

ML- CHAPTER

ML: is the field of study that gives the computer the ability to learn without explicitly program.

Taxonomy of ML

* Supervised learning : guided

* Unsupervised learning : unguided

* Reinforcement learning : learn and improve from previous knowledge. (feedback)

Regression vs Classification (Supervised)

continuous variable

discrete variable

clustering (unsupervised learning)

K-means clustering

Linear Regression

* the aim is to predict the price of other houses, as a function of the size of their living area
 ∵ given a training set ; we aim to learn a function $h: x \rightarrow y$

x : the space of input value

y : the space of output value

h : hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

θ_i : are the parameters (weights)

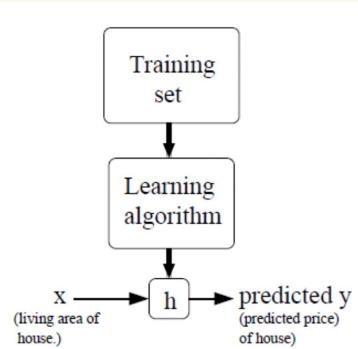
i : number of features.

for problem with two features : hypothesis will be :

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

for multiple input features : (multivariate linear regression) :

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x$$



Cost function : it takes an average difference of all the results of the hypothesis with input x 's and the actual output y 's.

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

difference b/w the predicted value and the actual value

it's called:

- * square error function or Mean square error (MSE) or least square error

Gradient Descent.

- * take the derivative of the cost function
- * it will give us the direction to move towards
- * go the opposite direction
- * the size of each step is determined by a parameter, which is called learning rate.
- * gradient descent algorithms
- * start with some "initial guess" for θ
- * repeat the convergence:

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
 - * where $j = 0, 1$ represent feature index number
 - * α is the learning rate
 - * at each iteration j , one should simultaneously update the parameters $\theta_0, \theta_1, \dots, \theta_d$
- * Batch gradient descent : the parameter are updated based on all example in the dataset.
- * Stochastic gradient descent : (incremental gradient descent) : the parameter are updated according to the gradient descent with respect to each training example individually.
- * Learning Rate: α : the size of each step is determined by it.
 - * if α is too small, gradient descent can be slow
 - * if α is too large, gradient descent can over shoot the minimum.

Debugging gradient descent.

- Make a plot with *number of iterations* on the x-axis.
- Plot the cost function, $J(\theta)$ over the number of iterations of gradient descent.
- If $J(\theta)$ ever increases, then you probably need to decrease α .

Automatic convergence test.

- Declare convergence if $J(\theta)$ decreases by less than E in one iteration, where E is some small value such as 10^{-3} .
- However in practice it's difficult to choose this threshold value.

The Normal Equation

- * is to minimize J explicitly by taking the derivative of J with respect to θ^i 's and setting them to zero.

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y})$$

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

Gradient descent vs. normal equations

Gradient Descent	Normal Equation *
Need to choose learning rate α	No need to choose learning rate α
Needs many iterations	No need to iterate
Works well when n is large	Slow if n is very large (in practice, n exceeds 10,000)
$O(kn^2)$	$O(n^3)$, need to calculate inverse of $X^T X$
Feature scaling is important	no need to do feature scaling
May find suboptimal solution	Finds unique optimal solution

Scikit learn uses normal equations in Linear regression implementation

Chapter (3)

Binary classification

Binary means 0 or 1

- * target variable can only take two values: 0 or 1
- * Binary classification problem is called **logistic function** or **sigmoid function**.

$$g'(z) = \frac{d}{dz} \frac{1}{1+e^{-z}} = g(z)(1-g(z))$$

* probability:

$$P(Y=1/x; \theta) + P(Y=0/x; \theta) = 1$$

Likely hood of parameters:

$$L(\theta) = P(\vec{y} | x; \theta)$$

Objective Function (log likelyhood)

$$J(\theta) = \log L(\theta) = \sum_{i=1} y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

Gradient ascent

$$\theta = \theta + \alpha \nabla_{\theta} l(\theta)$$

- * to maximise the likelihood we can use gradient ascent.

stochastic gradient ascent rule:

$$\theta_j = \theta_j + \alpha (y^i - h_{\theta}(x^i)) x_j^i$$

Newton's Method

: is for finding the roots of non linear equations $f(\theta) = 0$.

$$\theta = \theta - \frac{f(\theta)}{f'(\theta)}$$

Newton-Raphson method

- * generalisation of newton's method to multidimensional setting

$$\theta = \theta - H^{-1} \nabla_{\theta} l(\theta)$$

Hessian vector of partial derivatives

fisher scoring method

- * when Newton's method is applied to maximise the logistic regression log likelihood function $\ell(\theta)$, the resulting method is also called fisher scoring.

