Data Science: state of the art

Valentin Kuznetsov, Cornell University

SOSC 2019

Who am I?

- Theoretical Physicists (neutrino oscillations) at Irkutsk Univ & JINR
- Particle Physicists (tracking, silicon detectors) at CERN
- PhD in Physics (theory + experiment) at JINR
- Computing in HEP at JINR, CERN, Fermilab, Cornell
 - * HEP experiments: NOMAD, D0, Cleo-c, CMS
- Data Scientists at Cornell University
 - Data management, data discovery, services
 - BigData, Analytics, Monitoring, Machine Learning

Introduction



DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE



hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at

WHAT TO READ NEXT

Big Data: The Management Revolution

5 Essential Principles for Understanding Analytics

Data Scientists Don't Scale

VIEW MORE FROM THE

October 2012 Issue







Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.







DATA & AI LANDSCAPE 2019

é

INFRASTRUCTURE	ANALYTICS & MACHINE INTELLIGENCE	APPLICATIONS – ENTERPRISE		
HADOOP ON-PREMISE cloudera' Hortonworks MAPR Pivotal. IBM InfoSphere' jethro HADOOP IN THE CLOUD STREAMING / IN-MEMORY AWS ■ Microsoft Azure Cogogle Cloud Microsoft Azure Cogogle Cloud Microsoft Azure Microsoft Az	Abricks hfluent TA° kx Abricks ATTIV/O Datameer incorta. inter ana. Mode ENDOR SiSU switchboard Starburst Attrivector and Starburst MathWorks•	SALES MARKETING - B2B NSIDESALES.COM peoplea Clari A aviso tact.ai [] TROPS fuse[machines] clearbit Gi NOTCH MICO Clari A aviso tact.ai [] TROPS fuse[machines] clearbit Gi NOTCH MICO NOTCH MICO NOTC		
Nosol DATABASES NewSol DATABASES GRAPH DBs MPP DBs CLOUD EDW S Oracle Microsof Azure Pivotal Microsof Azure Pivotal Microsof Azure Micr	BI PLATFORMS VISUALIZATION MACHINE LEARNING	HUMAN CAPITAL Hue - Twe		
DATA TRANSFORMATION Italend Opentaho alteryx O TRIFACTA StreamSets UNIFI DATA INTEGRATION SAP Data Services Openation DATA INTEGRATION SAP Data Services Openation DATA GOVERNANCE Informatica OR Data Services Openation Sap Data Services Openation Mana O dataworld	Computer vision Microsoft Azure Microsoft Azure Micro	APVERTISING Apvertising Apvertising Apprexis Apprexi		
STORAGE CUSTER SVCS Compared Caudi Microsoft Azure Compared paragrage Compared pa	SEARCH COG ANALYTICS Splunk> elasticsearch SOCIAL ANALYTICS Splunk> elasticsearch WEB / MOBILE / COMMERCE ANALYTICS splunk> elasticsearch arapheore wy THIC control elagolia COVEO Hootsuite: spinkir sumologic WEB / MOBILE / COMMERCE ANALYTICS arapheore control elagolia COVEO Sumologic control serveras simple reach control single reach bitly similarWeb control simple reach custora	HEALTHCARE HEALTH		
CROSS-INFRASTRUCTURE/ANALYTICS CONSCRAME CAMPAGE Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Control of Sass 1010DATA VMWARC TIBCS TERADATA ORACLE INterpret Source Co				
FRAMEWORKS QUERY / DATA FLOW DATA ACCESS & DATABASES D Spack	CHESTRATION & MGMT Image: Streaming & militike intervention interventintervention intervention interventinterventi	RNING / DEEP LEARNING RNING / DEEP LEARNING Readword RNING / DEEP LEARNING RNING / DEEP LEARNING RNING / DEEP LEARNING RNING / DEEP LEARNING Readword RNING / DEEP LEARNING RNING / DEEP LEARNING REAdword RNING / DEEP LEARNING RNING		
HEALTH WALIDIC Practice fusion	DATA SOURCES & APIs JTERS D DOW JONES PLAID SCORE PLAID SCORE PL	ATTON INTELLIGENCE REQUARE S		

