ARTIFICIAL INTELLIGENCE
A program that can sense, reason, act, and adapt

MACHINE LEARNING
Algorithms whose performance improve as they are exposed to more data over time

DEEP LEARNING
Subset of machine learning in which multilayered neural networks learn from vast amounts of data

http://www.prowesscorp.com/whats-the-difference-between-artificial-intelligence-ai-machine-learning-and-deep-learning/

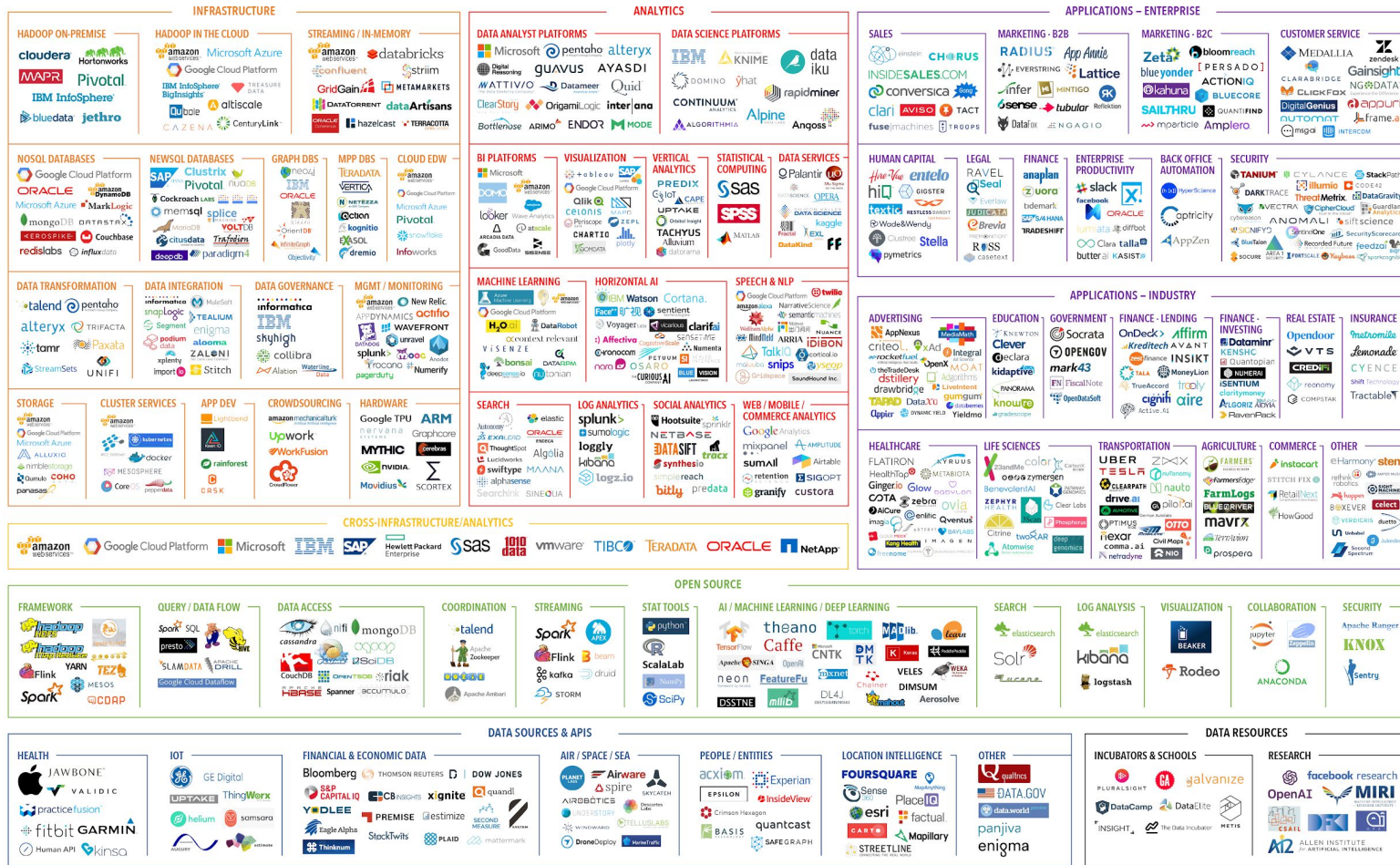THE AI ERA IS HERE

Ignited by a new computing model,
GPU Deep Learning

https://www.slideshare.net/NVIDIA/the-ai-era-ignited-by-gpu-deep-learning

# Trends



Legend: ● Big Data ● Machine Learning ● Data Science ● Tektonik

BIG DATA LANDSCAPE 2017

# BIG DATA & AI LANDSCAPE 2018

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# Data Scientist: The Sexiest Job of the 21st Century

