

# Machine Learning Process

## Question

- Classification: Is this A or B?
- Regression: How much, or how many of these?
- Anomaly Detection: Is this anomalous?
- Clustering: How can these elements be grouped?
- Reinforcement Learning: What should I do now?

## Direction

- SaaS - Pre-built Machine Learning models
  - Google Cloud
    - Vision API
    - Speech API
    - Jobs API
    - Video Intelligence API
    - Language API
    - Translation API
  - AWS
    - Rekognition
    - Lex
    - Polly
    - ... many others
  - ... many others
- Data Science and Applied Machine Learning
  - Google Cloud
    - ML Engine
  - AWS
    - Amazon Machine Learning
    - Tools: Jupyter / Datalab / Zeppelin
    - ... many others
  - ... many others
- Machine Learning Research
  - Tensorflow
  - MXNet
  - Torch
  - ... many others

## Data

- Find
- Collect
- Explore
- Clean Features
- Impute Features
- Engineer Features
- Select Features
- Encode Features
- Build Datasets

Machine Learning is math. In specific, performing Linear Algebra on Matrices. Our data values must be numeric.

## Model

Select Algorithm based on question and data available

## Cost Function

- The cost function will provide a measure of how far my algorithm and its parameters are from accurately representing my training data.
- Sometimes referred to as Cost or Loss function when the goal is to minimise it, or Objective function when the goal is to maximise it.

## Optimization

Having selected a cost function, we need a method to minimise the Cost function, or maximise the Objective function. Typically this is done by Gradient Descent or Stochastic Gradient Descent.

## Tuning

Different Algorithms have different Hyperparameters, which will affect the algorithms performance. There are multiple methods for Hyperparameter Tuning, such as Grid and Random search.

## Results and Benchmarking

- Analyse the performance of each algorithms and discuss results.
- Are the results good enough for production?
- Is the ML algorithm training and inference completing in a reasonable timeframe?

## Scaling

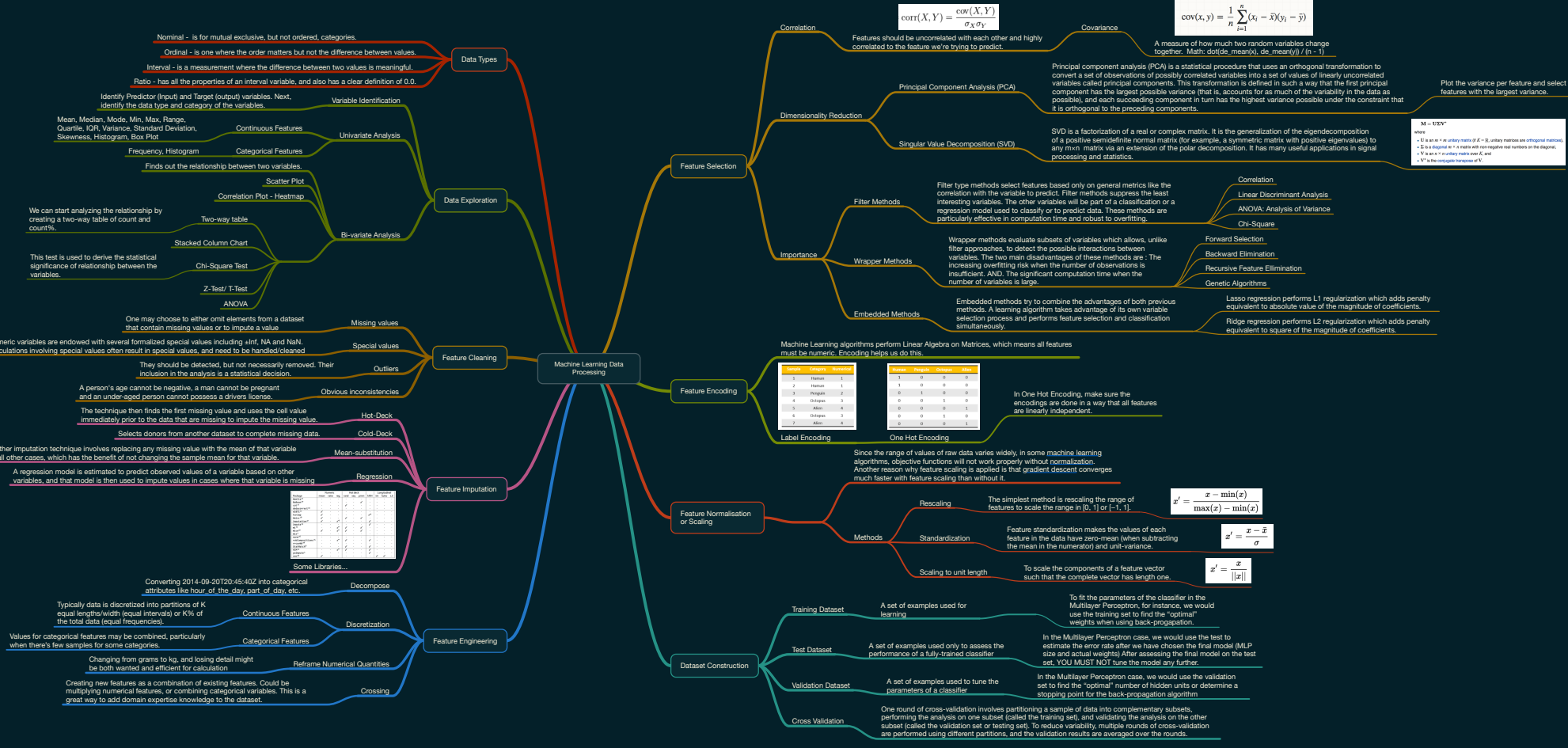
How does my algorithm scales for both training and inference?