July 16, 2019 - FINAL 2019 VERSION

© Matt Turck (@mattturck), Lisa Xu (@lisaxu92), & FirstMark (@firstmarkcap) mattturck.com/data2019

FIRSTMARK

You need to know (my bare minimum)

- Math, statistics, algorithms, be able to read scientific paper
- Programming languages: C/C++, Python, R, Go, etc.
- Shell scripting and unix tools: bash, sed, awk, etc.
- How to build / install packages / tools
 - from source code: make, autoconf, environment, tar, etc.
 - from package management tools: rpm, yum, apt, dpkg, pip, anaconda, and/or build your favorite Linux distribution
- Versioning tools: git, gitlab, bitbucket, etc.
- Compilers, linkers, structure of libraries, object files, etc.
- Statistical and visualization tools: R, MatLab, Pandas, NumPy, SciPy, matplotlib, etc.
- ML tools: Scikit-Learn, R, TensorFlow, Keras, xgboost, etc.

You need to know, cont'd

- Platforms: AWS, Microsoft Azure, Google Cloud, etc.
- Cloud infrastructures: Docker, Kubernetes, etc.
- BigData tools: Hadoop, Spark, HDFS, HDF5, etc.
- Databases: ORACLE, MySQL, SQLite, NoSQL, GraphDB, MongoDB, CouchDB, etc.
- * Monitoring: ElasticSearch, Kibana, Grafana, Prometheus, etc.
- Streaming: Spark, Kafka, Storm, etc.
- Collaboration: Jupyter, Zeppelin, Anaconda, SWAN, etc.
- Search: ElasticSearch, Lucene, Solr, etc.
- Lexical analysis & NLP: lexer, tokenizer, scanner, etc.
- Read, write, and ask questions about everything

Salary Growth Forecast for IT Jobs 2016-2017 (US)



Engineering challenges of 21st century National Academy of Engineering

- Advance personalized learning
- Make solar energy economical
- Enhance virtual reality
- Reverse-engineer the brain
- Engineer better medicine
- Advance Health informatics
- Restore and improve urban infrastructure

- Provide access to clean water
- Secure Cyberspace
- Prevent Nuclear terror
- Manage the Nitrogen cycle
- Develop carbon sequestration methods
- Engineer tools of scientific discovery

Problem statement



Acquire the data



Process the data



Understand the data



16

Actions

- Increase revenue
- Reduce operational costs
- Understand behavior
- Find anomalies
- Identify strategy
- Discover new features



Data Science



Scope of Data Science

- Data Exploration and Preparation
- Data Representation and Transformation
- Computing with Data
- Data Modeling
- Data Visualization and Presentation
- Science about Data Science



Data Exploration and Preparation

- Exploration
 - understand your data and perform Exploratory Data Analysis
- Preparation
 - perform data cleaning and understand anomalies and various artifacts



Data Representation and Transformation

- Data sources
 - Adata today comes from variety of sources, from home made txt files to SQL and noSQL databases, data streams, etc. Data Scientists need to know the structure, transformation, and algorithms to deal with modern data
- Data Transformation includes data cleaning, pruning, normalization, standardization, etc.
- Data Representation
 - use different mathematical structures for data representation

$$x' = rac{x-ar{x}}{\sigma}$$
 $x' = rac{x-\min(x)}{\max(x)-\min(x)}$

Categorical vs Numerical One-hot encoding Leave-one out encoding Word embedding

Computing with data

- Every Data Scientists should know several programming languages
 - scripting, specialized
 languages, general purpose
 languages
 - be fluent with cluster and cloud computing to run jobs over massive datasets
 - organize workflows, from scripting to reproducible notebooks



Data Visualization and Presentation

- Be able to tell the story
 - use histograms, scatterplots, time series plots, heat maps, etc.
 - use variety of visualization frameworks suitable for different needs



Data Modeling

Generative
 modeling, in which
 one proposes a
 stochastic model
 that could have
 generated the data
 (domain of
 Academic statistics)