Tree-Based method and Ensemble Learning

- * **Tree based** :- method for classification and regression.
This involve stratifying or segmenting the predictor space into a number of simple regions (rectangles).
- * Simple \rightarrow fast \rightarrow useful
- * works well with both categorical and quantitative features.

Terminology for decision Trees

- * terminal nodes
- * leaf Nodes $h(x)$
- * are typically upside down
- * internal nodes

Classification of Decision Tree

- * Decision tree : are trained in greedy, recursive fashion, downward from the root.
- * left subtree is only those point satisfying $x_j < v$ (smaller)
- * right subtree is only those point satisfying $x_j \geq v$ (bigger)

How to choose best split

- step 1 : Try all splits. (All features and all splits within a feature).
 step 2 : for a set s let $J(s)$ be the cost of s .
 step 3 : Choose split that minimise $J(s) + J(s_l)$; or
 step 4 : choose split that minimise the weighted average:

$$\frac{|s_l|J(s_l) + |s_r|J(s_r)}{|s| + |s_r|}$$

How to choose cost $J(s)$ Misclassification rate

- * label s with the class C that label the most example in s .
- * $J(s)$: number of example in s not in class C

Entropy

- * we can use a cost function based on the information theory by measuring the entropy.
- * let y be random class variable, and suppose $P(y=c) = p_c$.
- * the surprise of y being class C is: $-\log_2 p_c$. (nonnegative)
 - * event with probability 1 gives us zero surprise.
 - * event with probability 0 gives us infinite surprise.

$$H(S) = - \sum_c p_c \log_2 p_c$$

[α Not (-) some \leq c $p_c \log_2 p_c$?]

Information Gain

- * choose split that maximise the information gain:

$$I(S) = H(S) - H_{\text{after}}$$

where :

$H(S)$: entropy

H_{after} : is the average entropy after split.

$$H_{\text{after}} = \frac{|S_L| H(S_L) + |S_R| H(S_R)}{|S_L| + |S_R|}$$

Gini impurity

- * another way to asses the quality of a split is Gini Impurity (strictly concave)
- * it measures how often a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in subset.

$$G(S) = \sum_c p_c (1-p_c)$$

($G(S) = \sum c p_C \text{ of one minus } p_C$)

Stopping Criteria

- * deep decision tree = overfitting.
- * shallow decision tree = under fitting.
- * heuristic to consider to decide when to stop splitting:
 - * fixed depth:
 - * node purity
 - * Information gain Criteria.

Regression decision tree

- * the goal is to find the regions (boxes) R_1, \dots, R_j that minimize the residual sum of squares (RSS), given by:

$$\sum_{j=1}^J (y_i - \hat{y}_{R_j})^2$$

Where:

\hat{y}_{R_j} is the mean response for the training observation within the j^{th} box.

Advantage of Trees

- * Trees are very easy to explain to people.
- * trees can be displayed graphically.
- * can easily handle qualitative predictors

Disadvantage of trees

- * trees generally do not have the same level of predictive accuracy.
- * they can be non-robust.

weak learners

Ensemble learning

- * Ensemble method :- is an approach that combines many simple "building block" model in order to obtain a single and potentially very powerful model.
- * Bagging * Random forest * Boosting.

- * **Bagging:** is a general-purpose procedure for reducing the variance of a statistical learning method.

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

out-of-bag error estimation

- * each bagged tree, make use of $\frac{2}{3}$ of the observations. So; the remaining $\frac{1}{3}$ of the observation not used to fit the given bagged tree are referred to as the out-of-bag (OOB) observation.

Random forest: provides an improvement over bagged trees by decorrelating the trees.

Boosting:

- * Bagging involves creating multiple copies of original training data set, using the bootstrap, fitting a separate decision tree to each copy and then combining all the trees in order to create a single predictive model. Boosting works in similar way, but each tree is grown using information from previous grown tree.

