

Introduction to statistical learning

1. Introduction

V. Lefieux

June 2018



Big data

Data analytics

Data science

Statistical learning

Practical
informations

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

Statistical learning

Practical informations

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

Statistical learning

Practical informations

Data everywhere

Really "Big Data" – Data volume is increasing exponentially

By the year 2020, the digital universe will reach 44 zettabytes – that's a 10-fold increase from 2013.



Source: IDC's Digital Universe Study, April 2014, sponsored by EMC

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Data everywhere

Big data

Data analytics

Data science

Statistical learning

Practical
informations

- ▶ **Before:**
 - ▶ structured data,
 - ▶ generated by companies and organizations,
 - ▶ regular but not so frequent updates (e.g monthly).
- ▶ **Now:**
 - ▶ unstructured data,
 - ▶ generated by users,
 - ▶ real time data.

Some data generated by companies and organization



Big data

Data analytics

Data science

Statistical learning

Practical
informations

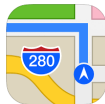
Some data generated by users

amazon

Google



LinkedIn



Big data

Data analytics

Data science

Statistical learning

Practical
informations

And now health data



Big data

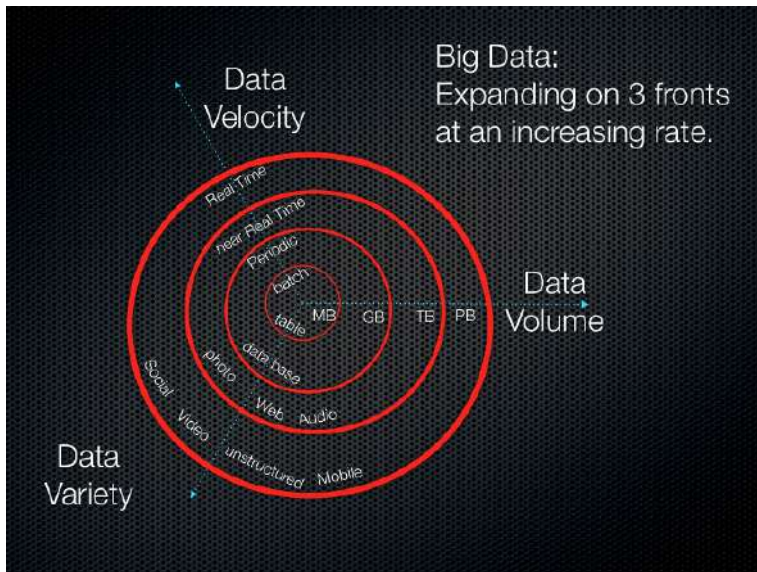
Data analytics

Data science

Statistical learning

Practical informations

3 V ?



Big data

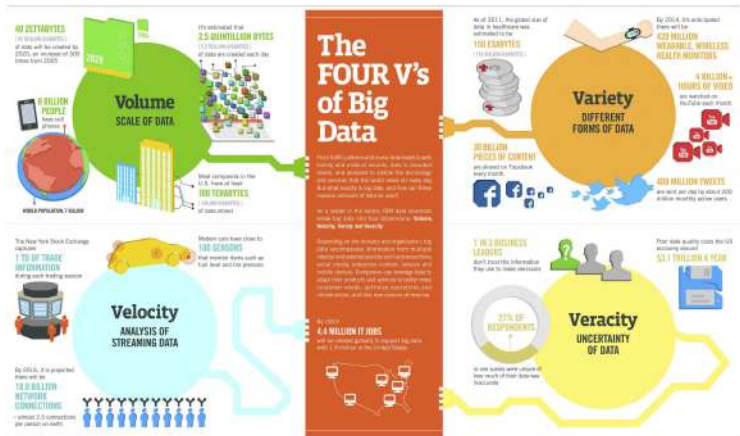
Data analytics

Data science

Statistical learning