Predictive
 modeling, in which
 one constructs
 methods which
 predict well over
 some given data
 universe (domain of
 Machine Learning)



Science about Data Science

 Data scientists are doing science about data science when they identify commonly occurring analysis or processing workflows



Does increase in AUC bring any value? Simple vs complex model Operational cost Efficiency

Common Task Framework (CTF)

- How to choose best model/approach
 - Provide publicly available datasets
 - A set of enrolled competitors whose common task is to infer a class prediction rule from the training data
 - A scoring referee, to which competitors can submit their training rule. The referee runs the prediction rule against a test dataset not available to competitors, and objectively and automatically reports the score achieved by the submitted rule.
- kaggle.com is an example of CTF
- DIY solutions:
 - split data into training, validation and test set
 - use different models
 - perform cross-validation, ensembles, etc.

Scientific analyses

- A broad collection of technical activities is not a science
 - to do science we must have continually evolving, evidence based approach
- Meta-data analysis: study data analyses on a given topic, does published analyses succeeded and can be improved (reproduced)
- Cross-study analysis: use common datasets to validate studies of different groups by fitting proposed models on different set of datasets and validating them on another set and measure concordance of predictive results
 - study individual models across datasets
 - study individual datasets across models
- Cross-workflow analysis: study effect of different analysis workflows on prediction outcome

Data, Algorithms, Techniques

Engineering Effort for Effective ML

From "Hidden Technical Debt in Machine Learning Systems",
 D. Sculley at al. (Google), paper at NIPS 2015



Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

Data pre-processing

- Most of the time will be spend in this step
- Data clean-up, data transformation, feature engineering
 - data transformation
 - scaling and normalization
 - encoding, aggregation features, log-transformation (to remove outliers)
 - data visualization, exploration
 - data augmentation, imputing, bucketing, binning, feature interactions
 - dimensionality reduction
- Your programming skills will be required here: R, Python, Databases, etc.

Types of data



How to better represent dates: categorical or numerical?

Data transformation

- Data transformation and aggregation: log, sum of values
- Scaling: a technique to scale data to a given range [0,1] or any other range
- Normalization/Standardization: a technique to scale data to mean with zero and and unit-variance
- Augmentation: a technique to create additional data based on input sample which slightly differ from it, e.g. image rotation, flip, scale, crop, etc.
- Bucketing/Binning: a technique to place similar values into buckets/bins

 $x' = rac{x - \min(x)}{\max(x) - \min(x)}$

 $x'=rac{x-ar{x}}{\sigma}$



32

One-hot-encoding

- It is a technique to handle
 "categorical" data
- It represents categorical column as vector of words
- You need to define word vector for full set of data (train + test datasets)
- Issues with NULL or missing data
 - delete rows with missing data
 - impute data for missing values

"One-Hot" refers to a state in electrical engineering where all of the bits in a circuit are 0, except a single bit with a value of 1 (it is said to be "hot").

Rome Paris
Rome =
$$[1, 0, 0, 0, 0, 0, 0, ..., 0]$$

Paris = $[0, 1, 0, 0, 0, 0, ..., 0]$
Italy = $[0, 0, 1, 0, 0, 0, ..., 0]$
France = $[0, 0, 0, 1, 0, 0, ..., 0]$

Leave-one-out encoding

- Use mean of all values within the same category except given row
- Add random noise
- Replace categorical value with leave-one-out times noise
- The test categorical values always represented as mean and no noise
- This technique may complement one-hot encoding

Split	UserID	Y	mean_y	random	new_Y
Train	A1	0			
Train	A1	1			
Train	A1	1			
Train	A1	0			
Test	A1	-			
Test	A1	-			
Train	A2	0			

We'll show how to add new categorial encoding features: mean_y: average value for given user random: random factor for given user new_Y: new encoding feature for given user Ref