Adaboost: adaptive boosting; is a method for binary classification.

Boosting for Regression

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x)$$

Tuning parameters for boosting

- * Number of trees (B)
- * Shrinkage parameter (λ)
- * number of splits (d)

Support Vector Machine (SVM)

- * SVM :- is a supervised ML algorithm used for both classification and regression tasks.
- * Hyperplane :- is the decision boundary used to separate data points of different classes in a feature space.
- * Maximum margin classifier :- is the biggest gap or margin b/w two classes.
- * Non-separable data :- is when data are not separable by a linear boundary.
- * Support Vector :- are the closest data points to the hyperplane.
- * Margin :- refers to the distance b/w support vector and the hyperplane.
- * Kernel :- is a mathematical function used in SVM to map input data (low D) to higher dimensional feature space.
- * Hard margin :- refers to the maximum-margin hyperplane that perfectly separates the data points of different classes without any misclassifications.
- * Soft margin :- when a data contain outliers or is not perfectly separable, SVM uses the soft margin technique which introduce **slack variable**.
- * C parameter :- is a regularisation term that balances margin maximisation and the penalty for misclassification.
- * Hinge loss :- is a common loss function in SVM. it penalise misclassification point or margin violation.
- * Dual problem :- it involves solving for the Lagrange multipliers associated with the support vectors.
- * Slack variable :- tell us where the observation is located, relative to hyperplane and relative to the margin.

Kernel function in SVM

- 1- Linear : $K(w, b) = w^T x + b$
 - 2- Polynomial : $K(w, x) = (\gamma w^T x + b)^N$
 - 3- Gaussian RBF : $K(w, x) = \exp(-\gamma \|x_i - x_j\|^2)$
 - 4- Sigmoid : $K(x_i, x_j) = \tanh(\alpha x_i^T x_j + b)$
- Radial Kernel

- * it Controls variance by squashing down most dimensions事儿而已。

How Radial Kernel works

if a given test observation $x^* = (x_1^*, \dots, x_p^*)^T$ is far from a training observation x_i in terms of Euclidean distance,

then $\sum_{j=1}^p (x_j^* - x_{ij})^2$ will be large,

and so $K(x_i, x_{i'}) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2)$ will be tiny.

This means that, x_i will play virtually no role in $f(x)$.

Advantages of Kernels

- * Computational

SVM in More than 2 Classes

- * One Versus all (OVA): each class versus the rest.
- * One Versus one (ONO): classify x^* to the class that wins the most pairwise competitions.

Different b/w SVM and logistic regression

- SVM tries to find the "best" margin that separates the classes and this reduce the risk of error in data, while LR does not.
- SVM works well with Structured and semi-structured data like text and images while LR works with already identified independent variables.
- SVM is based on geometrical properties of the data while LR is based on statistical approaches.
- The risk of overfitting is less in SVM, while LR is vulnerable to overfitting

which to use : SVM or LR?

- * When classes are (nearly) Separable, use SVM.
- * If you wish to estimate probabilities, use LR
- * for Nonlinear boundaries, Kernel, use SVM.

How does SVM algorithm works

1. Problem Setup:	2. Linear Separability:	3. Margin Maximization:
<p>SVM is given labeled training data: (x_i, y_i), where:</p> <ul style="list-style-type: none"> x_i represents the features of a data point. $y_i \in \{-1, 1\}$ indicates the class label (for binary classification). The goal is to classify new, unseen data points correctly. 	<p>Hyperplane: A linear decision boundary in the form $w \cdot x + b = 0$, where:</p> <ul style="list-style-type: none"> w is the weight vector (defining the orientation of the hyperplane). b is the bias term (shifts the hyperplane). x is the feature vector. Data points satisfy: <ul style="list-style-type: none"> $w \cdot x + b > 0$: Class +1 $w \cdot x + b < 0$: Class -1 	<p>SVM aims to find the hyperplane that maximizes the margin, which is the distance between the hyperplane and the nearest data points (support vectors).</p> <p>Margin = $\frac{2}{\ w\ }$, where $\ w\$ is the norm of w.</p> <p>Optimization Problem:</p> <ul style="list-style-type: none"> Minimize: $\frac{1}{2} \ w\ ^2$ (to maximize the margin). Subject to: <ul style="list-style-type: none"> $y_i(w \cdot x_i + b) \geq 1$ for all i. <p>This ensures all points are classified correctly and lie on the correct side of the margin.</p>

4. Non-Linearly Separable Data:
If data isn't linearly separable, SVM uses the slack variable ξ to allow some misclassification:
<ul style="list-style-type: none"> $y_i(w \cdot x_i + b) \geq 1 - \xi_i$
Add a penalty term $C \sum \xi_i$ to the objective function, where C is the regularization parameter.