## Deployment and Operationalisation

- How can feature manipulation be done for training and inference in real-time?
- How to make sure that the algorithm is retrained periodically and deployed into production?
- How will the ML algorithms be integrated with other systems?

## Infrastructure

- Can the infrastructure running the machine learning process scale?
- How is access to the ML algorithm provided? REST API? SDK?
- Is the infrastructure adapter to the algorithm we are running? Should GPU's be considered rather than CPUs?



# Machine Learning Concepts

- Motivation**
  - Prediction**: When we are interested mainly in the predicted variable as a result of the inputs, but not on the each way of the inputs affect the prediction. In a real estate example, Prediction would answer the question of: Is my house over or under valued? Non-linear models are very good at these sort of predictions, but not great for inference because the models are much less interpretable.
  - Inference**: When we are interested in the way each one of the inputs affect the prediction. In a real estate example, Prediction would answer the question of: How much would my house cost if I had a view of the sea? Linear models are more suited for inference because the models themselves are easier to understand than their non-linear counterparts.
- Confusion Matrix**

	Actual Positive	Actual Negative
Predicted Positive	True Positives (TP)	False Positives (FP)
Predicted Negative	False Negatives (FN)	True Negatives (TN)

Fraction of correct predictions, not reliable as skewed when the data set is unbalanced (that is, when the number of samples in different classes vary greatly).
- Accuracy**
  - Precision**:  $TP / (TP + FP)$  which tells us what proportion of persons we diagnosed as having cancer actually had cancer. In other words, proportion of TP to the set of positive cancer diagnoses. This is given by the rightmost column in the confusion matrix.  
Out of all the examples the classifier labelled as positive, what fraction were correct?
  - Recall**:  $TP / (TP + FN)$  which tells us what proportion of persons that actually had cancer were diagnosed by which having cancer. In other words, proportion of TP to the set of true cancer cases. This is given by the bottom row in the confusion matrix.  
Out of all the positive examples there were, what fraction did the classifier pick up?
  - f1 score**: Harmonic Mean of Precision and Recall:  $(2 * p * r) / (p + r)$
- Performance Analysis**
  - ROC Curve - Receiver Operating Characteristics**

Bias refers to the amount of error that is introduced by approximating a real-life problem, which may be extremely complicated, by a simple model. If Bias is high, and/or if the algorithm performs poorly even on your training data, try adding more features, or a more flexible model.

