

Companion to Machine Learning

Rohan Kumar

Contents

0.1	What is Machine Learning	3
0.2	Applications of Machine Learning	3
0.3	Types of Machine Learning	3
1	Data Analysis	7
2	Linear Models	7
2.1	Linear Regression	7
2.1.1	Formulation	7
2.1.2	Simple Regression	7
2.1.3	Multivariable Regression	7
2.1.4	Cost Function	8
2.1.5	Gradient Descent Solution	8
2.1.6	Normal Equation Solution	8
2.2	Generalized Linear Models	8
2.3	Regularization	9

Introduction

0.1 What is Machine Learning

Machine Learning is the field of study that gives computers the ability to learn from data without being explicitly programmed. This is good for problems that require a lot of fine-tuning or for which using a traditional approach yields no good solution. Machine Learning's data dependency allows it to adapt to new data and gain insight for complex problems and large amounts of data.

0.2 Applications of Machine Learning

Machine Learning can be used for a range of tasks and can be seen used in:

- Analyzing images of products on a production line to automatically classify them (Convolutional Neural Net)
- Forecasting company revenue based on performance metrics (Regression or Neural Net)
- Automatically classifying news articles (NLP using Recurrent Neural Networks)
- Summarizing long documents automatically (Natural Language Processing)
- Building intelligent bot for a game (Reinforcement Learning)

0.3 Types of Machine Learning

Supervised Learning

Hands-On Machine Learning

In supervised learning, the training set you feed to the algorithm includes the desired solutions, called labels. (e.g determining if an email is spam would be trained a dataset of example emails labelled as spam or not spam.)

Some commonly used supervised learning algorithms are:

- k-Nearest Neighbors
- Linear Regression
- Logistic Regression
- Support Vector Machines (SVMs)
- Decision Trees and Random Forests
- Neural Networks

Unsupervised Learning

Hands-On Machine Learning

In unsupervised learning, the training data is unlabeled and the system tries to learn without guidance. The system will try and automatically draw inferences and conclusions about the data and group it as such. (e.g. having a lot of data about blog visitors. Using a clustering algorithm we can group and detect similar visitors).

Some important unsupervised learning algorithms are:

- Clustering
 - K-Means
 - DBSCAN
 - Hierarchical Cluster Analysis
- Anomaly detection and novelty detection
 - One-class SVM
 - Isolation Forest
- Visualization and dimensionality reduction
 - Principal Component Analysis (PCA)
 - Kernel PCA
 - Locally Linear Embedding (LLE)
 - t-Distributed Stochastic Neighbor Embedding (t-SNE)
- Association rule learning
 - Apriori
 - Eclat

Semisupervised Learning

Hands-On Machine Learning

Labelling can be very time-consuming and costly, often there will be plenty of unlabelled and a few labelled instances. Algorithms that deal with data that is partially labeled is called semi-supervised learning. A good example of this is Google Photos. Google clusters and groups your photos based on facial recognition (unsupervised) and then you can label one photo and it will be able to label every picture like that (supervised). Most semi-supervised learning algorithms are combinations of unsupervised and supervised algorithms.

Reinforcement Learning

Hands-On Machine Learning

Reinforcement Learning is a learning algorithm based on a reward system. The learning system, called an agent, can observe the environment, select and perform actions, and get rewards in return (or penalties in the form of negative rewards). It will then learn by itself

what the best strategy, called a policy, to get the most reward over time. A policy defines what action the agent should choose when it is in a given situation.

Sources

Throughout this compendium, each piece of information will be formatted as such.

Name / Description of fact	Source
----------------------------	--------

Information about fact.	
-------------------------	--

The location which currently contains “Source” could potentially be filled with a variety of sources. Here is how to find the source based off the shortened form.

- **Hands-On Machine Learning** refers to Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd Edition by Aurélien Géron

1 Data Analysis

TODO:

- Data Model
- Overfitting/Underfitting
- Feature Normalization
- Types of Loss Functions

2 Linear Models

TODO:

- Gradient Descent
- Normalization
- Logistic Regression

2.1 Linear Regression

2.1.1 Formulation

Linear Regression is a supervised machine learning algorithm where the predicted output is continuous and has a constant slope. Our main objective is to generate a line that minimizes the distance from the line to all of data points. This is essentially minimizing the error and maximizing our prediction accuracy.

2.1.2 Simple Regression

A simple two variable linear regression uses the slope-intercept form, where m and b are the variables our algorithm will try to "learn". x represents our input data and y represents the prediction.

$$y = mx + b$$

2.1.3 Multivariable Regression

Often times there are more than one feature in the data and we need a more complex multi-variable linear equation as our hypothesis. We can represent our hypothesis with the follow multi-variable linear equation, where \mathbf{w} are the weights and \mathbf{x} is the input data.

$$\begin{aligned} h_{\mathbf{w}}(\mathbf{x}) &= w_0x_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n \\ &= \mathbf{w}^T \mathbf{x} \end{aligned}$$

2.1.4 Cost Function

To predict based on a dataset we first need to learn the weights that minimize the mean squared error (euclidean loss) of our hypothesis. We can define the following to be our cost function to minimize with m being the number of data points and i being the i^{th} training example.

