

CHAPTER 1: MACHINE LEARNING BASICS

Theory

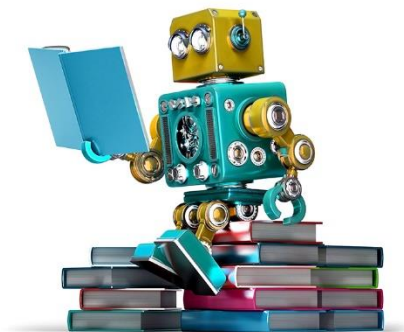
This chapter explains the basic concepts of machine learning and deep learning such as supervised and unsupervised, reinforcement learning, and linear regression, binary and multi-label classification and so on.

Basic concepts

- What is Machine Learning?
- What is learning?
 - o Supervised
 - o Unsupervised
 - o Reinforcement
- What is regression?
- What is classification?

Machine Learning

Today machine learning and deep learning is becoming a widespread area of artificial intelligence, including spam filtering, auto driving, face recognition, and more.



So what is machine learning?

Machine learning is a field of computer science that gives computers the ability to learn without being explicitly programmed.

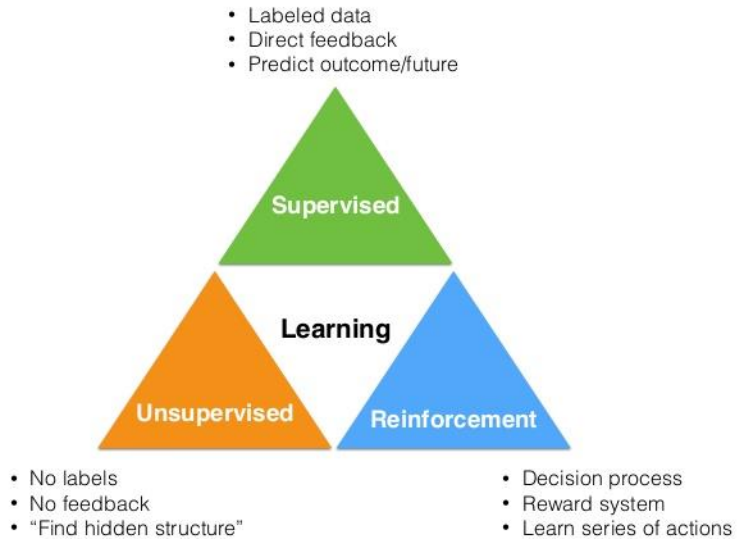
Arthur Samuel, an American pioneer in the field of computer gaming and artificial intelligence, coined the term "Machine Learning" in 1959 while at IBM. Evolved from the study of pattern recognition and computational learning theory in artificial intelligence, machine learning explores the study and construction of algorithms that can learn from and make predictions on data— such algorithms overcome following strictly static program instructions by making data-driven predictions or decisions, through building a model from sample inputs. Machine learning is employed in a range of computing tasks where designing and programming explicit algorithms with good performance is difficult or infeasible; example applications include email filtering, detection of network intruders or malicious insiders working towards a data breach, optical character recognition(OCR), learning to rank, and computer vision.

Machine learning is closely related to (and often overlaps with) computational statistics, which also focuses on prediction-making through the use of computers. It has strong ties to mathematical optimization, which delivers methods, theory and application domains to the field. Machine learning is sometimes conflated with data mining, where the latter subfield focuses more on exploratory data analysis and is known as unsupervised learning. Machine learning can also be unsupervised and be used to learn and establish baseline behavioral profiles for various entities and then used to find meaningful anomalies.

Within the field of data analytics, machine learning is a method used to devise complex models and algorithms that lend themselves to prediction; in commercial use, this is known as predictive analytics. These analytical models allow researchers, data scientists, engineers, and analysts to "produce reliable, repeatable decisions and results" and uncover "hidden insights" through learning from historical relationships and trends in the data.

