

Machine Learning Foundations

Master the Definitions and Concepts

Jason Brownlee

**MACHINE
LEARNING
MASTERY**





Machine Learning Mastery

Web: <http://MachineLearningMastery.com>

Email: jason@MachineLearningMastery.com

Machine Learning Foundations

Master the Definitions and Concepts

by Jason Brownlee, PhD

Copyright © 2014 Jason Brownlee, All Rights Reserved.

Share this Guide

If you know someone who can benefit from this guide, just send them this link:

<http://MachineLearningMastery.com>

Table of Contents

[Introduction](#)

[This Path is Right For You!](#)

[Guide Overview](#)

[What is Machine Learning?](#)

[Definitions](#)

[Mitchell's Machine Learning](#)

[Elements of Statistical Learning](#)

[Pattern Recognition](#)

[An Algorithmic Perspective](#)

[Definition of Machine Learning](#)

[Summary](#)

[Resources](#)

[Action Steps](#)

[What Are the Key Concepts in Machine Learning?](#)

[Data](#)

[Learning](#)

[Modelling](#)

[Summary](#)

[Resources](#)

[Action Steps](#)

[What Problems Can Machine Learning Address?](#)

[Example Problems](#)

[Types of Problems](#)

[Summary](#)

[Resources](#)

[Action Steps](#)

[What Algorithms Does Machine Learning Provide?](#)

[Algorithm Similarity](#)

[Regression](#)

[Instance-based Methods](#)

[Regularization Methods](#)

[Decision Tree Learning](#)

[Bayesian](#)

[Kernel Methods](#)

[Clustering Methods](#)

[Association Rule Learning](#)

[Artificial Neural Networks](#)

[Deep Learning](#)

[Dimensionality Reduction](#)

[Ensemble Methods](#)

[Summary](#)

[Resources](#)

[Action Steps](#)

[What Other Fields Are Related to Machine Learning?](#)

[Foundations](#)

[Probability](#)

[Statistics](#)

[Artificial Intelligence](#)

[Progenitors](#)

[Computational Intelligence](#)

[Data Mining](#)

[Data Science](#)

[Summary](#)

[Resources](#)

[Action Steps](#)

[Summary](#)

[Related Guides](#)

[Beginning Weka:](#)

[Discover Applied Machine Learning](#)

[Small Projects Methodology:](#)

[Learn and Practice Applied Machine Learning](#)

Introduction

Welcome to the guide **Machine Learning Foundations**. You have taken the first step and decided that you want to learn machine learning bad enough that you are willing to invest in yourself.

I want to make sure you get the most out of that investment. I'm here to help. If you are struggling or need some advice, email me anytime at jason@MachineLearningMastery.com

This Path is Right For You!

In this guide you will learn what machine learning is, as well as the high-level concepts you need to know to get started in the field.

There are some concepts that once you are aware of will provide a solid foundation for your journey through the field of machine learning.

