# How to manage data-related projects and not fail (too often)?

Jakub Nowacki

# `whoami`

Machine Learning Engineer
Educator

I can code, I do maths

@jsnowacki
https://www.linkedin.com/in/jakubnowacki/
https://github.com/jsnowacki

# Why do 87% of data science projects never make it into production?

# 85% of big data projects fail, but your developers can help yours succeed

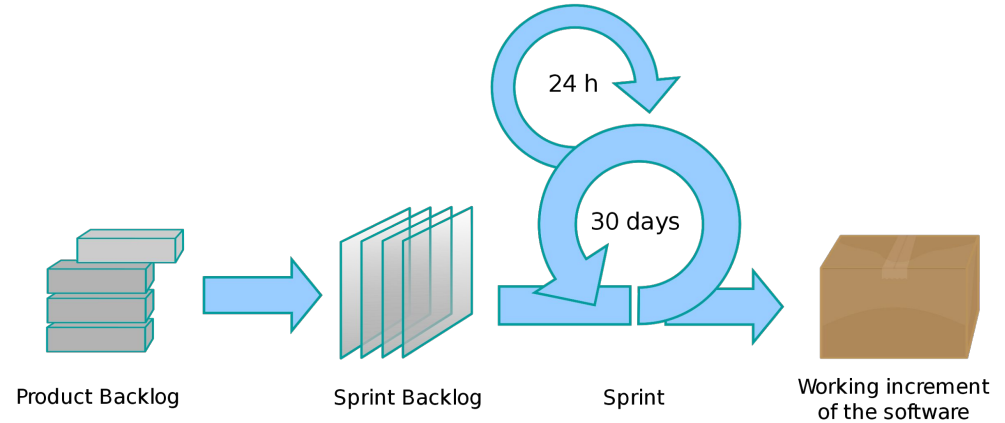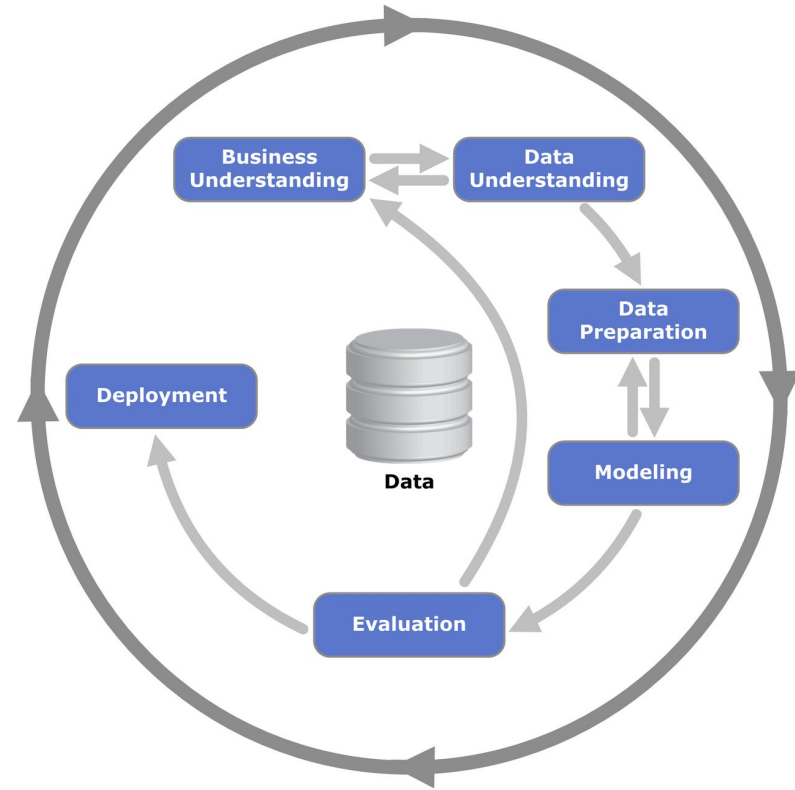# Looks familiar...?

- Why the process differ?
- Who are you anyway?
- What we cannot see?
- What are the costs?
- What to look out for?
- How the team differs?