# Keep Calm with Michael Jordan

Artificial Intelligence—The Revolution Hasn't Happened Yet

https://medium.com/@mijordan3/artificial-intelligence-the-revolution-hasnt-happened-yet-5e1d5812e1e7

# Keep Calm with Michael Jordan

Artificial Intelligence—The Revolution Hasn't Happened Yet

https://medium.com/@mijordan3/artificial-intelligence-the-revolution-hasnt-happened-yet-5e1d5812e1e7

# Keep Calm with Michael Jordan

Artificial Intelligence — The Revolution Hasn't Happened Yet

- We are currently building Machine Learning blocks
- Blocks are being put together in ad-hoc ways
- We are far from what we call "AI"

# Machine Learning Craftsmanship

Kelvin MOUTET

Machine Learning Engineer @ Prevision.io

# Software Craftsmanship

*As aspiring Software Craftsmen we are raising the bar of professional software development by practicing it and helping others learn the craft. Through this work we have come to value:*

*Not only working software ,*

*but also well-crafted software*

*Not only responding to change,*

*but also steadily adding value*

*Not only individuals and interactions,*

*but also a community of professionals*

*Not only customer collaboration,*

*but also productive partnerships*

*That is, in pursuit of the items on the left we have found the items on the right to be indispensable.*

*http://manifesto.softwarecraftsmanship.org/*

# PLAN

2 parts only !

# 1) Part 1

# Machine Learning Bare Necessities

- Data Scientists !

# Machine Learning Bare Necessities

- Data Scientists

- Tools !

# Machine Learning Bare Necessities

- Data Scientists

- Tools

- Data !

# Machine Learning Bare Necessities

- Data Scientists

- Tools

- Data

- BIG DATA !

Let's go !

# Machine Learning core steps



Get Data — 1

Clean, Prepare & Manipulate Data — 2

Train Model — 3

Test Data — 4

Improve — 5

# Go to Production !

# Machine Learning Infrastructure

# Welcome to AI ERA !!!

# Machine Learning "Success" Stories

# Why does it fail ?

# Software project failure !



The Top Six Reasons

- Incomplete requirements
- Lack of user involvement
- Lack of resources
- Unrealistic expectations
- Lack of executive support
- Changing requirements and specifications

www.sysgu.com
408-356-5793

# Machine Learning Project Failure (not exhaustive)

- Problem definition is bad or unclear (incomplete requirements)
- Scope is too large (unrealistic expectations)
- Overconfidence on data scientists (lack of resources)
- (Big) Data is too small or bad ! (lack of resources)
- True AI is not here (unrealistic expectations)
- Pessimistic on core business knowledge (lack of user involvement)
- Tools/ML Driven Development instead of Core Business Development
- Machine Learning algorithms understanding/coverage

1) ~~Part 1~~

2) Part 2

# 2) Machine Learning Craftsmanship

# Software Craftsmanship

*As aspiring Software Craftsmen we are raising the bar of professional software development by practicing it and helping others learn the craft. Through this work we have come to value:*

