

Big Data Engineering

Introduction to Machine Learning

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Contents

- Definitions and terminology
- The overall process
- Main techniques
- Algorithms and examples
- Big Data Machine Learning
- R and PMML
- Spark MLlib
- Introduction to the lab

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Definition of Machine Learning

- Algorithms that can learn from data

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Definition of Machine Learning

- Algorithms that can learn from data

Ok that was a circular definition 😊

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Definition take 2

- Algorithms that can analyse a set of data to find patterns and then make predictions when new data comes in

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Uses of Machine Learning

- Fraud Detection
 - Spam emails, fake reviews, credit card fraud
- Personalization
 - Recommendations
- Targeted Marketing
 - Predictive preferences, cross-selling
- Content Classification
 - Document classification, sentiment analysis
- Customer Support
 - Social media analysis
- Many others

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Learning phase



 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Usage phase



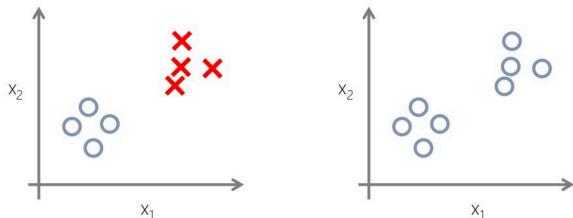
 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Terminology

- Sample
 - Some incoming data to be analysed
 - E.g. a JPG picture
- Feature
 - Some quantifiable data from the sample
 - E.g. colour, height, width, pixel data, etc
- Label
 - Some useful information about the sample that we wish to categorise:
 - E.g. looking at a picture this is a person
- Model
 - The output of some learning algorithm
 - The parameterization of an algorithm that can be run against new data

 © Paul Freudenthal 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Supervised vs Unsupervised



 © Paul Freudenthal 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Types of learning

- Supervised
 - The required labels are known
 - Aiming to find an algorithm that correctly identifies these
 - Iterative exploration and refinement
 - Useful for prediction
- Unsupervised
 - The labels are not known
 - The system identifies new classifications
 - Exploring the past, better understanding it
- Reinforcement
 - Learning as you go
 - E.g. learning to play chess while playing chess

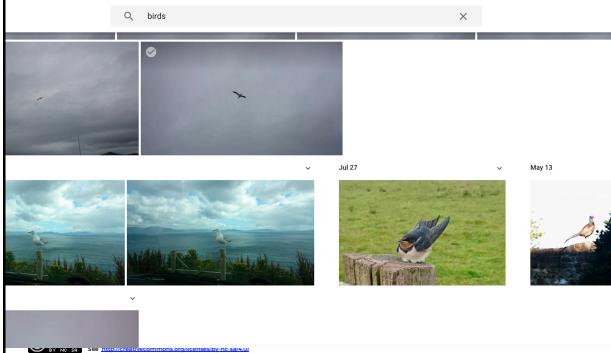
 © Paul Freudenthal 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Types of machine learning

- Classification
- Regression / Prediction
- Clustering
- Recommendation and Collaborative Filtering
- Frequent Pattern mining

 © Paul Freudenthal 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Classification



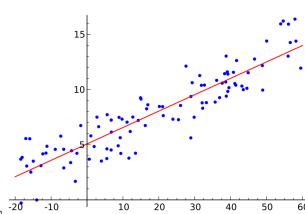
Classification

- Identifying a class into which this sample fits
 - E.g. look at a picture and decide if it contains a bird
 - A key part of artificial intelligence
 - Also deeply useful for making sense of big data

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Regression

- Applying a model based on previous data
 - Allows prediction of future state
- Many statistical techniques



© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Regression vs Classification

- Regression produces a real number or numbers
 - i.e. a continuously varying answer or answers
- Classification identifies a set or element of a set
 - E.g. False, Blue, Person, High-Value Customer

