Evaluation

1.1 Classification Metrics

1.1.1 Accuracy

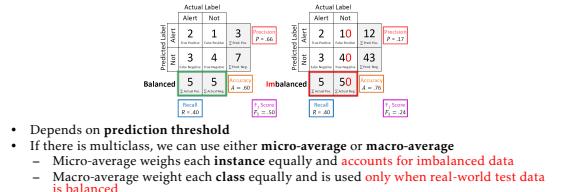
$$A = \frac{1}{m} \sum_{j=1}^m [\hat{y}_j = y_j]$$

- \hat{y}_j is the predicted value
- y_j is the ground truth value
- Interpreted as ratio of predicted elements that are equal to ground truth over total elements
- Accuracy other formula: $\frac{TP+TN}{TP+TN+FP+FN}$

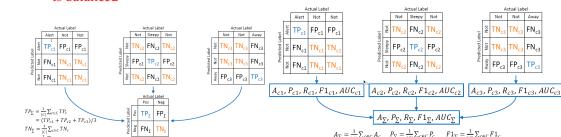
1.1.2 Confusion Matrix

		Actual Label		Precision
		Alert	Not	
Predicted Label	Alert	2 True Pos.	1 False Pos.	3 $P = \frac{TP}{TP+FP}$
	Not	3 False Neg.	4 True Neg.	7 $\Sigma \text{Pred. Neg.}$
		5 $\Sigma \text{Actual Pos.}$	5 $\Sigma \text{Actual Neg.}$	
				$F_1 \text{ Score} = \frac{2}{P + \frac{1}{R}}$
				Recall = $R = \frac{TP}{TP+FN}$

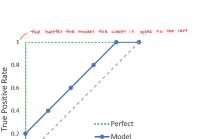
- Maximize recall to minimize false negatives
- Maximize precision to minimize false positive
- Only use F1 score if you can't justify between recall and precision
 - F1 score is not just a simple average but is instead more robust as it is less sensitive to extreme values
- Low metric score can be due to imbalanced data



- Depends on prediction threshold
- If there is multiclass, we can use either micro-average or macro-average
 - Micro-average weighs each instance equally and accounts for imbalanced data
 - Macro-average weight each class equally and is used only when real-world test data is balanced

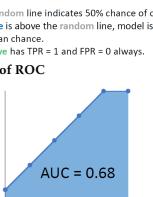


1.1.3 Receiver Operator Curve



- Diagonal random line indicates 50% chance of correctness.
- If ROC curve is above the random line, model is more accurate than chance.
- Perfect curve has TPR = 1 and FPR = 0 always.

1.1.4 Area Under Curve (AUC) of ROC



- AUC is a concise metric instead of a full figure.
- Concise metrics enable clearer comparisons.
- AUC > 0.5 means the model is better than chance.**
- AUC = 1 means model is very accurate.

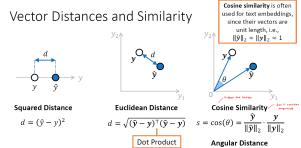
1.2 Regression Evaluation Metrics

1.2.1 Average Difference Metrics

$$\begin{aligned} \text{Mean Absolute Error (MAE)} &= \frac{1}{m} \sum_{j=1}^m |\hat{y}_j - y_j| \\ \text{Mean Squared Error (MSE)} &= \frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2 \\ \text{Root Mean Squared Error (RMSE)} &= \sqrt{\frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2} \end{aligned}$$

MSE and RMSE penalize larger differences more than MAE

1.2.2 Vector Distances and Similarity



1.2.3 Difference between MSE and Euclidean Distance

- MSE is the average of all the distances across the test set whereas euclidean distance is calculating the distance across feature vector

$$\begin{aligned} \text{Euclidean distance } d_j &= \sqrt{(y_j - y'_j)^2} \\ \text{MSE} &= \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^m (\hat{y}_j^{(i)} - y_j^{(i)})^2 \end{aligned}$$

2 Data Processing

2.1 Linear Separability

- In most real-world scenario, data tends to not be linearly separable (e.g. images, time-series data, text). Wrongly assuming that data features are linearly separable leads to irrelevant features being uninformative to train the model to discriminate between prediction labels
- Check using scatterplot (if only 2 features) and scatterplot matrix (if more than 2 features). We can also obtain a computational metric using linear SVM

