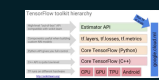
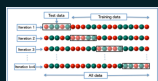
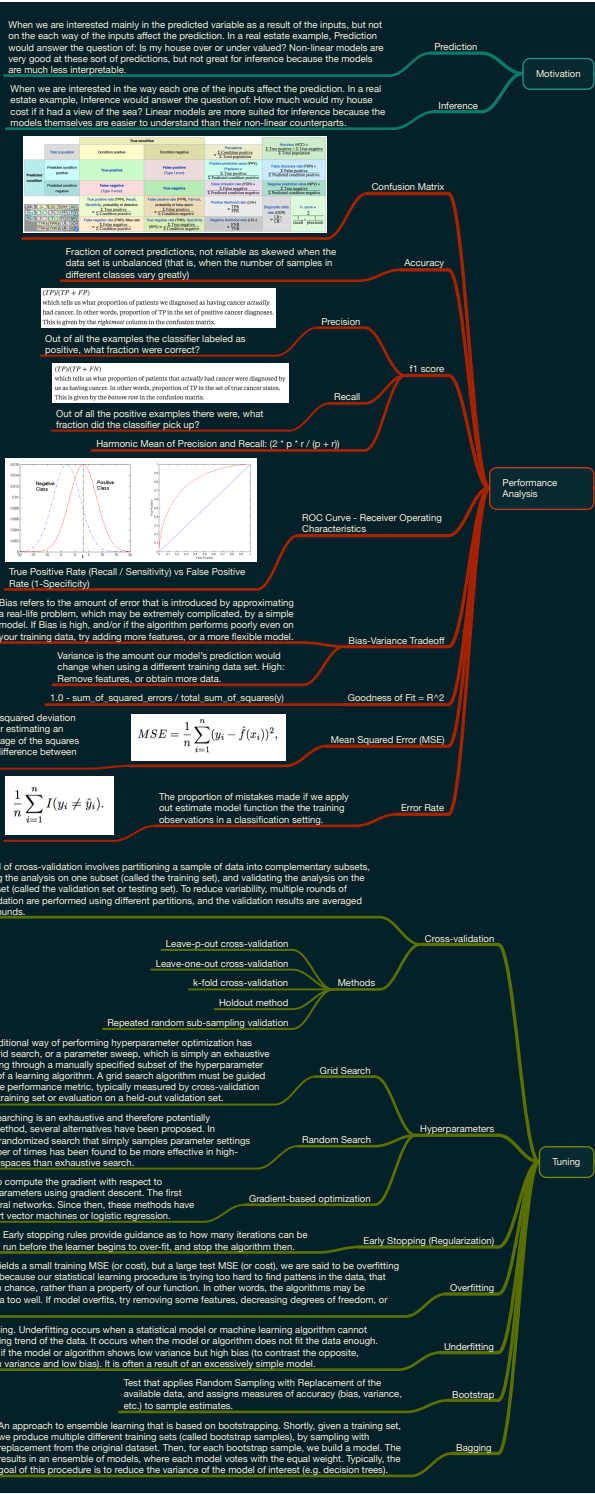
