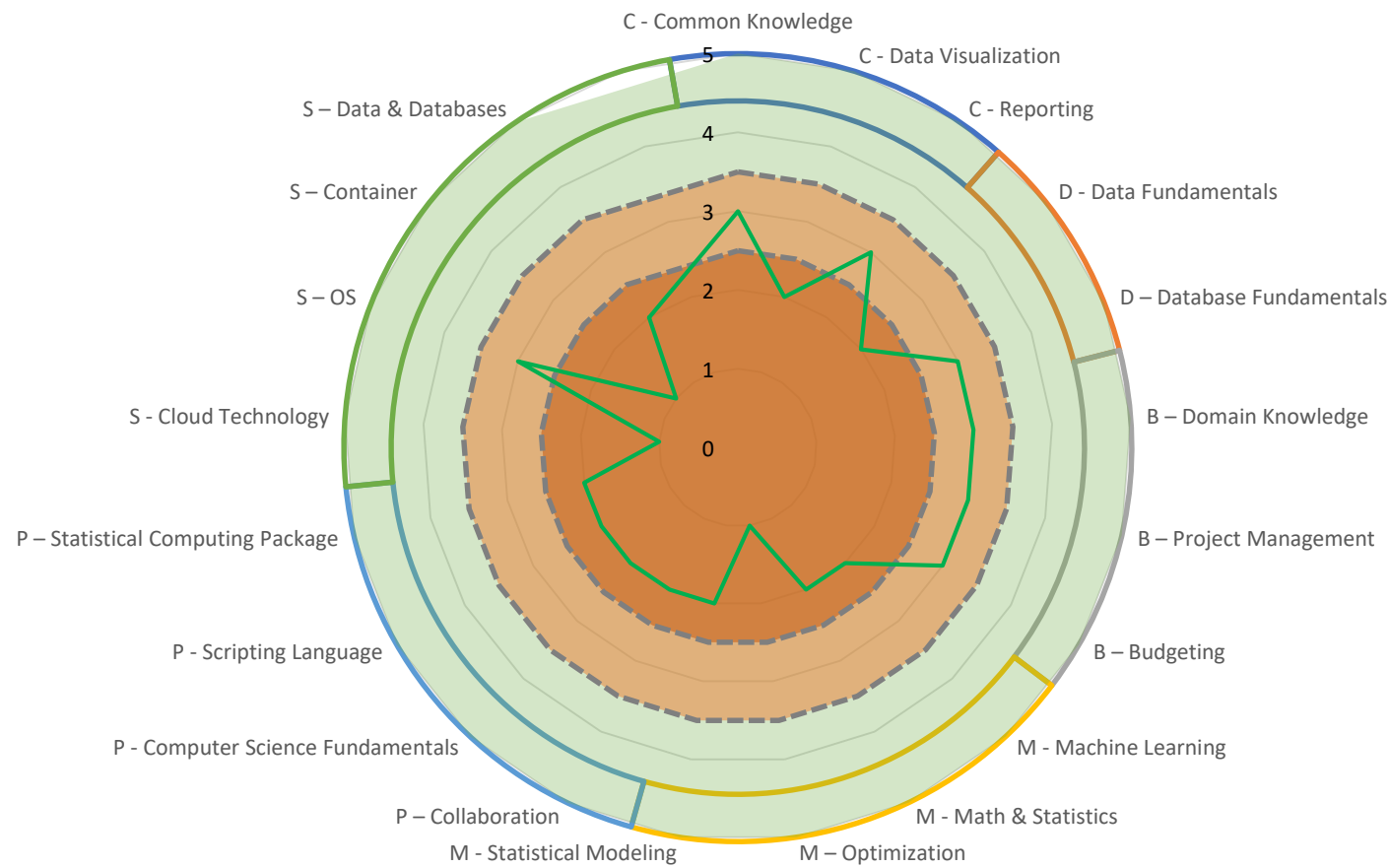


Data

DATA SCIENCE SKILLS

Uli



Data Types



- Data Types
- Data Structure
- Data Set
- Data Base

Types of Data

- Classification of Data
 - Boolean
 - Numeric
 - Integer
 - Floating
 - Composite (derived from more than one type)
 - Array
 - Record