2.1.1 Testing Linear Separability using Linear SVM

- Each $a(i)$ is the distance that the misclassified point j is from its correct margin
- Total violation: $\sum_{j=1}^n a(j)$
- Calculating the total violation indicates how linearly separable the data is in terms of its features
- Higher violation \Rightarrow Less linearly separable

2.1.2 Mitigation

- Find useful features through feature extraction (collecting new features)

- Transform features
 - Feature engineering ($x \rightarrow x^2$)
 - Change Basis Vectors (PCA, LDA)
 - Kernel trick
 - Features Learning (Neural Network)

2.1.3 Principal Component Analysis (PCA)

- PCA uses orthogonal transformation to convert a set of observations into a set of values of linearly uncorrelated variables called principal components
- Each principle component is a mix of features (e.g. PC1 can be 4 parts x_1 , 1 part x_2)
- We can then plot out something called a scree graph which tells you which Principle components are important and we can just keep those to reduce dimensionality
- Higher variance means more information
- Feature Normalization needs to be done

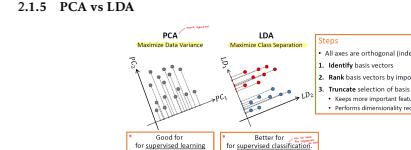
2.1.4 Linear Discriminant Analysis (LDA)

- LDA works by projecting the data onto a lower-dimensional space that maximizes the separation between the classes

$$F = \frac{|\mu_A - \mu_B|^2}{S_A + S_B}$$

- The F-test scores are tied to each linear discriminant and the higher F-test score means it separates classes better
- Discard LDA with low F-test to reduce dimensions

2.1.5 PCA vs LDA



2.2 Curse of Dimensionality

- Curse of dimensionality occurs when there are way more features than data points ($n \gg m$) as we are dealing with unstructured data (e.g. images, sensor data)
- Causes problems like overfitting since data are too sparse to inform about true decision boundary and model can easily fit to sparse data
- Also a problem for clustering (e.g. k-means) since distance are too similar
- Check using histogram of distances (check for variance)
 - Generally tedious to check various, general rule of thumb is for features to be $\frac{1}{5}$ of data points

2.2.1 Migrations

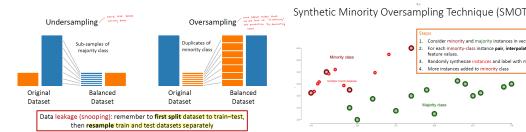
- Feature selection
 - Wrapper method (e.g. Recursive Feature Elimination, Elbow method)
 - Filter method (e.g. Only keep high info gain features, remove correlated features)
- Dimensionality Reduction
 - Linear Matrix Factorization: PCA, LDA
 - Auto Encoders
 - Feature Selection helps to enable faster model training and better interpretability

2.3 Imbalanced Data

- Values not evenly distributed in features due to rare events (e.g. rare cancer) or skewed data collection
- Leads to evaluation metrics being misleading to interpret and models overfitting to majority class
- Can be checked easily by using histogram or bar chart of features

2.3.1 Mitigation

- Collect more data instance
- Resample instances (undersampling, oversampling, SMOTE)
- Undersampling should only be used when there is not much loss of data (e.g. 55% majority class VS 45% minority class)
- If data is very largely imbalanced, then oversampling/SMOTE is preferred so that we don't lose too much data



3 Feature Engineering

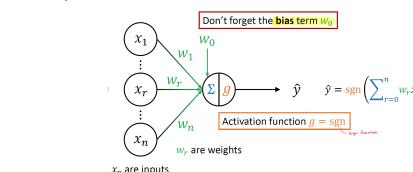
3.1 Feature Extraction and Engineering

- Process of transforming raw data to improve the accuracy of models
- Benefits:
 - Capture domain knowledge (e.g., periodicity or relationships between features)
 - Express non-linear relationships using linear models
 - Encode non-numeric features to be used as inputs to models

Tabular	Temporal	Image	Text
<ul style="list-style-type: none"> Domain-specific custom equation Features from domain experts Aggregation, difference, min, max 	<ul style="list-style-type: none"> Features from previous time step Aggregation statistics Linear regression Wave analysis features 	<ul style="list-style-type: none"> RGB image as 3D Image features from RGB histogram Shape features from edge detection Edge detection via Convolution 	<ul style="list-style-type: none"> Tokenization Stop words Bag-of-Words encoding

