

# Data Science in Action

Practical use cases that demonstrate  
how businesses generate value from data

21 March 2014 – SDS2014, Winterthur

Pivotal™

# Introduction to Pivotal Data Labs

## Our Team



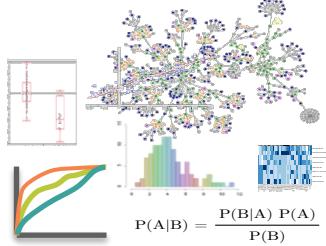
- High caliber global team of machine learning experts from a wide variety of quantitative backgrounds
- Equally capable in coding & statistics

## Our Tools



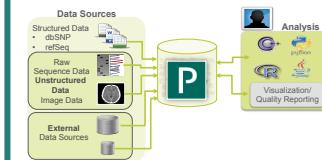
- Leading edge tools to implement machine learning collaboratively
- Have open-sourced several of our own tools for wide-spread use

## Our Methods



- Parallelized a wide variety of machine learning algorithms for optimum performance on the Pivotal platform
- Agile, test-driven, customer focused

## Our Process



- Analytical workflow aligned with business needs and optimized for speed
- Supports iterative and collaborative working

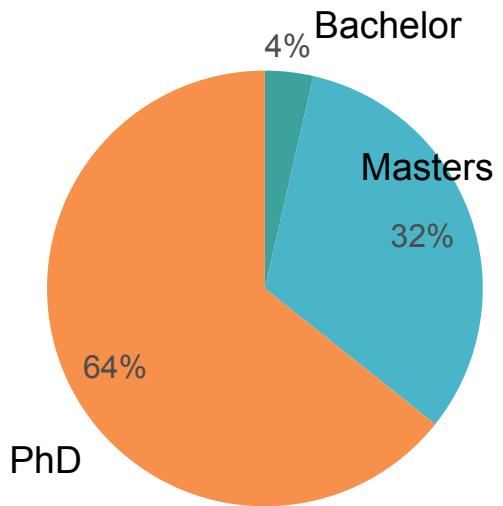
## Our Experience



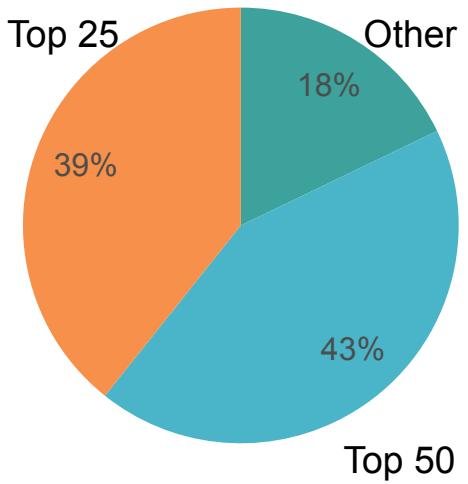
- More than 100 customer assignments carried out in the past 18 months
- Ensures quality and best practice in all our assignments

