

## 1 Linear Regression

- **Hypothesis function**

$$y = h(x) = \theta_0 + \theta_1 x$$

$$y = h(x) = \boldsymbol{\theta} \cdot \mathbf{x}$$

- **Loss function**

$$J(\theta_0, \theta_1) = \frac{1}{2} \frac{1}{m} \sum_{i=1}^m \left( h(x^{(i)}) - y_i \right)^2 = \frac{1}{2} \frac{1}{m} \sum_{i=1}^m \left( \theta_0 + \theta_1 x^{(i)} - y_i \right)^2$$

$$J(\theta_0, \theta_1, \dots, \theta_N) = \frac{1}{2} \frac{1}{m} \sum_{i=1}^m \left( \theta_0 + \sum_{k=1}^N \theta_k x_k^{(i)} - y_i \right)^2$$

- **Calculate the Gradient** gradient is the direction along which the loss function increases the most

$$\nabla_{\theta} J_{(0)} = \begin{pmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \end{pmatrix}_{(0)} \quad \nabla_{\theta} J_{(0)} = \begin{pmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_N} \end{pmatrix}_{(0)}$$

- **Update**

$$\boldsymbol{\theta}_{(1)} = \boldsymbol{\theta}_{(0)} - \alpha \nabla_{\theta} J_{(0)}$$

$$\boldsymbol{\theta}_{(n+1)} = \boldsymbol{\theta}_n - \alpha \frac{1}{m} \mathbf{X}^T (\mathbf{X} \boldsymbol{\theta}_n - \mathbf{y})$$

- **Repeat.**

$$\boldsymbol{\theta}_{(n+1)} = \boldsymbol{\theta}_{(n)} - \alpha \nabla_{\theta} J_n \quad n = 0, 1, 2, \dots$$

Stop if  $|J_{(n+1)} - J_n| \leq \text{tol}$  or  $n + 1 \geq \max$

### 1.1 Regularisation

A good machine learning model should have a good *bias-variance trade-off*. This can be achieved by minimising the generalisation error, which is the sum of the squared bias and the variance.

$$J(\theta_0, \theta_1) = \frac{1}{2} \frac{1}{m} \sum_{i=1}^m \left( h(x^{(i)}) - y_i \right)^2 + \lambda \sum_{j=1}^N \theta_j^2$$

If  $\lambda$  is small, large values of  $\sum_{j=1}^N \theta_j^2$  are allowed, meaning that the model contains too many parameters, thus it can **overfit**. If  $\lambda$  is large, only small values of  $\sum_{j=1}^N \theta_j^2$  are allowed. This means that the training will reduce the number of parameters, this it might underfit.

## 2 Logistic Regression

- **Hypothesis function**, models the probability that a given set of features belongs to class

$$h(\mathbf{x}) = \sigma(\boldsymbol{\theta} \cdot \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta} \cdot \mathbf{x})}$$

- **Loss function** (cross-entropy), entropy measures uncertainty on the possible values of a prediction

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left\{ y_i \ln[h(x^{(i)})] + (1 - y_i) \ln[1 - h(x^{(i)})] \right\}$$

- **Gradient descent** for logistic Regression:

Choose a starting point  $\Rightarrow$  Calculate the gradient  $\Rightarrow$  Update and repeat.

- If encounter **multi-class**, perform the same as many times as the number of categories of logistic regressions, each of which splits the data into two classes. End up with the same number of Hypothesis functions, select the highest probability.
- Metrics for classification

		Ground Truth		
		1	0	
Model	1	True Positive	False Positive	Precision
Prediction	0	False Negative	True Negative	
		Recall		

## 2.1 Supporting Vector Machines (SVMs)

In SVMs, we want to compute the optimal linear decision boundary, which is the **hyperplane** with the largest margin between the two classes.