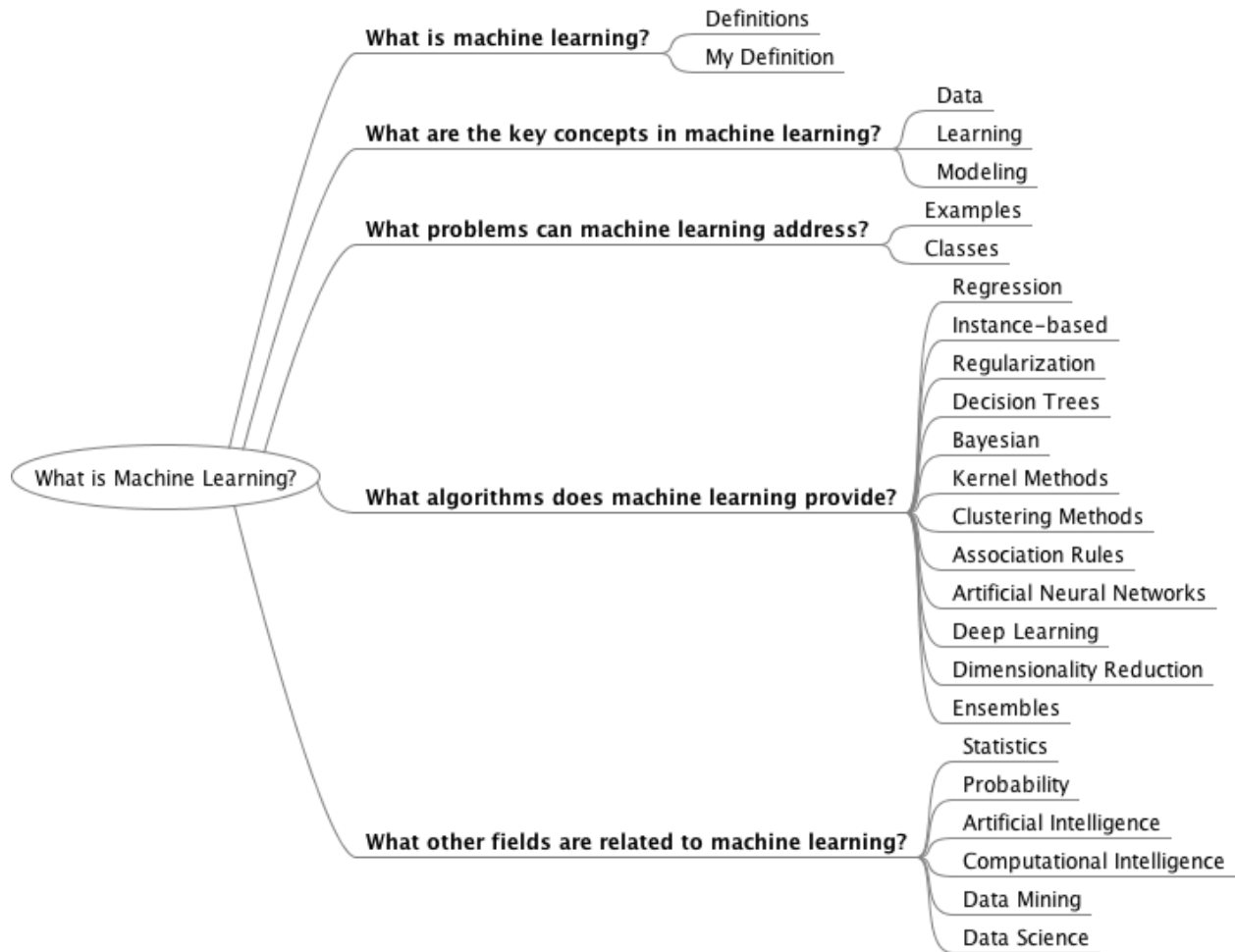
This guide is intended to point out those concepts to you including key definitions, types of problems and classes of algorithm.

Guide Overview

This self-study guide is broken down into 5 parts:

- What is machine learning
- What are the key concepts in machine learning
- What problems can machine learning address
- What methods does machine learning provide
- What other fields are related to machine learning

After completing this guide you will know what machine learning is and be able to explain it to friends and colleagues.



What is Machine Learning?

Given that you are interested in machine learning it is useful to you to have a working definition for machine learning. In this section you will explore authoritative definitions of machine learning from popular books in the field, and a handy one-sentence definition that I like to use.

Your goal with this section is to be able to describe machine learning in one or two sentences to anyone that asks.

Definitions

Let's start out by looking at four textbooks on Machine Learning that are commonly used in university level courses. These are our authoritative definitions and lay our foundation for deeper thought on the subject. I chose these four definitions to highlight some useful and varied perspectives on the field.

Mitchell's Machine Learning

Tom Mitchell in his book Machine Learning provides a definition in the opening line of the preface:

"The field of machine learning is concerned with the question of how to construct computer programs that automatically improve with experience."

I like this short and sweet definition and it is the basis for the programmers definition we come up with at the end of the section. Note the mention of "computer programs" and the reference to "automated improvement". Write programs that improve themselves, it's provocative!

In his introduction Mitchell provides a short formalism that you'll see much repeated in the field:

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E ."

Don't let the definition of terms scare you off, this is a very useful formalism. We can use this formalism as a template and put E , T , and P at the top of columns in a table and list out complex problems with less ambiguity. It could be used as a design tool to help us think clearly about what data to collect (E), what decisions the software needs to make (T) and how we will evaluate its results (P). This power is why it is often repeated as a standard definition.

Elements of Statistical Learning

The Elements of Statistical Learning: Data Mining, Inference, and Prediction was written by three Stanford statisticians and self-described as a statistical framework to organize their field of inquiry. In the preface is written:

“Vast amounts of data are being generated in many fields, and the statistician's job is to make sense of it all: to extract important patterns and trends, and to understand “what the data says”. We call this learning from data.”

I understand the job of a statistician is to use the tools of statistics to interpret data in the context of the domain. The authors seem to include all of the field of Machine Learning as aids in that pursuit. Interestingly, they chose to include "Data Mining" in the subtitle of the book.

Statisticians learn from data, but software does too and we learn from the things that the software learns.

Pattern Recognition

Bishop in the preface of his book Pattern Recognition and Machine Learning comments:

“Pattern recognition has its origins in engineering, whereas machine learning grew out of computer science. However, these activities can be viewed as two facets of the same field...”

Reading this, you get the impression that Bishop came at the field from an engineering perspective and later learned and leveraged the Computer Science take on the same methods. This is a mature approach and one we should emulate. More broadly, regardless of the field that lays claim to a method, if it suits our needs by getting us closer to an insight or a result by "learning from data", then we can decide to call it machine learning.