Algebraic functions

- `>>> 2 + 35`
- `>>> 7 - 52`
- `>>> 2*(3+1)8`
- `>>> 2 + 35`
- `>>> 7 - 52`
- `>>> 2*(3+1)8`
- `>>> 5/22.5`
- `>>> 2 + 35`
- `>>> 7 - 52`
- `>>> 2*(3+1)8`

Boolean functions

- Python can evaluate Boolean expressions
Boolean expressions evaluate to True or False
Boolean expressions often involve comparison operators <, >, ==, !=, <=, and >=

Excercises

- Translate the following into Python algebraic or Boolean expressions and then evaluate them:
- The difference between Annie's age (25) and Ellie's (21)
- The total of \$14.99, \$27.95, and \$19.83
- The area of a rectangle of length 20 and width 15
- 2 to the 10th power
- The minimum of 3, 1, 8, -2, 5, -3, and 0
- 3 equals 4-2
- The value of $17//5$ is 3
- The value of $17\%5$ is 3
- 284 is even284 is even and 284 is divisible by 3
- 284 is even or 284 is divisible by 3

Excercises

- Translate the following into Python algebraic or Boolean expressions and then evaluate them:
- The difference between Annie's age (25) and Ellie's (21)
- The total of \$14.99, \$27.95, and \$19.83
- The area of a rectangle of length 20 and width 15
- 2 to the 10th power
- The minimum of 3, 1, 8, -2, 5, -3, and 0
- 3 equals 4-2
- The value of $17//5$ is 3
- The value of $17\%5$ is 3
- 284 is even284 is even and 284 is divisible by 3
- 284 is even or 284 is divisible by 3

Variables

- Just as in algebra, a value can be assigned to a variable, such as x . When variable x appears inside an expression, it evaluates to its assigned value
- $x = 3$

Variables

- (Variable) names can contain these characters: a through z
- A through Z
- the underscore character _
- digits 0 through 9

Strings

In addition to number and Boolean values, Python support string values

'Hello, World!'

"Hello, World!"

A string value is represented as a sequence of characters enclosed within quotes

A string value can be assigned to a variable

String values can be manipulated using string operators and functions

Exercises

```
>>> s1 'good'>>> s2 'bad'>>> s3 'silly'
```

Write Python expressions involving strings s1, s2, and s3 that correspond to:

- 'll' appears in s3
- the blank space does not appear in s1
- the concatenation of s1, s2, and s3
- the blank space appears in the concatenation of s1, s2, and s3
- the concatenation of 10 copies of s3
- the total number of characters in the concatenation of s1, s2, and s3

- Data is any information you are collecting: numbers, statistics, measurements. It can also be words, observations, or other inputs.

	Continuous	Discrete		
	Quantitative data	Qualitative / Categorical / Attribute data		
Measurement	Units (example)	Ordinal (example)	Nominal (example)	Binary (example)
Time of day	Hours, minutes, seconds	1, 2, 3, etc.	N/A	a.m./p.m.
Date				
Cycle time				
Speed				
Brightness				
Temperature				
<Count data>				
Test scores				
Defects				
Defects				
Color				
Location				
Groups				
Anything				

Continuous Data

- has an infinite number of measurements depending on the resolution of the measurement system.
- There are no limits to the gaps between the measurements. It is data that can be expressed on an infinitely divisible scale.
- Examples:
 - Temperature
 - Height
 - ... ?

Discrete Data

- Data types that have a finite number of measurements and are based on counts. Data that can be sorted into distinct, countable, and in completely separate categories. The count value can not be divided further on an infinite scale with meaning.
- Example:
 - How many people can comfortably fit into an airplane? It doesn't make sense to say 129.7632213 people. It is either 129 or 130, in this case you would round down to 129. Attribute and discrete do not mean exactly the same when describing data, discrete has more than two outcomes.
 - ... ?