- Support hyperplanes (optimal hyperplane between two points)

$$\mathbf{w} \cdot \mathbf{x} + b = 1$$

$$\mathbf{w} \cdot \mathbf{x} + b = -1$$

How is the hyperplane with largest margin oriented can be found by solving an optimisation problem

Find  $\mathbf{w}, b$  such that  $\frac{2}{\|\mathbf{w}\|}$  is maximum

$$\text{or } \min_{(w,b)} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) \geq 1 \quad i = 1, 2, \dots, m$$

### 2.1.1 Regularisation

if the data is not linearly separable, we need to introduce regularisation to allow some misclassified data points by introducing a slack variable  $\varepsilon_i > 0$  such that: (i)  $\varepsilon_i = 0$  for correctly classified (ii)  $\varepsilon_i = |y_i - f(\mathbf{x}^{(i)})|$  otherwise.

$$\min_{(w,b,\varepsilon)} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \varepsilon_i \quad \text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) \geq 1 - \varepsilon_i \quad \varepsilon_i \geq 0$$

small C - soft (wide margin)  $C \rightarrow \infty$  - hard margin

### 2.1.2 Kernel Trick

For nonlinearly separable data, one option is to enrich the data with nonlinearities. This can be achieved by mapping the data into a nonlinear higher-dimensional space,  $\mathbf{x} \rightarrow \phi(\mathbf{x})$ , where  $\phi$  is a nonlinear map, for example, a polynomial  $(x_1^{(i)}, x_2^{(i)}) \rightarrow (x_1^{(i)}, x_2^{(i)}, x_1^{(i)2} + x_2^{(i)2})$ . Intuitively, by adding dimensions in the feature space, it is more likely that the data becomes linearly separated.

Kernel Trick decomposes the weight vector in the transformed feature space

$$\mathbf{w} = \sum_{j=1}^m \beta_j \phi(\mathbf{x}^{(j)})$$

where  $\beta_j$  become the new weights to find. In other words, we express the weight vector with coordinates defined in the transformed feature space.

$$f(\mathbf{x}^{(i)}) = \sum_{j=1}^m \beta_j \phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x}^{(i)}) + b$$

we resort to a theorem of functional analysis, which ensures that there exists a function,  $K$ , such that

$$K(\mathbf{x}^{(j)} \cdot \mathbf{x}^{(i)}) = \phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x}^{(i)})$$

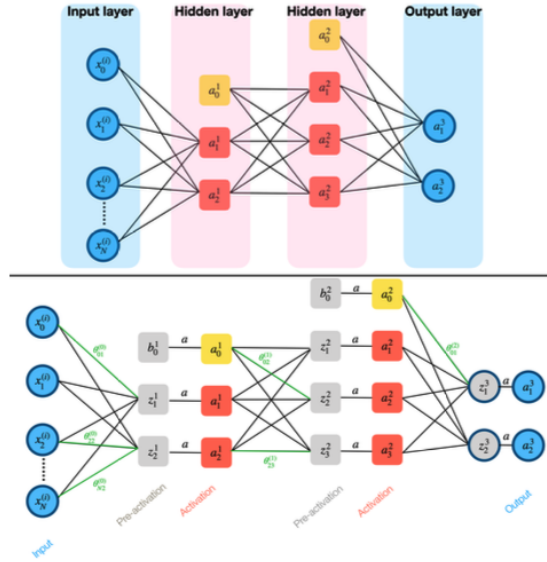
where  $K$  is the *kernel*.  $K$  suggested that we DO NOT need to increase the dimension of the feature space to include nonlinearities, we only need to apply the kernel to inner products in the original feature space. essentially finding optimal separating hyperplane without calculating anything about  $\phi(\mathbf{x}^{(i)})$ .

### 3 Feedforward Neural Network

*Neural Network generalise the principle of linearly combining simple nonlinear functions*

$$y \approx h(x : \theta)$$

function that maps input data "x" to output "y", with the function parameters represented by " $\theta$ ".



#### 3.1 Mathematical expression

- Input layer  $\rightarrow$  First hidden layer

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \Rightarrow \theta^{(0)} = \begin{pmatrix} \theta_{0,1}^{(0)} & \theta_{0,2}^{(0)} \\ \theta_{1,1}^{(0)} & \theta_{1,2}^{(0)} \\ \vdots & \vdots \\ \theta_{N,1}^{(0)} & \theta_{N,2}^{(0)} \end{pmatrix}, \quad z^1 = \begin{pmatrix} b_0^1 \\ \theta^{(0)T} x \end{pmatrix}, \quad a^1 = a\{z^1\}$$

where  $z$  is the *pre-activation*

- Second hidden layer

$$\theta^{(1)} = \begin{pmatrix} \theta_{0,1}^{(1)} & \theta_{0,2}^{(1)} & \theta_{0,3}^{(1)} \\ \theta_{1,1}^{(1)} & \theta_{1,2}^{(1)} & \theta_{1,3}^{(1)} \\ \theta_{2,1}^{(1)} & \theta_{2,2}^{(1)} & \theta_{2,3}^{(1)} \end{pmatrix}, \quad z^2 = \begin{pmatrix} b_0^2 \\ \theta^{(1)T} a^1 \end{pmatrix}, \quad a^2 = a\{z^2\}$$