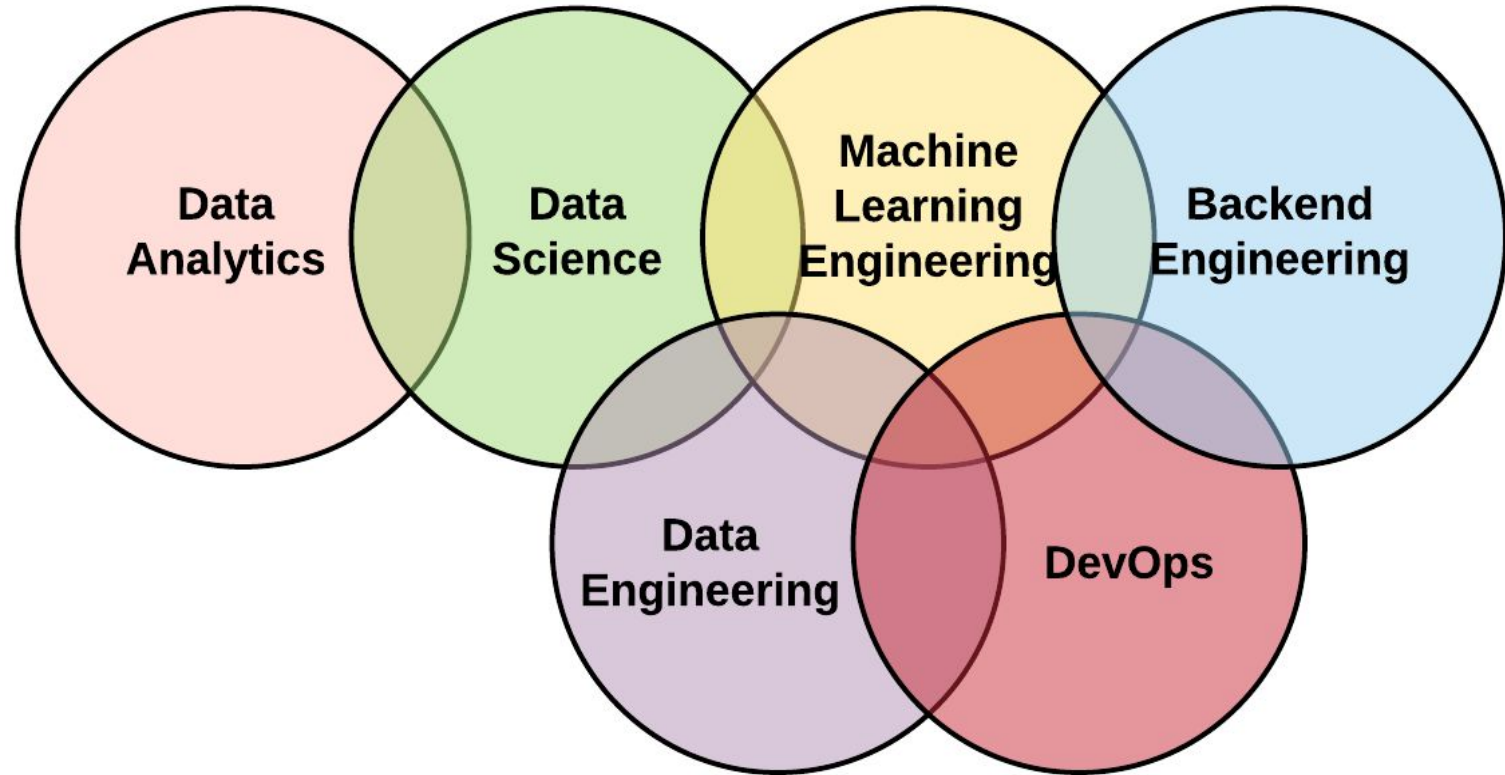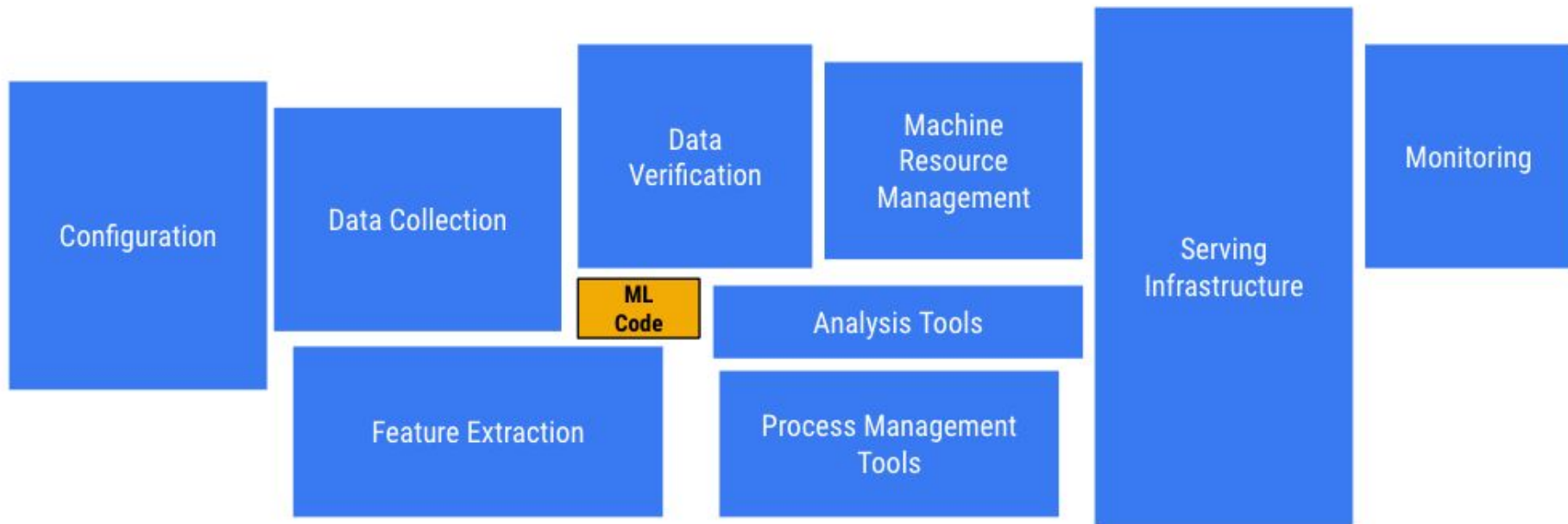# The process: CRISP DM vs SCRUM