An Algorithmic Perspective

Marsland adopts the Mitchell definition of Machine Learning in his book Machine Learning: An Algorithmic Perspective. I want to share a cogent note in the prologue of his book:

“One of the most interesting features of machine learning is that it lies on the boundary of several different academic disciplines, principally computer science, statistics, mathematics, and engineering. ...machine learning is usually studied as part of artificial intelligence, which puts it firmly into computer science ...understanding why these algorithms work requires a certain amount of statistical and mathematical sophistication that is often missing from computer science undergraduates.”

This is insightful and instructive. Firstly, he underscores the multidisciplinary nature of the field. We were getting a feeling for that from the above definition, but he draws a big red underline for us. Machine Learning draws from all manner of information sciences. Secondly, he underscores the danger of sticking to a given perspective too tightly. Specifically, the case of a the algorithmist who shies away from the mathematical inner workings of a method. No doubt, the counter case of the statistician that shies away from the practical concerns of implementation and deployment is just as limiting.

Definition of Machine Learning

So, let's see if we can use these pieces and construct a programmers definition of machine learning. A one-sentence definition of Machine Learning that I use is:

"Machine Learning is the training of a model from data that generalizes a decision against a performance measure."

Training a model suggests training examples. A model suggests state acquired through experience. Generalizes a decision suggests the capability to make a decision based on inputs and anticipating unseen inputs in the future for which a decision will be required. Finally, against a performance measure suggests a targeted need and directed quality to the model being prepared.

Summary

In this section you learned what machine learning is by reviewing authoritative definitions from popular textbooks in the field. You also learned how I think about machine learning and the definition I use when people ask me: what is machine learning?

Resources

Below are some resources that you might find useful for further reading.

- Tom Mitchell, [*Machine Learning*](#), 1997
- Trevor Hastie, Robert Tibshirani, Jerome Friedman, [*The Elements of Statistical Learning: Data Mining, Inference, and Prediction*](#), 2011
- Christopher Bishop, [*Pattern Recognition and Machine Learning*](#), 2007
- Stephen Marsland, [*Machine Learning: An Algorithmic Perspective*](#), 2009

Action Steps

Choose a definition (or come up with your own definition) of machine learning, one that you will feel confident to use when someone asks you what is machine learning.

Practice it and go and try it out on colleagues and friends on your blog, Facebook or Google+.

Share it with me.

Send an email to jason@MachineLearningMastery.com with the subject "My Definition".

What Are the Key Concepts in Machine Learning?

There are key concepts in machine learning that lay the foundation for understanding the field.

In this section you will learn the nomenclature (standard terms) that is used when describing data and datasets. You will also learn the concepts and terms used to describe learning and modelling from data that will provide a valuable intuition for your journey through the field of machine learning.

Data

Machine learning methods learn from examples. It is important to have good grasp of input data and the various terminology used when describing data. In this section you will learn the terminology used in machine learning when referring to data.

When I think of data, I think of rows and columns, like a database table or an Excel spreadsheet. This is a traditional structure for data and is what is common in the field of machine learning. Other data like images, videos, and text, so-called unstructured data is not considered at this time.

Instance: A single row of data is called an instance. It is an observation from the domain.

Feature: A single column of data is called a feature. It is a component of an observation and is also called an attribute of a data instance. Some features may be inputs to a model (the predictors) and others may be outputs or the features to be predicted.

Data Type: Features have a data type. They may be real or integer valued or may have a categorical or ordinal value. You can have strings, dates, times, and more complex types, but typically they are reduced to real or categorical values when working with traditional machine learning methods.

Datasets: A collection of instances is a dataset and when working with machine learning methods we typically need a few datasets for different purposes.

Training Dataset: A dataset that we feed into our machine learning algorithm to train our model.

Testing Dataset: A dataset that we use to validate the accuracy of our model but is not used to train the model. It may be called the validation dataset.

We may have to collect instances to form our datasets or we maybe given a finite dataset that we must split into sub-datasets.

Learning

Machine learning is indeed about automated learning with algorithms. In this section we will consider a few high-level concepts about learning.

Induction: Machine learning algorithms learn through a process called induction or inductive learning. Induction is a reasoning process that makes generalizations (a model) from specific information (training data).

Generalization: Generalization is required because the model that is prepared by a machine learning algorithm needs to make predictions or decisions based on specific data instances that were not seen during training.

Over-Learning: When a model learns the training data too closely and does not generalize, this is called over-learning. The result is poor performance on data other than the training dataset. This is also called over-fitting.

Under-Learning: When a model has not learned enough structure from the database because the learning process was terminated early, this is called under-learning. The result is good generalization but poor performance on all data, including the training dataset. This is also called under-fitting.

Online Learning: Online learning is when a method is updated with data instances from the domain as they become available. Online learning requires methods that are robust to noisy data but can produce models that are more intune with the current state of the domain.