4 Perceptron & Neural Networks

4.1 Perceptron Classification

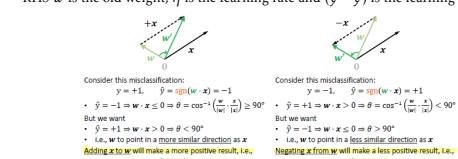


- Takes in a vector of input which are assigned weights
- Prediction class \hat{y} is given by applying the activation function g on the summation of the individual inputs multiplied by their weights $\rightarrow \hat{y} = g(w \cdot x) = g(w^T x)$

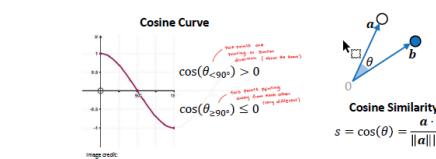
4.2 Perceptron Learning Algorithm (PLA)

- One of the simplest way to learn what weights to assign to inputs and is typically not used in modern ML approach. Never converges for non-linear or noisy data

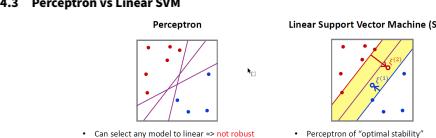
- Initialize weights w
 - Could be all zero, or random small values
- For each instance i with features $x^{(i)}$
 - Classify $\hat{y}^{(i)} = \text{sgn}(w^T x^{(i)})$
- Select one misclassified instance
 - Update weights: $w \leftarrow w + \eta(y - \hat{y})x$
- Iterate steps 2 to 3 until
 - Convergence (classification error < threshold), or
 - Maximum number of iterations
- The update of weight is given by $w \leftarrow w + \eta(y - \hat{y})x$, where LHS w is the new weight, RHS w is the old weight, η is the learning rate, and $(y - \hat{y})$ is the learning error



- Metric used to measure learning error is cosine similarity



4.3 Perceptron vs Linear SVM



- Can select any model to linear \Rightarrow not robust (uses different weights for different initializations)
- Cannot converge on non-linearly separable data

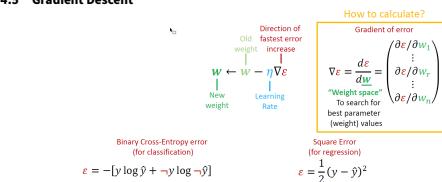
4.4 Activation Functions

Step

$$\text{sgn}(x) = \begin{cases} +1 & x > 0 \\ -1 & x \leq 0 \end{cases}$$

- 4 common activation functions are sign, sigmoid, tanh and ReLU
- All 4 activation functions are **differentiable**, allowing for gradient descent to be used to find minimum of error faster
- Sigmoid, tanh and ReLU are non-linear functions which allows for models to introduce **non-linearity** and create a non-linear, more complex decision boundary

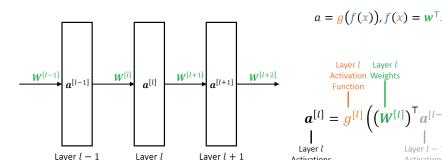
4.5 Gradient Descent



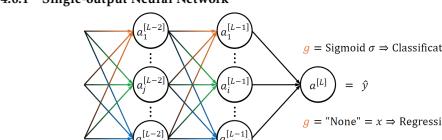
4.5.2 Summary of Gradient Calculation for Activation Functions

Error $e(y)$	Activation $\sigma(f)$	$\frac{de}{df}$	Weighted Sum $\frac{df}{dw}$	$\frac{de}{dw}$
$\frac{1}{2} (y - \hat{y})^2$	$\sigma(y - \hat{y})$	$-1 + e^{-x}$	$\frac{1}{e^{-x}} (1 - e^{-x})$	$w^T x$
$-\log \hat{y}$	$\text{softmax}(f)$	$\frac{-1}{\hat{y}} e^{-x}$	$\frac{1}{\hat{y}} e^{-x}$	$\frac{1}{\hat{y}} e^{-x}$
$\frac{1}{2} (y - \hat{y})^2$	$\text{ReLU}(f)$	$\text{ReLU}'(x)$	$\text{ReLU}'(x)$	$\text{ReLU}'(x) x$
$-\log \hat{y}$	$\text{Softmax}(f)$	$\frac{\partial}{\partial f_i} \text{softmax}(f)$	$\frac{\partial}{\partial f_i} \text{softmax}(f)$	$\frac{\partial}{\partial f_i} \text{softmax}(f) x$