https://en.wikipedia.org/wiki/Scrum_(software_development)

https://en.wikipedia.org/wiki/Cross-industry_standard_process_for_data_mining

# What is actually involved in data project?

# Hidden technical debt in ML systems

# DATA & AI LANDSCAPE 2019

## INFRASTRUCTURE

**HADOOP ON-PREMISE**
cloudera, Hortonworks, MAPR, Pivotal, IBM InfoSphere, jethro

**HADOOP IN THE CLOUD**
AWS, Microsoft Azure, Google Cloud, SAP Cloud Platform, IBM InfoSphere BigInsights, arm, Qubole, CAZENA

**STREAMING / IN-MEMORY**
Amazon Kinesis, databricks, SAP Cloud Platform, ORACLE, confluent, striim, hazelcast, GridGain, GIGASPACES, Wallaroo.labs, FASTDATA.io, kx

### NoSQL DATABASES
Google Cloud, AWS, ORACLE, Microsoft Azure, mongoDB, MarkLogic, Couchbase, DATASTAX, redislabs, AEROSPIKE, ArangoDB, SCYLLA

### NewSQL DATABASES
SAP, Clustrix, Pivotal, MemSQL, import.io, Cockroach Labs, VOLTDB, splice, imply, paradigm4, TiDB

### GRAPH DBs
neo4j, Amazon Neptune, Cambridge Semantics, Objectivity

### MPP DBs
TERADATA, VERTICA, IBM, ORACLE, Kognitio, Exasol, dremio, Yellowbrick

### CLOUD EDW
AWS, Google Cloud, Microsoft Azure, Pivotal, snowflake, Infoworks

### SERVERLESS
PULSAR, nuclio, Function Service

### DATA TRANSFORMATION
talend, pentaho, alteryx, TRIFACTA, tamr, Paxata, Infoworks, Fivetran, StreamSets, UNIFI

### DATA INTEGRATION
SAP Data Services, MuleSoft, TIBCO, imaplogic, enigma, Qlik Data Catalyst, Segment, ATTUNITY, ZALONI, import.io, SNOWPLOW, MATILLION