- Output layer

$$\theta^{(1)} = \begin{pmatrix} \theta_{0,1}^{(2)} & \theta_{0,2}^{(2)} \\ \theta_{1,1}^{(2)} & \theta_{1,2}^{(2)} \\ \theta_{2,1}^{(2)} & \theta_{2,2}^{(2)} \\ \theta_{3,1}^{(2)} & \theta_{3,2}^{(2)} \end{pmatrix}, \quad z^3 = \theta^{(2)T} a^2, \quad h = a^3 = a\{z^3\}$$

Or in generalised from

$$a^{L+1} = a\{\theta^{(L)T} a^L\}$$

### 3.2 Number of parameters

$$\begin{aligned} \text{n. of parameters} &= \underbrace{(s_0 + 1) \times s_1}_{(N+1) \times 2} + \underbrace{(s_1 + 1) \times s_2}_{3 \times 3} + \underbrace{(s_2 + 1) \times s_3}_{4 \times 2} = (N + 1) \times 2 + 17 \\ &= \sum_{i=0}^L (s_i + 1) \times s_{i+1} \end{aligned}$$

### 3.3 Training

Training dataset (60-80%), Validation (10-20%), and test (10-20%)

Initialise  $\Rightarrow$  Forward propagation  $\Rightarrow$  Backpropagation  $\Rightarrow$  Update  $\Rightarrow$  Repeat until convergence

## 4 Convolution Neural Network

Drawbacks from feed-forward neural networks

- textitfully connected layers not suitable for images: contains too many parameters
- need to *flatten* 2D data into 1D array — losing spatial information.

Generic process

- input  $\rightarrow$  filter (convolutional layer)  $\rightarrow$  pooling layer  $\rightarrow$  (repeat)  $\rightarrow$  flatten  $\rightarrow$  fully-connected layer
- Numbers of neurons in a convolutional layer

$$N_d = \frac{W_d - F_d + 2P_d}{S_d} + 1$$

where  $N_d$  is the number of neurons,  $W$  is the width of the input channel,  $F$  is the width of the filter,  $P$  is the width of the pooling, and  $S$  is the Stride number.

- Number of parameter in convolution layer  $l$

$$l = (F_x \times F_y \times \text{n. channels} + 1 \text{ bias}) \times K$$

where  $K$  is the number of filters (also =number of feature maps)

## 5 Sequence Modelling

When both the data and its order contain information to perform the task, the points in the dataset are dependent on the position of other points in the dataset, e.g., speech recognition, DNA, time series. We lose key information if we shuffle the data.

### 5.1 Sequence modelling with feedforward neural network

Consider time as **feature (inputs)**, and the physical variable as **label**, may be OK for short time series, but generally not a good approach due to:

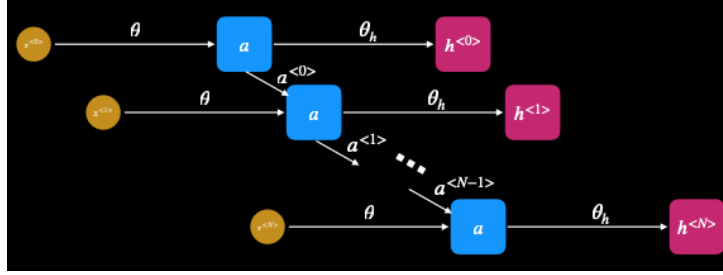
- do not learn correlation as opposed to dynamical system having temporal correlations, and patterns.
- large number of parameters
- needs the inputs of a fixed size.