Nominal Data

- The **lowest level** of data classification. A numerical label that represents a qualitative description. These numbers are labels or assignments of numbers that represent a category or classification. This is also referred to as categorical data usually of more than two categories and is a form of discrete data and should apply nonparametric test to analyze. The number assignment does not reflect that one category is better or worse than another.
- Example:
 - Gender
 - 1 = Male
 - 2 = Female

Ordinal Data

- The next level higher of data classification than nominal data. Numerical data where number is assigned to represent a qualitative description similar to nominal data. These are measures by only the rank order.
- However, these numbers can be arranged to represent worst to best or vice-versa. Ordinal data is a form of discrete data and should apply non-parametric test to analyze.
- Example:
 - Classifying households as low income, middle-income, and high income
-

Interval Data

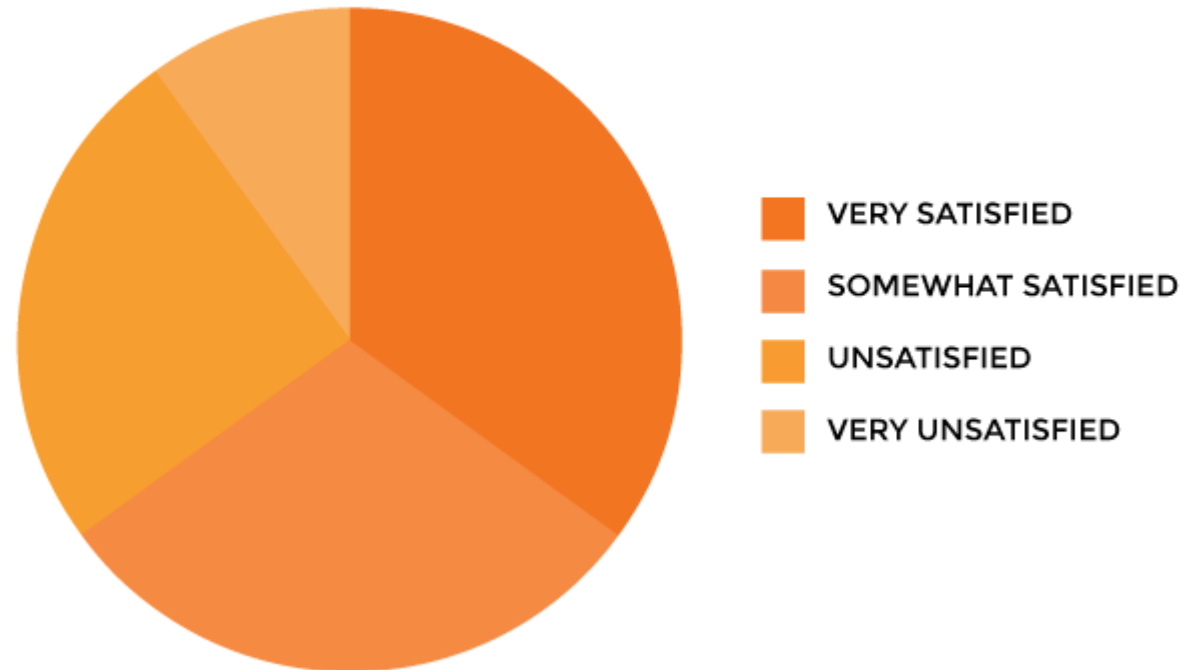
- The next higher level of data classification. Numerical data where the data can be arranged in a order and the differences between the values are meaningful but not necessarily a zero point. These are measures using equal intervals.
- Interval data can be both continuous and discrete. Zero degrees Fahrenheit does not mean it is the lowest point on the scale, it is just another point on the scale.
- Example:
 - Temperature
 - ...
-

Exercise

	Continuous	Discrete		
	Quantitative data	Qualitative / Categorical / Attribute data		
Measurement	Units (example)	Ordinal (example)	Nominal (example)	Binary (example)
Time of day	Hours, minutes, seconds	1, 2, 3, etc.	N/A	a.m./p.m.
Date	Month, date, year	Jan., Feb., Mar., etc.	N/A	Before / After
Cycle time	Hours, minutes, seconds, month, date, year	10, 20, 30, etc.	N/A	Before / After
Speed	Miles per hour/centimeters per second	10, 20, 30, etc.	N/A	Fast / Slow
Brightness	Lumens	Light, medium, dark	N/A	On / Off
Temperature	Degrees C or F	10, 20, 30, etc.	N/A	Hot / Cold
<Count data>	Number of things	10, 20, 30, etc.	N/A	Large / Small
Test scores	Percent, number correct	F, D, C, B, A	N/A	Pass / Fail
Defects	N/A	Number of cracks	N/A	Good / Bad
Defects	N/A	N/A	Oversized, missing	Good / Bad
Color	N/A	N/A	Red, blue, green	N/A
Location	N/A	N/A	East, West, South	Domestic / International
Groups	N/A	N/A	HR, Legal, IT	Exempt / Non-exempt
Anything	Percent	10, 20, 30, etc.	N/A	Above / Below