Offline Learning: Offline learning is when a method is created on pre-prepared data and is then used operationally on unobserved data. The training process can be controlled and can tuned carefully because the scope of the training data is known. The model is not updated after it has been prepared and performance may decrease if the domain changes.

Supervised Learning: This is a learning process for generalizing on problems where a prediction is required. A "teaching process" compares predictions by the model to known answers and makes corrections in the model.

Unsupervised Learning: This is a learning process for generalizing the structure in the data where no prediction is required. Natural structures are identified and exploited for relating instances to each other.

Modelling

The artefact created by a machine learning process could be considered a program in it's own right. I like to call it a model, as in, it is a model of the problem we are trying to solve.

Model Selection: We can think of the process of configuring and training the model as a model selection process. Each iteration we have a new model that we could choose to use or to modify. Even the choice of machine learning algorithm is part of that model selection process. Of all the possible models that exist for a problem, a given algorithm and algorithm configuration on the chosen training dataset will provide a finally selected model.

Inductive Bias: Bias is the limits imposed on the selected model. All models are biased which introduces error in the model, and by definition all models have error (they are generalizations from observations). Biases are introduced by the generalizations made in the model including the configuration of the model and the selection of the algorithm to generate the model. A machine learning method can create a model with a low or a high bias and tactics can be used to reduce the bias of a highly biased model.

Model Variance: Variance is how sensitive the model is to the data on which it was trained. A machine learning method can have a high or a low variance when creating a model on a dataset. A tactic to reduce the variance of a model is to run it multiple times on a dataset with different initial conditions and take the average accuracy as the models performance.

Bias-Variance Tradeoff: Model selection can be thought of as a the trade-off of the bias and variance. A low bias model will have a high variance and will need to be trained for a long time or many times to get a usable model. A high bias model will have a low variance and will train quickly, but suffer poor and limited performance.

Summary

This section provided a useful glossary of terms that you can refer back to anytime for a clear definition. It can also provide the basis for you to build up your own glossary of terms and expand the definitions as you learn more.

Resources

Below are some resources that you might find useful for further reading.

- Tom Mitchell, [The need for biases in learning generalizations](#), 1980
- [Understanding the Bias-Variance Tradeoff](#)

Action Steps

Select one term from the glossary that interests you the most and research it further. Write a few paragraphs to a page about the term including authoritative definitions and resources to learn more.

Share your definition. Post it to your blog, Facebook or Google+.

Share it with me.

Send an email to jason@MachineLearningMastery.com with the subject "My Defined Term".

What Problems Can Machine Learning Address?

As part of understanding what machine learning is, it is useful to understand the types of problems that machine learning can solve. This is valuable, because knowing the type of problem you are facing allows us to think about the data we need and the types of algorithms to try.

In this section we you will first look at some well known and understood examples of machine learning problems in the real world. You will then look at a taxonomy (naming system) for standard machine learning problems and learn how to identify a problem as one of these standard cases.

Example Problems

Machine Learning problems are abound. They make up core or difficult parts of the software you use on the web or on your desktop everyday. Think of the "do you want to follow" suggestions on twitter and the speech understanding in Apple's Siri.

Below are 10 examples of machine learning that really ground what machine learning is all about.

- **Spam Detection:** Given email in an inbox, identify those email messages that are spam and those that are not. Having a model of this problem would allow a program to leave non-spam emails in the inbox and move spam emails to a spam folder. We should all be familiar with this example.
- **Credit Card Fraud Detection:** Given credit card transactions for a customer in a month, identify those transactions that were made by the customer and those that were not. A program with a model of this decision could refund those transactions that were fraudulent.
- **Digit Recognition:** Given zip codes hand written on envelopes, identify the digit for each handwritten character. A model of this problem would allow a computer program to read and understand handwritten zip codes and sort envelopes by geographic region.
- **Speech Understanding:** Given an utterance from a user, identify the specific request made by the user. A model of this problem would allow a program to understand and make an attempt to fulfil that request. The iPhone with Siri has this capability.
- **Face Detection:** Given a digital photo album of many hundreds of digital photographs, identify those photos that include a given person. A model of this decision process would allow a program to organize photos by person. Some digital cameras and software like iPhoto has this capability.
- **Product Recommendation:** Given the purchase history for a customer and a large inventory of products, identify those products in which that customer will be interested and likely to purchase. A model of this decision process would allow a program to make recommendations to a customer and motivate product purchases. Amazon has this

capability. Also think of Facebook, Google+ and Facebook that recommend users for you to connect with.