### DATA GOVERNANCE
IBM, SailPoint, McAfee Skyhigh Security Cloud, collibra, dynatrace, SignalFx, Alation, druva, unravel, OKERA, HMMUTA, MANTA, data.world, zenysis, OpsRamp, MAGNITUDE

### MGMT / MONITORING
AWS, New Relic, actifio, rubrik, APPDYNAMICS, WAVEFRONT, VMware, splunk, Numerify, ScienceLogic, datadog, DATOS IO, pagerduty

### STORAGE
AWS, Microsoft Azure, PURE STORAGE, qumulo, wasabi, Cohesity

### CLUSTER SVCS
Amazon, packet, elastic

### DATA GENERATION & LABELLING
Upwork, appen, Mighty AI, scale, HIVE, Labelbox, DEFINED.ai, fiddler

### AI OPS
ALGORITHMIA, SPELL, comet, Vertex.ai, PG Storm, FLOYDHUB

### GPU DBs & CLOUD
NVIDIA, kinetica, SQREAM, brytlyt, BLAZINGDB, MapD, Moviidus

### HARDWARE
Google TPU, arm, intel, GRAPHCORE, MYTHIC, CEREBRAS, CORNAMI, VATHYS

## ANALYTICS & MACHINE INTELLIGENCE

### DATA ANALYST PLATFORMS
Microsoft, Pentaho, alteryx, Digital Reasoning, GUAVUS, AYASDI, ATTIVIO, Datameer, incorta, inter.ana, sisu, Switchboard, Starburst

### DATA SCIENCE PLATFORMS
IBM, databricks, data iku, DOMINO, rapidminer, TIBCO, ANACONDA, SAS, Alteryx, KNIME, MathWorks

### BI PLATFORMS
looker, AWS, DOMO, ARCADIA DATA, ThoughtSpot, ATSCALE, SiSense, Qlik, GoodData, Information Builders, birst, MicroStrategy, Keen IO

### VISUALIZATION
tableau, Microsoft Power BI, SAP, Google Cloud, Periscope Data, zepl, GEOMATA, plotly, CHARTIO

### MACHINE LEARNING
Azure Machine Learning, Google Cloud, AWS, DataRobot, H2O.ai, gamalon, ViSENZE, ELEMENT, deepsense.ai

### COMPUTER VISION
Microsoft Azure, Amazon Rekognition, clarifai, EVER AI, neurala, twentybn, UBIQUITY, YITU, trax, synthesia, DataGrid

### HORIZONTAL AI
IBM Watson, Cortana, sentient, Voyager, Affectiva, Numenta, naralogics, CURIOUS AI, OSARO, SHIELD AI, cogito, snips

### SPEECH & NLP
Google Cloud, twilio, semantic, Eigen Technologies, PRIMER, cresta, neuro, fortia, Polyai

### SEARCH
elasticsearch, ORACLE, ENDECA, algolia, coveo, Lucidworks, ATTIVIO, swiftype, EXALEAD, alphasense, MAANA, omni:us, SINEQUA

### LOG ANALYTICS
splunk, sumologic, solarwinds, TIMBER, kibana, logz.io

### SOCIAL ANALYTICS
Hootsuite, sprinkr, NETBASE, synthesio, tracx, simplereach, bitly, SimilarWeb

### WEB / MOBILE / COMMERCE ANALYTICS
Google Analytics, mixpanel, AMPLITUDE, Airtable, RESCI, SIGOPT, granify, custora

## APPLICATIONS – ENTERPRISE

### SALES
CHORUS, INSIDESALES.COM, peopleai, conversica, clari, aviso, tact.ai, TROOPS, fuse machines, Clearbit

### MARKETING – B2B
RADIUS, App Annie, MINTIGO, Lattice, sense, tubular, RK, KNOTCH, mrp

### MARKETING – B2C
zeta, bloomreach, SendGrid, EVERSTRING, BLUECORE, mparticle, Amplero, amperity, QUANTIFIND, Engagio, Lytics, PERSADO, remesh