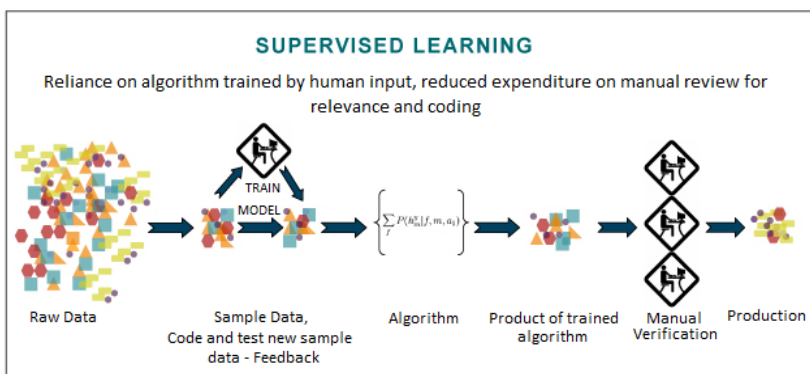
According to the Gartner hype cycle of 2016, machine learning is at its peak of inflated expectations. Effective machine learning is difficult because finding patterns is hard and often not enough training data is available; as a result, machine-learning programs often fail to deliver.

Machine Learning could be divided into supervised, unsupervised and reinforcement learning.



Supervised Learning

Supervised learning is the learning with labeled examples – training set.



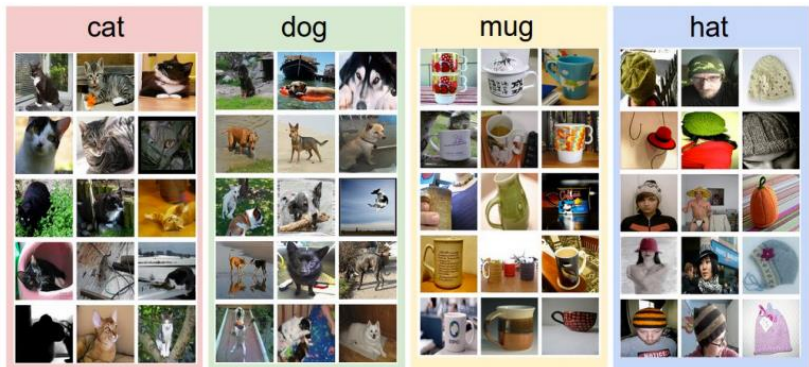
Supervised learning is the machine learning task of inferring a function from labeled training data. The training data consist of a set of training examples. In supervised learning, each example is a pair consisting of an input object

(typically a vector) and a desired output value (also called the supervisory signal). A supervised learning algorithm analyzes the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances. This requires the learning algorithm to generalize from the training data to unseen situations in a "reasonable" way.

The parallel task in human and animal psychology is often referred to as concept learning.

Supervised learning is the most common problem type in Machine Learning.

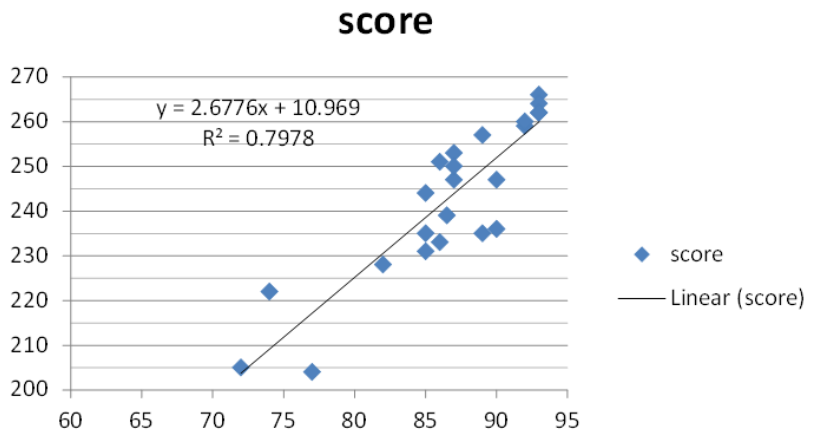
- Image labelling: learning from tagged images



- Email spam filter: learning from labeled (spam or ham) email



- Predicting exam score: learning from previous exam score and time spent



Types of Supervised Learning: Linear Regression

In statistics, linear regression is a linear approach for modeling the relationship between a scalar dependent variable y and one or more explanatory variables (or independent variables) denoted X . The case of one explanatory variable is called simple linear regression.

For more than one explanatory variable, the process is called multiple linear regression. (This term is distinct from multivariate linear regression, where multiple correlated dependent variables are predicted, rather than a single scalar variable.)