# Pivotal Data Science Team

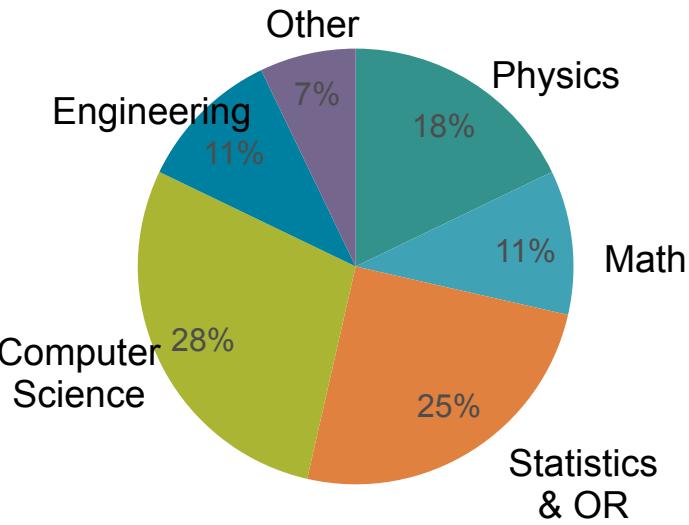
**By Degree**



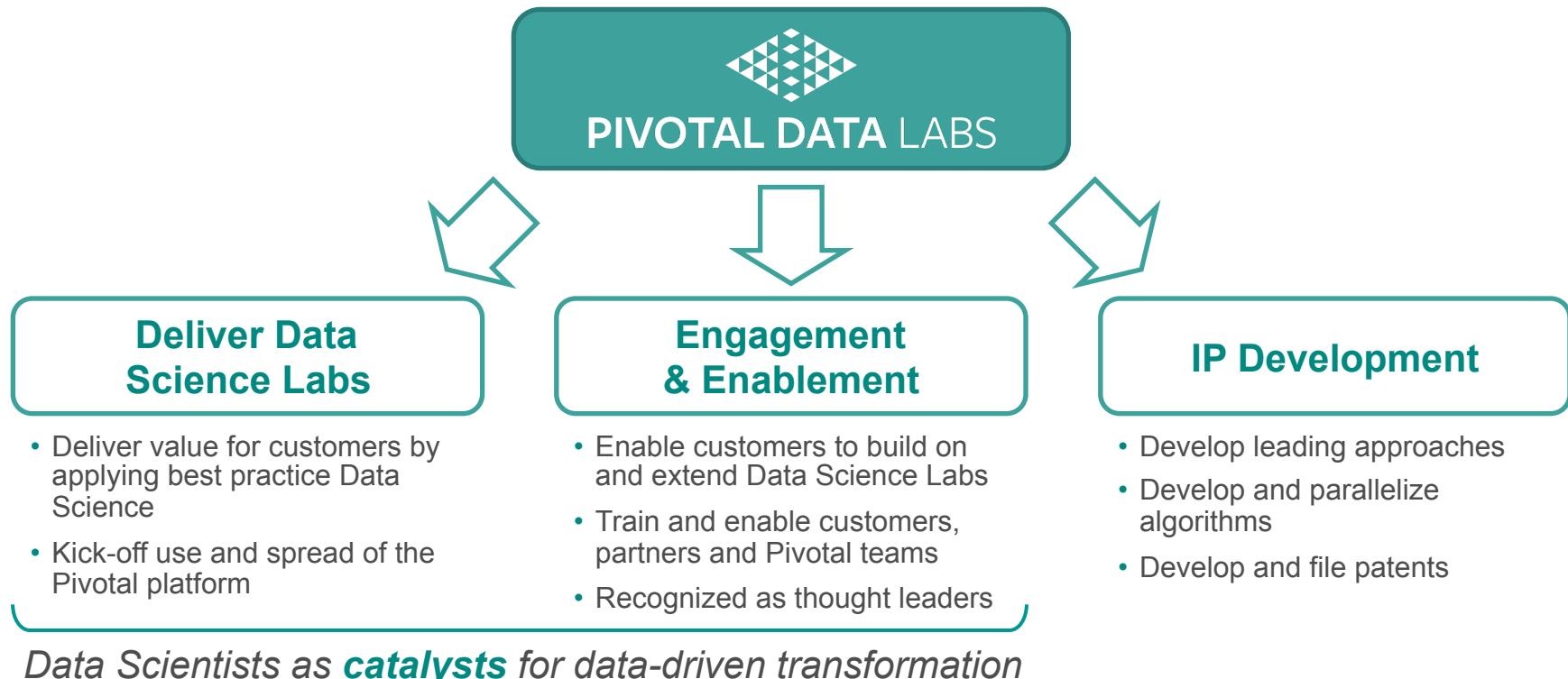
**By University**



**By Subject**



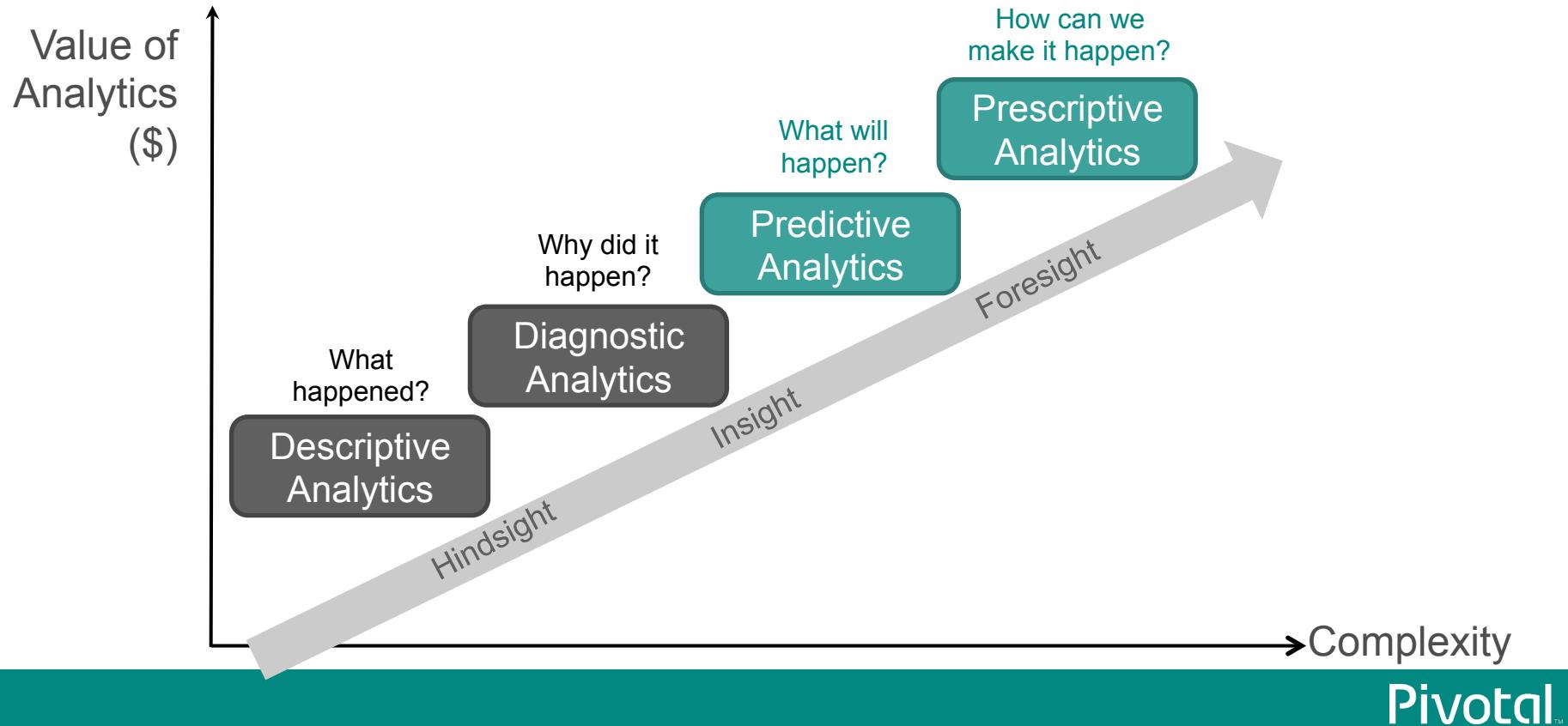
# What does the Pivotal Data Science team do?



# Data Science in Action

- The Value of Data Science
- The Practice of Data Science
- In-depth Use Case: Traffic Prediction
- Use Case Overviews
- Q&A

# What do we mean with Data Science?



# Big Data & Data Science

Decision = Data + Rules

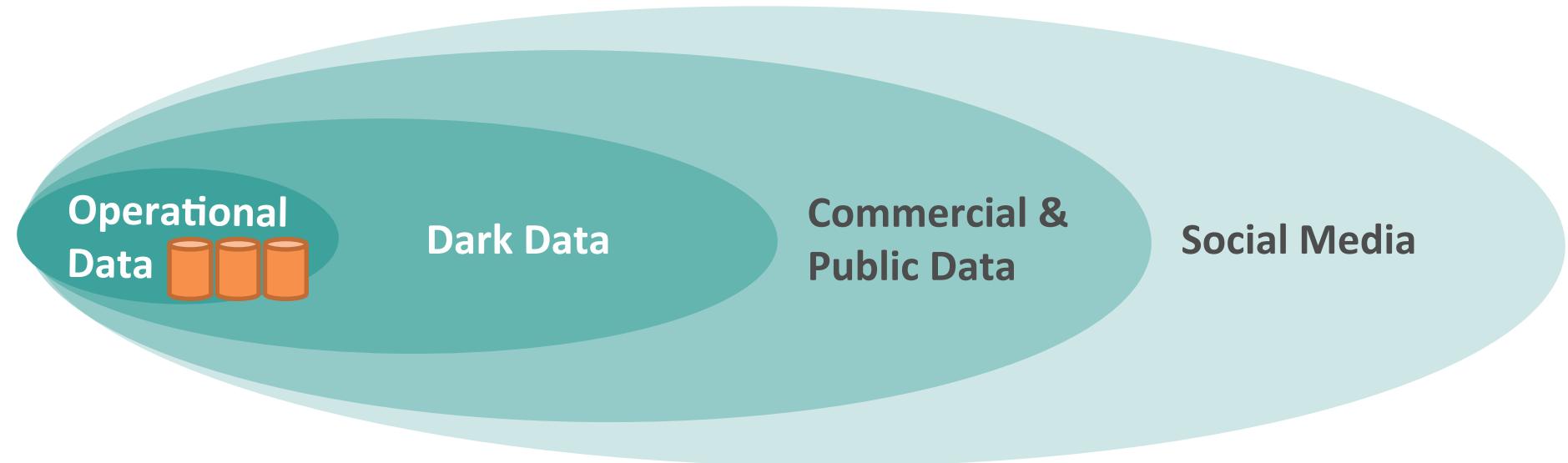


“Big Data”