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Bayes Theorem



$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

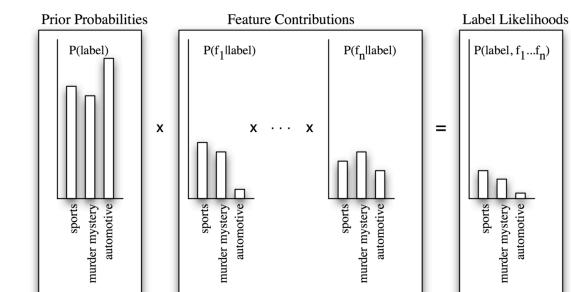
$P(A|B)$ is the probability of A given B
 $P(A)$ is the probability of A without regard to B

© Paul Frenzen 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Classification Algorithms

Naïve Bayes

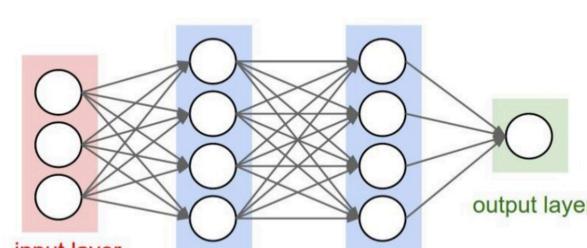
Prior Probabilities Feature Contributions Label Likelihoods



$P(\text{label})$ $P(f_1|\text{label}) \times \dots \times P(f_n|\text{label})$ $P(\text{label}, f_1 \dots f_n)$

sports murder mystery automotive

Deep Learning



input layer hidden layer 1 hidden layer 2 output layer

© Paul Frenzen 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Google TensorFlow

TensorFlow:
Large-Scale Machine Learning on Heterogeneous Distributed Systems
 (Preliminary White Paper, November 9, 2015)

Marin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Jeffrey Dean, Michael Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng
 Google Research*

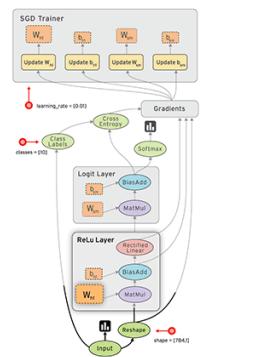
Abstract:
 TensorFlow [1] is an interface for expressing machine learning algorithms, and an implementation for executing such algorithms. A computation expressed using TensorFlow can be executed with little or no change on a wide variety of heterogeneous systems, ranging from mobile devices such as phones and tablets up to large-scale distributed systems of thousands of machines connected via commodity networking such as GPU cards. The system is flexible and can be used to express a wide variety of algorithms, including training and inference using deep neural network models, and it has been used for conducting research and for developing production systems that produce across more than a dozen areas of computer science and other fields, including speech recognition, computer vision, robotics, information retrieval, natural language processing, scientific information extraction, and computational drug discovery. This paper describes the TensorFlow interface and an implementation of that interface that we have built at Google. The TensorFlow API and a reference

* TensorFlow was developed by the Brain Team, a machine learning organization with members from the Google Brain team, more than 50 teams at Google and other Alphabet companies, and many others. TensorFlow is a trademark of Google LLC. TensorFlow has been deployed to tens of thousands of servers in Google Datacenters and is used in dozens of products, including Google Search [11], our advertising products, our speech recognition systems [50, 6, 46], Google Photos [43], Google Maps and StreetView [19], Google Translate [18], YouTube, and many others.

Based on our experience with DistBelief [2] and a model comparison of the desirable system properties and requirements for training and using neural networks, we have built TensorFlow, our end-to-end system for training and using deep learning models. TensorFlow takes computations described using a dataflow-like model and maps them onto a wide variety of different hardware platforms, ranging from running inference on mobile device platforms such as Android and iOS to modern

TensorFlow

- Announced and open sourced in Nov 2015
 - Strong adoption in the meantime
- CPU and GPU support with no coding changes
- Neural networks plus arbitrary dataflows
- www.tensorflow.org



© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Inspirobot.me



© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Spark MLLib's algorithms

Problem Type	Supported Methods
Binary Classification	linear SVMs, logistic regression, decision trees, random forests, gradient-boosted trees, naive Bayes
Multiclass Classification	decision trees, random forests, naive Bayes
Regression	linear least squares, Lasso, ridge regression, decision trees, random forests, gradient-boosted trees, isotonic regression

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Clustering

- Grouping items into clusters
 - Where items in a cluster are more similar to each other than to items in other clusters
- Basically creating the classifications from the data rather than applying them a priori

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

K-Means Clustering

Demonstration of the standard algorithm

1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).