Kernel Trick:

- For highly complex data, SVM uses kernels to transform the input space into a higher-dimensional space where a linear hyperplane can separate the classes.
- Examples of kernels:
 - Linear: $K(x, z) = x \cdot z$
 - Polynomial: $K(x, z) = (x \cdot z + c)^d$
 - RBF (Gaussian): $K(x, z) = \exp(-\|x - z\|^2 / (2\sigma^2))$

Limitations:

- Computationally expensive for large datasets.
- Sensitive to the choice of kernel and hyperparameters (C, σ).

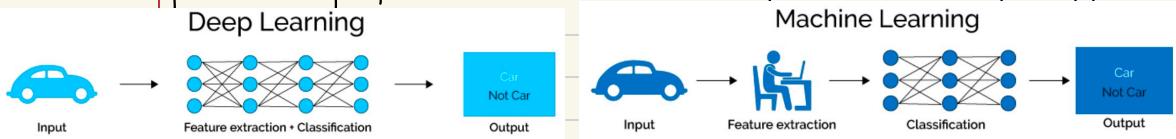
Advantages of SVM:
<ul style="list-style-type: none"> • Works well in high-dimensional spaces. • Effective for small datasets with a clear margin of separation. • Can handle non-linear data using kernels.
5. Predictions:
<ul style="list-style-type: none"> Once trained, the SVM predicts the class of a new data point x using: $f(x) = \text{sign}(w \cdot x + b)$

For non-linear kernels, $f(x)$ depends on the kernel function and the support vectors.

Deep Learning

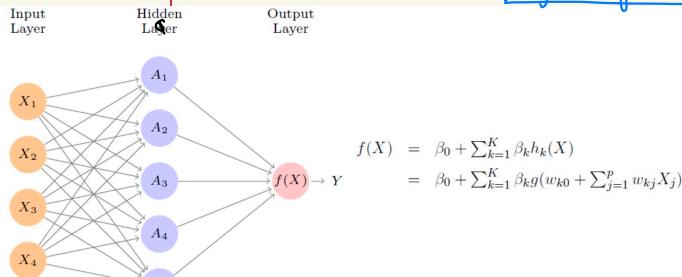
DL:- is a branch of ML that is Based on ANN architecture.
why DL is useful?

- * Can learn both Supervised and Unsupervised.
- * Effective end-to-end 'join System learning
- * utilise large amounts of training data
- * Learned features are easy to adapt and fast to learn.
- * provide very flexible, universal, learnable framework for representing world.



- * **Representational power:** NN with at least one hidden layer are universal approximators.
- * **perceptron:** is the basic processing element. It has input that may come from the environment or may be the output of other perceptrons.

Single layer Neural Network



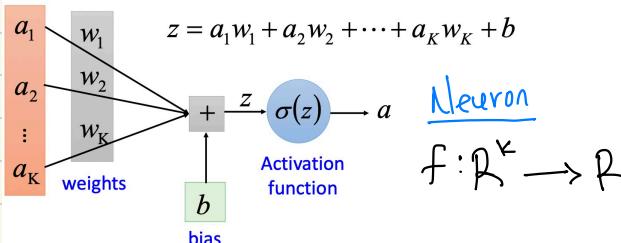
Steps:

- ① initialised with random weight.
- ② present a training instance ← Repeat
- ③ feed it through to get output from this steps till Okay
- ④ Compare with target output.
- ⑤ Adjust weight based on error
- ⑥ present a training instance

Some points

- * if $f(X)$ is a non-linear, a network with 1 hidden layer can, in theory, learn perfectly any classification problem.
 - * if $f(X)$ is linear, the NN can only draw straight decision boundary.
 - * NN use non-linear $f(X)$ so they can draw complex boundaries, but keep the data unchanged.
 - * SVM only draw straight lines, but they transform the data first in a way that makes that okay.
- Example:** Handwriting Digit Recognition.