Data  
Science

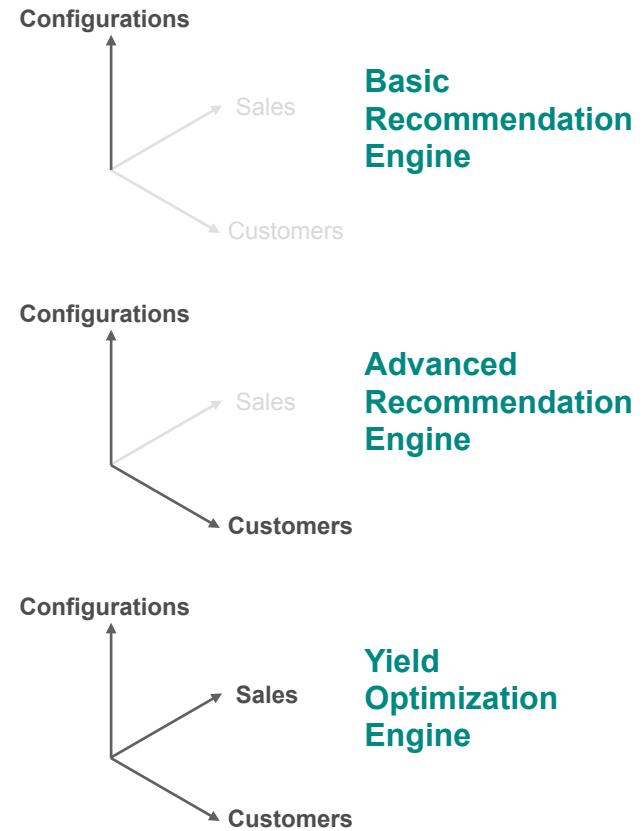
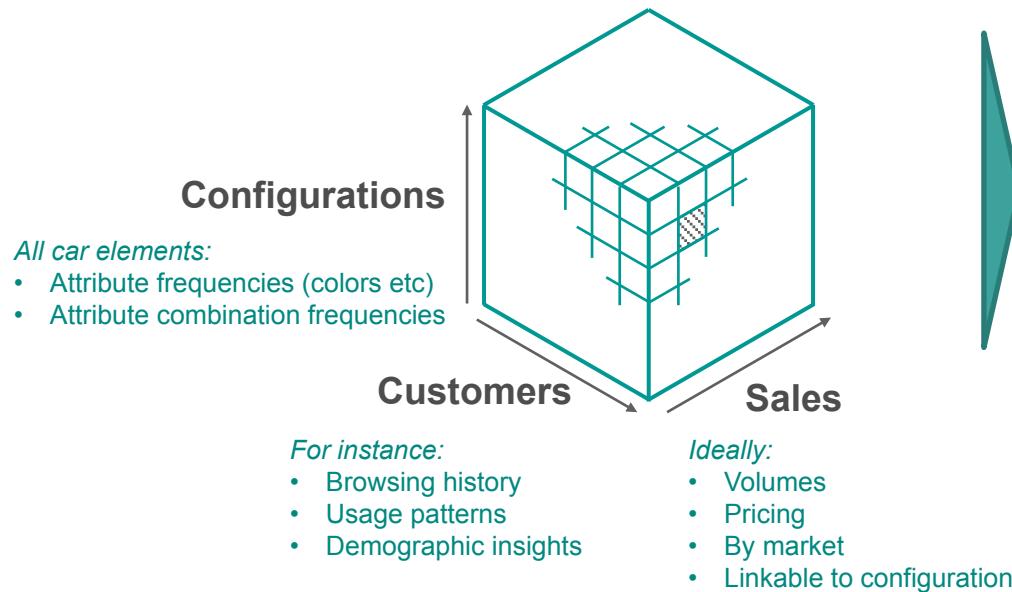
# “Big Data”



# Combining data sources: Example

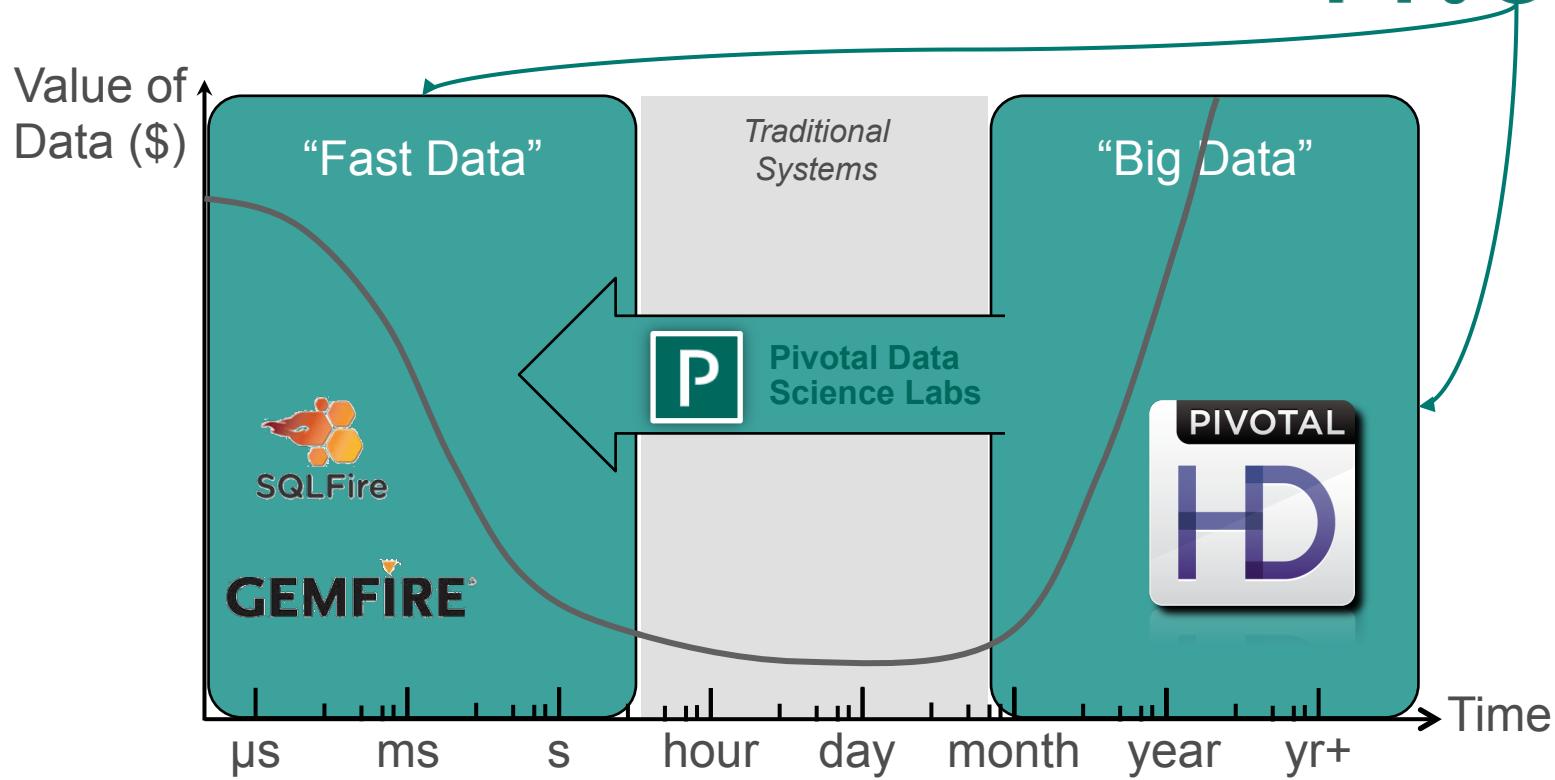
<p><b>IPSQ (Quality)</b> Owner: TS Production team Test flags from production line 1 year ~300GB</p>	<p><b>APDM</b> Owner: TS Production team Full vehicle history including IPST (technical), IPSL (logistics), IPSQ test flags and all test results. 30 years ~TBs</p>
<p><b>FASTA</b> Owner: Aftersales Dealership electronic tests Identifies early issues with cars &gt;25TB</p>	<p><b>IQS: Initial Quality Survey from JD Power</b> Owner: R&amp;D Survey responses from new owners after 90 days for approx 1700 vehicles Few thousand lines ~MB</p>
<p><b>Social Data</b> Owner: R&amp;D Pulling 500MB per day from Twitter</p>	<p><b>TQP</b> Owner: Supplier management PDFs of parts spec sheets ~ 500GB</p>