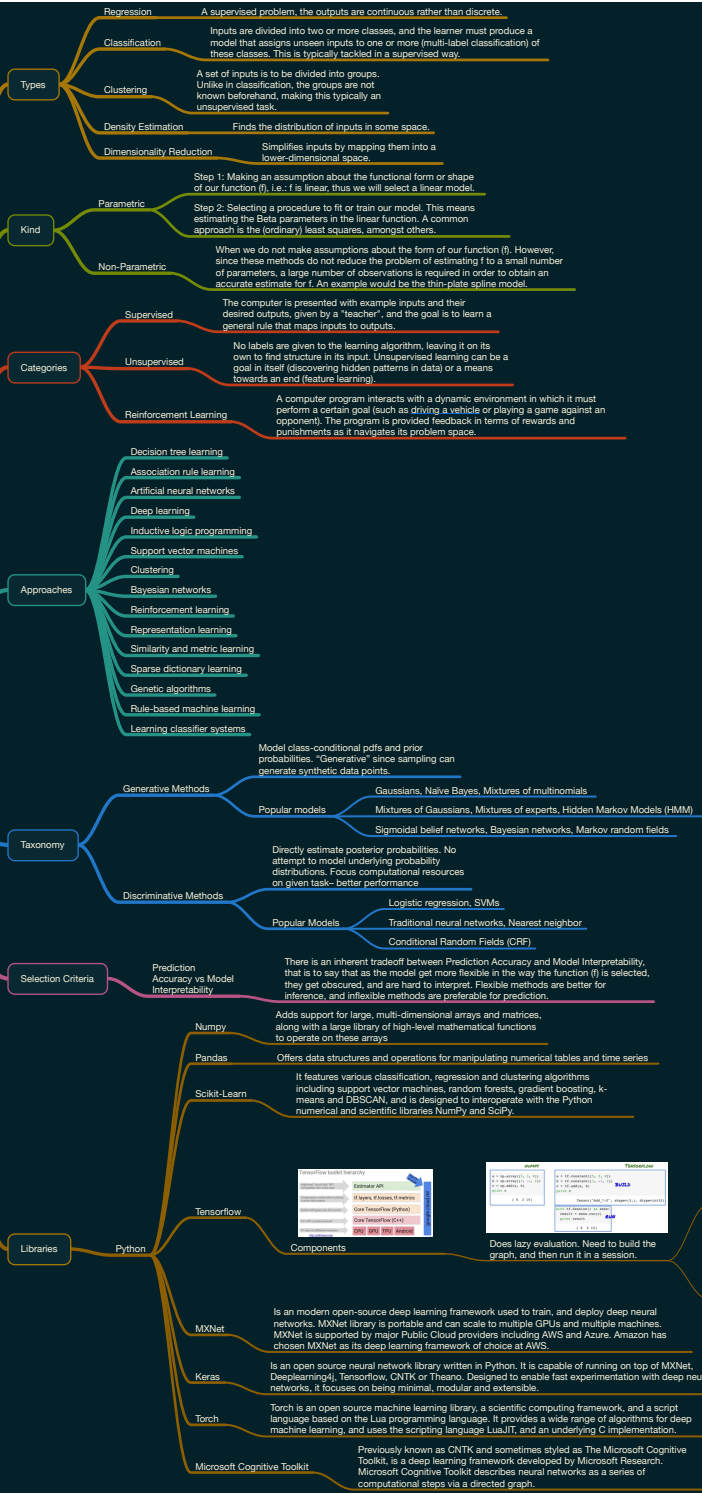
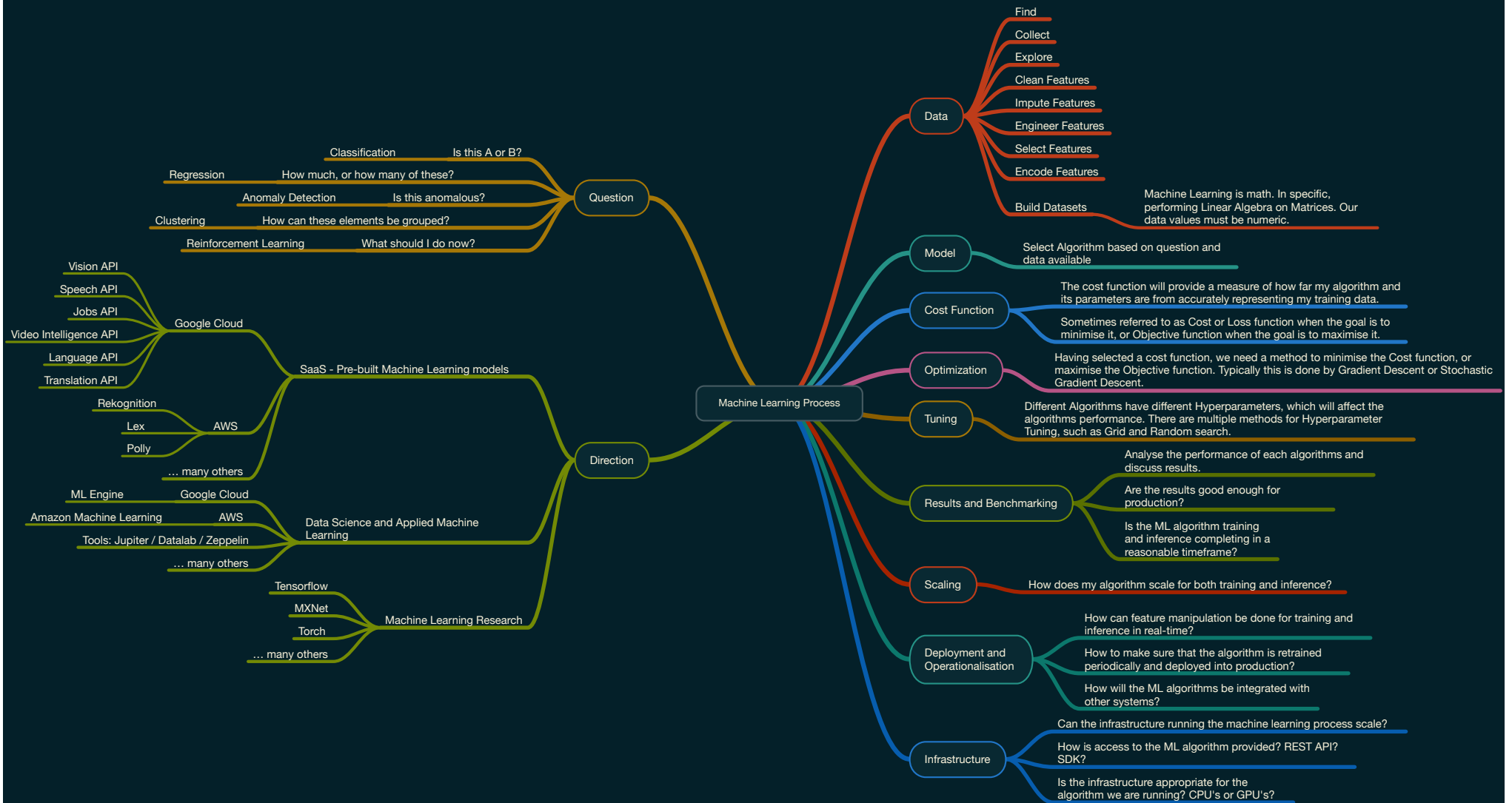