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2m} \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^i) - \mathbf{y}^i)^2 \\ &= \frac{1}{2m} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \end{aligned}$$

2.1.5 Gradient Descent Solution

Now to solve for \mathbf{w} we can use Gradient Descent and iteratively update \mathbf{w} until it converges. We get the slope of the cost function to be:

$$\frac{\partial J\mathbf{w}}{\partial \mathbf{w}_j} = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}^i - \mathbf{y}^i) \mathbf{x}_j^i$$

now applying a step α we can iteratively change \mathbf{w} until it reaches the global minima.

$$\mathbf{w}_j := \mathbf{w}_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\mathbf{w}}(\mathbf{x}^i) - \mathbf{y}^i)$$

2.1.6 Normal Equation Solution

The closed form solution to the linear system in \mathbf{w}

$$\frac{\partial J\mathbf{w}}{\partial \mathbf{w}_j} = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}^i - \mathbf{y}^i) \mathbf{x}_j^i = 0$$

writing this as a linear system in w we get $A\mathbf{w} = b$ where

$$A = \sum_{n=1}^N (\mathbf{x}_n \mathbf{x}_n^T) \text{ and } b = \sum_{n=1}^N (\mathbf{x}_n y_n)$$

so we can solve for $\mathbf{w} = \mathbf{A}^{-1}\mathbf{b}$ and get the following vectorized solution.

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

2.2 Generalized Linear Models

Often times our data won't be linear and it could be of a higher degree polynomial or a completely different distribution altogether. We can turn this non-linear problem into a linear regression problem by mapping the data to a different vector space using a basis function.

To demonstrate, let us consider a $N \times 1$ (feature) dataset. Let ϕ denote the polynomial basis function where $\phi_j(\mathbf{x}) = x^j$. Then we can express our hypothesis as:

$$h_{\mathbf{w}}(\mathbf{x}) = w_0\phi(x) + w_1\phi_1(x) + w_2\phi_2(x) + \dots + w_3\phi_3(x)$$

A dataset with 3 features with a polynomial basis would have a hypothesis as such

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1^2x_2 + w_6x_1x_2^2 + w_7x_1^2x_2^2 + w_8x_1^3 + w_9x_2^3$$

This can then be extrapolated to n -features. Some commonly used basis functions are:

- Polynomial: $\phi_j(\mathbf{x}) = x^j$
- Gaussian: $\phi_j(\mathbf{x}) = e^{(\frac{x-\mu_j}{2s^2})}$
- Sigmoid: $\phi_j(\mathbf{x}) = \sigma(\frac{x-\mu_j}{s})$
- Fourier Basis, Wavelets, etc ...

2.3 Regularization

Small outliers can drastically change our values of \mathbf{w} so rely on regularization to reduce overfitting. Polynomial models can be easily regularized by reducing the number of polynomial degrees. For a linear model, regularization is typically achieved by constraining the weights of the model. The regularization term should only be added to the cost function during training. Once the model is trained, the non-regularized cost should be used to measure the model's performance. The bias term w_0 is not regularized.

Ridge Regression

Ridge Regression (Tikhonov Regularization) is a regularized version of Linear regression with a regularization term of $\frac{\lambda}{2}\|\mathbf{w}\|_2^2$ (l_2 -norm) added to the cost function. This forces the learning algorithm to fit the data but also keep the model weights as small as possible. The hyperparameter λ controls how much you want to regularize the model.

$$J(\mathbf{w}) = MSE(\mathbf{w}) + \frac{\lambda}{2}\|\mathbf{w}\|_2^2$$

Lasso Regression

Least Absolute Shrinkage and Selection Operator Regression is another regularized version of Linear Regression, it adds a regularization term to the cost function but uses the l_1 norm of the weight vector instead of half the square of the l_2 norm.

$$J(\mathbf{w}) = MSE(\mathbf{w}) + \lambda \sum_{i=1}^n |w_i|$$

An important characteristic of Lasso Regression is that it tends to eliminate the weights of the least important features (i.e, set them to zero). Lasso Regression automatically performs feature selection and outputs a *sparse model*.

Elastic Net

Elastic Net is a middle ground between Ridge Regression and Lasso Regression. The regularization term is a simple mix of both Ridge and Lasso's regularization terms and you can control the mix ratio r . When $r = 0$, Elastic Net is equivalent to Ridge Regression, and when $r = 1$, it is equivalent to Lasso Regression.

$$J(\mathbf{w}) = MSE(\mathbf{w}) + \frac{(1-r)\lambda}{2} \|\mathbf{w}\|_2^2 + r\lambda \sum_{i=1}^n |w_i|$$

Early Stopping

Early Stopping is a different way to regularize iterative learning algorithms such as Gradient Descent. This method aims to stop training as soon as the validation error reaches a minimum. For all convex optimization problems there will be a global minima, once that global minima is reached the curve will start going up. This proposes to stop as soon as we reach the minimum.

Index

Early Stopping, 10

Elastic Net, 10

Lasso Regression, 9

Reinforcement Learning, 4

Ridge Regression, 9

Semisupervised Learning, 4

Supervised Learning, 3

Unsupervised Learning, 3