# Generating value from data: Car configurator example



# The value of data over time

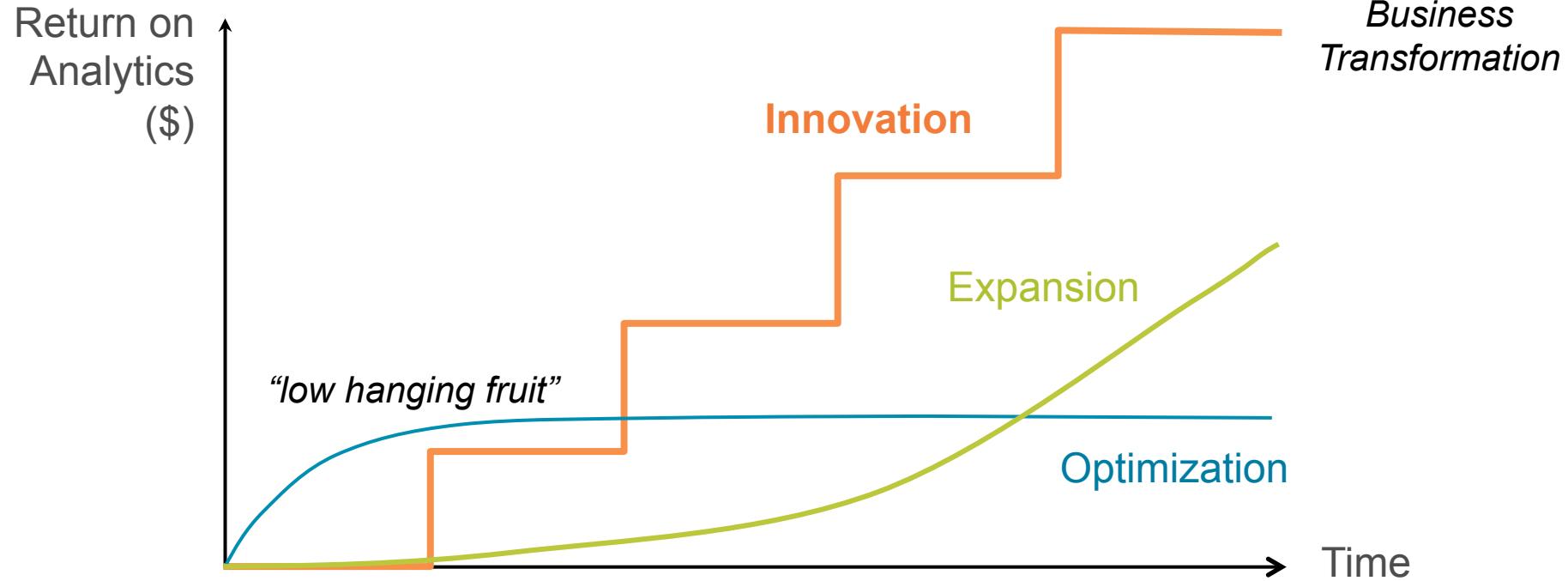
Pivotal™



# Data Science

- The use of statistical and **machine learning** techniques on **big multi-structured data** in a distributed computing environment to identify **correlations** and causal relationships, classify and predict events, identify patterns and anomalies, and infer probabilities, interest, and sentiment.
- In order to **drive automated low latency actions in response to events of interest**

# Why do Data Science?



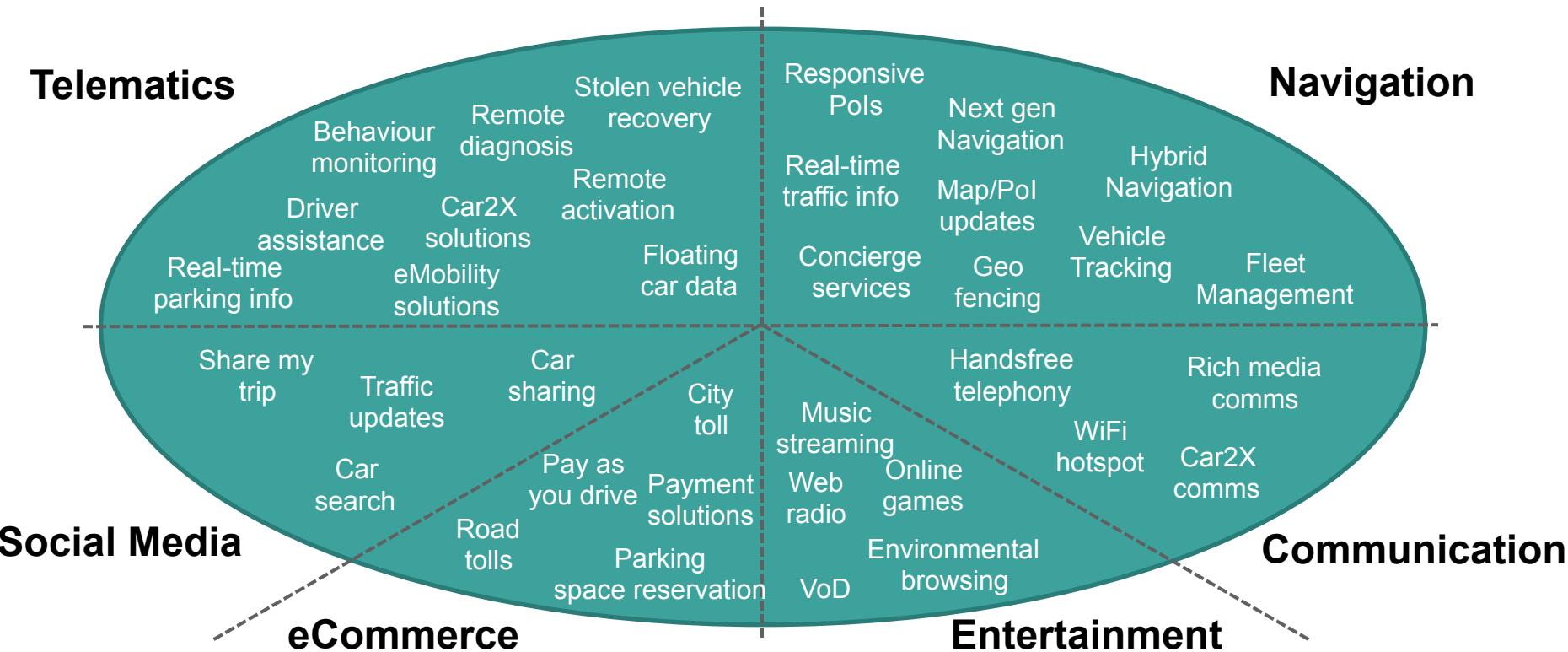
# What transforms businesses today?

- Digitization
- Internet of Things
- Pervasive Computing
- Pervasive Connectivity

# Example: major paradigm shifts in automotive

Genesis	Mass Production	Modern Manufacturing	Platform Strategy	What's Next?
1885	1908	1950s	1980s	2020
				
<b>Not a horse</b>	<b>Mass availability</b> <i>"You can have any color of car, provided it's black"</i>	<b>Brand proliferation</b> <i>You can have any color</i>	<b>Globalization</b> <i>You can have anything anywhere</i>	<b>Connected, autonomous vehicles</b>

# The Connected Car drives innovation



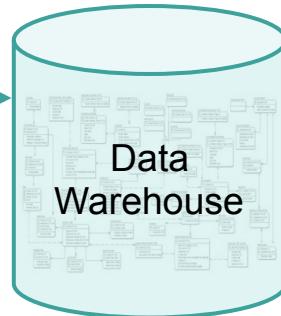
# Data Science in Action

- The Value of Data Science
- **The Practice of Data Science**
- In-depth Use Case: Traffic Prediction
- Use Case Overviews
- Q&A