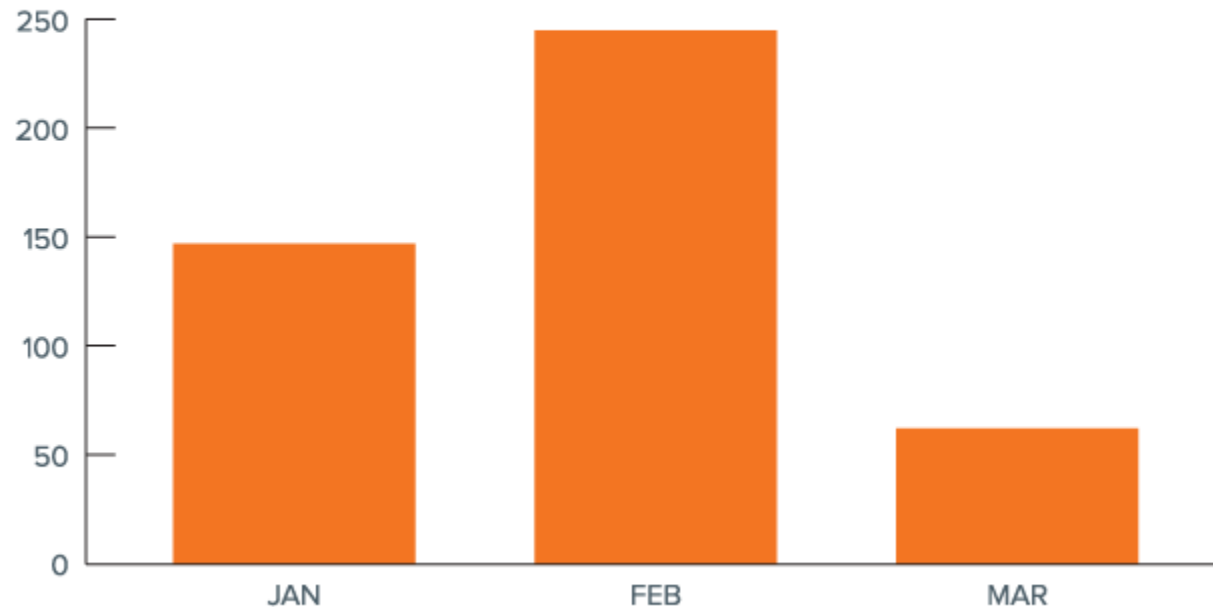
- **Cross-Sectional:** The sample of elements is measured only once. This shows you a snapshot of variables at a point in time (e.g., market survey).

CUSTOMER SATISFACTION



- **Longitudinal:** The data sample is measured repeatedly over time (e.g., stock prices, monthly sales data).

PAGE VIEWS, BY MONTH



What Makes a Data Set?

- A data set is comprised of variables; each individual data point—the thing that is measured or counted—is a **variable**. Each variable can be examined on its own or in relation to other variables to reveal insights, including:
- Mean:



What Makes a Data Set?

- Range:



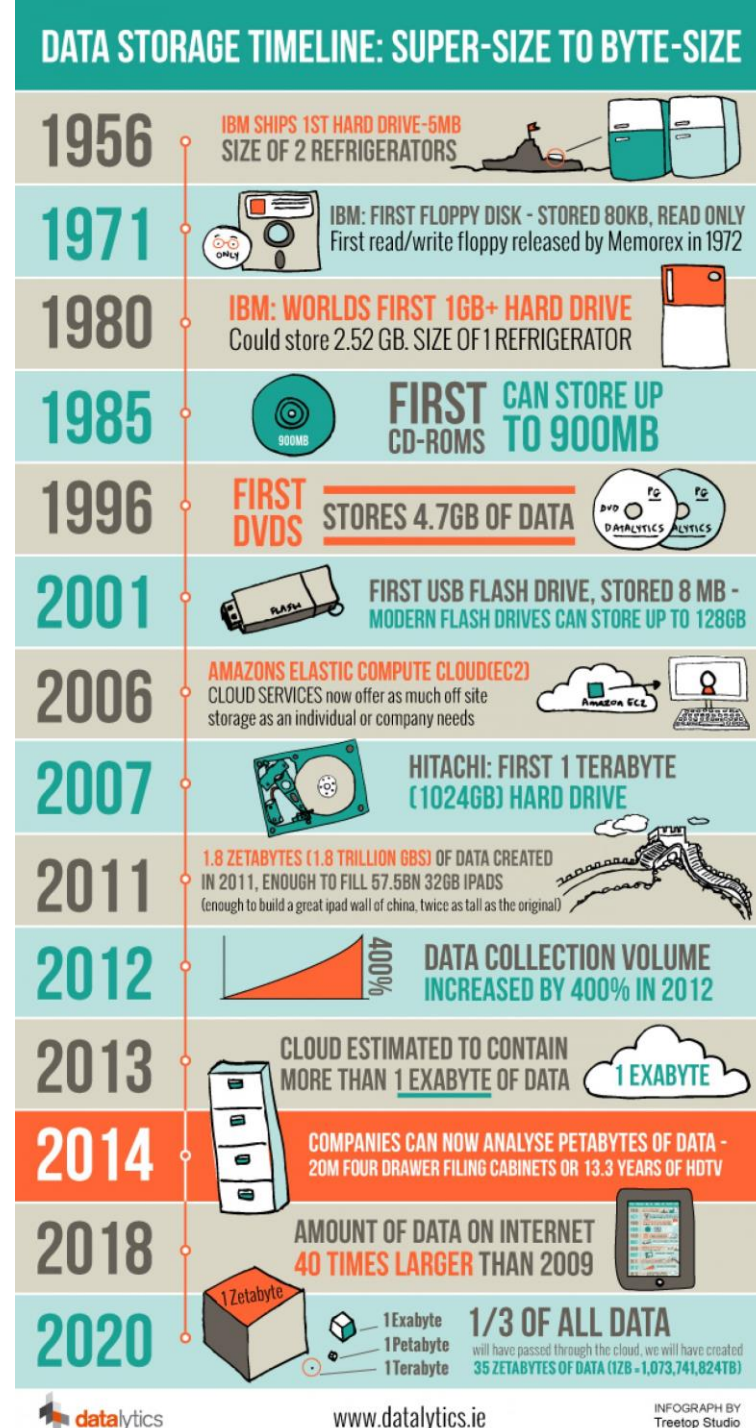
What Makes a Data Set?

- Range:



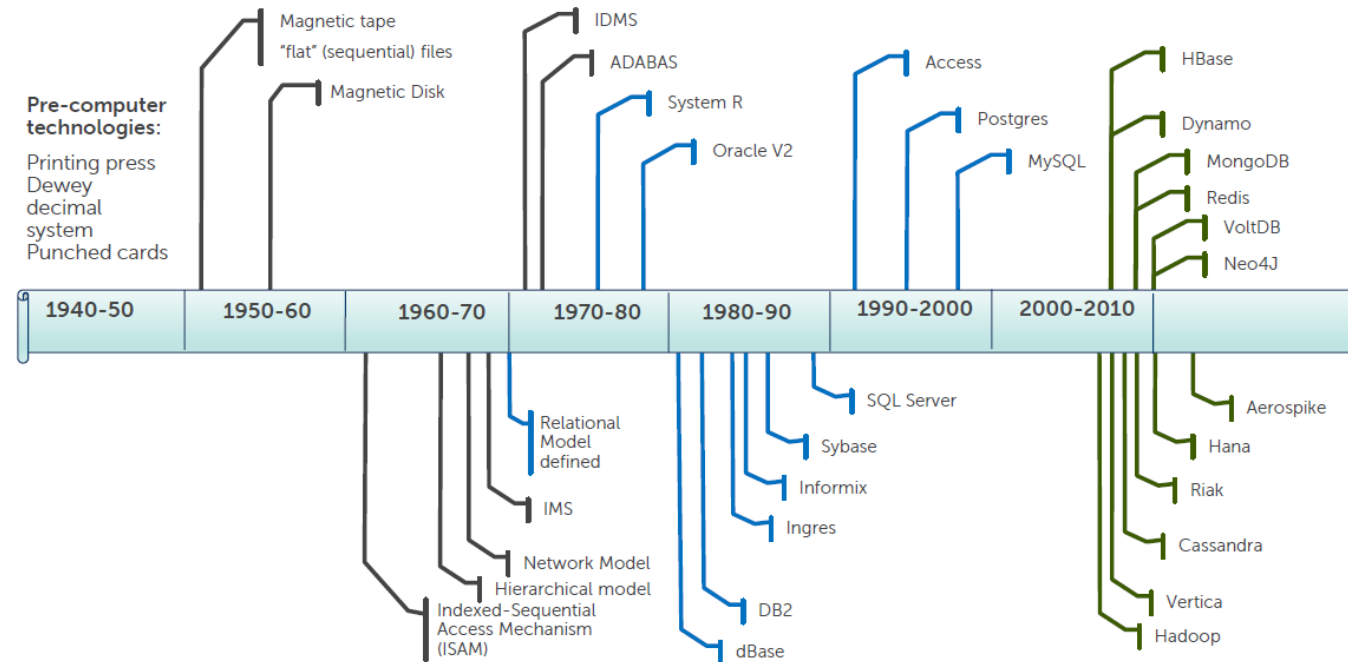
Data Storage

- Files
 - CSV
 - XML
 - JSON
- Place
 - HDD
 - Memory
 - Locally
 - Cloud
 - Database



History of databases

History of databases



BIG DATA LANDSCAPE 2017

INFRASTRUCTURE

ANALYTICS

APPLICATIONS – ENTERPRISE

The diagram illustrates the Hadoop ecosystem components, categorized into three main groups:

- HADOOP ON-PREMISE:** Includes Cloudera, Hortonworks, MapR, Pivotal, IBM InfoSphere, bluedata, and jethro.
- HADOOP IN THE CLOUD:** Includes Amazon Web Services, Microsoft Azure, Google Cloud Platform, IBM InfoSphere BigInsights, TREASURE DATA, altilscale, Doble, CAZENA, and CenturyLink.
- STREAMING / IN-MEMORY:** Includes Amazon Web Services, databricks, confluent, strim, GridGain, METAMARKETS, DATATORRENT, dataArtisans, ORACLE, hazelcast, and TERADATA.

The diagram is divided into two main sections by a horizontal line. The left section is titled 'DATA ANALYST PLATFORMS' and the right section is titled 'DATA SCIENCE PLATFORMS'. Each section contains a grid of logos for various platforms.

DATA ANALYST PLATFORMS:

- Microsoft
- pentaho
- alteryx
- Digital Reasoning
- guavus
- AYASDI
- MATTIVO
- Datameer
- Quid
- ClearStory
- OriginalLogic
- interana
- Bottlenose
- ARIMO
- ENDOR
- MODE

DATA SCIENCE PLATFORMS:

- IBM
- KNIME
- dataiku
- DOMINO
- yhat
- CONTINUUM ANALYTICS
- rapidminer
- ALGORITHMIA
- Alpine Analytics
- Anqoss

The collage displays logos for various database technologies, organized into five main categories:

- NO SQL DATABASES:** Includes Google Cloud Platform, Oracle, Amazon DynamoDB, Microsoft Azure, MarkLogic, mongoDB, DATASTAX, KERO SPIKE, Couchbase, redislabs, and influxdata.
- NEWSQL DATABASES:** Includes SAP HANA, Clustring, Pivotal, nuodb, Cockroach LABS, memsql, spice, MariaDB, YOLDB, citusdata, Trafalgar, and paradigim4.
- GRAPH DBs:** Includes neo4j, IBM ORACLE, OrientDB, InfiniteGraph, and OrientDB.
- MPP DBs:** Includes TERADATA, VERTICA, Netezza, Action, kognitio, EXSOL, and dremio.
- CLOUD EDW:** Includes Amazon Redshift, Google Cloud Platform, Microsoft Azure, Pivotal, Snowflake, and Infoworks.

