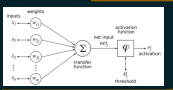


$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

The perceptron is an algorithm for learning a binary classifier: a function that maps its input x (a real-valued vector) to an output value $f(x)$ (a single binary value)

Perception



Unit (Neuron)

Comprised of multiple Real-Valued inputs. Each input must be linearly independent from each other.

Input Layer



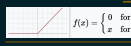
Hidden Layers

Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch is an estimate of the gradient over the training set, whose quality improves as the batch size increases. Second, computation over a batch can be much more efficient than in computations for individual examples, due to the parallelism afforded by the modern computing platforms.

Batch Normalization

Cost Function

Defines the output of that node given an input or set of inputs.



ReLU



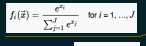
Sigmoid / Logistic



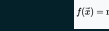
Tanh



Softplus



Softmax



Maxout

Leaky ReLU, PReLU, RReLU, ELU, SELU, and others.

Is a method used in artificial neural networks to calculate the error contribution of each neuron after a batch of data. It calculates the gradient of the loss function. It is commonly used in the gradient descent optimisation algorithm. It is also called backward propagation of errors, because error is calculated at the output and distributed back through the network layers.

Backpropagation

Learning Rate

Weight Initialization

Is a regularisation technique for reducing overfitting in neural networks by preventing complex co-adaptations on training sets. It is a very efficient way of performing model averaging with neural networks. The term "dropout" refers to dropping out units (both hidden and visible) in a neural network.

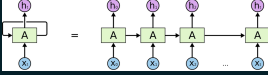
Dropout

Output Layer



They have applications in image and video recognition, recommender systems and natural language processing.

Pooling
Convolution
Subsampling

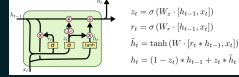


Is a class of artificial neural network where connections between units form a directed cycle. This allows it to exhibit dynamic temporal behavior. Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs.



Is a kind of deep neural network created by applying the same set of weights recursively over a structure, to produce a structured prediction over variable-size input structures, or a scalar prediction on k , by traversing a given structure in topological order.

RNNs have been successful for instance in learning sequence and tree structures in natural language processing, many phrase and sentence continuous representations based on word embedding.

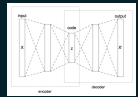


It is a type of recurrent (RNN), allowing data to flow both forwards and backwards within the network.

An LSTM is well-suited to learn from experience to classify, process and predict time series given time lags of unknown size and bound between important events. Relative inaccessibility to gap length gives an advantage to LSTM over alternative (RNNs, hidden Markov models and other sequence learning methods in numerous applications.



GANs are a class of artificial intelligence algorithms used in unsupervised machine learning, implemented by a system of two neural networks contesting with each other in a zero-sum game framework.



The aim of an auto-encoder is to learn a representation (encoding) for a set of data, typically for the purpose of dimensionality reduction. Recently, the auto-encoder concept has become more widely used for learning generative models of data.

Is an artificial neural network used for unsupervised learning of efficient codings.

Neural Networks

Machine Learning Models

Regression

$$\hat{y} = \beta_0 + \beta_1 x + \epsilon$$
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad \text{data } i = 1, \dots, n$$

Linear Regression

Is a flexible generalization of ordinary linear regression that allows for response variables that have error distribution models other than a normal distribution. The GLM generalizes linear regression by allowing the linear model to be related to the response variable via a link function and by allowing the magnitude of the variance of each measurement to be a function of its predicted value.

Generalised Linear Models (GLMs)

Identity $\mu = \mathbf{X}\beta$

Inverse $\mu = (\mathbf{X}\beta)^{-1}$

Logit $\mu = \frac{\exp(\mathbf{X}\beta)}{1 + \exp(\mathbf{X}\beta)} = \frac{1}{1 + \exp(-\mathbf{X}\beta)}$

Cost Function is found via Maximum Likelihood Estimation

Locally Estimated Scatterplot Smoothing (LOESS)

Ridge Regression

Least Absolute Shrinkage and Selection Operator (LASSO)

$$p(\mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

Logistic Regression



Bayesian

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

Naive Bayes

Multinomial Naive Bayes

Bayesian Belief Network (BBN)

$$\hat{y} = \underset{k}{\operatorname{argmax}} \frac{p(C_k)}{n} \prod_{i=1}^n p(x_i | C_k)$$

Naive Bayes Classifier: We neglect the denominator as we calculate for every class and pick the max of the numerator.

Dimensionality Reduction

Principal Component Analysis (PCA)

Partial Least Squares Regression (PLSR)

Principal Component Regression (PCR)

Partial Least Squares Discriminant Analysis

Quadratic Discriminant Analysis (QDA)

Linear Discriminant Analysis (LDA)

Instance Based

k-nearest Neighbour (KNN)

Learning Vector Quantization (LVQ)

Self-Organising Map (SOM)

Locally Weighted Learning (LWL)

Decision Tree

Random Forest

Classification and Regression Tree (CART)

Gradient Boosting Machines (GBM)

Conditional Decision Trees

Gradient Boosted Regression Trees (GBRT)

Clustering

Hierarchical Clustering

Linkage

complete

single

average

centroid

Dissimilarity Measure

Euclidean

Euclidean distance or Euclidean metric is the "ordinary" straight-line distance between two points in Euclidean space.

Manhattan

The distance between two points measured along axes at right angles.

k-Means

How many clusters do we select?

k-Medians

Fuzzy C-Means

Self-Organising Maps (SOM)

Expectation Maximization

DBSCAN

Data Structure Metrics

Dunn Index

Connectivity

Distance With

Non-overlapping APRI

Average Distance AD

Average Distance Between Means ADM

Figure of Merit FOM

Validation

Stability Metrics