# Traditional Analytics Process

## Data Sources

- Structured Data
  - Sensors
  - Flight recordings



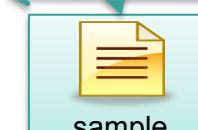
## Analysis



## Visualization/ Quality Reporting

Scheduling Analysis

Operations Analysis



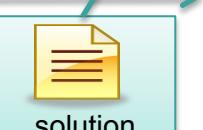
In-memory  
statistics tool



forecast



In-memory  
optimization tool



Pivotal™

Time-to-Insights

# Augmenting an analytical architecture

## Data Sources

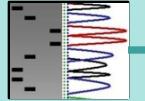
### Structured Data

- Sensors
- Flight recordings



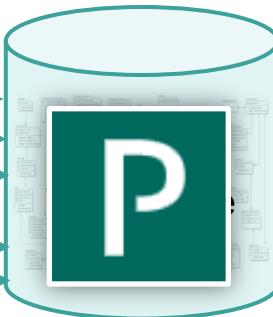
### Unstructured Data

- Image Data
- Geolocation Data
- Voice Transcription



### External Data Sources

- Weather Data
- Open Gov Data



## Analysis



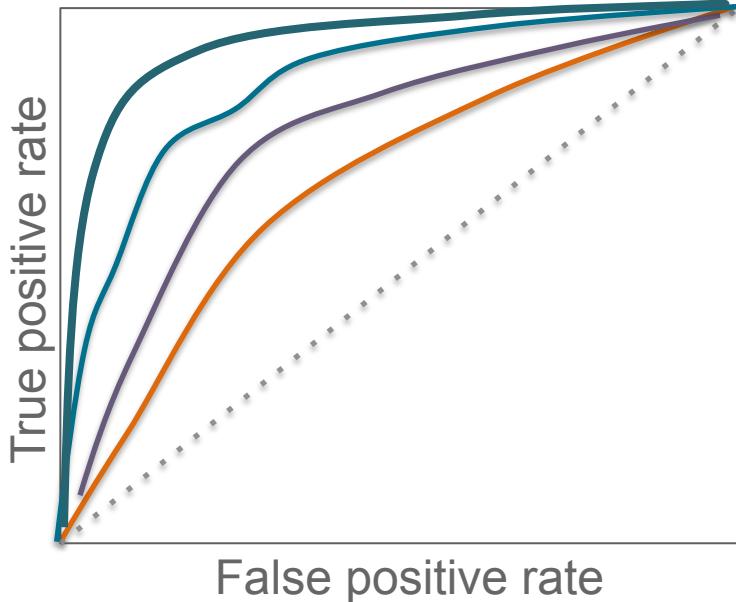
## Visualization/ Quality Reporting

## Benefits of a new architecture:

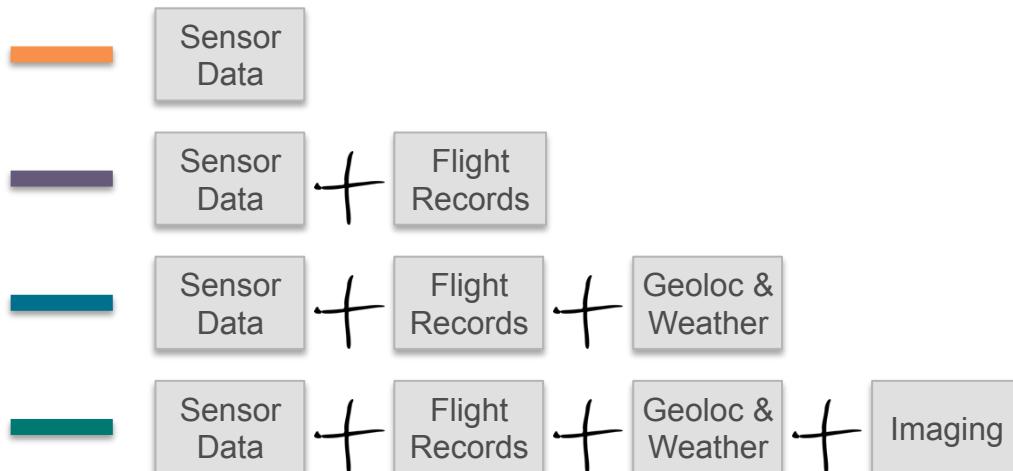
- **Eliminates** data movement
- Enables **rapid** data re-processing
- **Seamless** integration of additional external resources into analyses

# Machine Learning and Big Data

*Getting the whole picture improves predictive power*



- More data from different sources
- Provides a more complete view
- Improves statistics and inference



# The main types of use cases in practice

## Data Mining

- Categorize types (segmentation)
- Categorize behaviors/usage
- Identify co-occurrences and associations
- Identify anomalies
- Identify attitudes
- Resolve entities

## Predicting Behavior

- Predict churn likelihood
- Predict cross/up sales potential
- Predict fraud/waste/abuse likelihood
- Predict performance
- Predict reliability/quality
- Make a “recommendation”

## Optimization

- Optimize processes
- Optimize process parameters
- Optimize asset allocation



# As Data Scientists, what do we want?

Infrastructure  
Independent

Fast &  
Scalable

Schema  
Free

Real  
Time

Easy to  
Use

- Open source
- PaaS
- In-database analytics
- MPP
- Hadoop
- In-memory data grids embedded into the platform
- SQL, not Java
- Faster than Hive

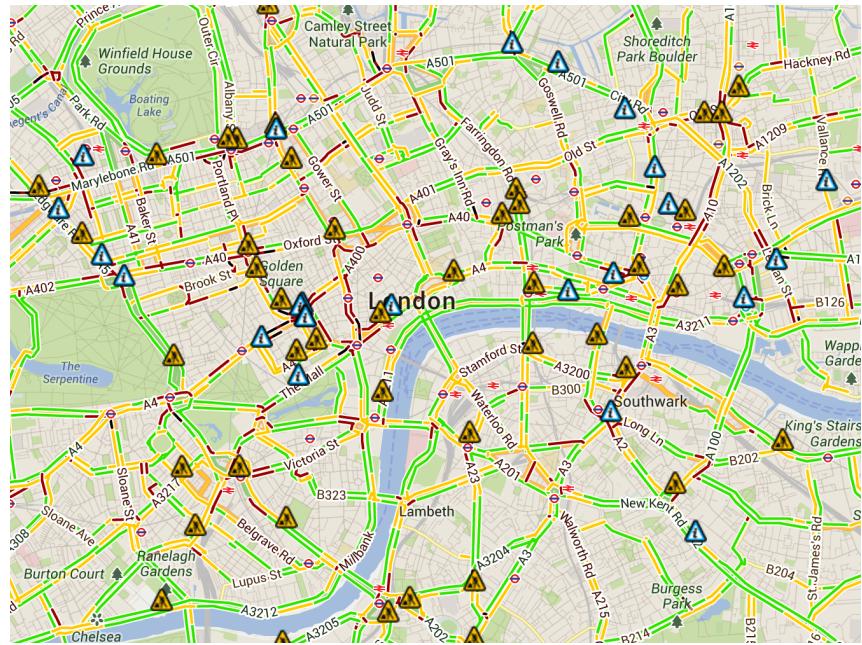
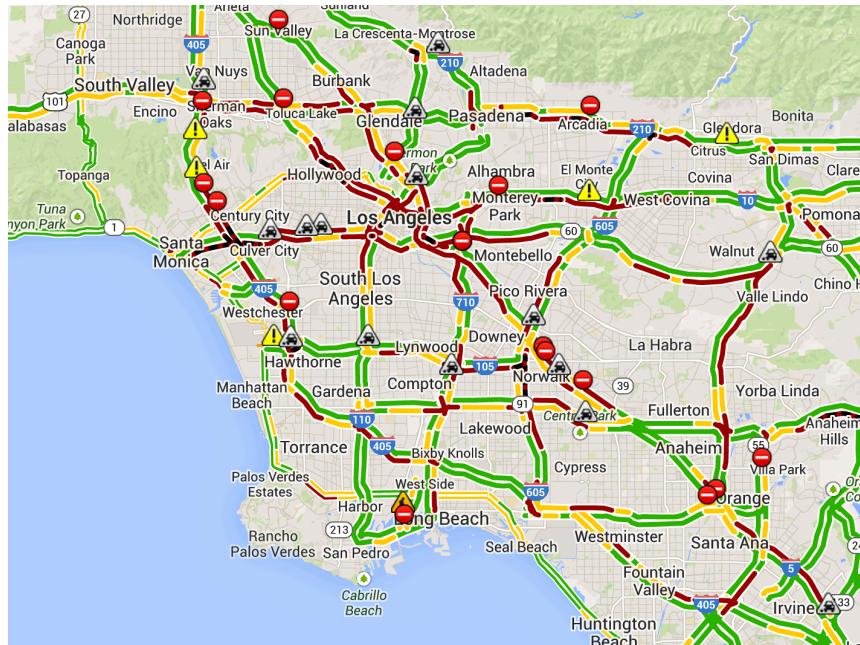
# Data Science in Action

- The Value of Data Science
- The Practice of Data Science
- **In-depth Use Case: Traffic Prediction**
- Use Case Overviews
- Q&A

An aerial photograph of a multi-lane highway during rush hour. The road is packed with numerous cars, creating a dense, dark grey texture. The highway curves slightly to the left, with a green embankment on the left side and some buildings visible. The right side shows a steep hillside covered in greenery.