Word embedding

- A way to capture multi-dimensional relationships between categories
 - e.g. Sun and Sat may have similar effect while other days may be treated independently
 - you define a dimension of word vector upfront
 - it projects categorical variables into another phase space, e.g. days may be sunny or rainy, season or off season; all of these features are hidden from original data representation
- Use NN or other ML algorithms to train the model to find best representation of embedded variables

real\hidder	Dog Age Cat
puppy	[0.9, 1.0, 0.0]
dog	[1.0, 0.2, 0.0]
kitten	[0.0, 0.1, 0.9]
cat	[0.0, 1.0, 1.0]
	puppy dog kitten cat

Ref

Data visualization

- Graphical representation may reveal important features of the data
 - find correlations, identify range, etc.
- Identify features which may require transformations, e.g. see outliers or skewness in data
- It helps to identify a strategy how to deal with different features







ML algorithm

- Inputs: X, e.g. timestamp, price, color, size, etc.
- ✤ Features: X, transformed inputs
- Labels: y (stay vs leave)
- Weights: W (matrix)
- Activation function: φ (step function, e.g. sigmoid)
- * Predictions: $z = \phi(W^T X)$ yields (-1,1)
- * Cost function: J(**W**), e.g. $\sum (y_i z_i)^2/2$
- Algorithm: minimizes cost function & find best separation



	ТҮРЕ		DESCRIPTION	ADVANTAGES	DISADVANTAGES
ear		Linear regression	The "best fit" line through all data points. Predictions are numerical.	Easy to understand you clearly see what the biggest drivers of the model are.	 X Sometimes too simple to capture complex relationships between variables. X Tendency for the model to "overfit".
Line		Logistic regression	The adaptation of linear regression to problems of classification (e.g., yes/no questions, groups, etc.)	Also easy to understand.	 X Sometimes too simple to capture complex relationships between variables. X Tendency for the model to "overfit".
		Decision tree	A graph that uses a branching method to match all possible outcomes of a decision.	Easy to understand and implement.	X Not often used on its own for prediction because it's also often too simple and not powerful enough for complex data.
Tree-based		Random Forest	Takes the average of many decision trees, each of which is made with a sample of the data. Each tree is weaker than a full decision tree, but by combining them we get better overall performance.	A sort of "wisdom of the crowd". Tends to result in very high quality models. Fast to train.	 X Can be slow to output predictions relative to other algorithms. X Not easy to understand predictions.
	Ŷ	Gradient Boosting	Uses even weaker decision trees, that are increasingly focused on "hard" examples.	High-performing.	 X A small change in the feature set or training set can create radical changes in the model. X Not easy to understand predictions.
Neural networks		Neural networks	Mimics the behavior of the brain. Neural networks are interconnected neurons that pass messages to each other. Deep learning uses several layers of neural networks put one after the other.	Can handle extremely complex tasks - no other algorithm comes close in image recognition.	 X Very, very slow to train, because they have so many layers. Require a lot of power. X Almost impossible to understand predictions.

40

<u>Ref 1</u>

<u>Ref 2</u>

Loss functions



<u>Ref 1</u>

Regularization

- <u>Ref 2</u>
- One of the major aspects of training the model is overfitting, when ML model tries too hard to capture the noise in your training dataset
- * **Regularization** term is an addition to loss function which helps generalize the model. It helps to learn simpler model, induce models to be sparse, introduce group structure into learning problem $\min_{f} \sum_{i=1}^{n} V(f(x_i), y_i) + \lambda R(f)$
 - * L1 or Lasso regularization adds penalty which is a sums of the absolute values of weights $Min(\sum_{i=1}^{n} (y_i w_i x_i)^2 + p \sum_{i=1}^{n} |w_i|) \qquad MSE+L1$
 - * L2 or Ridge regularization adds penalty which is a sums of the squared values of weights $Min(\sum_{i=1}^{n} (y_i w_i x_i)^2 + p \sum_{i=1}^{n} (w_i)^2) \qquad MSE+L2$
- Dropout is a term introduced in NN context where hidden nodes are dropped randomly and allow model to generalize better
- Early Stopping is time regularization technique which stop training based on given criteria