## 5.2 Recurrent Neural Network

Store the previous time-step and feed it to the next time step as the hidden state.

$$a^{<N>} = a(\theta_a^{<N-1>} + \theta x^{<N>} + b)$$

$$h^{<N>} = \theta_h a^{<N>} + b_h$$



- The weights and biases are shared across the time-step
- **time is not a feature**
- the loss function is a sum of the loss functions of every time step by applying chain rule across different time steps.

$$J = \sum_{i=0}^N J^{<N>}$$

$$\frac{\partial J}{\partial a^{<i>}} = \frac{\partial J}{\partial a^{<N>}} \frac{\partial a^{<N>}}{\partial a^{<N-1>}} \dots \frac{\partial a^{<i+1>}}{\partial a^{<i>}} = \frac{\partial J}{\partial a^{<N>}} \underbrace{\theta_a \theta_a \dots \theta_a}_{N-i \text{ times}} = \frac{\partial J}{\partial a^{<N>}} \theta_a^{N-1}$$

- Unstable gradient occurs from above equation, when  $\theta_a$  is small, the gradient will become smaller and leads to **vanishing gradient** problem, make it difficult for the network to remember information from previous time steps. This can be solved by using ReLU activation functions due to its unbounded positive range.
- If  $\theta_a$  is large, this leads to exploding gradient, which can be solved by clip (rescale) the gradient
- Recurrent neural network are not ideal for learning temporal correlations in long time series due to short-term memory.

### NOTE:

- 3D tensor with shape (batch\_size, sequence\_length, input\_dim).
- **example** refers to a sequence of input data, where each element in the sequence is a time step, e.g., audio clip, where each time step is a small segment of the audio waveform.
  - In recurrent networks examples have to be input one at a time so that data in each forward propagation is still, effectively, one-dimensional
- **feature** refer to the characteristics of each time step or input element in the sequence, e.g. spectral representations of the audio waveform.

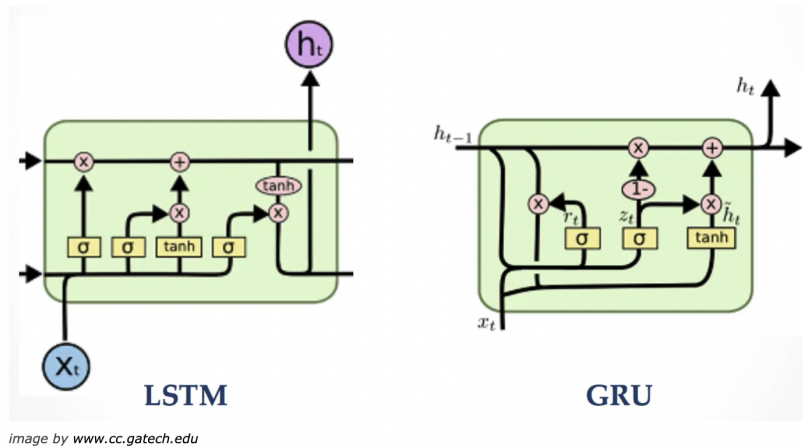
## 5.3 Long-short time memory units LSTMs

Introducing a *cell state*  $C$  to convey (store) information from the far past – long term memory, with the combination of the same short-memory past.

- Consists of three gates to control the status of the states (forget, input, output). Each gate is controlled by a **sigmoid** activation function, due to the extreme nature of being either 1 or 0.
- **Forget gate.** This determines what percentage of the Long-term memory is remembered.
- **Candidate state** combines the short-term memory and the input to create a **potential long-term memory**  $g$ . By using **tanh** to restrain the data in the range of -1 to 1.
- **Input gate** determines what percentage of the potential long-term memory is added to the long term memory  $C$ .
- **Output gate** calculates the potential long term memory from the cell state using a tanh function. Then, a sigmoid function is also used to determine what percentage of this potential long term memory is pass on as an output of the LSTM unit.

## 5.4 Gated recurrent units GRUs

- got rid of the cell state, and stored both long and short term memory in the **hidden state** and consists of 2 gates only.
- **Update gate** controls both the *input* and *forget* gate in determining what to forget and what to add on.
- **Reset gate** decides how much past information to forget.
- **Candidate state** contains past information.
- GRUs are simpler and with less parameters, making them easier to train than LSTMs



### NOTE:

A GRU layer takes as input a 3D tensor of shape  $(\mathbf{B}, \mathbf{T}, \mathbf{F})$ , where  $\mathbf{B}$  is the batch size,  $\mathbf{T}$  is the number of timestep in each sequence, and  $\mathbf{F}$  is the number of features in each timestep. It processes each timestep in the input sequence independently and generates an output for each timestep. Therefore, the output of the GRU layer will also be a 3D tensor of shape  $(\mathbf{B}, \mathbf{T}, \mathbf{N})$ , where  $\mathbf{N}$  is the number of neurons in the layer.