Elements of Neural Network



- * **Neuron:** the basic unit that receive inputs each neuron is governed by a threshold and an activation function.
- * **Connections:** link between neurons that carry information, regulated by weight and bias.
- * **Weight & biases:** These parameter determines the influence of a connection.
- * **propagation function:** Mechanism that help process and transfer data across layers of neurons.
- * **learning Rule:** the method that adjust the weight and bias overtime to improve accuracy.

Soft Max Layer

- * is an activation function commonly used in NN for multiclassification problems.
- * **Soft max function:** is a mathematical function that convert a vector of raw prediction scores (often called logits) from the NN to probabilities.

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Activations functions

- * **Sigmoid function:** takes a real-valued number and squashes it into the range b/w 0 and 1.
- * **Tanh function:** takes a real-valued number and squashes it into the range b/w 1 and -1.
- * **ReLU (Rectified Linear Unit):** takes a real-valued number and thresholds it at zero.
- * **Linear function:** means that the output signal is proportional to the input signal to the neuron.

Training NNs

1. **Data preporcessing:** help convergence during training
 - * mean subtraction: to obtain zero-centred data
 - * Normalization:
 - * divide each feature by its standard deviation.
 - * OR scale the data within the range [0,1] or [-1,1].
2. **To train a NN:** Set the parameters θ such that for a training subset of images, the corresponding elements in the predicted output have maximum value.
3. Define a loss function $L(\theta)$ that calculate the difference (error) b/w model prediction and the true label.
4. for training set of N images, calculate the total loss overall all image
 - * find the optimal parameter θ that minimise the total loss $L(\theta)$

Loss functions

classification task

* Training example: ground truth class label y_i :

* out put layer: Softmax activation

* loss function: cross entropy.

Regression task

* Training example: ground truth output y_i

* output layer: Linear or Sigmoid activation

* loss function: MSE / MAE

* Gradient descent Algorithm: is an optimisation algorithm that minimise the cost function of a neural network modelling during training.

* Back propagation: is an algorithm that back-propagate the error from the output nodes to the input nodes. it compute the gradient descent of the loss function with respect to the network weight.

* Batch gradient descent: to update the model parameter values like weight and bias, the entire training data set is used to compute the gradient and update the parameter at each iteration.

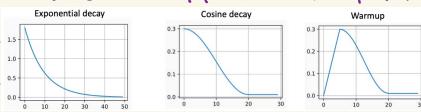
* Stochastic gradient descent (SGD): only one training example is used to compute the gradient and update the parameter at each iteration.

* Mini-batch gradient descent: a small batch of training example is used to compute the gradient and update the parameter at each iteration.

* Momentum-based gradient descent: it incorporates information from the previous weight updates to help the algorithm converge more quickly to the optimal solution.

* Adam: it combine the benefits of momentum, Adagrad, the learning rate is adaptively adjusted for each parameter based on the moving average of the gradient.

- * **Learning rate**: the gradient tell us the direction in which the loss the loss has the steepest rate of increase, but it does not tell us how far along the opposite direction we should step.
- * **Learning rate scheduling**: is applied to change the value of the learning rate during training.
 - * **Annealing**: is reducing the learning rate over time (area learning rate decay)
 - * **Warmup**: is gradually increasing the learning rate initially, and afterward let it cool down until the end of the training.



Generalization

- * **Underfitting**: the model is too simple to capture the data complexity.
high bias, low variance.
- * **Overfitting**: when the model does not make accurate predictions on testing data.
high variance low bias
- * **Weight decay**: is a term that penalises large weights added to the loss function.
- * **Drop out**: randomly drop unit (along with their connections) during training.
- * **Early stopping**: stop when the validation accuracy (or loss) has not improved after n epochs.
- * **Batch normalisation layers**: they calculate the mean μ and variance σ^2 of a batch data, and normalised the data x to a zero mean and unit variance.

Hyper parameter tuning

- * number of layers, and number of neurons per layers
- * initial learning rate
- * learning rate decay schedule
- * Optimiser type
- * loss function
- * Grid Search: check all values in range with a step value
- * Random search: Randomly sample values for the parameter.
- * Bayesian hyper parameter optimization: is an active area of research.

K-fold cross validation :- for hyper parameter tuning is common when the size of the training data is small.