*Not only working software ,*

  *but also well-crafted software*

*Not only responding to change,*

  *but also steadily adding value*

*Not only individuals and interactions,*

  *but also a community of professionals*

*Not only customer collaboration,*

  *but also productive partnerships*

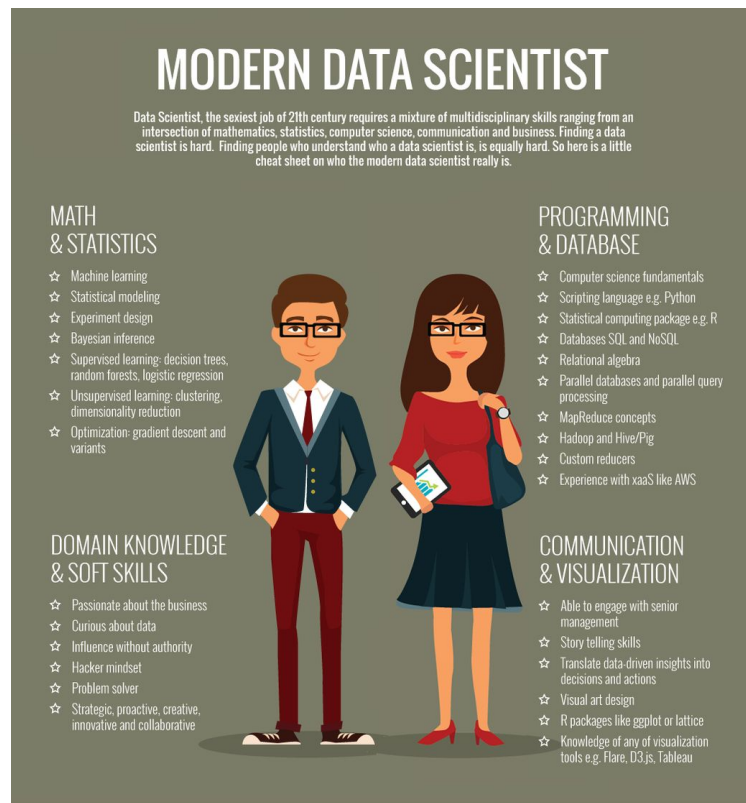*That is, in pursuit of the items on the left we have found the items on the right to be indispensable.*
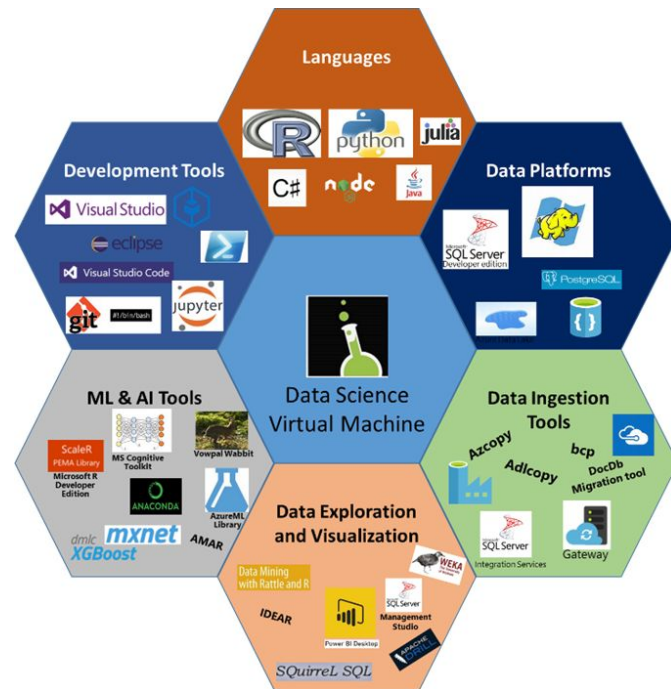
# Machine Learning Craftsmanship

As aspiring Software Craftsmen we are raising the bar of professional software development by practicing it and helping others learn the craft. Through this work we have come to value:

Not only working software,

but also well-crafted software

Not only responding to change,

but also steadily adding value

Not only individuals and interactions,

but also a community of professionals

Not only customer collaboration,

but also productive partnerships

That is, in pursuit of the items on the left we have found the items on the right to be indispensable.

# Machine Learning Craftsmanship

Mindset !

# A good way to start a Machine Learning project

# Examples !

# Machine Learning problem solving

Inspired by those two articles :

https://towardsdatascience.com/1-year-doing-data-science-in-the-real-world-54f49b591991
by **Jonny Brooks-Bartlett**, Data scientist Deliveroo

https://towardsdatascience.com/first-create-a-common-sense-baseline-e66dbf8a8a47 by
**Rama Ramakrishnan**, Senior Vice President, Salesforce.

# Identify your problem then build value !

# A good way to start a ML project

# Customers targeting

Send an advertising of your latest new products, at most, to 100,000 customers from the database.

How do you solve this ?

# Customers targeting

Pick the most **loyal** customers

      a)    Shop a lot (frequency) and b) spends (money) a lot over a period (year, month)

      c)    Recently shopped with you (number of day, week since last transaction)

Decile the customer databases, sort and pick the 100 000 first !

This is a core business tools called **Recency-Frequency-Monetary (RFM) heuristic**

# Product recommendations

Build a personalized product recommendation zone that will be displayed on the home page.

How do you do it ?

# Product recommendations

Transform the problem to : What's the simplest recommandation you can show visitors?

- Top selling products over all or by categories on a time window
- It works for all visitors (with or without informations)
- Personalisation with tweak (remember categories visited)

This is a good baseline before doing ML based methods

# What next ?

# Machine Learning common issues !

# Machine Learning algorithms issues

[Data] Cold Start / Warm Start
[Data] Conceptual Drift / Covariate shift

[Model] No free lunch Theorem against you

[Data/Model] Reproduce
[Data/Model] Data Bias / Leak of information
[Data/Model] Cross validation and test validation

# Machine Learning project issues

[Model Selection] Usability (under core  business constraints)

[Model selection] Core business metrics vs ML metrics !

[Model evaluation] A/B testing

[Production] Automatise update mode

[Production] Monitoring (speed, failure, user feedback)

# Conclusion

❏ Talk with core business, as software engineers do, and build a common sense baseline !

❏ Fear over engineering and useless Machine Learning

❏ Keep learning :
  - Mooc, tutorials, books, articles, talks
  - Forum, blogs, twitter, people

# Machine Learning Craftsmanship

Kelvin MOUTET

Machine Learning Engineer @ Prevision.io