Data Science recipe

- Understand your data: preprocessing, cleaning, augmentation, onehot-encoding
- Categorize the problem: classification, regression, clustering, dimensionality reduction
- Choose the language and toolkit: R, Python, Hadoop+Spark, ML providers
- Choose the right technique: trees, bagging, stacking, boosting, (rank | weight) averaging, NNets
- Start coding using your favorite ML framework and visualization tools

Techniques



- exactly once, and gets to be in a validation set set *k-1* times
- This significantly reduces bias as we are using most of the data for fitting, and also significantly reduces variance as most of the data is also being used in validation set
- Stratified K-Fold Cross Validation deals with imbalanced data, each fold contains the same percentage of samples of each class
- Leave-P-Out Cross Validation: leaves p data points out of train set, i.e. if we have n data points, then n-p points used for training and p points for validation



Machine Learning



Search Contents

riance problems



Plotting
 validatio
 variance

High bia
 fit the tra

incre para addi

 High-vai overfittir

> to co and : mod

Ensembles

<u>Ref 1</u> <u>Ref 2</u> <u>Ref 3</u>

All models are wrong, but some are useful (George Box)

Sometimes intentionally built weak models are good blending candidates

- Bagging
 - * building multiple models (typically of the same type) from different subsamples of the training dataset
- Boosting
 - building multiple models (typically of the same type) each of which learns to fix the predictions errors
 of a prior model in the chain
- Stacking
 - building multiple models (typically of the different types) and supervisor model that learns how to best combine the predictions of the primary model
- Weighting | Blending
 - combine multiple models into single prediction using different weight functions

Diversity is a key: use different un-correlated models, e.g. GBM, RF, SVM, NN



: voting



better to fight over-fitting better to get lower errors

Bagging vs Boosting

Ref

Similarities

Both are ensemble methods to get N learners from one

Generate several training sets by random sampling

Make final decision by averaging N learners or taking majority of them



bagging

Differences

boosting

build independently for Bagging, and Boosting tries to add new models that do well where previous models fail

Boosting weights the data to scale in favor of most difficult cases

Bagging: equally weighted average Boosting: weighted average, more weight to those who perform better on training set

Stacking

- Stacking (also called metaensembling) is a model ensembling technique used to combine information from multiple predictive models to generate a new model
- Usually outperform individual models used in ensemble, e.g. GBM+RF+NN
- Most effective when base models are independent
- May be applied at multiple level, e.g. stacking first set, then second set, etc.

Consider datasets A,B,C. Target variable (y) is known for A,B.

Ref

Technical tricks

Use one set of features (text) for simple model 1, and use numerical features and model1 prediction * for model 2, etc.



- ✤ Use chained models: build stand-alone model for G, then used in next model, e.g. F=>G=>B=>A
- ✤ Feature engineering:
 - one-hot-encoding, leave-one-out, word embedding and add them to original data set
 - split days into years, months, dates and threat them as categorical variables
 - aggregate values, e.g. sum all numerical values in a row and/or use its mean/median
 - handle missing values, e.g. apply mean across column or even apply additional training to find their values

Tools and frameworks





learn Classification Regression Clustering

Dimensionality reduction Model selection Preprocessing



DataFrame data.table ggplot xgboost NeuralNetwork Trees, Bagging



ML for "standard" use-cases

- In most cases you may rely on R or Python eco-system. In Python <u>scikit-learn</u> is de-facto standard, in R all ML tools are available through 3rd party packages via install.packages(<pkg>)
- Majority of DataScientists in kaggle competition use <u>xgboost</u>, the distributed gradient boosting library (both R and Python APIs are available) based on parallel tree boosting algorithm (aka GBDR, GBM)
- Less known libraries are:
 - Weka is Waikato Environment for Knowledge Analysis is a suite of machine learning software written in Java, developed at the University of Waikato, New Zealand (GUI environment)
 - <u>StackNet</u> is a computational, scalable and analytical Meta modeling framework (developed by toplevel kaggle competitor Kaza-Nova and used in many competition to won first places). Written in Java and uses uses Wolpert's stacked generalization to improve accuracy of ML models. The network is built iteratively one layer at a time (using stacked generalization), each of which uses the final target as its target.
 - h2o Open Source Fast Scalable Machine Learning Platform For Smarter Applications (Deep Learning, Gradient Boosting, Random Forest, Generalized Linear Modeling (Logistic Regression, Elastic Net), K-Means, PCA, Stacked Ensembles, Automatic Machine Learning (AutoML)