In linear regression, the relationships are modeled using linear predictor functions whose unknown model parameters are estimated from the data. Such models are called linear models. Most commonly, the conditional mean of y given the value of X is assumed to be an affine function of X ; less commonly, the median or some other quantile of the conditional distribution of y given X is expressed as a linear function of X . Like all forms of regression analysis, linear regression focuses on the conditional probability distribution of y given X , rather than on the joint probability distribution of y and X , which is the domain of multivariate analysis.

Linear regression was the first type of regression analysis to be studied rigorously, and to be used extensively in practical applications. This is

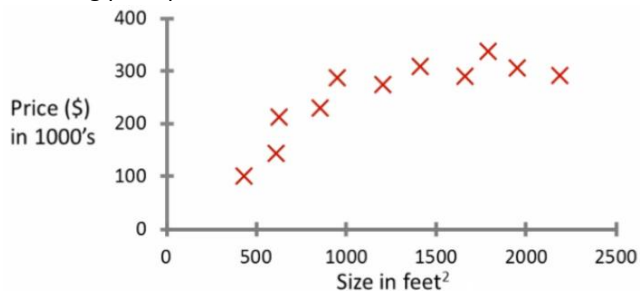
because models which depend linearly on their unknown parameters are easier to fit than models which are non-linearly related to their parameters and because the statistical properties of the resulting estimators are easier to determine.

Examples:

- Predicting final exam score based on time spent

x (hours)	y (score)
10	90
9	80
3	50
2	30

- Housing price prediction



Types of Supervised Learning: Binary classification

Binary or binomial classification is the task of classifying the elements of a given set into two groups (predicting which group each one belongs to) on the basis of a classification rule. Contexts requiring a decision as to whether or not an item has some qualitative property, some specified characteristic, or some typical binary classification include:

- Medical testing to determine if a patient has certain disease or not
 - the classification property is the presence of the disease.

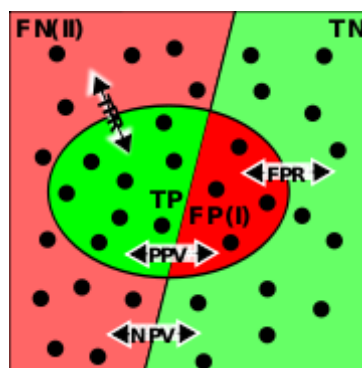
- A "pass or fail" test method or quality control in factories, i.e. deciding if a specification has or has not been met – a Go/no go classification.
- Information retrieval, namely deciding whether a page or an article should be in the result set of a search or not – the classification property is the relevance of the article, or the usefulness to the user.

Binary classification is dichotomization applied to practical purposes, and in many practical binary classification problems, the two groups are not symmetric – rather than overall accuracy, the relative proportion of different types of errors is of interest. For example, in medical testing, a false positive (detecting a disease when it is not present) is considered differently from a false negative (not detecting a disease when it is present)

Evaluation of binary classifiers

There are many metrics that can be used to measure the performance of a classifier or predictor; different fields have different preferences for specific metrics due to different goals.

For example, in medicine sensitivity and specificity are often used, while in information retrieval precision and recall are preferred. An important distinction is between metrics that are independent on the prevalence (how often each category occurs in the population), and metrics that depend on the prevalence – both types are useful, but they have very different properties.



Given a classification of a specific data set, there are four basic data: the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). These can be arranged into a 2x2 contingency table,

with columns corresponding to actual value – condition positive (CP) or condition negative (CN) – and rows corresponding to classification value – test outcome positive or test outcome negative. There are eight basic ratios that one can compute from this table, which come in four complementary pairs (each pair summing to 1). These are obtained by dividing each of the four numbers by the sum of its row or column, yielding eight numbers, which can be referred to generically in the form "true positive row ratio" or "false negative column ratio", though there are conventional terms. There are thus two pairs of column ratios and two pairs of row ratios, and one can summarize these with four numbers by choosing one ratio from each pair – the other four numbers are the complements.

The column ratios are True Positive Rate (TPR, aka Sensitivity or recall), with complement the False Negative Rate (FNR); and True Negative Rate (TNR, aka Specificity, SPC), with complement False Positive Rate (FPR). These are the proportion of the population with the condition (resp., without the condition) for which the test is correct (or, complementarily, for which the test is incorrect); these are independent of prevalence.