- **Medical Diagnosis:** Given the symptoms exhibited in a patient and a database of anonymized patient records, predict whether the patient is likely to have an illness. A model of this decision problem could be used by a program to provide decision support to medical professionals.
- **Stock Trading:** Given the current and past price movements for a stock, determine whether the stock should be bought, held or sold. A model of this decision problem could provide decision support to financial analysts.
- **Customer Segmentation:** Given the pattern of behaviour by a user during a trial period and the past behaviours of all users, identify those users that will convert to the paid version of the product and those that will not. A model of this decision problem would allow a program to trigger customer interventions to persuade the customer to convert early or better engage in a limited trial.
- **Shape Detection:** Given a user hand drawing a shape on a touch screen and a database of known shapes, determine which shape the user was trying to draw. A model of this decision would allow a program to show the platonic version of that shape the user drew to make crisp diagrams. The Instaviz iPhone app does this.

These 10 examples give you a good sense of what a machine learning problem looks like. There is a corpus of historic examples, there is a decision that needs to be modelled and a business or domain benefit to having that decision modelled and efficaciously made automatically.

Some of these problems are some of the hardest problems in Artificial Intelligence, such as Natural Language Processing and Machine Vision (doing things that humans do easily). Others are still difficult, but are classic examples of machine learning such as spam detection and credit card fraud detection.

Types of Problems

Reading through the list of example machine learning problems above, I'm sure you can start to see similarities. This is a valuable skill, because being good at extracting the essence of a problem will allow you to think effectively about what data you need and what types of algorithms you should try.

There are common classes of problem in Machine Learning. The problem classes below are archetypes for most of the problems we refer to when we are doing Machine Learning.

- **Classification:** Data is labelled meaning it is assigned a class, for example spam/non-spam or fraud/non-fraud. The decision being modelled is to assign labels to new unlabelled pieces of data. This can be thought of as a discrimination problem, modelling the differences or similarities between groups.

- **Regression:** Data is labelled with a real value (think floating point) rather than a label. Examples that are easy to understand are time series data like the price of a stock over time. The decision being modelled is the relationship s between inputs and outputs.
- **Clustering:** Data is not labelled, but can be divided into groups based on similarity and other measures of natural structure in the data. An example from the above list would be organising pictures by faces without names, where the human user has to assign names to groups, like iPhoto on the Mac.
- **Rule Extraction:** Data is used as the basis for the extraction of propositional rules (antecedent/consequent or if-then). Such rules may, but are typically not directed, meaning that the methods discover statistically supportable relationships between attributes in the data, not necessarily involving something that is being predicted. An example is the discovery of the relationship between the purchase of beer and diapers (this is data mining folk-law, true or not, it's illustrative of the desire and opportunity).

When you think a problem is a machine learning problem (a decision problem that needs to be modelled from data), think next of what type of problem you could phrase it as easily or what type of outcome the client or requirement is asking for and work backwards.

Summary

In this section you explored examples of common machine learning problems. They were not problems specific to machine learning, but they are problems that machine learning methods can address. You also learned about four main classes of machine learning problems, specifically classification, regression, clustering and rule extraction.

Resources

Below are some resources that you might find useful for further reading.

- [What are the Top 10 problems in Machine Learning for 2013?](#)
- [AI-Complete](#)

Action Steps

The action steps for this section are to reinforce the definition of machine learning in conjunction with the sample problems.

Describe 5 additional complex problems that you think machine learning methods would suit and define the E , T , and P terms from Tom Mitchell's definition. Think about some of your interactions with online and offline software in the last week. I'm sure you could easily guess at another ten or twenty examples of machine learning you have directly or indirectly used.

Share your problem definitions. Post it to your blog, Facebook or Google+.

Share it with me.

Send an email to jason@MachineLearningMastery.com with the subject "My Complex

Problems".

What Algorithms Does Machine Learning Provide?

The study of machine learning is the study of machine learning algorithms. There are so many algorithms available. The difficulty is that there are classes of method and there are extensions to methods and it quickly becomes very difficult to determine what constitutes a canonical algorithm.

In this section you will take a tour of the popular machine learning algorithms and learn how to categorize algorithms by their mechanism and behaviour.

Algorithm Similarity

There are lots of ways to categorize machine learning algorithms. Algorithm similarity is the easiest way to get started.

Algorithms are universally presented in groups by similarity in terms of function or form. For example, tree based methods, and neural network inspired methods. This is a useful grouping method, but it is not perfect.

There are still algorithms that could just as easily fit into multiple categories such as Learning Vector Quantization that is both a neural network inspired method and an instance-based method. There are also categories that have the same name that describes the problem and the class of algorithm such as Regression and Clustering. As such, you will see variations on the way algorithms are grouped depending on the source you check. Like machine learning algorithms themselves, there is no perfect model, just a 'good enough' model.

In the rest of this section we will go through a list many of the popular machine learning algorithms grouped the way I think is the most intuitive. It is not exhaustive in either the groups or the algorithms, but I think it is representative and will be useful for you to get an idea of the lay of the land.

Regression

Regression is concerned with modelling the relationship between variables that is iteratively refined using a measure of error in the predictions made by the model. Regression methods are a work horse of statistics and have been co-opted into statistical machine learning. This may be confusing because we can use regression to refer to the class of problem and the class of algorithm. Really, regression is a process.