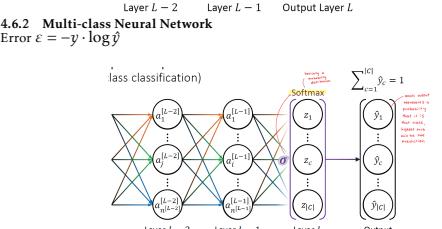
4.6 Artificial Neural Network (ANN)



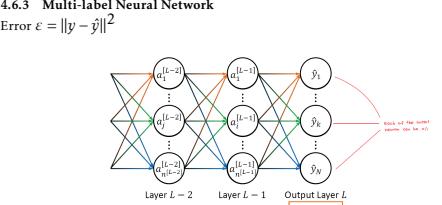
4.6.1 Single-output Neural Network



4.6.2 Multi-class Neural Network



4.6.3 Multi-label Neural Network



4.4 Sigmoid vs Softmax

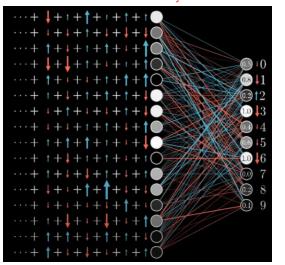
Sigmoid vs. Softmax

Sigmoid	Multiple Sigmoids	Softmax
• Is prediction true?	• Which predictions are true?	• What class is most probably true?
• For binary classification	• For multi-label classification	• For multi-class classification
• $\sigma(x) = \frac{e^x}{1+e^x}$	• $\sigma(x_k) = \frac{e^{x_k}}{\sum_{i=1}^K e^{x_i}}$	• $\sigma(x_k) = \frac{e^{x_k}}{\sum_{i=1}^K e^{x_i}} = \frac{e^{x_k}}{e^x + e^{-x}}$

5 Backpropagation

5.1 What is Backpropagation

- Used to efficiently compute the gradient of the weighted input of each layer from the last layer of the neural network to the front
 - avoids recomputing gradients
 - doesn't compute unnecessary intermediate values
 - computes gradient of each layer
- Idea for backpropagation is that the error of the prediction is propagated backwards through network, leading to adjustment of the weights and biases to minimize error
- Backprop relies on the fact that each function/layer is differentiable



5.2 Important Math Notations

$$\begin{aligned}\frac{d\hat{y}}{dW^{[l]}} &= a^{[l-1]}(\delta^{[l]})^T \\ \delta^{[l]} &= \left(\frac{dg^{[l]}}{df^{[l]}} \right) (W^{[l+1]}\delta^{[l+1]}) \\ \frac{de}{dW^{[l]}} &= \frac{de}{d\hat{y}} \frac{d\hat{y}}{dW^{[l]}}\end{aligned}$$

5.3 Backpropagation Visualization

$$\begin{aligned}(x^{[0]}, W^{[1]} f^{[1]} g^{[1]} \dots f^{[l-2]} g^{[l-2]} f^{[l-1]} g^{[l-1]} f^{[l]} g^{[l]}) &= \hat{y} \\ (a^{[l-1]}, W^{[l]} f^{[l]} g^{[l]}) &= \hat{y} \\ \frac{d\hat{y}}{dW^{[l]}} &= \frac{\partial \hat{y}}{\partial W^{[l]}} \\ (a^{[l-2]}, W^{[l-1]} f^{[l-1]} g^{[l-1]} f^{[l-2]} l^{[l-2]} g^{[l-2]} f^{[l-1]} g^{[l-1]} f^{[l]} g^{[l]}) &= \hat{y} \\ \frac{d\hat{y}}{dW^{[l-1]}} &= \frac{\partial \hat{y}}{\partial W^{[l-1]}} \\ (a^{[l-3]}, W^{[l-2]} f^{[l-2]} g^{[l-2]} f^{[l-3]} l^{[l-3]} g^{[l-3]} f^{[l-1]} g^{[l-1]} f^{[l]} g^{[l]}) &= \hat{y} \\ \frac{d\hat{y}}{dW^{[l-2]}} &= \frac{\partial \hat{y}}{\partial W^{[l-2]}}\end{aligned}$$