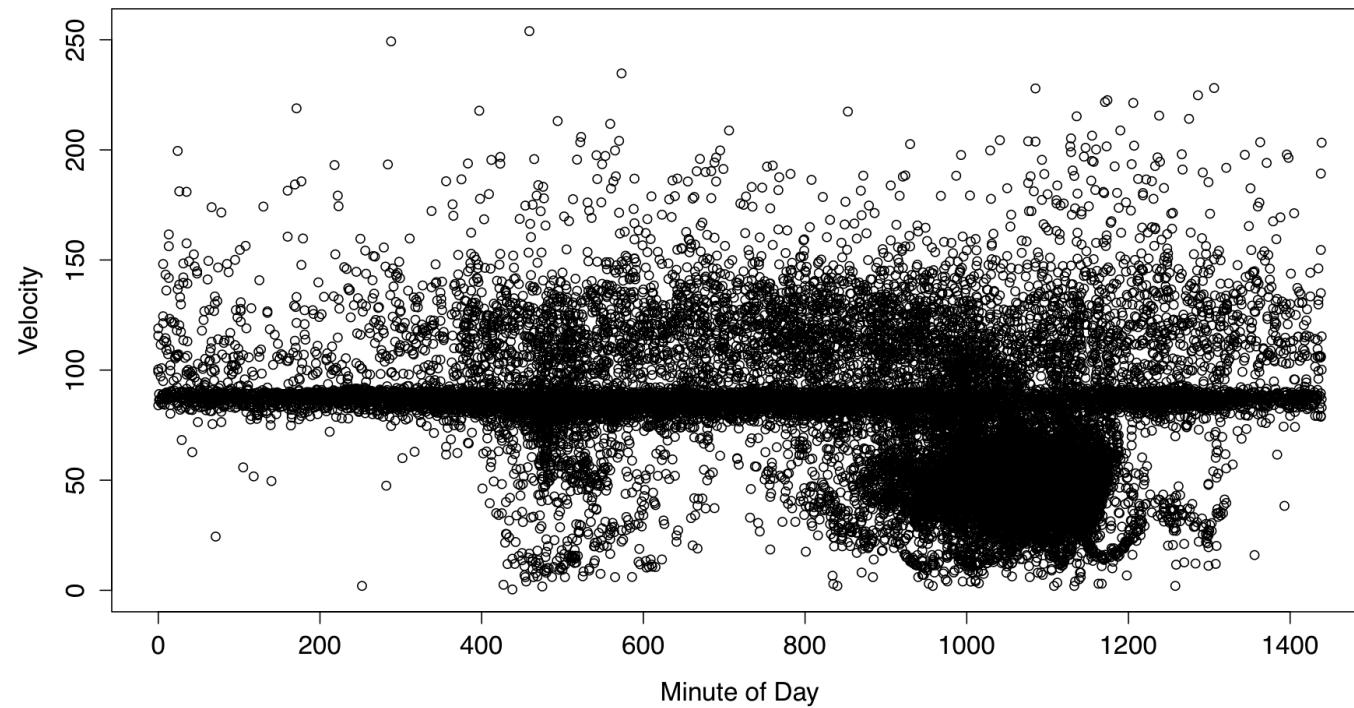
What does  
traffic data  
look like?

...like this?