Some example algorithms are:

- Ordinary Least Squares Linear Regression
- Logistic Regression
- Stepwise Regression
- Multivariate Adaptive Regression Splines (MARS)

- Locally Estimated Scatterplot Smoothing (LOESS)

Instance-based Methods

Instance based learning model a decision problem with instances or examples of training data that are deemed important or required to the model. Such methods typically build up a database of example data and compare new data to the database using a similarity measure in order to find the best match and make a prediction. For this reason, instance-based methods are also called winner-take all methods as well as case-based and memory-based learning. Focus is put on representation of the stored instances and similarity measures used between instances.

- k-Nearest Neighbour (kNN)
- Learning Vector Quantization (LVQ)
- Self-Organizing Map (SOM)

Regularization Methods

An extension made to another method (typically regression methods) that penalizes models based on their complexity, favoring simpler models that are also better at generalizing. I have listed Regularization methods here because they are popular, powerful and generally simple modifications made to other methods.

- Ridge Regression
- Least Absolute Shrinkage and Selection Operator (LASSO)
- Elastic Net

Decision Tree Learning

Decision tree methods construct a model of decisions made based on actual values of attributes in the data. Decisions fork in tree structures until a prediction decision is made for a given record. Decision trees are trained on data for classification and regression problems.

- Classification and Regression Tree (CART)
- Iterative Dichotomiser 3 (ID3)
- C4.5
- Chi-squared Automatic Interaction Detection (CHAID)
- Decision Stump
- Random Forest
- Multivariate Adaptive Regression Splines (MARS)
- Gradient Boosting Machines (GBM)

Bayesian

Bayesian methods are those that explicitly apply Bayes' Theorem for problems such as classification and regression.

- Naive Bayes
- Averaged One-Dependence Estimators (AODE)
- Bayesian Belief Network (BBN)

Kernel Methods

Kernel Methods are best known for the popular method Support Vector Machines which is really a constellation of methods in and of itself. Kernel Methods are concerned with mapping input data into a higher dimensional vector space where some classification or regression problems are easier to model.

- Support Vector Machines (SVM)
- Radial Basis Function (RBF)
- Linear Discriminate Analysis (LDA)

Clustering Methods

Clustering, like regression describes the class of problem and the class of methods. Clustering methods are typically organized by the modelling approaches such as centroid-based and hierarchical. All methods are concerned with using the inherent structures in the data to best organize the data into groups of maximum similarity.

- k-Means
- Expectation Maximisation (EM)

Association Rule Learning

Association rule learning are methods that extract rules that best explain observed relationships between variables in data. These rules can discover important and commercially useful associations in large multidimensional datasets that can be exploited by an organisation.

- Apriori algorithm
- Eclat algorithm

Artificial Neural Networks

Artificial Neural Networks are models that are inspired by the structure of biological neural networks. They are a class of pattern matching method that are commonly used for regression

and classification problems but are really an enormous subfield comprised of hundreds of algorithms and variations for all manner of problem types. Some of the classically popular methods include (I have separated Deep Learning from this category):

- Perceptron
- Back-Propagation
- Hopfield Network
- Self-Organizing Map (SOM)
- Learning Vector Quantization (LVQ)

Deep Learning

Deep Learning methods are a modern update to Artificial Neural Networks that exploit abundant cheap computation. They are concerned with building much larger and more complex neural networks, and many methods are concerned with semi-supervised learning problems where large datasets contain very little labelled data.

- Restricted Boltzmann Machine (RBM)
- Deep Belief Networks (DBN)
- Convolutional Network
- Stacked Auto-encoders

Dimensionality Reduction

Like clustering methods, Dimensionality Reduction methods seek and exploit the inherent structure in the data, but in this case in an unsupervised manner or order to summarise or describe data using less information. This can be useful to visualize high-dimensional data or to simplify data which can then be used in a supervised learning method.

- Principal Component Analysis (PCA)
- Partial Least Squares Regression (PLS)
- Sammon Mapping
- Multidimensional Scaling (MDS)
- Projection Pursuit

Ensemble Methods

Ensemble methods are models composed of multiple weaker models that are independently trained and whose predictions are combined in some way to make the overall prediction. Much effort is put into what types of weak learners to combine and the ways in which to combine them. This is a very powerful class of techniques and as such is very popular.

- Bagging
- Bootstrapped Aggregation (Boosting)
- AdaBoost
- Stacked Generalization (blending)
- Gradient Boosting Machines (GBM)
- Random Forest

Summary

In this section you took a tour of popular machine learning algorithms grouped by algorithm similarity (function and form). You learned that there are a lot of algorithms that can be used under the guise of machine learning, including standard methods from the field of statistics such as linear regression.