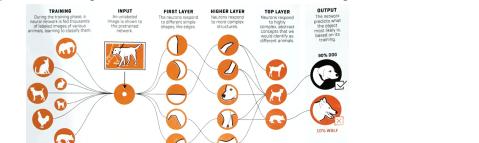
6 Convolutional Neural Network

6.1 Deep Neural Network

- Basically a neural net that has ≤ 3 hidden layers
- Layers can be convolution, pooling, concatenation, dropout, fully connected or softmax
- Use to predict complex target functions in the real world with a lot of parameters
- Requires large amount of data to increase performance and avoid curse of dimensionality

6.2 Convolutional Neural Network

- Suitable for image processing
 - Image classification - detect face emotions
 - Object detection - self-driving cars
 - Image segmentation - cancer cell detection
- Convolutional Neural Networks layers are analogous to applying different filters and kernels where it gradually builds up in complexity
 - First layer could be simple edge shape detection
 - Higher level can have more complex abstractions like shape of a tail, shape of a leg
 - Top layer would be a general abstraction of classes that could be predicted



6.3 Kernel Size, Stride, Padding

$$\begin{array}{c} \text{Input: } \begin{matrix} 0 & 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 3 & 4 & 5 & 0 \\ 0 & 6 & 7 & 8 & 0 \end{matrix} \\ \text{Kernel: } \begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix} \\ \text{Output: } \begin{matrix} 0 & 3 & 8 & 4 \\ 9 & 16 & 25 & 10 \\ 21 & 37 & 43 & 16 \\ 6 & 7 & 8 & 0 \end{matrix} \end{array}$$

Figure 1: Padding

Figure 2: Stride

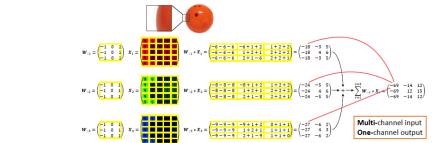
6.3.1 Calculation for Output Dimension

$$\text{dim}_y = \left\lfloor \left(\frac{h_x + h_p - h_k}{s} \right) + 1 \right\rfloor \times \left\lfloor \left(\frac{w_x + w_p - w_k}{s} \right) + 1 \right\rfloor$$

- h_x, w_x is the row x col of the input
- h_p, w_p is row x col of padding
- h_k, w_k is row x col of kernel
- h_s, w_s is row x col of stride

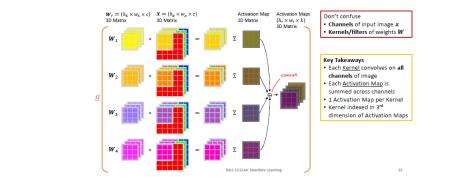
6.4 Multi-Channel Convolution

- Have 1 kernel for each channel (R, G, B) of the image, do convolution and then add up the resulting 2d matrix



6.5 Multiple Convolution Kernels

- Each kernel now convolves on all channels of the image
- Can think of each kernel looking for specific "feature" of the image (e.g. which part was more yellow, pink, orange etc.)
- Results in an activation map that basically tells you which part of the image each feature lies more in
- Concat the activation maps together in final output



6.6 Multi-Channel Multi-Kernel Convolution

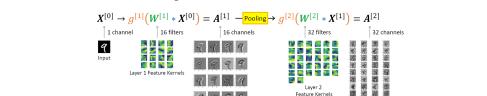
- Intuition is that multi-channel is used to detect color-specific features whereas multi-kernel can be used to detect different kinds of features (e.g. edges, textures etc.)
- In a multi-channel multi-kernel convolution operation, each kernel is applied to the corresponding channel of the input data, and the results are summed to produce the final output feature map. This process allows the network to detect complex patterns in the input data by combining the features detected by the different kernels