The collage is organized into six vertical columns, each representing a different category of data science tools:

- BI PLATFORMS:** Includes logos for Microsoft, Tableau, SAP, Amazon Web Services, Oracle, Qlik, Celonis, Looker, Wave Analytics, Alteryx, Arcadia Data, GoodData, and Sisense.
- VISUALIZATION:** Includes logos for Google Cloud Platform, Microsoft, Olik, Celonis, Periscope, Chartio, and Geopointa.
- VERTICAL ANALYTICS:** Includes logos for Predix, Capex, Uptake, Oracle, Tachyus, Allium, and Datorama.
- STATISTICAL COMPUTING:** Includes logos for SAS, SPSS, and MATLAB.
- DATA SERVICES:** Includes logos for Palantir, IBM, and Data Science.
- DATA SCIENCE:** Includes logos for Kaggle, Excel, DataKind, and FICO.

HUMAN CAPITAL

- HereVue
- entelo
- hiQ
- DIGSTER
- textio
- RESTLESS BANDIT
- Wade&Wendy
- Clustrio
- Stella
- pymetrics

LEGAL

- RAVEL
- Seal
- Everlaw
- JUDICATA
- Brevia
- RESS
- casevity

FINANCE

- anaplan
- Worx
- tidemark
- 500 S4 HANA
- TRADESHIFT
- lumatra
- diffbot
- Clara Talita
- butter ai
- kasisto

ENTERPRISE PRODUCTIVITY

- slack
- facebook
- ORACLE
- AppZen

BACK OFFICE AUTOMATION

- HyperScience
- capricity

SECURITY

- TANIUM
- CYCLANCE
- StackPath
- DARKTRACE
- illumio
- CODE42
- ThreatMetrix
- DataGravity
- VECTRA
- CipherCloud
- Guardian
- cyberzen
- ANOMALI
- siftscience
- SINEIFY
- OneOneOne
- SecurityScorecard
- BlueTalon
- Recorded Future
- feedai
- ASAFI
- FortiScale
- Yumenu
- ASAFI

The diagram is divided into four colored rectangular sections, each representing a different category of data integration tools. Each section has a title in bold, uppercase letters and a list of companies with their logos.

- DATA TRANSFORMATION** (Orange background):
 - talend
 - pentaho
 - alteryx
 - TRIFACTA
 - tamr
 - Paxata
 - StreamSets
 - UNIFI
- DATA INTEGRATION** (Yellow background):
 - informatica
 - snoplogic
 - Segment
 - TEALUM
 - enigma
 - alooma
 - podium data
 - ZALONI
 - xpienty
 - import
 - Stitch
- DATA GOVERNANCE** (Green background):
 - informatica
 - IBM
 - skyhigh
 - collibra
 - Alation
 - Waterline Data
- MGMT / MONITORING** (Blue background):
 - amazon
 - New Relic
 - APPROVIMATICS
 - oetfio
 - WAVEFRONT
 - unravel
 - DATADOG
 - splunk
 - ooc
 - Avrocast
 - procon
 - pagerduty
 - Numerify

MACHINE LEARNING

- Google Cloud Platform
- H₂O
- DataRobot
- VISEN
- bonsai
- deepsense
- nitorian

HORIZONTAL AI

- IBM Watson Cortana
- Face++ 旷视 sentient
- Voyager
- Affective Cognitive Scale
- Groconcom
- PETRUM
- SARO
- CURIOUS AI
- BLUTE VISION

SPEECH & NLP

- Google Cloud Platform
- Narrative Science
- semanticrhythms
- ARRIA
- iBIBO
- Talkia
- cortical.io
- snips
- yseon
- Gradspace Soundhound

The image displays a collection of logos for various cloud and data services, organized into five columns:


- STORAGE:** Includes logos for Amazon Web Services, Google Cloud Platform, Microsoft Azure, Alluxio, nimbustorage, Cumulo, and panadas.
- CLUSTER SERVICES:** Includes logos for Amazon EMR, Eucalyptus, Kubernetes, Docker, Mesosphere, CoreOS, and Pepperdata.
- APP DEV:** Includes logos for Lightbend, Amazon Mechanical Turk, Upwork, WorkFusion, Rainforest, and CRSK.
- CROWDSOURCING:** Includes logos for Amazon Mechanical Turk, Upwork, WorkFusion, and CrowdPower.
- HARDWARE:** Includes logos for Google TPU, ARM, Nervana Systems, Graphcore, Mythic, Corsair, NVIDIA, and Movidius.