### CUSTOMER EXPERIENCE / SERVICE
qualtrics, MEDALLIA, SurveyMonkey, User Testing, CLARABRIDGE, zendesk, Kustomer, afiniti, Gainsight, pendo, HEAP, Amplitude, Watson Assistant, DigitalGenius, ASAPP, ada, AUTOMAT

### ENTERPRISE PRODUCTIVITY
slack, x.ai, ORACLE, GURU, lumiata, DIFFBOT, clara, talla, Kasisto

### HUMAN CAPITAL
Hire Vue, hiQ, pymetrics, mya, Allyo, textio, Wade&Wendy, Stella, entelo, RESTLESS BANDIT, beamery

### LEGAL
RAVEL, Everlaw, kira, JUDICATA, text IQ, Comply Advantage, IRONCLAD, ROSS, casetext

### REGTECH & COMPLIANCE
Seal, BigID, Zuora, FBREVIA, INTOSUM, DATA REPUBLIC

### FINANCE
Anaplan, SAP, TRADESHIFT, VIDADO, AppZen, WorkFusion, workato, ANTWORKS, ALKYMI

### BACK OFFICE AUTOMATION & RPA
UiPath, blueprism, Catalytic, KRYON

### SECURITY
TANIUM, CYLANCE, zscaler, StackPath, illumio, CODE42, CipherCloud, DARKTRACE, ANOMALI, VECTRA, DATAVISOR, sift science, Secureworks, SOCURE, Recorded Future, feedzai, Cybereason, SIGHT, bitglass, BLUEHEXAGON, Semmle, ARMORBLOX

## APPLICATIONS – INDUSTRY

### ADVERTISING
AppNexus, criteo, xAd, integral, ORACLE MOAT, theTradeDesk, dstillery, TAPAD, dataxu, gumgum, AppNexus

### EDUCATION
Udacity, KNEWTON, Clever, kidaptive, PANORAMA, gradescope

### REAL ESTATE
REDFIN, Opendoor, VTS, CREDIFI, GEOPHY, reonomy, SmartProcure, COMPSTAK, SPACEMAKER, SKYLINE, OpenSpace

### GOV'T
OPENGOV, mark43, Palantir, Dataminr, LiveStories, Quid, PRIMER, FORGE

### INTELLIGENCE
KENSHO, Quantopian, ADDEPAR, NUMERAI, iSENTIUM, ALGORIZ, 100Credit, aire, cignifi

### FINANCE – INVESTING
metromile, ROOT, zesty.ai

### FINANCE – LENDING
affirm, Kreditech, AVANT, TALA, BLEARBANC, upstart, WeLab, Wecash, MoneyLion, ognh

### INSURANCE
metromile, Lemonade, Hippo, Shift Technology, ROOT, zesty.ai, CAPE

### HEALTHCARE
flatiron, Clover, XYRUS, HealthTap, METABIOTA, Gingerio, Glow, babylon, 3Med, zebra, Flatiron, oviia, TEMPUS, AiCure, insitro, notable, RECURSION, citizen, prognos, Qventus, ATERYS, IMAGEN, PAIGE, DATAVANT, innaccer, LeanTaaS

### LIFE SCIENCES
color, verily, WuXiNextCODE, Clear Labs, nference, DNAnexus, CYTERA, BAYLABS, twoXAR, deep genomics, Atomwise

### TRANSPORTATION
UBER, TESLA, WAYMO, ZOOX, CLEARPATH, CRUISE, NURO, Aurora, nauto, AIMOTIVE, NIO, OPTIMUS, moovit, netradyne, Civic Maps, INRIX, Kodiak, comma.ai

### AGRICULTURE
FARMERS, ondeck, Granular, JOHN DEERE, BLUE RIVER, Farmers Edge, AgroStar, FarmLogs, TARANIS, GAMAYA, Terravion, prospera