6.7 Fully Connected vs Convolution Layers

- Fully Connected (FC) layers
 - Each activation is a 1D vector
 - Each activation is a 2D vector
 - Layer of multiple activations is a 3D vector
- Convolutional (Conv) layers
 - Each activation is a 2D vector
 - Each activation is a 3D vector
 - Each activation is a 4D vector
 - Each activation is a 20D vector
- Activation function can be the same
- Activation stores activation of each neuron i separately
- Activation stores activation map of each kernel k separately

6.8 Convolutional Layers

- Typically starts off with relatively low number of kernels to detect basic features, increase number of kernels in deeper layers to get more descriptive power as we combine more and more simple features



- Uses pooling layers to downsample feature maps which helps to train later kernels to detect higher-level features
- Also helps in reducing dimensionality
- Typical methods are max-pool (most common), average pool, sum-pool

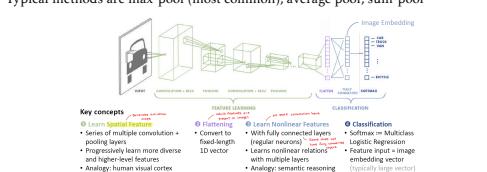


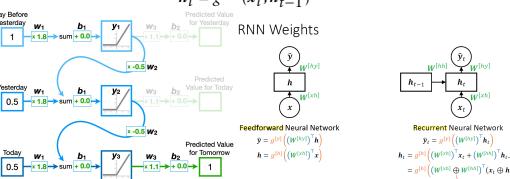
Figure 1: Padding

Figure 2: Stride

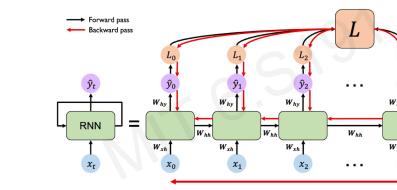
7 Recurrent Neural Network

- Suitable for text and time series data
- Used in applications like speech recognition, music generation, sentiment classification, DNA sequence analysis, machine translation
- Output of the hidden layers of the previous time step is fed into the training of the hidden layer of the next time step
 - Notion of using past data to predict things in the future
 - We are using the output of the hidden layer of the previous time step and not just the input data of the previous time step as hidden layer contains a more "compressed" representation of the input data with mostly useful information
- Weights for training each time step remains the same for efficiency and to ensure that model can predict independent of time

$$\begin{aligned}\hat{y}_t &= g^{[y]}(h_t) \\ h_t &= g^{[h]}(x_t, h_{t-1})\end{aligned}$$



7.1 Backpropagation Through Time

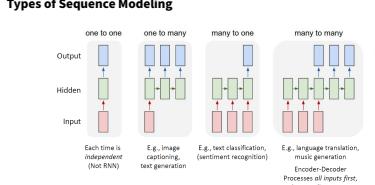


- In typical feed-forward neural nets, the backpropagation only occurs through one of the verticals in the diagram
- However, in RNN, since there is a notion of time data, we also have to propagate the weight information through the horizontal (time) axis
- Extremely slow since we have to complete the entire sequence for all the time steps before we can do any gradient/weight updates
 - Transformers are a way to speed them up by calculating the backprop in parallel without needing to depend on the next time step

7.2 Training RNN vs Prediction using RNN

- Find "best" explainer g that minimizes $\xi(x)$
 - "Exact" vs "Simple"
 - Locally-weighted error loss function: $\ell(f, g, \pi_x) = \sum_{i \in N(x)} \pi_i(x^{(i)}) (f(x^{(i)}) - g(x^{(i)}))^2$
 - Sparsity regularization: $c^{-\|f(x^{(i)})\|_1^2}$
- Model-agnostic means that LIME can be used for explaining any models
- Note that it only gives explanation for a local area
 - If area is defined big enough then can explain "globally"
- 8.4 Gradient-Weighted Class Activation Mapping (GRAD-CAM)
 - Used when explaining images as it has too many features which is unsuitable for LIME
 - Explains using something called saliency map which is basically a heat map
 - Grad-CAM aims to apply idea of feature importance to activation maps
- During training, we will always want to feed in the actual label as input for each of the time step
- During prediction, the input of the next timestep will be the predicted label in the previous timestep

7.3 Types of Sequence Modeling



7.4 RNN Training Issues

7.4.1 Overfitting

- Recap: overfitting is when the performance of validation < training
- RNN/Deep Learning models can overfit when there are too many parameters and model is more expressive than needed
- Mitigation:
 - Gather more data and get larger training dataset
 - Implement a dropout layer - randomly drop out neurons during batch training so that you disallow propagation through these neurons
 - Note that all neurons still used for prediction

