# Arquitectura Big Data
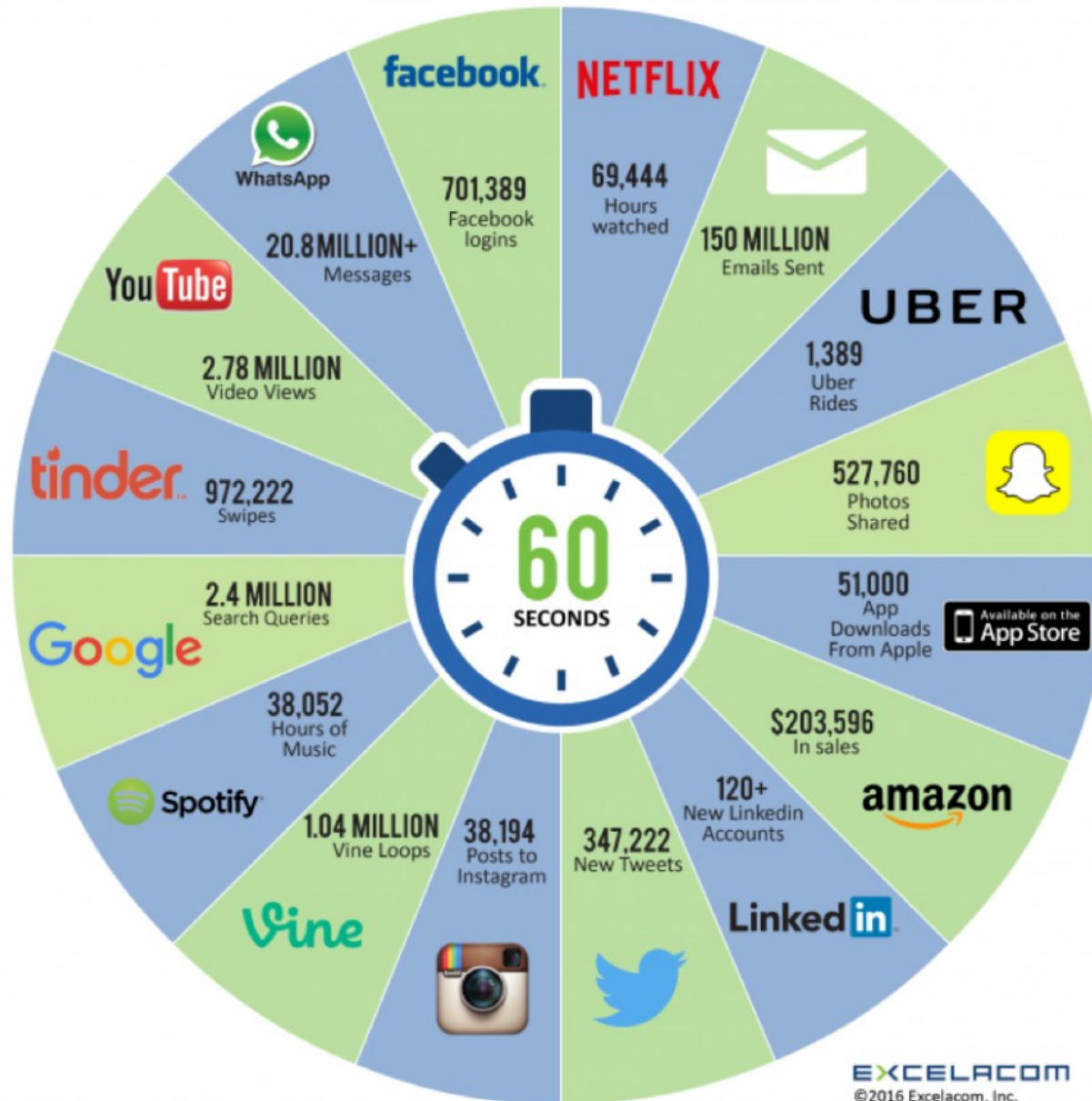
Angela Devia

ing.angela.devia@gmail.com

Científicas de Datos

www.cientificasdedatos.com

# 2016
## What happens in an INTERNET MINUTE?

**NETFLIX**
69,444 Hours watched

150 MILLION Emails Sent

**facebook**
701,389 Facebook logins

**WhatsApp**
20.8 MILLION+ Messages

**UBER**
1,389 Uber Rides

527,760 Photos Shared

**YouTube**
2.78 MILLION Video Views

**tinder**
972,222 Swipes

51,000 App Downloads From Apple
Available on the App Store

## 60 SECONDS

2.4 MILLION Search Queries
**Google**

$203,596 In sales
**amazon**

38,052 Hours of Music
**Spotify**

120+ New Linkedin Accounts
**Linked in**

1.04 MILLION Vine Loops
**Vine**

38,194 Posts to Instagram

347,222 New Tweets

# The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

## Volume
### SCALE OF DATA

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

2005
2020

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

## Velocity
### ANALYSIS OF STREAMING DATA

The New York Stock Exchange captures
**1 TB OF TRADE INFORMATION**
during each trading session

Modern cars have close to