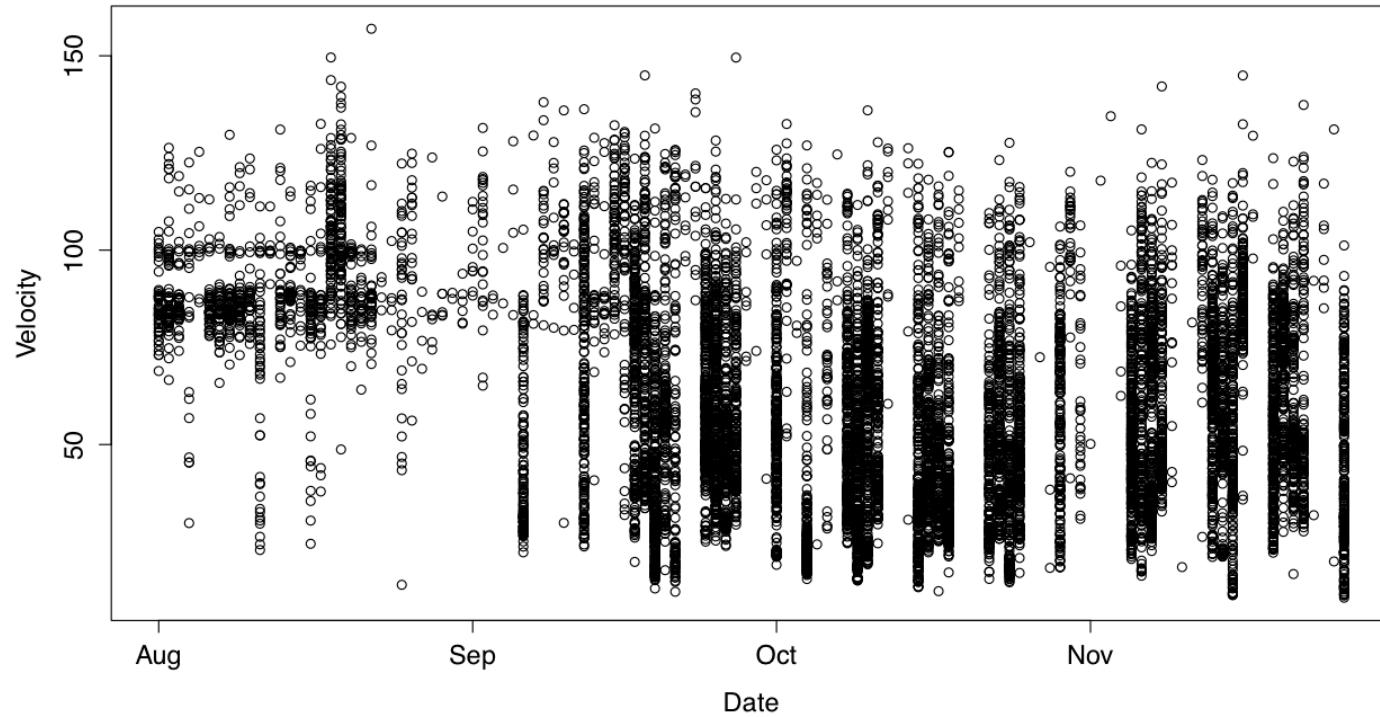




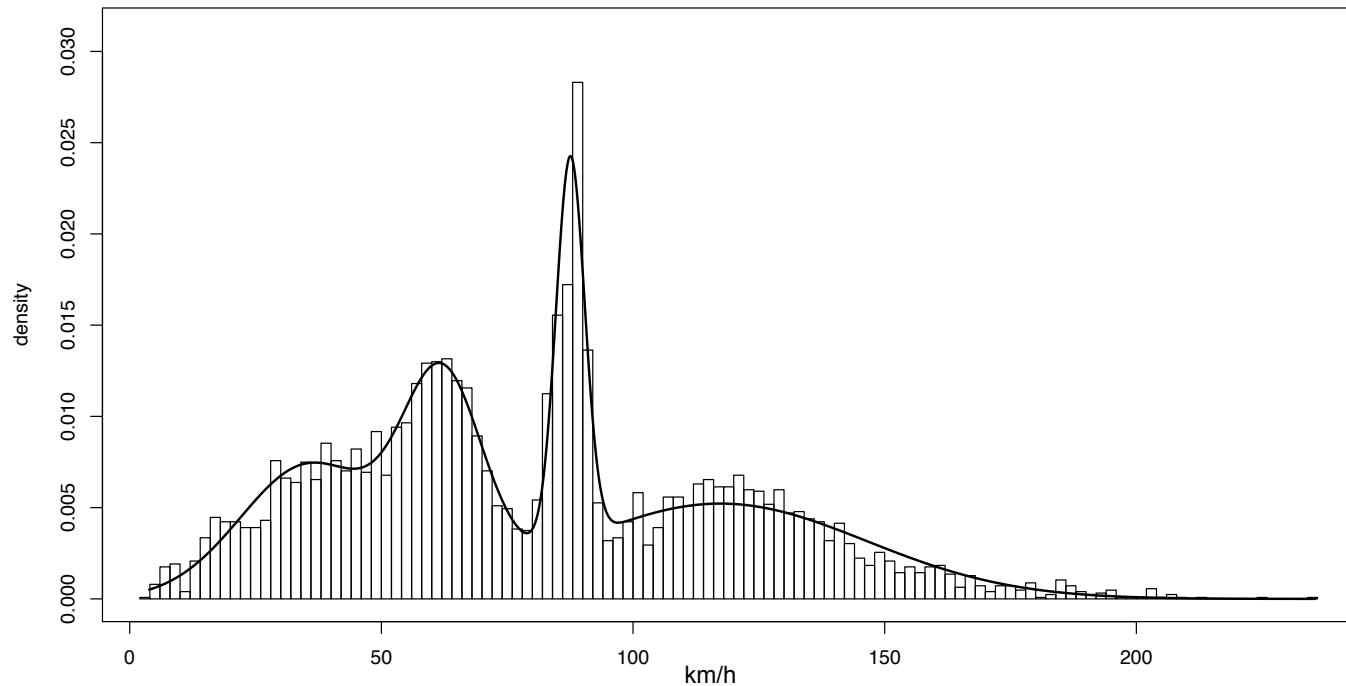
# Velocity by Time of Day



# Distribution of Velocity over Time

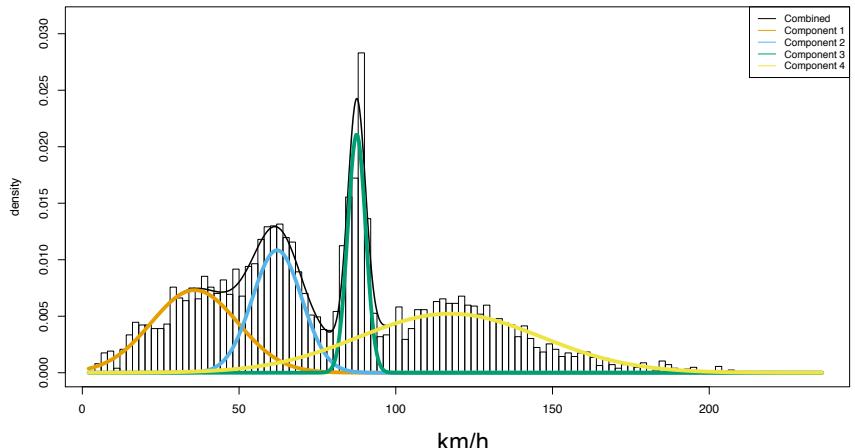
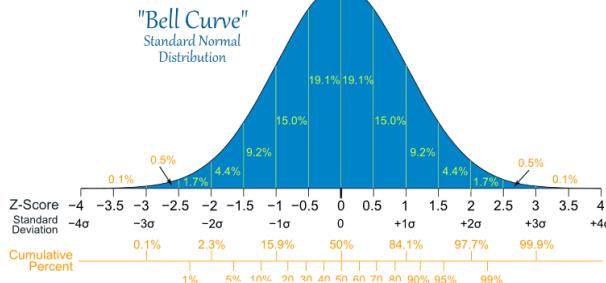


# Velocity Distribution

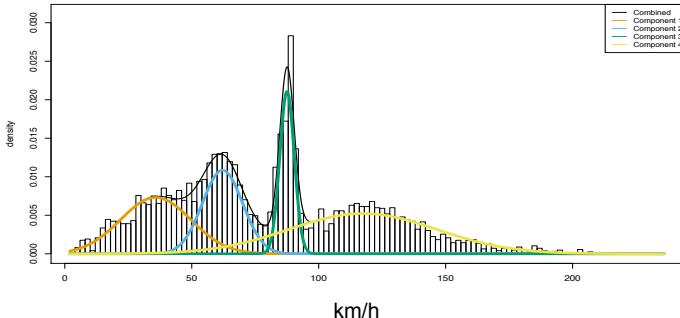
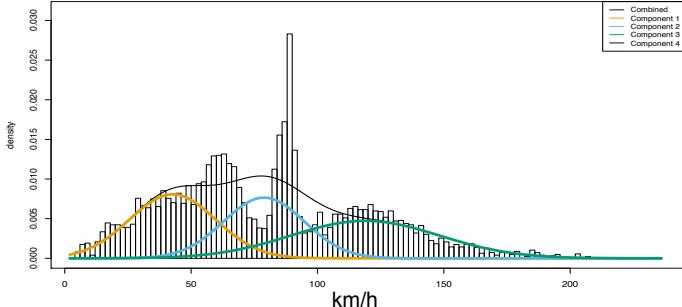
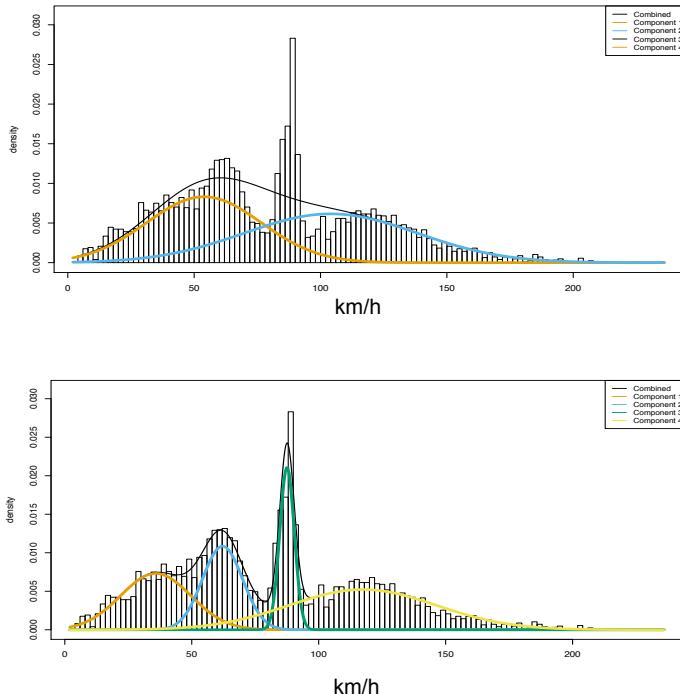
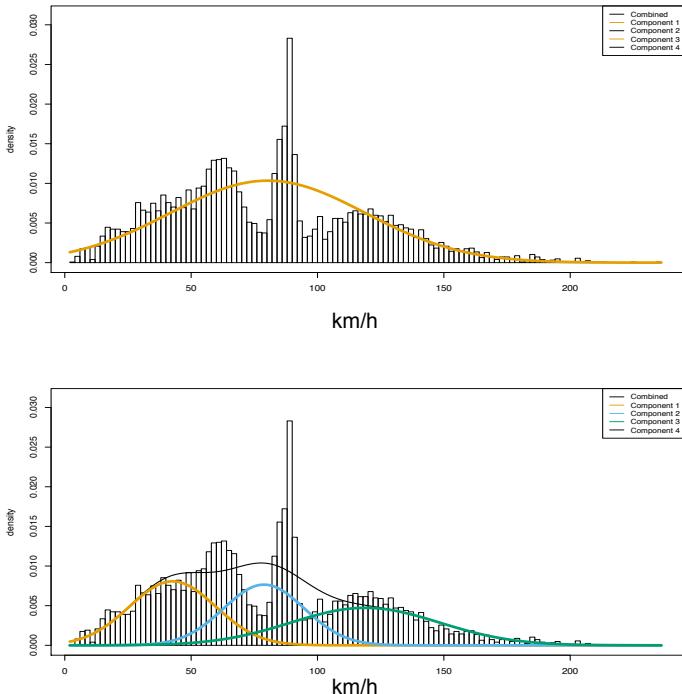