7.4.2 Saturating/Vanishing Gradient Problem

Saturating Gradient Problem

Saturating Gradient Problem due to activation functions
Mitigate with ReLU activation function

When $a > 0$, $\frac{da}{dx} = 0$ (saturating)

When $a < 0$, $\frac{da}{dx} \rightarrow \infty$ (vanishing)

When $a = 0$, $\frac{da}{dx}$ is finite

With ReLU gradient is always 1 (for $x > 0$)

Can change sign for $x < 0$

• More importance for x_i

• Directly supportive (positive) or opposing (negative)

• Direction indicates increasing or decreasing influence

Vanishing Gradient Problem

When $a < 0$, $\frac{da}{dx} \rightarrow \infty$ (vanishing)

With gradient update much slower

• If $a < 0$, $\frac{da}{dx} \rightarrow \infty$ (vanishing)

• If $a > 0$, $\frac{da}{dx} \rightarrow 0$ (saturating)

E.g., $0.5^{10} \approx 0.0003$

8 Explainable AI (XAI)

8.1 Feature Importance

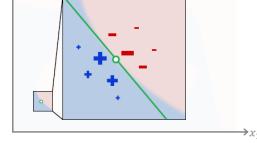
- Want to explain what features are important for prediction and in what ways the features influenced the prediction
- Implementations:
 - Weights for linear/logistic regression
 - Surrogate Weights from LIME
 - Saliency Maps (heatmap) of CNN

8.2 Interpreting Logistic and Linear Regression

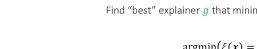
8.3 Local Interpretable Model-Agnostic Explanations (LIME)

8.3.1 Naive Approaches to Describe Non-linear Decision Boundary

- Use linear model to explain a complicated non-linear model
 - Simple
 - Too much error between explanation function and actual prediction function - lots of false positives and false negatives if we use simple linear model
- Use gradients to explain non-linear model
 - Calculates steepness for each feature x_j
 - More importance for x_j
 - Detects supportive (positive) or opposing (negative)
 - Direction indicates increasing or decreasing influence
- Difficult to remember, since gradients are different for each instance (point)



- Choose an instance x to explain
- Focus on a local region
- Get training set as the neighbors of x
- Train a surrogate linear model $g(x) = w \cdot x + \sum_{j=0}^n w_j x_j$, where the weights are higher the closer they are to x



- Find "best" explainer g that minimizes $\xi(x)$
 - "Exact" vs "Simple"
 - Locally-weighted error loss function: $\ell(f, g, \pi_x) = \sum_{i \in N(x)} \pi_i(x^{(i)}) (f(x^{(i)}) - g(x^{(i)}))^2$
 - Sparsity regularization: $c^{-\|f(x^{(i)})\|_1^2}$
- Model-agnostic means that LIME can be used for explaining any models
- Note that it only gives explanation for a local area
 - If area is defined big enough then can explain "globally"

8.4 Gradient-Weighted Class Activation Mapping (GRAD-CAM)

- Used when explaining images as it has too many features which is unsuitable for LIME
- Explains using something called saliency map which is basically a heat map
- Grad-CAM aims to apply idea of feature importance to activation maps
 - Compute Activation Maps $A^{[L]}$ of last conv layer L via Forward Propagation
 - Choose class label c to explain about (e.g., predict "go", "car")
 - Filter prediction \hat{y} to be about class c
 - Given: $\hat{y} = \hat{y}_{(0)} \hat{y}_{(1)} \dots \hat{y}_{(C)}$, $\hat{e}(c) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix}$, then $\hat{y}^c = \hat{y} \circ \hat{e}(c) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix}$
 - To generate explanation only for that class c
 - Compute importance weight $\alpha_i^{(c)}$ for Activation Map $A_i^{[L]}$
 - Backprop from $\hat{y}^{(c)}$ to get gradients (relative to activations) at last conv layer
 - Compute weighted sum with ReLU to get Class Activation Map
 - Compute weighted sum with ReLU to get Class Activation Map