The list of algorithms and classes of algorithms listed in this section can be used as a reference when you encounter an algorithm and want an idea of what other algorithms it is like. It can also provide the basis of your own algorithm glossary that you can continue to add to as you learn about more algorithm and start to describe what they do.

Resources

Below are some resources that you might find useful for further reading.

- [Top 10 Algorithms in Data Mining](#)
- [List of Machine Learning algorithms on Wikipedia](#)

Action Steps

Pick one algorithm that you are curious about and to research it. Answer the following 5 questions about your chosen algorithm:

- What is the algorithm name and what other names and abbreviations is the algorithm known by?
- What is the general information processing strategy of the algorithm (i.e. what are the inputs, outputs and generally what is it trying to achieve)?
 - If possible draw a really high-level flow diagram of the algorithm
- What are 5-10 heuristics or rules of thumb to keep in mind when using the algorithm
- What are one or two libraries or tools that provide an implementation of the algorithm you could use (directly or as an inspiration) if you wanted to integrate it into an application?
- What are 5-10 resources you can use when you need to know more specific details on how the algorithm works?

Share your algorithm description. Post it to your blog, Facebook or Google+.

Share it with me.

Send an email to jason@MachineLearningMastery.com with the subject "My Algorithm".

What Other Fields Are Related to Machine Learning?

Machine Learning is a multidisciplinary field and it can be very confusing when you are getting started to differentiate machine learning from the closely related fields of Artificial Intelligence and Data Mining.

In this section you will learn about those fields that are related to machine learning. Specifically, you will learn about the boundaries of the field by learning how machine learning builds on fields of mathematics and artificial intelligence and is used within fields such as data mining and data science.

Foundations

Machine Learning is built on the field of Mathematics and Computer Science. Specifically, machine learning methods are best described using linear and matrix algebra and their behaviours are best understood using the tools of probability and statistics.

In this section you will consider the fields of Statistics, Probability and Artificial Intelligence that represent the foundational subjects for machine learning.

Probability

The field of probability theory is the study of characterising the likelihood of random events. Probability theory is a branch of mathematics and provides the basis for the field of statistics.

Machine learning methods are often described in the language of probability and there are methods that directly employ probability theories such as Bayes' Theorem.

Statistics

The field of statistics is the study of methods to collect, analyze, describe and present data. Statistics is a branch of mathematics. The field is concerned with questions like what does the data mean.

Machine learning can be well understood in a statistical framework where learning from training data is taken as a modelling of the structures and relationships in the data. As such, statistical modelling methods are adopted in machine learning but machine learning includes more than statistical modelling methods.

Artificial Intelligence

The field of artificial intelligence is the study and construction of computational systems that do things that humans can do or that do things that we think are intelligent. For example humans can move around an environment, understand what they see and understand language they read and hear, and we have corresponding subfields of robotics, computer vision and natural language processing. A grand master chess champion is considered intelligent, and so chess playing intelligent systems are created. Artificial Intelligence is a branch of computer science. The field is concerned with questions of what is intelligence and how to create intelligences.

Learning is a feature of an intelligent system. As such, Machine Learning is considered a branch of artificial intelligence concerned with the study and construction of systems that are capable of learning.

Progenitors

Algorithms that can learn from data to describe the data and predict outcomes for unseen data are useful for addressing complex problems. As such, machine learning methods are used in applied computer science fields such as Data Mining and Data Science. Additionally, there are related fields of Artificial Intelligence that study intelligent methods that also learn from data and their environment. Examples include Computational Intelligence and Metaheuristics.

In this section you will review the related fields of Computational Intelligence, Data Mining and Data Science and learn how machine learning methods applied.

Computational Intelligence

The field of computational intelligence is concerned with the study and construction of systems that are easy to specify but result in complex emergent behaviours. Many computational intelligence systems are inspired by natural systems such as evolution, the immune system and the nervous system for subfields such as evolutionary computing, artificial immune systems and artificial neural networks. Computational Intelligence is a branch of artificial intelligence. The field is concerned with questions of explaining how complex emergent behaviours are derived from simple rules and what problems they are best suited to address.

Many computational intelligence systems learn from interactions with their environment and as such have been adopted as machine learning methods.

Data Mining

The field of data mining is the study and construction of systems that discover interesting relationships from large data sets. As such data mining spans both the storage and maintenance of data and the process of making discoveries in the data. Data mining is a process and is also known as knowledge discovery in databases (KDD). Data Mining is a subfield of computer science. The field is concerned with questions of what relationships are interesting and how to best discover them.

Machine learning provides a set of tools used in the data mining process for learning relationships in data that provide the basis of discovery.

Data Science

The field of Data Science is concerned with the practicality of solving complex problems using data. Data science is a subfield of computer science. Data science is the application of the data mining process and the use of machine learning methods in a specific domain. A data scientist is a practitioner of data science.