### COMMERCE
Instacart, FAIRE, STITCH FIX, Retailnext, CELECT

### INDUSTRIAL
AVEVA, SIEMENS, PREDIX, UPTAKE, SCORTEX, TACHYUS

### OTHER
pharmacy, stem, Amper, ByteDance, Toppr, celect, SOJERN, Electric, ZINIER, Spoke

## CROSS-INFRASTRUCTURE/ANALYTICS
AWS, Google Cloud, Microsoft, IBM, SAP, Hewlett Packard Enterprise, SAS, 1010DATA, vmware, TIBCO, TERADATA, ORACLE, NetApp, syncsort, MAPR, cloudera
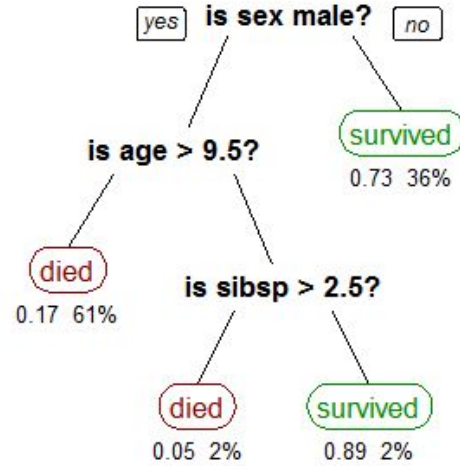
## OPEN SOURCE

### FRAMEWORKS
Spark, Flink, MESOS, CDAP, RedHat, HELIX

### QUERY / DATA FLOW
Spark SQL, HIVE, presto, APACHE DRILL, SLAMDATA, GraphQL, Flink

### DATA ACCESS & DATABASES
mongoDB, redis, cassandra, Zookeeper, Cockroach Labs, druid, CouchDB, riak, HBASE, Cloud Spanner, SciDB

### ORCHESTRATION & MGMT
talend, python, Apache Ambari, Apache Airflow, MESOS, etcd, Kong, Apache RocketMQ

### STREAMING & MESSAGING
Spark, nifi, Beam, Flink, kafka, STORM

### STAT TOOLS & LANGUAGES
Scala, SciPy, Julia

### AI OPS & INFRA
mlflow, Kubeflow, Studio, DVC, SELDON, Polyaxon

### AI / MACHINE LEARNING / DEEP LEARNING
TensorFlow, Microsoft Cognitive Toolkit, Caffe, OpenAI, DMTK, theano, Apache SINGA, DIMSUM, FeatureFu, mxnet, VELES, ONNX, Ludwig, PyTorch, neon, DSSTNE, DL4J, MAHOUT, Aerosolve, fast.ai, mlr

### SEARCH
elasticsearch, Solr, Lucene

### LOGGING & MONITORING
elasticsearch, kibana, SENTRY, logstash, Prometheus, fluentd, fluentd, Grafana

### VISUALIZATION
matplotlib, TensorBoard, seaborn, Bokeh

### COLLABORATION
BeakerX, Jupyter, ANACONDA

### SECURITY
Apache Ranger, KNOX, Sentry, accumulo

## DATA SOURCES & APIs

### HEALTH
VALIDIC, practicefusion, fitbit, GARMIN, HUMAN API, kinsa, MIMO

### IOT
GE Digital, UPTAKE, thingworx, samsara, helium

### FINANCIAL & ECONOMIC DATA
Bloomberg, THOMSON REUTERS, DOW JONES, S&P CAPITAL IQ, CB INSIGHTS, PLAID, SECOND MEASURE, estimize, PREMISE, Quandl, Eagle Alpha, StockTwits, xignite, Thinknum, earnest, predata

### AIR / SPACE / SEA
Orbital Insight, GE Digital, AIRBOTICS, spire, kespry, PRECISIONHAWK, pitney bowes, DroneDeploy, WINDWARD, MarineTraffic, RS Metrics

### PEOPLE / ENTITIES
acxiom, experian, EPSILON, InsideView, Crimson Hexagon, tellus labs, BASIS, Quantcast, SAFEGRAPH

### LOCATION INTELLIGENCE
FOURSQUARE, mapAnything, mapbox, HEXAGON, sense360, PlaceIQ, esri, factual, CARTO, Mapillary, StreetMap, cuebiq, Radar, OpenStreetMap