Practical
informations

4 V's



Big data

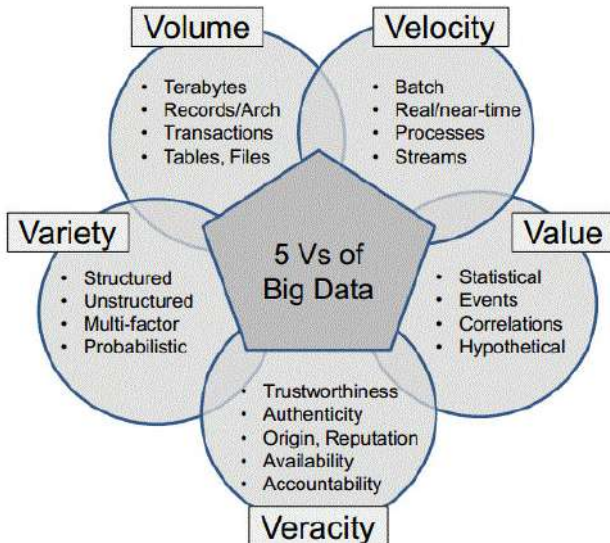
Data analytics

Data science

Statistical learning

Practical informations

5 V ?



Big data

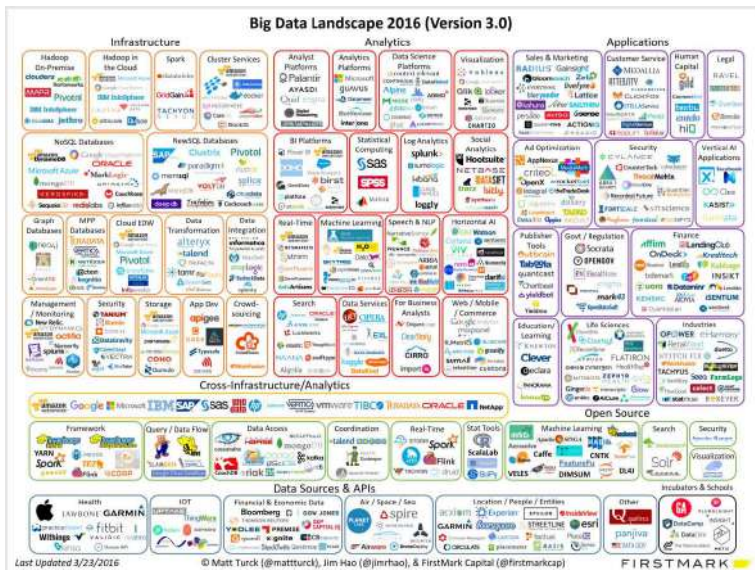
Data analytics

Data science

Statistical learning

Practical informations

A landscape



Big data

Data analytics

Data science

Statistical learning

Practical informations

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

Statistical learning

Practical informations

Data analytics

Big data

Data analytics

Data science

Statistical learning

Practical
informations

A process:

- ▶ collecting,
- ▶ organizing (cleaning and storing),
- ▶ analyzing,
- ▶ visualizing

large sets of data.

An objective: discover useful information to improve business decisions.

A new idea ?

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Four major influences act on data analysis today:

- ▶ *The formal theories of statistics.*
- ▶ *Accelerating developments in computers and display devices.*
- ▶ *The challenge, in many fields, of more and ever larger bodies of data.*
- ▶ *The emphasis on quantification in an ever wider variety of disciplines.*

Not so new !

Data analysis and statistics: an expository overview

J. W. Tukey and M. B. Wilk

1966

Four major influences act on data analysis today:

- ▶ *The formal theories of statistics.*
- ▶ *Accelerating developments in computers and display devices.*
- ▶ *The challenge, in many fields, of more and ever larger bodies of data.*
- ▶ *The emphasis on quantification in an ever wider variety of disciplines.*