Variance is the amount our model's prediction would change when using a different training data set. High: Remove features, or obtain more data.

Goodness of Fit =  $R^2/2$   
 $1.0 - \text{sum of squared errors} / \text{total sum of squares}$
  - Bias-Variance Tradeoff**
  - Mean Squared Error (MSE)**:  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$   
The mean squared error (MSE) or mean squared deviation (MSD) of an estimator (of a procedure for estimating an unobserved quantity) measures the average of the squares of the errors or deviations—that is, the difference between the estimator and what is estimated.
  - Error Rate**:  $\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$   
The proportion of mistakes made if we apply out estimate model function to the training observations in a classification setting.
- Cross-validation**
  - Methods**: Leave-p-out cross-validation, Leave-one-out cross-validation, k-fold cross-validation, Holdout method, Repeated random sub-sampling validation.
  - Grid Search**: The traditional way of performing hyperparameter optimization has been grid search, or a parameter sweep, which is simply an exhaustive searching through a manually specified subset of the hyperparameter space of a learning algorithm. A grid search algorithm must be guided by some performance metric, typically measured by cross-validation on the training set or evaluation on a held-out validation set.
  - Random Search**: Since grid searching is an exhaustive and therefore potentially expensive method, several alternatives have been proposed. In particular, a randomized search that simply samples parameter settings a fixed number of times has been found to be more effective in high-dimensional spaces than exhaustive search.
  - Hyperparameters**
    - Gradient-based optimization**: For specific learning algorithms, it is possible to compute the gradient with respect to hyperparameters and then optimize the hyperparameters using gradient descent. The first usage of these techniques was focused on neural networks. Since then, these methods have been extended to other models such as support vector machines or logistic regression.
    - Early Stopping (Regularization)**: Early stopping rates provides guidance as to how many iterations can be run before the learner begins to overfit, and stop the algorithm then.
    - Overfitting**: When a given method yields a small training MSE (or cost), but a large test MSE (or cost), we are said to be overfitting the data. This happens because our statistical learning procedure is trying too hard to find patterns in the data, that might be due to random chance, rather than a property of our function. In other words, the algorithms may be learning the training data too well. If model underfits, try removing some features, decreasing degrees of freedom, or adding more data.
    - Underfitting**: Opposite of Overfitting. Underfitting occurs when a statistical model or machine learning algorithm cannot capture the underlying trend of the data. It occurs when the model or algorithm does not fit the data enough. Underfitting occurs if the model or algorithm shows low variance but high bias to contrast the opposite, overfitting from high variance and low bias. It is often a result of an excessively simple model.
  - Tuning**
    - Bootstrap**: Test that applies Random Sampling with Replacement of the available data, and assigns measures of accuracy (bias, variance, etc.) to sample estimates.
    - Bagging**: An approach to ensemble learning that is based on bootstrapping. Shortly, given a training set, we produce multiple different training sets (called bootstrap samples), by sampling with replacement from the original dataset. Then, for each bootstrap sample, we build a model. The results in an ensemble of models, where each model votes with the equal weight. Typically, the goal of this procedure is to reduce the variance of the model of interest (e.g. decision trees).