## 6 Unsupervised Learning

Methods that do not require labelled dataset. **Clustering** groups the instances based on their similarity without class labels.

### 6.1 k-Means clustering (Lloyd-Forgy algorithm)

Given a set of vector-valued data,  $\mathbf{x}^{(i)}$ , with the goal of grouping  $m$  observation into  $k$  clusters. k-means algorithm is iterative and is described as following:

- **Choose a metric (distance function).** This quantitatively defines the notion of *similarity* by measuring the distance between points in the *feature space*. The distance are usually calculated using the squared Euclidean distance  $\|\cdot\|^2$ .
- **Choose number of clusters,  $k$**
- **Choose the centroids** These are the baricentres (means),  $\mu_j$ , of the clusters, representing the positions of the clusters.
- **Compute the distance and Assign** of each instance with respect to each centroid,  $\|\mathbf{x}^{(i)} - \mu_j\|^2$  for  $i = 1, 2, \dots, k$ .
- **Compute the baricentre** of each centroid,  $\mu_{j,new} = 1/N_j \sum_{j=1}^{N_j} \mathbf{x}^{(j)}$ , where  $N_j$  is the number of points,  $\mathbf{x}^{(i)}$ , that belong in cluster  $j$ . Therefore,  $\sum_{i=1}^k N_i = m$ .
- **Update centroid.** The new centroid is the baricentre from the previous new mean.
- Compute the new distance, under the mean distance from the centroid stops changing.

Although k-means algorithm presents a fast and effective method to cluster data points, it has some limits: (i) though the algorithm is guaranteed to converge, it might not be the globally optimal solution, mainly because it depends on the initial centroid initialisation and numbers of cluster chosen. (ii) it does not perform well when clusters have varying sizes.

There are several approaches to be taken to optimise the solution: (i) different centroid initialisation, (ii) number of clusters, a simple method is to plot the variance and the number of cluster, and choose the number of clusters at the inflection point (elbow). This is called the *elbow method*. (iii) Feature scaling, scale the input features before running k-means.

## 6.2 Density-based spatial clustering of applications with noise (DBSCAN)

A clustering algorithm that can handle nested clusters, and identifying nested algorithms at higher dimensions where we can't perceive by eye. DBSCAN groups points that are closely-packed, identifies the points that are isolated as outliers.

- **Choose a metric (distance function)** — Euclidean distance,  $\|\cdot\|^2$
- **Choose the vicinity,  $\varepsilon$** , to other points, i.e., the neighbourhood  $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2 \leq \varepsilon$ . Choose  $N$  to define the number of points that are needed to classify a neighbourhood as high density.
- **Select a instance and Count** how many instances are located with distance  $\varepsilon$  from it.
  - If an instance has at least  $N$  in its  $\varepsilon$ , then it is considered as a *core instance*. All instances in  $\varepsilon$  of a core instance belong to the same cluster. This may also include other core instances.
  - If fewer than  $N$  belong in  $\varepsilon$ , it is identify as *Non-core* instance.
- **Classify** anomalies that are not core instances.

DBSCAN provides the advantage of not needing to specify the number of clusters, can find clusters with complicated shapes and can also identifies outliers. However, it does not cluster effectively when there are large differences in densities.

## 6.3 Hierarchical clustering

Hierarchical clustering does not required user input cluster numbers, it does it by inspecting the result. This section will be focusing on using **Bottom-up**, in which each instance defines a cluster a iteration 0. The algorithm then pairs and merge clusters.

- **Choose a metric (distance function)** — Euclidean distance,  $\|\cdot\|^2$ .
- **Compute** the distance between all the instances.
- **Merge** the closest pair into a new cluster, which is located midway its original instance.

- **Repeat** Step 1-2 with  $m - 1$  instances.
- **Stop** when left with one cluster.

## 6.4 Elliptic envelope for anomaly detection

Using Gaussian density distribution as an approach to identify outliers.

### 6.4.1 Multivariate Gaussian distribution

A column vector of random variable (feature),  $\mathbf{X} = [X_1, X_2, \dots, X_N]^T$  is Gaussian distributed.