The row ratios are Positive Predictive Value (PPV, aka precision), with complement the False Discovery Rate (FDR); and Negative Predictive Value (NPV), with complement the False Omission Rate (FOR). These are the proportion of the population with a given test result for which the test is correct (or, complementarily, for which the test is incorrect); these depend on prevalence.

In diagnostic testing, the main ratios used are the true column ratios – True Positive Rate and True Negative Rate – where they are known as sensitivity and specificity. In informational retrieval, the main ratios are the true positive ratios (row and column) – Positive Predictive Value and True Positive Rate – where they are known as precision and recall.

One can take ratios of a complementary pair of ratios, yielding four likelihood ratios (two column ratio of ratios, two row ratio of ratios). This is primarily done for the column (condition) ratios, yielding likelihood ratios in diagnostic testing. Taking the ratio of one of these groups of ratios yields a final ratio, the diagnostic odds ratio (DOR). This can also be defined directly as $(TP \times TN) / (FP \times FN) = (TP / FN) / (FP / TN)$; this has a useful interpretation – as an odds ratio – and is prevalence-independent.

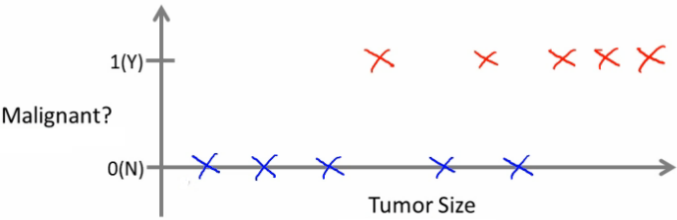
There are a number of other metrics, most simply the accuracy or Fraction Correct (FC), which measures the fraction of all instances that are correctly categorized; the complement is the Fraction Incorrect (FIC). The F-score combines precision and recall into one number via a choice of weighing, most simply equal weighing, as the balanced F-score (F1 score). Some metrics come from regression coefficients: the markedness and the informedness, and their geometric mean, the Matthews correlation coefficient. Other metrics include Youden's J statistic, the uncertainty coefficient, the Phi coefficient, and Cohen's kappa.

Examples:

- Pass/non-pass based on time spent

x (hours)	y (pass/fail)
10	P
9	P
3	F
2	F

- Decide breast cancer as malignant or benign based on tumor size



Types of Supervised Learning: Multi-label classification

In machine learning, multi-label classification and the strongly related problem of multi-output classification are variants of the classification problem where multiple labels may be assigned to each instance. Multi-label classification is a generalization of multiclass classification, which is the single-label problem of categorizing instances into precisely one of more

than two classes; in the multi-label problem there is no constraint on how many of the classes the instance can be assigned to.

Formally, multi-label classification is the problem of finding a model that maps inputs x to binary vectors y (assigning a value of 0 or 1 for each element (label) in y).

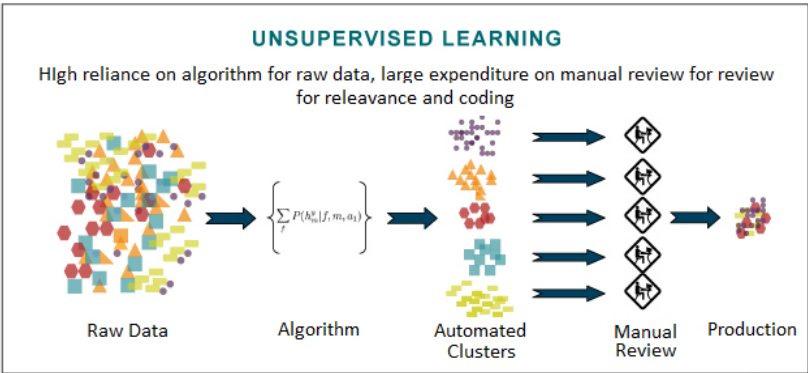
Examples:

- Letter grade (A, B, C, E and F) based on time spent

x (hours)	y (grade)
10	A
9	B
3	D
2	F

Unsupervised Learning

Unsupervised machine learning is the machine learning task of inferring a function to describe hidden structure from "unlabeled" data (a classification or categorization is not included in the observations).