Machine Learning Concepts





Mathematics

Basic Operations: Addition, Multiplication, Subtraction

Transformations

Exponents and Logarithms

Differentiation Chain Rule

Integration

Linear Algebra

Calculus

Statistics

Probability

Optimization

Regression

Density Estimation

Mathematics

Basic Operations: Addition, Multiplication, Subtraction

Transformations

Exponents and Logarithms

Differentiation Chain Rule

Integration

Linear Algebra

Calculus

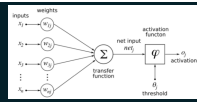
Statistics

Probability

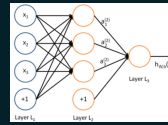
Optimization

Regression

Density Estimation



A unit often refers to the activation function in a layer by which the inputs are transformed via a nonlinear activation function (for example by the logistic sigmoid function). Usually, a unit has several incoming connections and several outgoing connections.



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

Layers other than the input and output layers. A layer is the highest-level building block in deep learning. A layer is a container that usually receives weighted input, transforms it with a set of mostly non-linear functions and then passes these values as output to the next layer.

With SGD, the training proceeds in steps, and at each step we consider a mini-batch $x_{i:m}$ of size m . The mini-batch is used to approximate the gradient of the loss function with respect to the parameters.

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.

However, if you actually try that, the weights will change far too much each iteration, which will make them "overcorrect" and the loss will actually increase/diverge. So in practice, people usually multiply each derivative by a small value called the "learning rate" before they subtract it from its corresponding weight.

Neural networks are often trained by gradient descent on the weights. This means at each iteration we use backpropagation to calculate the derivative of the loss function with respect to each weight and subtract it from that weight.