### OTHER
DATA.GOV, IMAGENET, LabelMe, CRUX, id:grafitti.io

## DATA RESOURCES

### DATA SERVICES
OPERA, GENERAL ASSEMBLY, DATA SCIENCE, PLURALSIGHT, galvanize, DataCamp, DataElite, fractal, kaggle, INSIGHT, DataKind, The Data Incubator, EXL, INNOPLEXUS

### INCUBATORS & SCHOOLS
GA, galvanize, DataCamp, DataElite, kaggle, DataKind, METIS

### RESEARCH
OpenAI, facebook research, MIRI, MILA, VECTOR INSTITUTE, CSAIL, AI2, ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

---

**FIRSTMARK** — EARLY STAGE VENTURE CAPITAL

# Man vs machine: let's automate!


https://en.wikipedia.org/wiki/Factory



is sex male?  yes / no

is age > 9.5?

survived
0.73 36%

died
0.17 61%

is sibsp > 2.5?

died
0.05 2%

survived
0.89 2%

https://en.wikipedia.org/wiki/Decision_tree_learning

# It ain't cheap!

"We train XLNet-Large on 512 TPU v3 chips for 500K steps with an Adam optimizer, linear learning rate decay and a batch size of 2048, which takes about 2.5 days."

https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/

# There is no data!



Image classification
Easiest classes
red fox (100)  hen-of-the-woods (100)  ibex (100)  goldfinch (100)  flat-coated retriever (100)
tiger (100)  hamster (100)  porcupine (100)  stingray (100)  Blenheim spaniel (100)