Big data

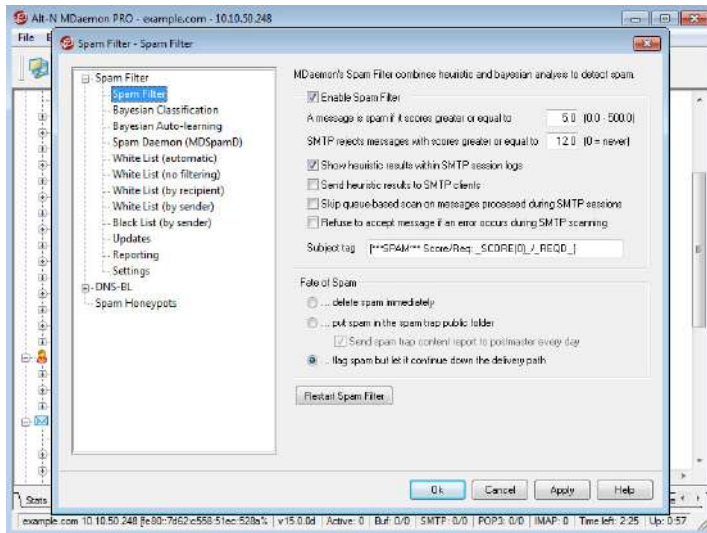
Data analytics

Data science

Statistical learning

Practical
informations

Spam filter



Big data

Data analytics

Data science

Statistical learning

Practical informations

White Papers & Reports



Moving Up the Digital Marketing Maturity with Big Data Analytics

Much of the big data evolution is being driven by the vast commercial possibilities of the Internet. The scope and pace of change suggests the advent of an even more fundamental shift, focusing the broader advertising community on a series of critical questions: How can data most effectively be used to address core business needs? Where can this data be sourced? How can it be shared, optimized and enhanced for analysis and monetization?

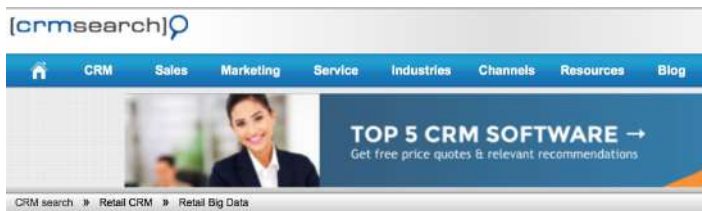
Are you asking yourself these same questions? We hope so.


This paper will achieve the following objectives:

1. History of media and its impact to marketing
2. Critical marketing imperatives for data-driven digital marketing
3. Four critical optimizations
4. Steps to success in the digital marketing big data journey
5. Summarize concepts and steps digital marketers need to embrace if they are to harness, manage and successfully exploit big data



Customer relationship management (CRM)



[crmsearch] 

[CRM](#) [Sales](#) [Marketing](#) [Service](#) [Industries](#) [Channels](#) [Resources](#) [Blog](#)

TOP 5 CRM SOFTWARE →
Get free price quotes & relevant recommendations

CRM search » Retail CRM » Retail Big Data



Big Data in Retail Examples

By [Chuck Schaeffer](#)

5 Retail Big Data Examples with Big Paybacks

Big data is delivering some big results for retailers.

- ▶ *Hotel chain uses big data to increase bookings.*
- ▶ *Pizza chain earns more dough in bad weather.*
- ▶ *Music distributor applies big data for demand planning.*
- ▶ *Financial services company scores new clients.*
- ▶ *Retailer creates pregnancy detection model.*

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Smart grids



And smart cities.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

nature International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video | For Authors

News & Comment | News | 2015 | October | Article

NATURE | NEWS

Genome researchers raise alarm over big data

Storing and processing genome data will exceed the computing challenges of running YouTube and Twitter, biologists warn.

Erika Check Hayden

07 July 2015

Rights & Permissions

The computing resources needed to handle genome data will soon exceed those of Twitter and YouTube, says a team of biologists and computer scientists who are worried that their discipline is not geared up to cope with the coming genomics flood.

Other computing experts say that such a comparison with other 'big data' areas is not convincing and a little glib. But they agree that the computing needs of genomics will be enormous as sequencing costs drop and ever more genomes are analysed.

Refugee trauma

The mental-health crisis among migrants

The refugees and migrants surging into Europe are suffering very high levels of psychiatric disorders. Researchers are struggling to help.

Like Share 293k people like this. Be the first of your friends.