Like data mining, machine learning provides a set of tools used in data science for learning relationships in data in order to characterise data or make predictions.

Summary

In this section you understood machine learning by defining the foundations upon which it was built. You learned high-level definitions for the fields related to machine learning and learned how machine learning is used in those fields.

Resources

Below are some resources that you might find useful for further reading.

- Leo Breiman, [*Statistical Modeling: The Two Cultures*](#), 2001
- Stuart Russell and Peter Norvig, [*Artificial Intelligence: A Modern Approach*](#), 2009
- Andries Engelbrecht, [*Computational Intelligence: An Introduction*](#), 2007

Action Steps

The action step for this section is to pick one of the foundation or progenitor fields of machine learning described in this section and to collect 3-5 authoritative definitions for that field.

You can collect these definitions from books. There are many freely available books and websites on each of these topics that you can use. There are also many excellent books on Amazon and Google Books on these subjects that you can use. These services allow you to search and browse within the books before purchasing and can be used to collect the quotes you need to complete the action steps for this section.

Collect the quotes and their source and explain what you think they mean.

Share your definitions. Post it to your blog, Facebook or Google+.

Share it with me.

Send an email to jason@MachineLearningMastery.com with the subject "My Domain".

Summary

You now have a grasp of what machine learning is and the confidence to explain it to friends and colleagues.

You have also learned the key concepts of machine learning including data, learning and modelling terms. Finally, you learned how to describe a complex problem as a machine learning problem and to research an algorithm in which you are interested.

Through the action steps you learned the following tactics that you can take forward in your journey through machine learning:

- Collect authoritative definitions for terms from multiple sources and craft or select a definition you feel comfortable and confident to use.
- Keep a glossary of terms used in the field and expand the definition of a terms you find interesting or ambiguous.
- Define problems in terms of Tom Mitchell's formalism to quickly have a grasp of how they may be considered in a machine learning context.
- Breakdown a description of a machine learning algorithm in way that is useful and build up a library of machine learning algorithm descriptions.
- Perspective on the field of machine learning is key both for the foundational fields that describe the techniques and the progenitor fields that apply the techniques.

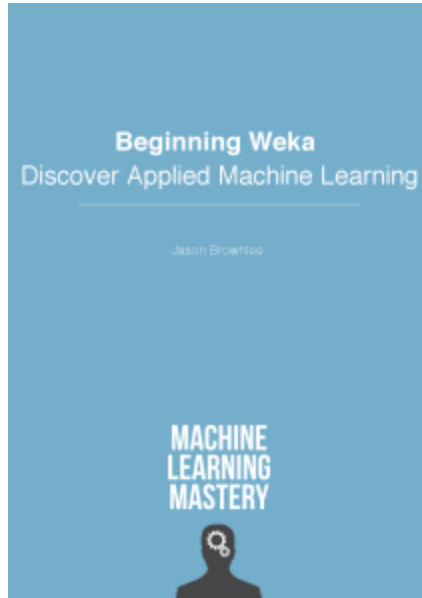
Send Feedback

I'm very interest to hear what you thought of this guide and whether you have any further questions.

Please send me an email with your comments to jason@MachineLearningMastery.com
I look forward to hearing from you.

Related Guides

This section lists related guide on machine learning that you might find interesting.

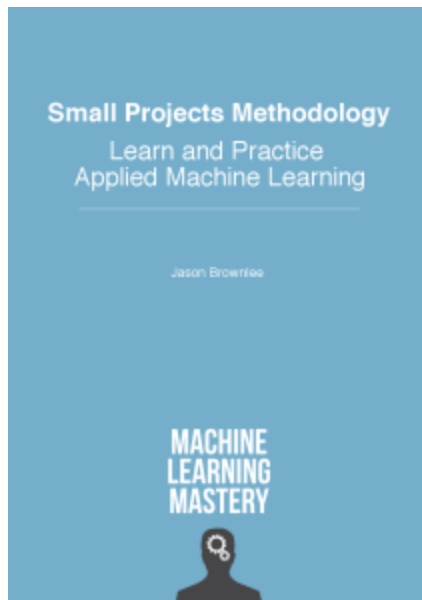


Beginning Weka: Discover Applied Machine Learning

If you are interested in learning the process of applied machine learning with tutorials and practical case studies, you might be interested in Beginning Weka.

It includes a 6-part guide on the process of applied machine learning, a 5-part getting started guide for Weka and 3 case studies (predicting the onset of diabetes, free electron structure in the ionosphere and the recurrence of breast cancer).

You can learn more at:
<http://BeginningWeka.com>



Small Projects Methodology: Learn and Practice Applied Machine Learning

If you are interested in some ideas for self-study projects, you might be interested in the Small Projects Methodology.

It provides 4 classes of project ideas with tactics and strategies you can use, as well as 90 ideas of projects that you could study

You can learn more at
<http://SmallProjectsMethodology.com>