Hardest classes
muzzle (71)  hatchet (68)  water bottle (68)  velvet (68)  loupe (66)
hook (66)  spotlight (66)  ladle (65)  restaurant (64)  letter opener (59)
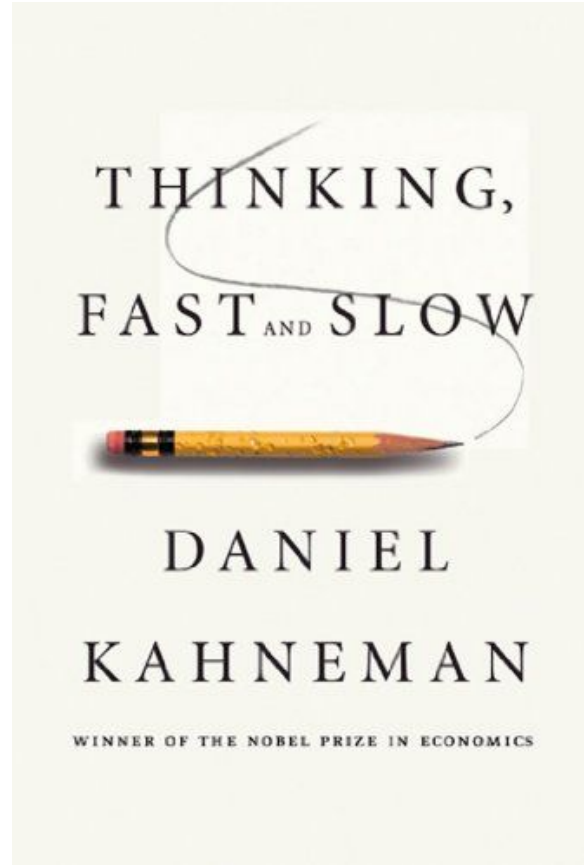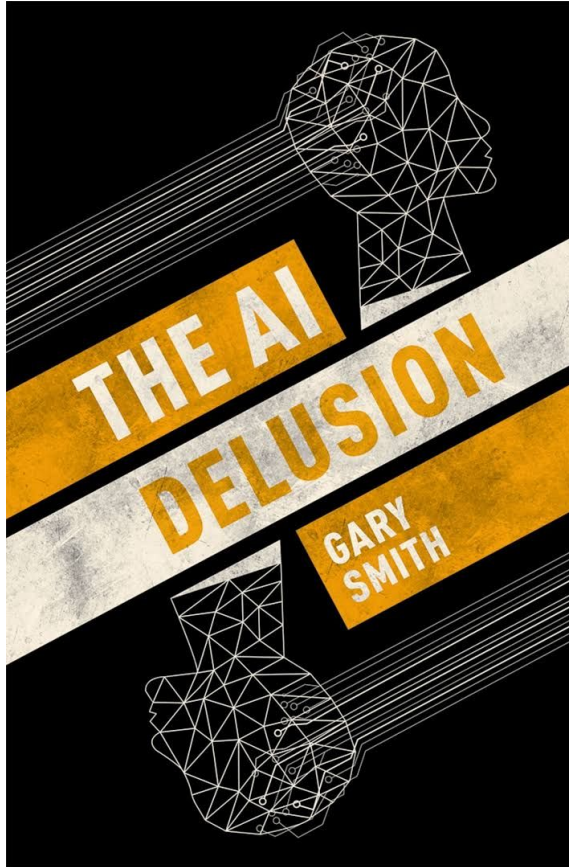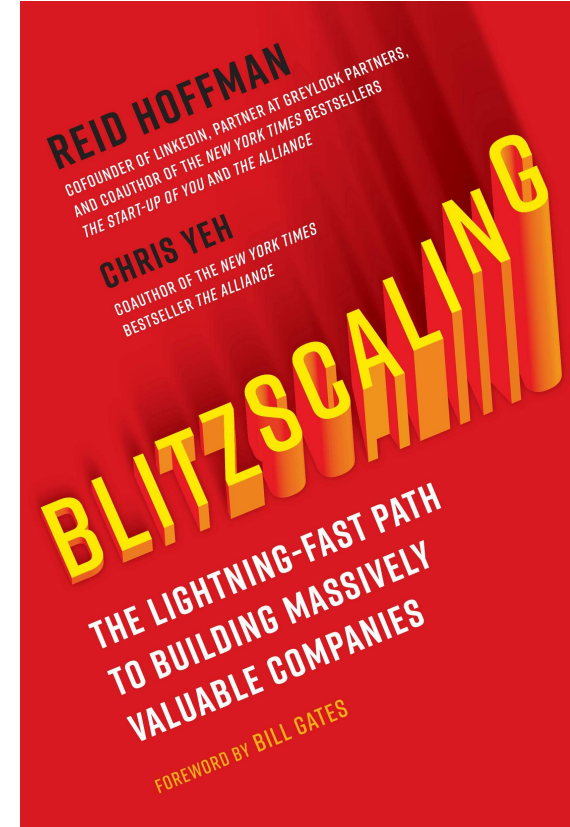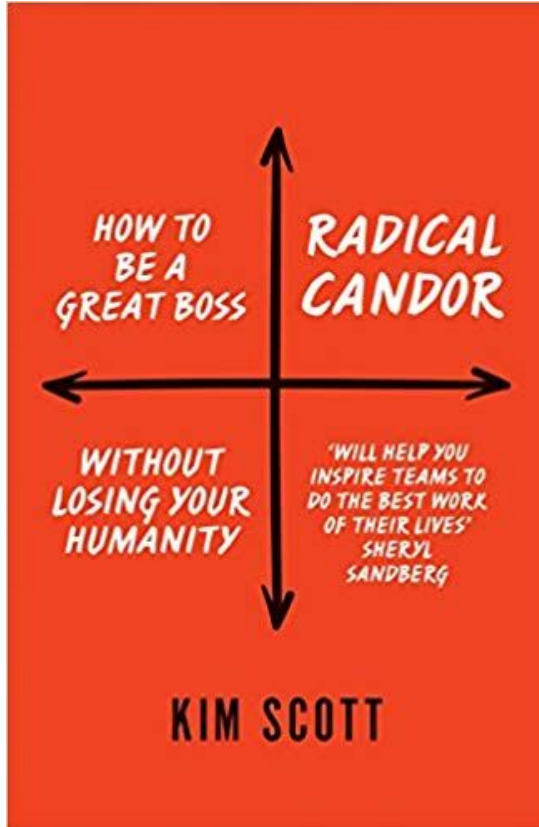


WIKIPEDIA
The Free Encyclopedia

Try it first!

# Be careful!

# Different type of teams

# Thank you!

Questions?