- Illustration of a 5-fold cross-validation



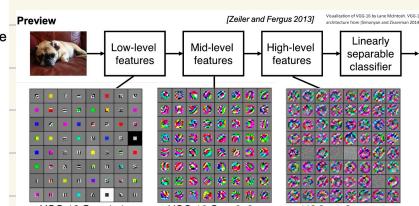
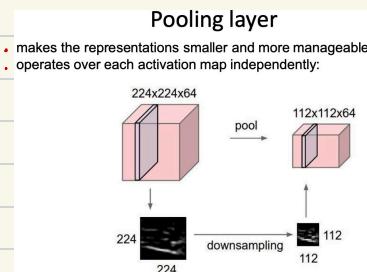
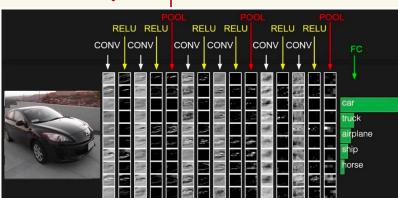
- Ensemble learning: is training multiple classifier separately and combining their prediction.

Convolutional Neural Network (CNNs)

- They were primarily designed for image data. CNN use convolutional operators for extracting data features.

How CNN works

- The CNN build up an image in a hierarchical fashion.
- Edges and shapes (Local features) are recognised and pieced together to form more complex shapes (Compound features).
- The hierarchical construction is achieved using convolution and pooling layers.
- Convolution layer** :- when convolutional filters are scanned over the image, they capture useful features. (e.g. edge detection).
- pooling layer** :- reduce the spatial size of the feature map.
 - Max pooling** :- report the maximum output within a rectangular neighbourhood.
 - Average pooling** :- report average output of a rectangular neighbourhood.
- fully connected layers** :- Contain neurons that connect to the entire input volume, as in ordinary neural networks.



Unsupervised Learning (Clustering) + PCT

- * clustering : is a technique for finding similar group in data, called clusters.
- * clustering algorithms :
 - * partitional clustering
 - * Hierarchical clustering
- * a distance function
- * clustering quality
 - * inter-cluster distance \Rightarrow maximised / minimized

K-means clustering

- * is a partitional clustering. for dividing data points into K clusters. based on their features, with each other cluster represented by the mean of its data points.

K-means algorithm

- 1 Randomly choose K data point to be the initial Centroids, cluster centres.
- 2 assign each data point to the closest centroid.
- 3 Re-compute the Centroids using the current cluster memberships.
- 4 if 9 convergence criterion is not met, go to (2)

- 1 no (or minimum) re-assignments of data points to different clusters.
- 2 no (or minimum) change of centroids
- 3 minimum decrease in the sum of squared error (SSE)

$$SSE = \sum_{j=1}^k \sum_{x \in C_j} \text{dist}(x, m_j)^2$$

C_j = the jth cluster, m_j = is the centroid of cluster C_j and $\text{dist}(x, m)$ = is the distance b/w data point x and centroid m_j .

strengths of k-means

- * simple
- * efficient
- * both K and t are small

Kmeans is the most popular clustering algorithm

Note: it terminate at local optimum if SSE is used. The Global optimum is hard to find due to complexity.

Weaknesses of K-means

- * The user needs to specify the K.
- * The algorithm is sensitive to outliers.
- * The algorithm is only applicable if the mean is defined.

Common ways to represent clusters

- * Compute the radius and standard deviation of the cluster to determine its spread.
- * Using classification model
- * use frequent value to represent cluster

Hierarchical clustering

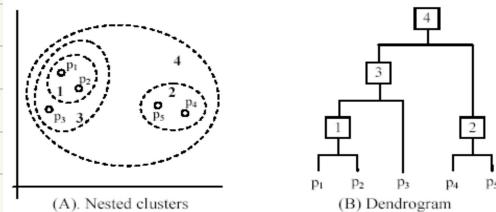
- * It build a hierarchy of clusters either from **bottom up (Agglomerative)** or from the **top down (divisive)**, where each data point starts in its cluster and pair of clusters are merge or split recursively.

Types of hierarchical clustering

- * **Agglomerative (bottom up) clustering:** it builds the **dendrogram (tree)** from bottom level and merge the most similar pair of clusters.
- * **Divisive (top down) clustering:** it start with all data points in one cluster the root, split the root into a set of child cluster. each child cluster is recursively divided further.