**100 SENSORS**
that monitor items such as fuel level and tire pressure

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
– almost 2.5 connections per person on earth

## Variety
### DIFFERENT FORMS OF DATA

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

By 2014, it's anticipated there will be
**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**
are watched on YouTube each month

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

## Veracity
### UNCERTAINTY OF DATA

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

**27% OF RESPONDENTS**
in one survey were unsure of how much of their data was inaccurate

IBM

# THE FLOOD OF BIG DATA

DRIVING MARKETING EFFECTIVENESS BY MANAGING

**FACEBOOK**

**EMAIL**

## BIG DATA = BIG OPPORTUNITY

**35** ZETTABYTES OF DATA GENERATED ANNUALLY BY 2020[7]

**60%** GROWTH IN STRUCTURED AND UNSTRUCTURED DATA ANNUALLY[8]

**80% GROWTH** IN UNSTRUCTURED DATA[10]

**2.7** ZETTABYTES OF DATA EXIST IN THE DIGITAL UNIVERSE[9]

**5** EXABYTES OF DATA GENERATED EVERY TWO DAYS[11]

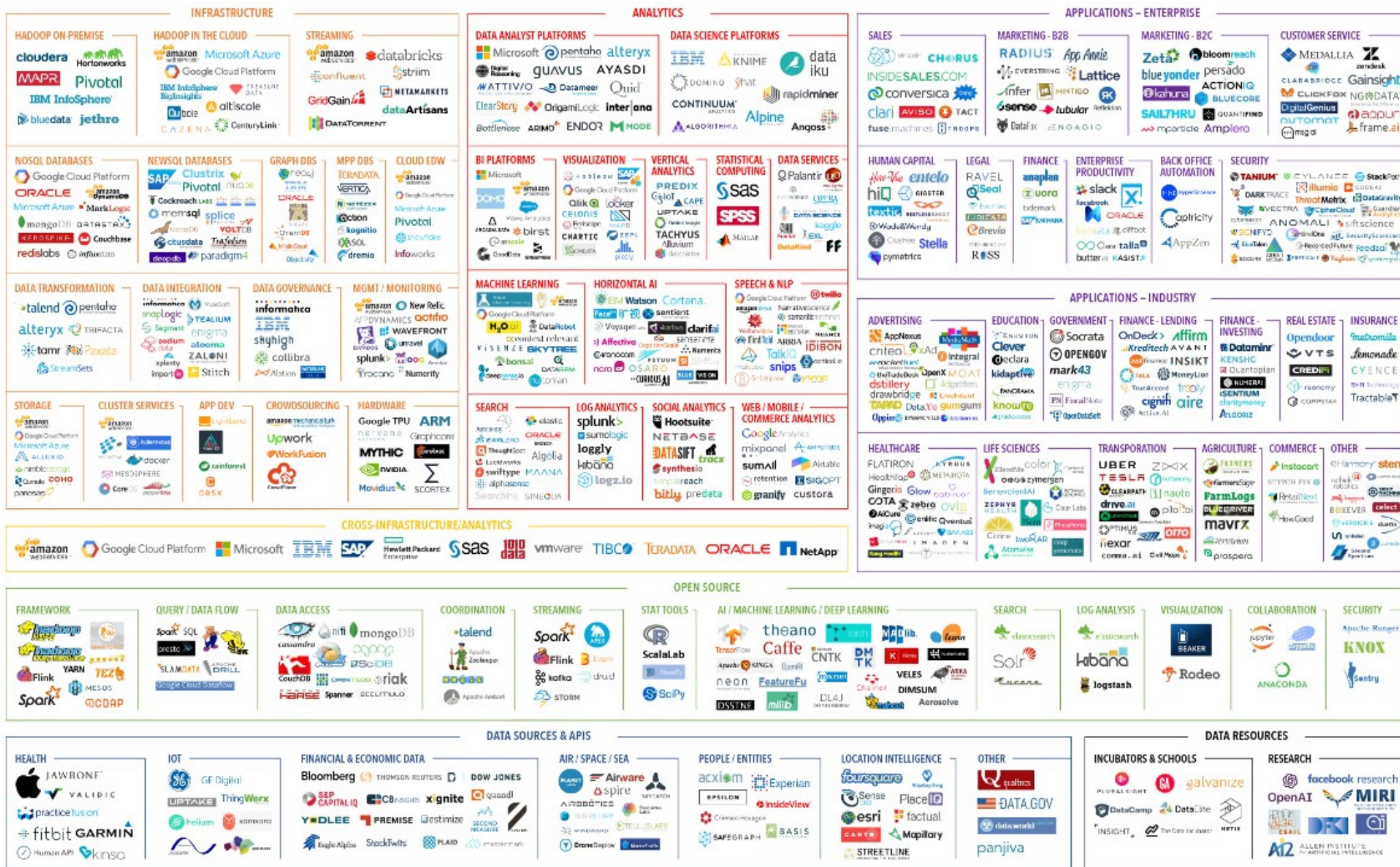## CAPITALIZING ON THIS OPPORTUNITY WILL REQUIRE:

RULES-DRIVEN INTEGRATION OF DISPARATE DATA

IMPROVED OPERATING INFRASTRUCTURES

NETWORK OF DATA-CENTRIC TECHNOLOGY AND PARTNERS

MARKETING DATA GOVERNANCE

# BIG DATA LANDSCAPE 2017



A full-page infographic titled "Big Data Landscape 2017" organizing companies into categories: Infrastructure, Analytics, Applications – Enterprise, Applications – Industry, Cross-Infrastructure/Analytics, Open Source, Data Sources & APIs, and Data Resources.

# Extended Relational Reference Architecture

# The
# Data Science Process

**Ask** an interesting question.

What is the scientific **goal?**
What would you do if you had all the **data?**
What do you want to **predict** or **estimate?**

**GET** the data.

How were the data **sampled?**
Which data are **relevant?**
Are there **privacy** issues?

**EXPLORE** the data.

**Plot** the data.
Are there **anomalies?**
Are there **patterns?**

**MODEL** the data.

**Build** a model.
**Fit** the model.
**Validate** the model.

**Communicate** and **visualize** the results.

What did we **learn?**
Do the results make **sense?**
Can we tell a **story?**

**2016 Magic Quadrant for Advanced Analytics Platforms**

CHALLENGERS | LEADERS

SAS
IBM
KNIME
RapidMiner
Dell
SAP
Angoss

Microsoft
Alteryx
FICO
Predixion Software
Alpine Data
Lavastorm
Megaputer
Prognoz
Accenture

NICHE PLAYERS | VISIONARIES

ABILITY TO EXECUTE
COMPLETENESS OF VISION

As of February 2016

**2017 Magic Quadrant for Data Science Platforms**

CHALLENGERS | LEADERS

RapidMiner
IBM
SAS
KNIME
MathWorks
NEW
Quest
Alteryx
Angoss
Microsoft
SAP
FICO
H2O.ai
NEW
Dataiku
NEW
Teradata
NEW
Domino Data Lab
NEW
Alpine Data