9 Unsupervised Learning

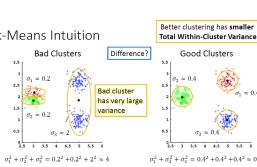
9.1 K-Means Clustering

9.1.1 Clustering Applications

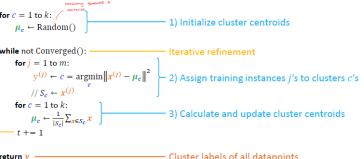
- Customer segmentation, Image Segmentation, Behavior Segmentation, Color Quantization (reduce image size)

9.1.2 Intuition of K-Means

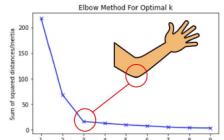
- Group clusters and try to minimize within-cluster variance
- More variance \rightarrow more dissimilarity between members of the same cluster
- Metric for variance is Within-Cluster Sum-of-Squares (WCSS): $L = \argmin_S \sum_{c=1}^k \sum_{x \in S_c} \|x - \mu_c\|^2$
- WCSS is the sum of the squared distances between each data point in a cluster and its centroid. It's calculated for each cluster and then summed up for all clusters.



9.1.3 k-means Clustering Algorithm



9.1.4 Elbow Method



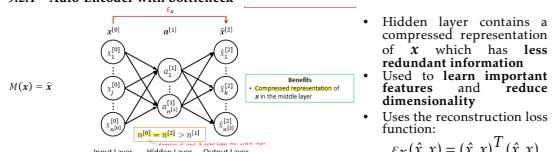
9.1.5 k-means VS KNN

	k-means Clustering	k-Nearest Neighbors (KNN)
Learning Paradigm	Unsupervised	Supervised
Purpose	Group neighbors	Label based on neighbors
k is ...	Number of clusters	Number of neighbors
Distance metric	Only squared Euclidean (to match Variance)	Any distance metric (e.g., Euclidean, Manhattan, Cosine)
Measures distance between	Training set x points and cluster centroids	Test set x points and training set neighbors
Need model training?	Yes	No

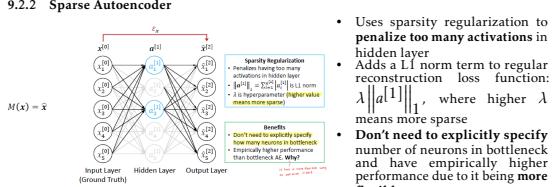
9.2 Autoencoders

- Autoencoders are designed to efficiently **compress** and **encode** data
- Primarily used for **dimensionality reduction** or **feature learning**

9.2.1 Auto-Encoder with bottleneck



9.2.2 Sparse Autoencoder

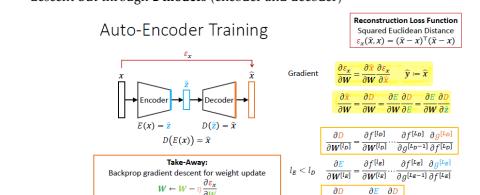


9.3 Comparison with PCA

- Both have ability to reduce dimensionality
- Autoencoders with linear activation/linear activation is similar to PCA but with non-orthogonal latent features
- Autoencoders performs better than PCA with non-linear activation since it can be more expressive

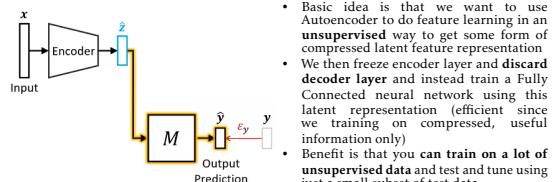
9.4 How to train Autoencoders

- Basically the same as normal neural networks using **backpropagation** with gradient descent but through 2 models (encoder and decoder)



9.5 Applications of Autoencoder

9.5.1 Feature Representation Learning



9.5.2 Anomaly Detection

- Idea is to train autoencoder which is very good at reconstruction **regular** instances
 - Only train on "good" instances
- When an irregular input comes in and autoencoder tries to reconstruct it, **error rate will be high** and we can catch it

9.5.3 Denoising Autoencoder for Robust Learning

- Idea is that you perturb dataset with noise and pass it to autoencoder to train
- Since autoencoders are good at removing redundant information, the goal is that autoencoder will only pick up on the non-noisy part of data and ignore noise
- Helps to get models which are **resistant** and **robust** to noise