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_t(\theta)$$

Simplest recipe: keep it fixed and use the same for all parameters.

$$\alpha = \frac{\alpha_0 \tau}{\max(t, \tau)}$$

Reduce by 0.5 when validation error stops improving

Reduction by $O(1/\tau)$ because of theoretical convergence guarantees, with hyper-parameters ϵ_0 and τ and t is iteration numbers

Better yet: No hand-set learning of rates by using AdaGrad

Better results by allowing learning rates to decrease Options:

The implementation for weights might simply drawing values from a normal distribution with zero mean, and unit standard deviation. It is also possible to use small numbers drawn from a uniform distribution, but this seems to have relatively little impact on the final performance in practice.

Thus, you still want the weights to be very close to zero, but not identically zero. In this way, you can random these neurons to small numbers which are very close to zero, and it is treated as symmetry breaking. The idea is that the neurons are all random and unique in the beginning, so they will compute distinct updates and integrate themselves as diverse parts of the neural network.

In the ideal situation, with proper data normalization it is reasonable to assume that approximately half of the weights will be positive and half of them will be negative. A reasonable-sounding idea then might be to set all the initial weights to zero, which you expect to be the "best guess" in expectation.

All Zero Initialization

Initialization with Small Random Numbers

Calibrating the Variances

One problem with the above suggestion is that the distribution of the outputs from a randomly initialized neuron has a variance that grows with the number of inputs. It turns out that you can normalize the variance of each neuron's output to 1, by scaling its weight vector by the square root of its fan-in (i.e., its number of inputs)

$$J = \max(0, 1 - s + \kappa_s)$$

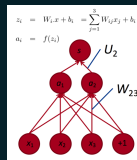
$s = U^T f(Wx + b)$
 $s_0 = U^T f(Wx_0 + b)$

$$s = U^T f(Wx + b) \quad x \in \mathbb{R}^{20 \times 1}, W \in \mathbb{R}^{8 \times 20}, U \in \mathbb{R}^{8 \times 1}$$

$$s = U^T a$$
$$a = f(s)$$
$$z = Wx + b$$

Neural Network taking 4 dimension vector representation of words

$$\frac{\partial a}{\partial x_j} = \sum_{i=1}^K \frac{\partial a}{\partial u_i} \frac{\partial u_i}{\partial x_j} = \sum_{i=1}^K \frac{\partial (U^T a)}{\partial u_i} \frac{\partial u_i}{\partial x_j} = \sum_{i=1}^K U_i \frac{\partial f(Wx + b)}{\partial x_j} = \sum_{i=1}^K U_i f'(Wx + b) \frac{\partial W_{ij} x}{\partial x_j} = \sum_{i=1}^K U_i W_{ij}$$
$$W_{ij} = \sum_{j=1}^n \frac{\partial f}{\partial x_j} W_{ij}$$



In this method, we reuse partial derivatives computed for higher layers in lower layers, for efficiency.

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

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

ReLU

$$f(x) = \frac{1}{1 + e^{-x}}$$

Sigmoid / Logistic

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Binary

$$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

Tanh

$$f(x) = \ln(1 + e^x)$$

Softplus

$$f_i(\vec{x}) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}} \quad \text{for } i = 1, \dots, J$$

Softmax

$$f(\vec{x}) = \max_i x_i$$

Maxout

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

Neural Networks

Machine Learning Models

Regression

$$Y = \beta_0 + \beta_1 X + \epsilon$$
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (\text{data: } i = 1 \dots n)$$

Linear Regression

Generalised Linear Models (GLMs)

Identity

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

Link Function

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

Inverse

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

Logit

Cost Function is found via Maximum Likelihood Estimation

Locally Estimated Scatterplot Smoothing (LOESS)

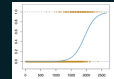
Ridge Regression

Least Absolute Shrinkage and Selection Operator (LASSO)

Logistic Regression

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

Logistic Function



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)

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

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)

Random Forest

Classification and Regression Tree (CART)

Gradient Boosting Machines (GBM)

Conditional Decision Trees

Gradient Boosted Regression Trees (GBRT)

Decision Tree

Clustering

Hierarchical Clustering

Linkage

complete

single

average

centroid

Disimilarity 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 Structures Metrics

Dunn Index

Connectivity

Silhouette Width

Stability Metrics

Non-overlap APN

Average Distance AD

Average Distance Between Means ADM

Figure of Merit FOM