# Find Velocity Groups

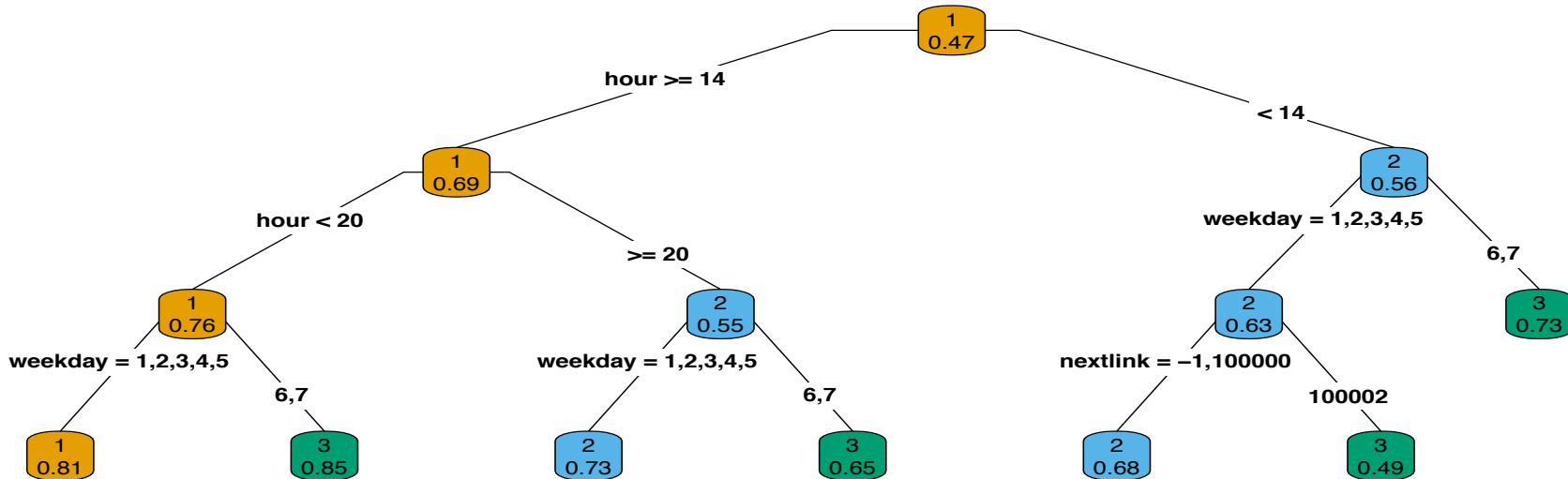
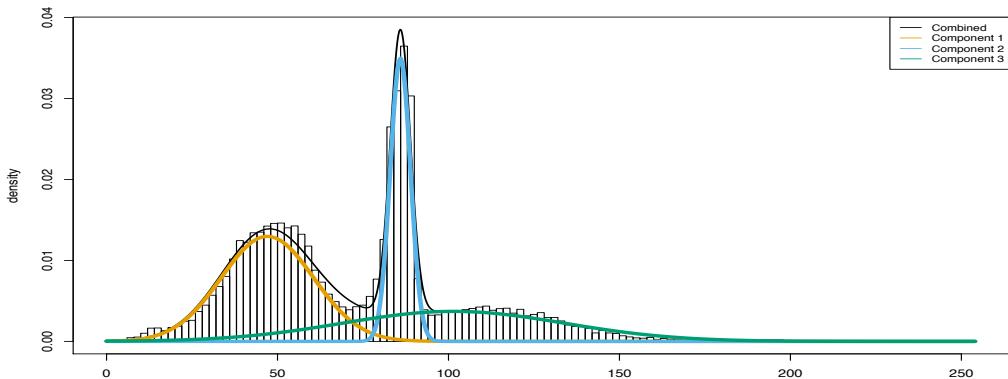
- Velocity distributions can be **fit well** with Gaussians
- An ‘overlay’ of multiple Gaussians is called **Gaussian Mixture Model**
- GMM fitting of the velocity distribution is done by **Expectation-Maximization** algorithm
- Shapes and positions of Gaussians **determine velocity groups**



# Gaussian Mixture Model

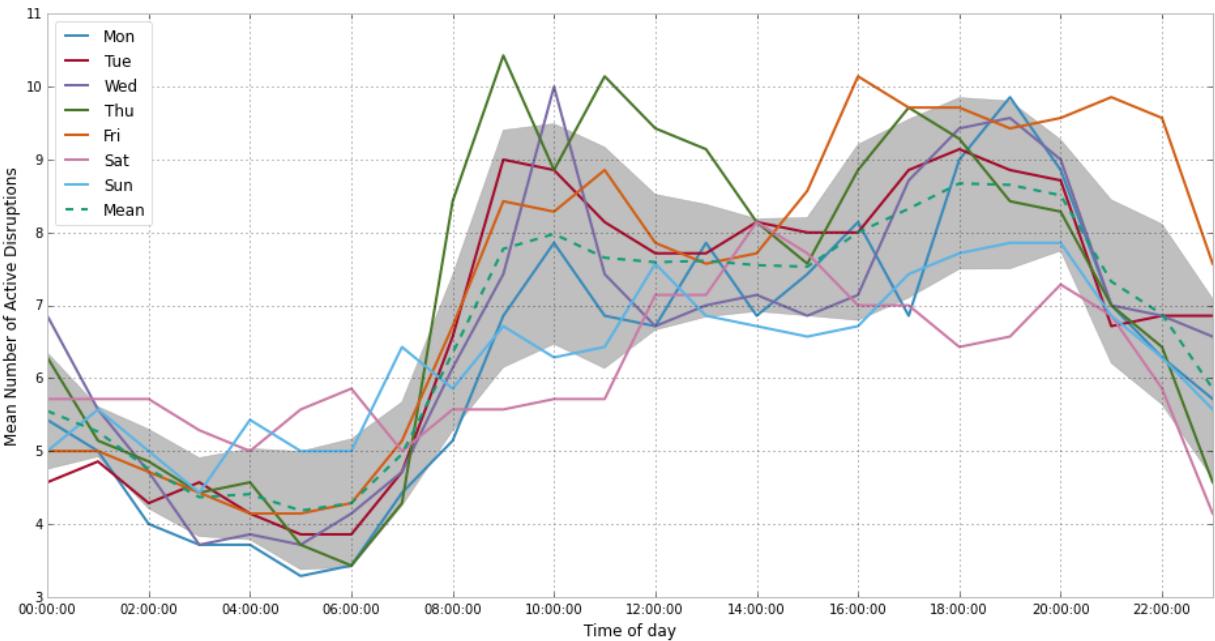


# Decision Trees Example



# Sneak Peek at our TfL Data Demo

- Used the freely accessible **TfL data** for a demo
- Shows # of **active disruptions** over different days in London



- Rush hour effects visible
- Nights are more quiet, but more disruptions on weekend nights

# Data Science in Action

- The Value of Data Science
- The Practice of Data Science
- In-depth Use Case: Traffic Prediction
- **Use Case Overviews – *not presented***
- Q&A

# Pivotal

BUILT FOR THE SPEED OF BUSINESS