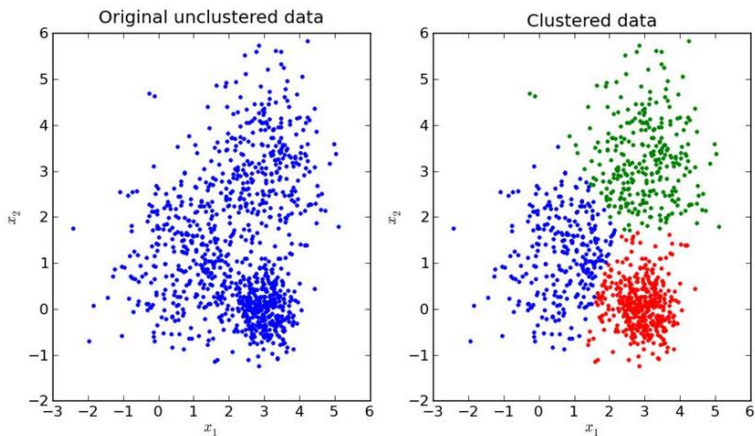
Since the examples given to the learner are unlabeled, there is no evaluation of the accuracy of the structure that is output by the relevant algorithm—

which is one way of distinguishing unsupervised learning from supervised learning and reinforcement learning.

A central case of unsupervised learning is the problem of density estimation in statistics, though unsupervised learning encompasses many other problems (and solutions) involving summarizing and explaining key features of the data.

Approaches to unsupervised learning include:

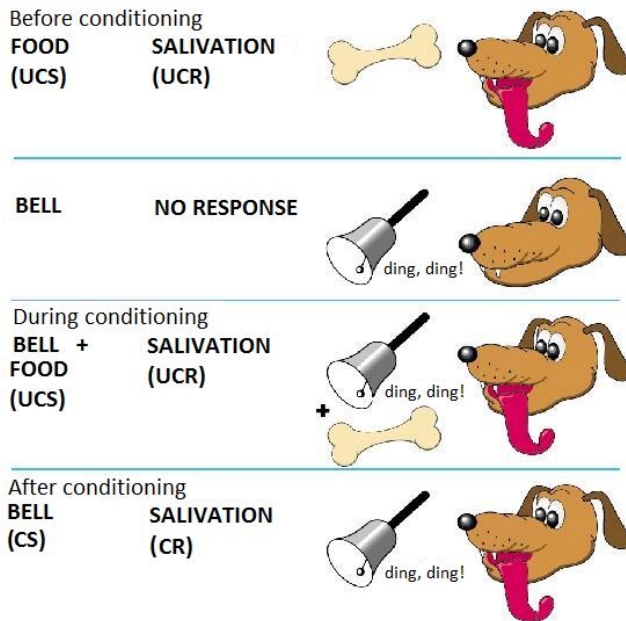
- Clustering
 - K-means
 - Mixture models
 - Hierarchical clustering



- Anomaly detection



- Neural Networks
 - o Hebbian Learning
 - o Generative Adversarial Networks

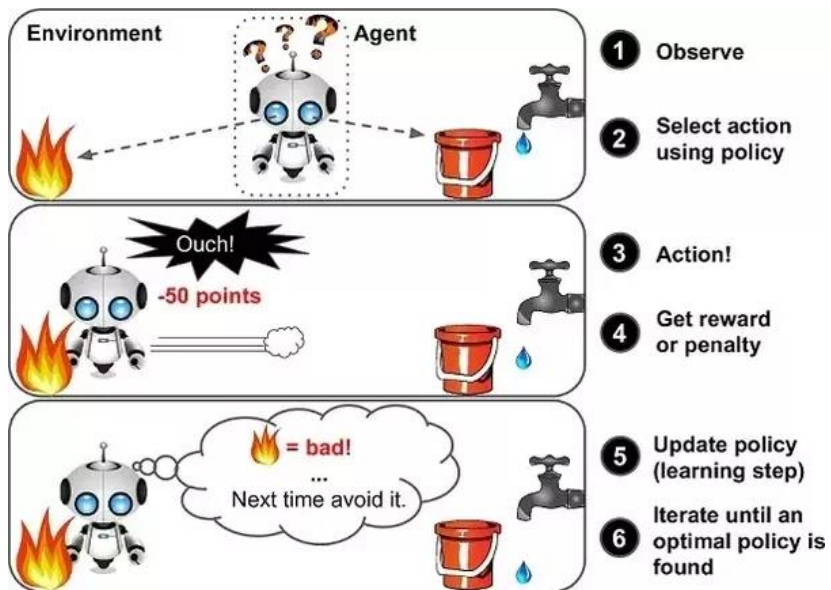


Reinforcement Learning

Reinforcement learning(RL) is an area of machine learning inspired by behaviorist psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