- Types**
  - Regression**: A supervised problem, the outputs are continuous rather than discrete.
  - Classification**: Inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (multi-label classification) of these classes. This is typically tackled in a supervised way.
  - Clustering**: A set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task.
  - Density Estimation**: Finds the distribution of inputs in some space.
  - Dimensionality Reduction**: Simplifies inputs by mapping them into a lower-dimensional space.
- Kind**
  - Parametric**
    - Step 1: Making an assumption about the functional form or shape of our function (f), i.e.: f is linear, thus we will select a linear model.
    - Step 2: Selecting a procedure to fit or train our model. This means estimating the Beta parameters in the linear function. A common approach is the (ordinary) least squares, amongst others.
  - Non-Parametric**: When we do not make assumptions about the form of our function (f). However, since these methods do not reduce the problem of estimating f to a small number of parameters, a large number of observations is required in order to obtain an accurate estimate for f. An example would be the thin-plate spline model.
- Categories**
  - Supervised**: The computer is presented with example inputs and their desired outputs, given by a "teacher", and the goal is to learn a general rule that maps inputs to outputs.
  - Unsupervised**: No labels are given to the learning algorithm, leaving it on its own to find structure in its input. Unsupervised learning can be a goal in itself (discovering hidden patterns in data) or a means towards an end (feature learning).
  - Reinforcement Learning**: A computer program interacts with a dynamic environment in which it must perform a certain goal (such as driving a vehicle or playing a game against an opponent). The program is provided feedback in terms of rewards and punishments as it navigates its problem space.
- Approaches**
  - Decision tree learning
  - Association rule learning
  - Artificial neural networks
  - Deep learning
  - Inductive logic programming
  - Support vector machines
  - Clustering
  - Bayesian networks
  - Reinforcement learning
  - Representation learning
  - Similarity and metric learning
  - Sparse dictionary learning
  - Genetic algorithms
  - Rule-based machine learning
  - Learning classifier systems
- Taxonomy**
  - Generative Methods**
    - Model class-conditional pdfs and prior probabilities. "Generative" since sampling can generate synthetic data points.
    - Popular models
      - Gaussians, Naïve Bayes, Mixtures of multinomials
      - Mixtures of Gaussians, Mixtures of experts, Hidden Markov Models (HMM)
      - Sigmoidal belief networks, Bayesian networks, Markov random fields
  - Discriminative Methods**
    - Directly estimate posterior probabilities. No attempt to model underlying probability distributions. Focus computational resources on given task-better performance.
    - Popular Models
      - Logistic regression, SVMs
      - Traditional neural networks, Nearest neighbor
      - Conditional Random Fields (CRF)
- Selection Criteria**
  - Prediction Accuracy vs Model Interpretability**: There is an inherent tradeoff between Prediction Accuracy and Model Interpretability, that is to say that as the model get more flexible in the way the function (f) is selected, they get obscured, and are hard to interpret. Flexible methods are better for inference, and inflexible methods are preferable for prediction.
- Libraries**
  - Python**
    - Numpy**: Adds support for large, multi-dimensional arrays and matrices, along with a large library of high-level mathematical functions to operate on these arrays.
    - Pandas**: Offers data structures and operations for manipulating numerical tables and time series.
    - Scikit-Learn**: It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy.
    - Tensorflow**
      - Components**
        - Does lazy evaluation. Need to build the graph, and then run it in a session.
    - MXNet**: Is an modern open-source deep learning framework used to train, and deploy deep neural networks. MXNet library is portable and can scale to multiple GPUs and multiple machines. MXNet is supported by major Public Cloud providers including AWS and Azure. Amazon has chosen MXNet as its deep learning framework of choice at AWS.
    - Keras**: Is an open source neural network library written in Python. It is capable of running on top of MXNet, TensorFlow, CNTK or Theano. Designed to enable fast experimentation with deep neural networks, it focuses on being minimal, modular and extensible.
    - Torch**: Torch is an open source machine learning library, a scientific computing framework, and a script language based on the Lua programming language. It provides a wide range of algorithms for deep machine learning, and uses the scripting language LuaJIT, and an underlying C implementation.
    - Microsoft Cognitive Toolkit**: Previously known as CNTK and sometimes styled as The Microsoft Cognitive Toolkit, is a deep learning framework developed by Microsoft Research. Microsoft Cognitive Toolkit describes neural networks as a series of computational steps via a directed graph.