NICHE PLAYERS | VISIONARIES

ABILITY TO EXECUTE
COMPLETENESS OF VISION

As of February 2017

# IBM Watson

**TARGETS**
This project has **1 target**
Edit

**ANALYSIS DETAIL**
**131** input fields were evaluated.
**120** were potentially useful.

120

**TOP FIELD ASSOCIATIONS**
**62 strong associations** were found between fields.
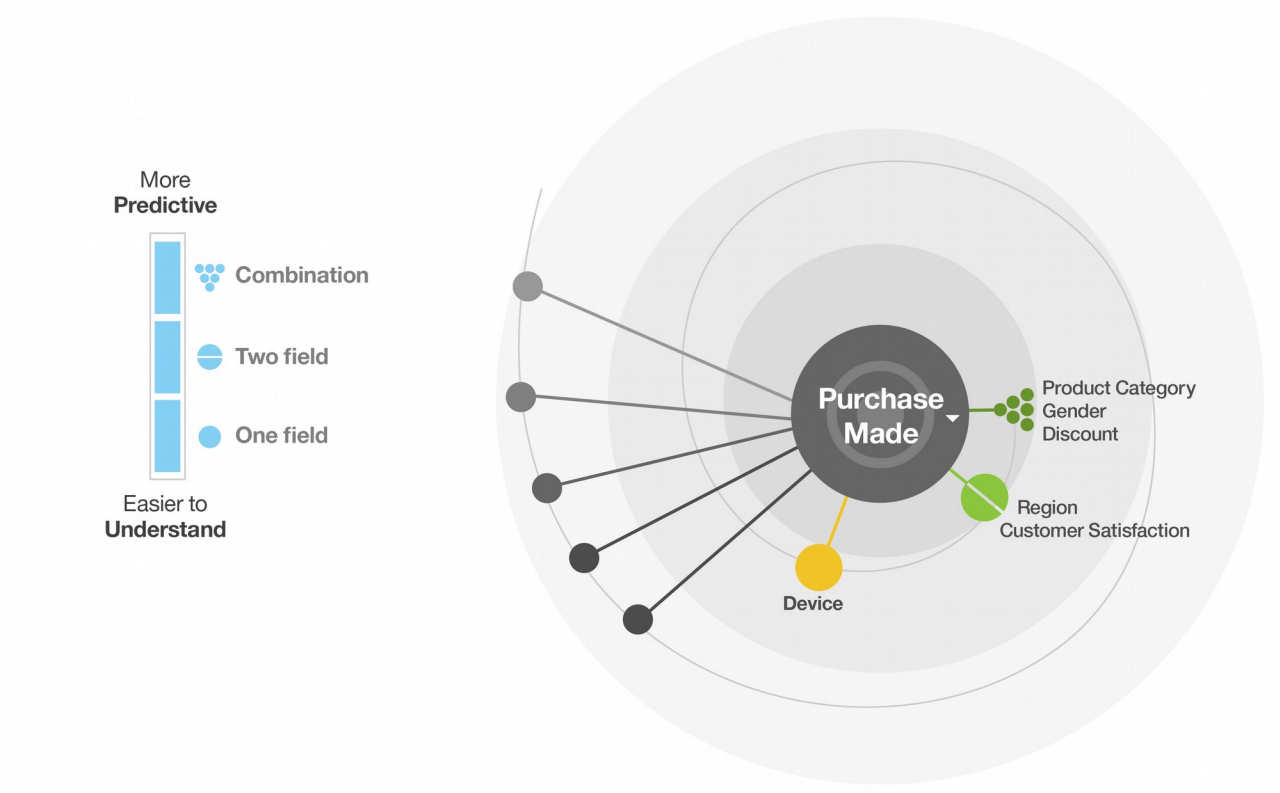View

**SALES**
A model with high predictive strength using **20 inputs** was found.

20

**FAVORITES**
One insights has been marked as a favorite.

# What predicts **Purchase Made**?

**2 strong predictors** and **1 moderate predictor** have been found and are shown below.

More
**Predictive**

Combination

Two field

One field

Easier to
**Understand**

Purchase Made

Product Category
Gender
Discount

Region
Customer Satisfaction

Device

## What influences **Purchase Made**?

● ● ○

**1**
Product Category, Gender, and **Discount** predict **Purchase Made**

Product Category
Gender  Discount  Time
Sentiment  Device  Region

**2**
**Region** and **Customer Satisfaction** predict **Purchase Made**

**3**
**Device** predicts **Purchase Made**

## What else is interesting about this?

● ○

**Average Sale Amount** differs across **Time**

**High Fashion** and **Region** are strongly associated.

**Discount** and **Time on Site** are associated.

Web Sales Data

⊕

⏱ Date  ⏱ Time  abc Source  abc Device  abc Product_Category  abc Purchase Made?  # Sale Amount  abc Country  # Discount

# IBM   resources

- [https://www.ibm.com/analytics/us/en/industry/government/](https://www.ibm.com/analytics/us/en/industry/government/)


- **IBM Certified Data Architect - Big Data**

[http://www-03.ibm.com/certify/certs/50001701.shtml](http://www-03.ibm.com/certify/certs/50001701.shtml)


- [https://datascience.ibm.com/](https://datascience.ibm.com/)

# Felicitaciones por tu interés en los datos y cómo usarlos 

Ángela Devia - ing.angela.devia@gmail.com

## Gracias