2. k clusters are created by associating every observation with the nearest mean. The partitions here represent the Voronoi diagram generated by the means.

3. The centroid of each of the k clusters becomes the new mean.

4. Steps 2 and 3 are repeated until convergence has been reached.

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution NonCommercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

MLLib's clustering

- K-means
- Gaussian mixture
- Power iteration clustering (PIC)
- Latent Dirichlet allocation (LDA)
- Streaming k-means

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution NonCommercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Recommendation and Collaborative Filtering

- Given a user's interaction with items, what else are they likely to prefer

Large-scale Parallel Collaborative Filtering for the Netflix Prize

Yunhong Zhou, Dennis Wilkinson, Robert Schreiber and Rong Pan
HP Labs, 1501 Page Mill Rd, Palo Alto, CA, 94304
{yunhong.zhou, dennis.wilkinson, rob.schreiber, rong.pan}@hp.com

Abstract. Many recommendation systems suggest items to users by utilizing techniques of collaborative filtering (CF) based on historical records of items that the user has viewed or explicitly rated. Two major problems that most CF approaches have to resolve are scalability and sparseness of the user profiles. In this paper, we describe Alternating-Least-Squares with Weighted- λ -Regularization (ALS-WR), a parallel algorithm that we designed for the Netflix Prize, a large-scale collaborative filtering challenge. We use parallel Matlab on a Linux cluster

Frequent Pattern Mining

Related topics: [@cassiomeio](#) Your last post was about a hangout. These are the topics you relate to hangout: tonight, movies, board game, friends, awesome, photos, bar and NBA.

Related words for: "hangout"

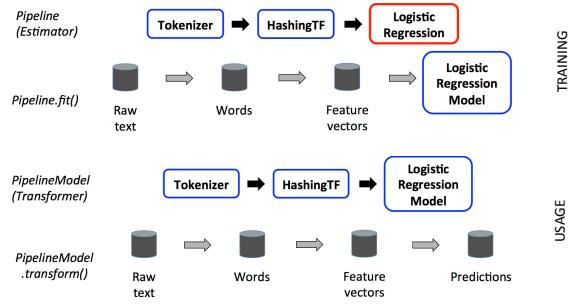
Word	Percentage
tonight	23.81 %
photos	23.81 %
board game	9.52 %
awesome	9.52 %
friends	14.29 %
movies	9.52 %
:	9.52 %

MLLib FPM

- Frequent pattern mining
 - FP-growth
 - association rules
 - PrefixSpan

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Spark MLLib Pipelines



Big Data ML

- Obviously we can learn more insights with more data
- Many examples
 - Netflix competition
 - Google, Facebook, Twitter etc are all doing big data ML
- Obviously we want the right algorithms:
 - E.g. Kmeans++ is a parallelizable version of Kmeans
- MLLib and Mahout come pre-built with these

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Amazon Machine Learning

Powerful machine learning technology

Based on Amazon's battle-hardened internal systems



Not just the algorithms:

- Smart data transformations
- Input data and model quality alerts
- Built-in industry best practices

Grows with your needs

- Train on up to 100 GB of data
- Generate billions of predictions
- Obtain predictions in batches or real-time

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License
See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Predictive Model Markup Language (PMML)

- An XML language for sharing models from machine learning
- Supported by R and other packages
- Mahout has no support
- Spark can export but not import
 - And only from Scala, not Python



© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Recap

- Machine Learning is a powerful way of gaining insight and value from big data
 - Recommendation
 - Classification and prediction
 - Clustering and understanding
- Many coding and deployment options
- Built into Spark, Hadoop and AWS

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Questions?

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution Non-Commercial ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>