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|_{\Sigma}^2}{2}\right), \quad \boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

where  $\boldsymbol{\mu}$  is the column-vector of the means and  $\Sigma$  is the **covariance matrix**, which can be computed from a sample,  $m$ , as (usually presented as matrix). Covariance matrix is **symmetric** and can give insights into whether the variables are correlated (only diagonal values, or the principal axes of the ellipsoid is not parallel to the feature axis), or uncorrelated (non-diagonal). It also represents the eigenvalues, which are non-negative.

$$\Sigma = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^T$$

And finally the norm that measures the distance between the mean and a feature is

$$\|\mathbf{x} - \boldsymbol{\mu}\|_{\Sigma}^2 = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

The inverse of the covariance matrix,  $\Sigma^{-1}$ , is also known as the *precision* matrix.

### 6.4.2 Elliptic envelope algorithm

elliptic envelope algorithm is a distance-based anomaly detection method that can be summarised as

1. **Assume** the datapoints to be Gaussian distributed
2. **Compute** mean and covariance (covariance matrix also defines the shape of the ellipse)
3. **Compute** the distance,  $D$ , between testpoints and the mean
4. **Choose** the decision boundary,  $c$  (the range of the elliptic)
5. **If**  $D > c$ , the test point is an outlier (anomaly).

## 7 Decision Tree

Decision tree contains three main components: (i) **Internal nodes** that provides classification criteria, (ii) **Branches** that listed out the features that satisfies the previous criteria, and (iii) **End nodes (leaves)** that is the final result for classification.

Noteworthy that when classifying between binary classes, the criteria can be set for e.g.,  $x_1 = 3$  where equal sign could be used. But when dealing with numerical values, **inequalities** must be used,  $x \leq 3$ .

### 7.1 Impurity Measures

Decision tree is constructed in the learning process node by node. The following nodes, or *children nodes*, for the partition are constructed with minimum impurity. Two impurity measure will be discussed here — *information gain*, *Gini diversity index*

- Information gain



$x_1$	3	7	3	4	3	6	5	6
$x_2$	5	6	3	8	9	5	8	4
label	pass	fail	pass	fail	fail	pass	fail	fail

$x_1$	3	6	3	6	7	4	5	3
$x_2$	3	4	5	5	6	8	8	9
label	pass	fail	pass	pass	fail	fail	fail	fail

- information gain tries to find a suitable separation line in the features that can classify data in a more effective way. e.g. by shuffling the order of the feature: In this case, we can place a separation line at  $x_2 \geq 6$  to better separate the fail by one node.
- **Entropy** can be used to evaluate the difficulty of predicting category.

$$H = - \sum_{i=1}^n (p_i \log_2 p_i)$$

The aim is to find features with **highest information gain** (this feature will be split into children nodes to further). **IG** is **highest** when splitting the feature makes the biggest reduction in entropy.

$$IG(feature) = H(\text{label}) - H(\text{label}, feature)$$

- Gini Impurity Criterion

- How often a randomly chosen element from the set is incorrectly labelled. Higher GI, higher chance of misclassification

$$GI = 1 - \sum_{i=0}^J P_i^2$$

- Split with the lowest weighted FI for the child trees to achieve largest reduction in GI.
- Shallow trees has less variance but higher bias. Deep trees has less bias but more prone to overfitting.

**NOTE:**

For a classification model, high variance  $\rightarrow$  overfitting, high bias  $\rightarrow$  underfitting.

- If the model classifies only yellow helicopters as a helicopter.
  - It is overfitting (high variance). This is because the model may be too complex and, thus, it fits features that are not relevant to the classification (the colour).
- If the model classifies any round object as a helicopter.
  - It is underfitting (high bias). This is because the model is unable to learn more complex feature, such as the bladed, the helicopter tail, etc.

## 8 Principal Component Analysis (PCA)

finds an optimal direction (principle axis of the ellipsoid in gaussian distribution) of maximum variance. The maximum variance in the direction  $\hat{w}_1$  is given by the maximum eigenvalue,  $\lambda_1$ , of the covariance matrix. Geometrically, this direction is the principal axis (eigenvector) of the ellipsoid, which is the direction along which the ellipsoid is most elongated.

The difference between PCA and linear regression is:

- The sensitivity on how you define input and output is large ( $y = ax + b$  or  $x = ay + b$ ), provides very different result
- linear regression is sensitive to outliers
- The errors in PCA are orthogonal to the principal axis.