Neural network frameworks

- Torch is an open source machine learning library, a scientific computing framework, and a script language based on the Lua programming language.
- Theano is a numerical computation library for Python that allows you to define, optimize, and evaluate mathematical expressions involving multi-dimensional arrays efficiently. In Theano, computations are expressed using a NumPy-esque syntax and compiled to run efficiently on either CPU or GPU architectures.
- * <u>Caffe</u> is a deep learning framework (C++ and Python) made with expression, speed, and modularity in mind.
- <u>TensorFlow</u> is an open-source software library (C++, Python, Go) for data-flow programming across a range of tasks. It is a symbolic math library, and is also used for machine learning applications such as neural networks.
- <u>PyTorch</u> is a deep learning framework for fast, flexible experimentation. It is Tensors and Dynamic neural networks in Python with strong GPU acceleration.
- * Keras is a high-level neural networks API, written in Python and capable of running on top of TensorFlow, CNTK, or Theano
- * <u>fast.ai</u> library simplify training ML/DL learning using best DataScientists and ML practitioner practices.
- * Apache MXNet framework (Python and R) is a modern deep learning framework
- <u>onnx.ai</u> is an Open Neural Network exchange format which allows to import and export Neural Network models from/to different frameworks

Visualization of Neural Networks

- TensorFlow playground: provides an intuitive web based interface to train Neural Networks for a given dataset
- <u>ConvNetJS</u> is a Javascript library for training Deep Learning models (Neural Networks) entirely in your browser
- * LSTMVis visual analysis for Recurrent Neural Networks
- <u>Netron</u> is a visualizer for Deep Learning and machine learning models
- * <u>Ann-visualizer</u>, is a python library for visualizing Artificial Neural Networks
- Keras-vis is a high-level toolkit for visualizing and debugging your trained keras neural net models
- VisualDL is an open-source cross-framework web dashboard that richly visualizes the performance and data flowing through your neural network training

ML for Big Data

- Some datasets can't be trained with standard ML tools since they are too big to fit into memory, therefore you can't use "standard" tools like scikit-learn or R
- Gradient Boosting Algorithm (<u>GBM</u>) is a ML technique which produces a prediction model in a form of ensemble of weak prediction models, typically decision trees
 - Boosting is an ensemble technique in which the predictors are not made independently, but sequentially. Therefore a large dataset can be learned in "chunks" with GBM
- * <u>Vowpal Wabbit</u> is online learning algorithm designed to deal with tera-features datasets
- Spark ML Big Data platform (<u>MLlib</u>), Spark is a technique to deal and process large datasets using Hadoop platform which now has a set of ML algorithms available as a part of platform

Courses

- <u>kaggle.com</u> is a place to do data science projects, it is your ULTIMATE source of knowledge in DataScience, ML, DL and AI
- fast.ai provides cutting edge about deep learning
- * Google TensorFlow Development Summit new ideas and practical implication of TF
- Machine Learning A-Z: Hands on Python & R In Data Science covers machine learning workflows
- * Scala and Spark for Big Data and Machine Learning covers Big Data technology
- Building Neural Network from scratch: github and blog
- Machine Learning courses ranked by user reviews

Resources

- * How to get started with ML
- * Choosing the right ML algorithm
- * Colah's blog
- * Stacking Made Easy
- * Gradient Descend Optimization
- * ML, Python and Math Cheat Sheets
- * Data Science interview questions
- * Neural Network zoo
- * Large Scale Deep-Learning with TensorFlow
- Learning Machine Learning
- * Cheat Sheet for AI, ML, NN, BigData Salary history and career path of a Data Scientist



The Story

BBCFOUR

https://www.youtube.com/watch?feature=player_embedded&v=jbkSRLYSojo