At the bottom right, the SCORTEX logo is also visible.

The diagram is organized into four columns, each representing a different type of analytics tool:






- SEARCH**
 - Autonomy
 - EXALENGE
 - ThoughtSpot
 - Lucidworks
 - swiftype
 - alphasense
 - Searchkick
 - SINEQUA
- LOG ANALYTICS**
 - splunk
 - sumologic
 - loggly
 - kibana
 - logz.io
- SOCIAL ANALYTICS**
 - Hootsuite
 - NETBASE
 - DATASIFT
 - synthesio
 - simplereach
 - bitly
 - predata
- WEB / MOBILE / COMMERCE ANALYTICS**
 - Google Analytics
 - mixpanel
 - sumall
 - retention
 - granify
 - AMPLITUDE
 - Airtable
 - SIGOPT
 - custora






































































































































CROSS-INFRASTRUCTURE/ANALYTICS

amazon web services™ Google Cloud Platform Microsoft IBM SAP Hewlett Packard Enterprise sas 1010 data vmware TIBCO TERADATA ORACLE NetApp

OPEN SOURCE

FRAMEWORK	QUERY / DATA FLOW	DATA ACCESS	COORDINATION	STREAMING	STAT TOOLS	AI / MACHINE LEARNING / DEEP LEARNING	SEARCH	LOG ANALYSIS	VISUALIZATION	COLLABORATION	SECURITY
 Hadoop  MapReduce  Flink  YARN  Mesos  Spark  CDAP	 Spark SQL  Presto  SLiM DATA  Google Cloud Dataflow	 Cassandra  MongoDB  CouchDB  HBase  Spanner  Accumulo	 Talend  Apache Zookeeper  Apache Ambari	 Spark  Flink  Kafka  Storm	 Python  ScalaLab  NumPy  SciPy	 TensorFlow  Apache MXNet  neon  DSSNE  Caffe  CNTK  FeatureFu  DL4J  Theano  MxNet  Keras  PyTorch  Vespa  Weka  Chainer  Aerosolve	 Elasticsearch  Solr  Lucene	 Elasticsearch  Kibana  Logstash	 Beaker  Rodeo	 Jupyter  ANACONDA	 Apache Ranger  KNOX  Sentry

DATA SOURCES & APIS

HEALTH

Apple, Jawbone, Validic, Practicefusion, Fitbit, Garmin, Human API, Kinsa

IOT

GE Digital, Uptake, ThingWorx, Helium, Samsara, Eagle Alpha, Radian6

FINANCIAL & ECONOMIC DATA

Bloomberg, Thomson Reuters, Dow Jones, S&P Capital IQ, CB Insights, Xignite, Quandl, Yodlee, Premise, Estimize, Second Measure, StockTwits, Plaid, Mattermark, Thinknum

AIR / SPACE / SEA

Planet Labs, Airware, Spire, SkyCatcher, AeroBotics, Inmarsat, Inmarsat, Tellus Labs, Inmarsat, DroneDeploy, MarineTraffic

PEOPLE / ENTITIES

Axiom, Experian, Epsilon, InsideView, Crimson Hexagon, Basis, Quantcast, SafeGraph

LOCATION INTELLIGENCE

Foursquare, Sense, PlaceIQ, Esri, Factual, Carto, Mapillary, StreetLine

OTHER

Qualtrics, Data.gov, Data.world, Panjiva, Enigma

DATA RESOURCES

INCUBATORS & SCHOOLS

- PLURALSIGHT
- GA
- galvanize
- DataCamp
- DataElite
- INsIGHT
- The Data Incubator
- METIS

RESEARCH

- facebook research
- OpenAI
- MIRI
- CSAIL
- DFK
- CI
- ALLEN INSTITUTE FOR ARTIFICIAL INTELLIGENCE