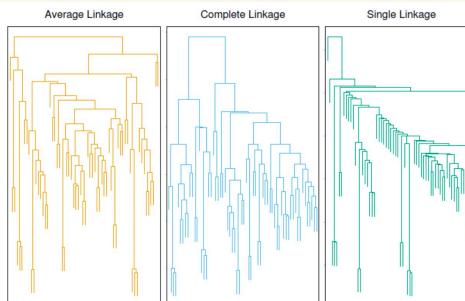


An example: working of the algorithm



Messuring the distance of two clusters

1. **Single link:** is distance b/w two **closest** data points, in the two clusters
2. **Complete link:** is the distance b/w to **furthest** data points in the two clusters.
3. **Average link:** the distance b/w two clusters is the average distance of all pair wise distances b/w the data points in two clusters.
4. **Centroid method:** the distance b/w two clusters is the distance b/w their centroids.



Average and complete linkage tend to yield more balanced clusters

Distance functions

○ If $h = 1$, it is the **Manhattan distance**

$$dist(\mathbf{x}_i, \mathbf{x}_j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ir} - x_{jr}|$$

○ **Chebychev distance:** one wants to define two data points as "different" if they are different on any one of the attributes.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, \dots, |x_{ir} - x_{jr}|)$$

○ **Weighted Euclidean distance**

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_r(x_{ir} - x_{jr})^2}$$

○ If $h = 2$, it is the **Euclidean distance**

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2}$$

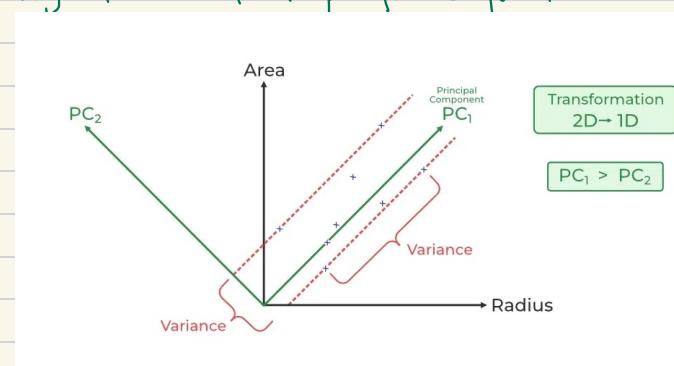
○ **Squared Euclidean distance:** to place progressively greater weight on data points that are further apart.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2$$

Principle Component analysis (PCA)

* **PCA:** is an unsupervised learning algorithm techniques used to examine the interrelation among a set of variables. The main goal is to reduce the dimensionality of data set while preserving the most important pattern or relationship b/w the variables without any prior knowledge of the target variable.

* **The first principle component:** Of a set of features, is the normalised linear combination of the features. The first PC capture the most variation in the data, but the second PC capture the maximum variance that is orthogonal to the first principle component and so on.



Need to Read & watch more tutorials on PCA to fully understand what it means.

Generative learning algorithm

- * **GILA**: are algorithms that try to learn mapping directly from the space of input x to the label.

Bayes theorem for classification

$$p(Y=k|x=x) = \frac{p(x=x|Y=k)p(Y=k)}{p(x=x)}$$

- * it classifies an observation x to the class for which $p(Y=k|x=x)$ is largest, has the lowest possible error rate out of all classifier.
- why discriminant analysis?

- * is popular when we have more than two response classes.
- * is more stable than the logistic regression model.

discriminant function

- * **discriminant score**:

$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

LDA when $p > 1$

- * extend the LDA classifier to be case of multiple predictors.
- * Draw Multivariate Gaussian distribution

Quadratic Discriminant Analysis

- * it provides an alternative quadratic discriminant analysis approach like LDA, QDA classifier result from assuming that the observation from each class are drawn from a Gaussian distribution, and plugin estimates for the parameters into Bayes theorem in order to perform prediction.

Types of errors

- * **false positive rate**: the fraction of negative example that are classified as positive.
- * **false negative rate**: the fraction of positive example that are classified as negative

Summaries

- * **logistic regression**: is very popular for classification, especially when $K=2$.
- * **LDA**: is useful when n is small, or the classes are well separated, and Gaussian assumption are reasonable. Also when $K > 2$.
- * **Naive Bayes**: is useful when p is very large.