The problem, due to its generality, is studied in many other disciplines, such as game theory, control theory, operations research, information theory, simulation-based optimization, multi-agent systems, swarm intelligence, statistics and genetic algorithms. In the operations research and control literature, the field where reinforcement learning methods are studied is called approximate dynamic programming. The problem has been studied in the theory of optimal control, though most studies are concerned with the existence of optimal solutions and their characterization, and not with the learning or approximation aspects. In economics and game theory,

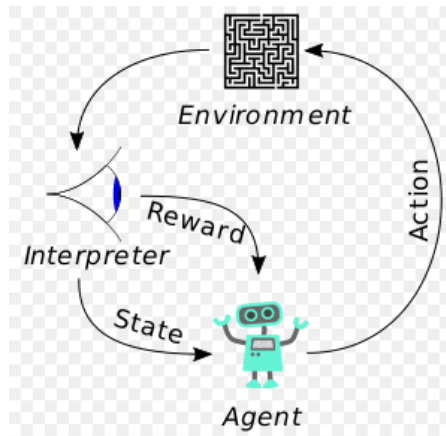
reinforcement learning may be used to explain how equilibrium may arise under bounded rationality.



In machine learning, the environment is typically formulated as a Markov decision process (MDP), as many reinforcement learning algorithms for this context utilize dynamic programming techniques. The main difference between the classical techniques and reinforcement learning algorithms is that the latter do not need knowledge about the MDP and they target large MDPs where exact methods become infeasible.

Reinforcement learning differs from standard supervised learning in that correct input/output pairs are never presented, nor sub-optimal actions explicitly corrected. Instead the focus is on on-line performance, which involves finding a balance between exploration (of uncharted territory) and exploitation (of current knowledge). The exploration vs. exploitation trade-off in reinforcement learning has been most thoroughly studied through the multi-armed bandit problem and in finite MDPs.

The typical framing of a Reinforcement Learning (RL) scenario is as follows.



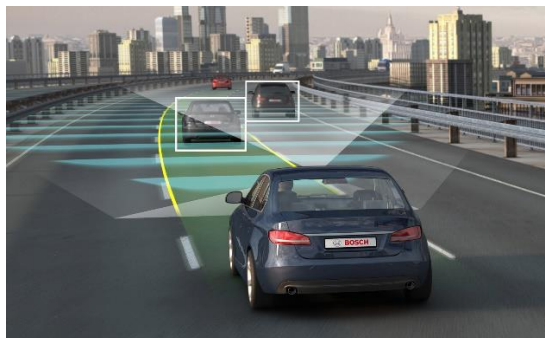
An agent takes actions in an environment, which is interpreted into a reward and a representation of the state, which are fed back into the agent.

Reinforcement learning could be used in various area including:

- Game tactics such as chess and go



- Autonomous vehicle



- Decision making



AIM

Our goals are to understand and utilize the following contents through this learning.

- Basic understanding of machine learning
 - o Supervised learning
 - o Unsupervised learning
 - o Reinforcement learning
- Solve your problems using machine learning tools
 - o Tensorflow
 - o Python

LAB EXERCISE 1

“There are no activities required for this lab”

SUMMARY

Machine learning is a field of computer science that gives computers the ability to learn without being explicitly programmed.

Supervised learning is the learning with labeled examples – training set.

Unsupervised machine learning is the machine learning task of inferring a function to describe hidden structure from "unlabeled" data.

Reinforcement learning(RL) is an area of machine learning inspired by behaviorist psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

REFERENCES

- https://en.wikipedia.org/wiki/machine_learning
- <http://www.holehouse.org/mlclass>
- <https://en.wikipedia.org/wiki/regression>

INDEX

Theory	1
Basic concepts	1
Machine Learning	1
Supervised Learning	3
Types of Supervised Learning: Linear Regression	5
Types of Supervised Learning: Binary classification	6
Types of Supervised Learning: Multi-label classification	9
Unsupervised Learning	10
Reinforcement Learning	12
AIM	16
LAB EXERCISE 1	17
SUMMARY	18
REFERENCES	19