Enterprise Big Data

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

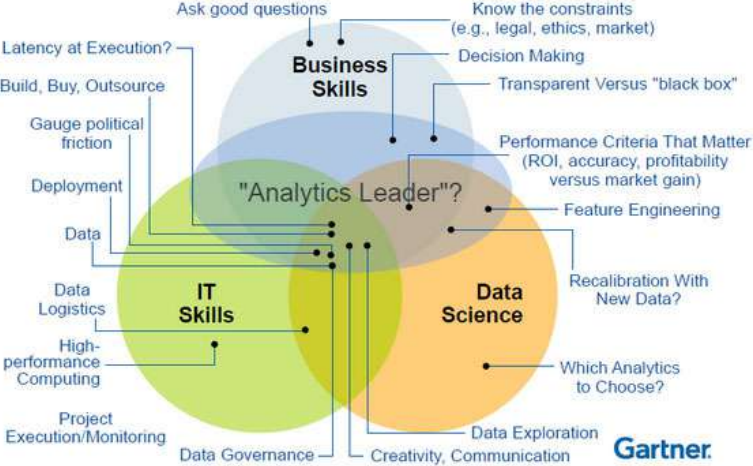
Statistical learning

Practical informations

Data scientist skills

- Big data
- Data analytics
- Data science
- Statistical learning
- Practical informations

Driving the Success of Data Science Solutions: Skills, Roles and Responsibilities ...



Gartner.

Some definitions: Data science



*Data science is an interdisciplinary field about processes and systems to extract knowledge or insights from data in various forms, either structured or unstructured, which is a continuation of some of the data analysis fields such as **statistics**, **machine learning**, data mining, and predictive analytics, similar to Knowledge Discovery in Databases (KDD).*

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

Statistical learning

Practical informations

Some definitions: Machine learning



Machine learning is a field of computer science that often uses statistical techniques to give computers the ability to “learn” (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Some definitions: Statistical learning



WIKIPEDIA
The Free Encyclopedia

Statistical learning theory is a framework for machine learning drawing from the fields of statistics and functional analysis. Statistical learning theory deals with the problem of finding a predictive function based on data. Statistical learning theory has led to successful applications in fields such as computer vision, speech recognition, bioinformatics and baseball.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Statistical learning vs Machine learning

- ▶ Machine learning, from Artificial Intelligence: large scale applications, prediction accuracy.
- ▶ Statistical learning, from Statistics: interpretability, precision, uncertainty, inference.
- ▶ For some statisticians: statistical learning is a mathematical formalisation of the machine learning.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Some concepts: online/offline learning

- ▶ **Online learning** (real-time): under time constraints.

Some examples:

- ▶ Personalized advertising.
 - ▶ Personalized healthcare.
 - ▶ Navigation & transit tools.
 - ▶ Autonomous cars.
 - ▶ Load curve forecasts.
 - ▶ Weather forecasts.
-
- ▶ **Offline learning** (batch).

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Some concepts: supervised/unsupervised learning

- ▶ **Supervised learning:**

Infer (predict) a function/relationship from labeled training data (e.g. classification, regression).

- ▶ **Unsupervised learning:**

Find “structure” in unlabeled data (e.g. clustering). Even if it is more subjective than supervised learning, it can be useful as a pre-processing step for supervised learning.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Supervised learning

Big data

Data analytics

Data science

Statistical learning

Practical
informations

There are many different paradigms, including:

- ▶ Parametric statistics (linear or non-linear).
- ▶ Non-parametric statistics (local estimation methods, e.g smoothing kernel methods, k -nearest neighbors).
- ▶ Tree based methods.
- ▶ Support Vector Machines.
- ▶ Deep learning.

Some key points

- ▶ Trade-off between prediction accuracy and interpretability.
- ▶ Avoid over-fitting.
- ▶ Parsimonious model vs (full) black box: “less is more”.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Table of contents

Big data

Data analytics

Data science

Statistical learning

Practical informations

Big data

Data analytics

Data science

Statistical learning

Practical informations

Outline

- ▶ Introduction.
- ▶ Unsupervised learning: PCA & clustering.
- ▶ Supervised learning:
 - ▶ Cross validation & bootstrap.
 - ▶ Reminders on linear regression & logistic regression.
 - ▶ Tree based methods.
 - ▶ Support Vector Machines.

Big data

Data analytics

Data science

Statistical learning

Practical
informations

Software tools

Big data

Data analytics

Data science

Statistical learning

**Practical
informations**

