

ODTUG  
Kscope19   
SEATTLE, WASHINGTON • JUNE 23-27

PLEASE FILL OUT  
YOUR EVALUATIONS

SEATTLE

 Washington State  
Convention Center

# Is It Corked? Wine Machine Learning Predictions with OAC

Francesco Tisiot  
BI Tech Lead at Rittman Mead



A close-up, high-contrast photograph of an espresso machine's spout pouring a stream of golden-brown espresso into a clear glass cup. The machine is dark and sleek, with the coffee stream being the primary light source in the scene.

# Francesco Tisiot

## BI Tech Lead at Rittman Mead



Verona, Italy



Rittman Mead Blog



10 Years Experience in BI/Analytics



[francesco.tisiot@rittmanmead.com](mailto:francesco.tisiot@rittmanmead.com)



@FTisiot



Oracle ACE

# About Rittman Mead

Rittman Mead is a **data and analytics company** who specialise in data visualisation, predictive analytics, enterprise reporting and data engineering.

We use our skill, experience and know-how to work with organisations across the world to interpret their data. We enable the business, the consumers, the data providers and IT to work towards a common goal, **delivering innovative and cost-effective solutions** based on our core values of thought leadership, hard work and honesty.

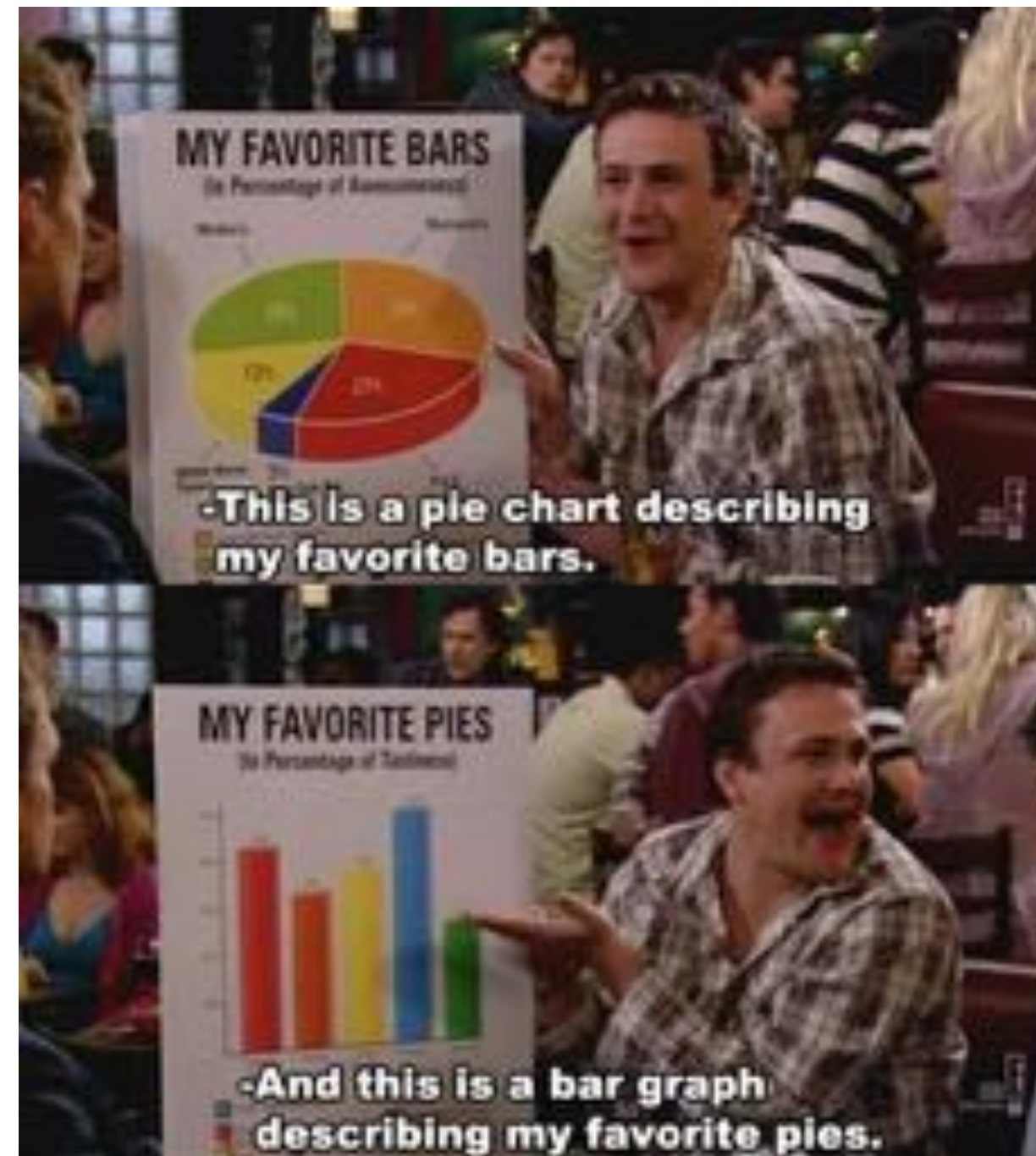
We work across **multiple verticals** on projects that range from mature, large scale implementations to proofs of concept and can provide skills in **development, architecture, delivery, training and support**.

# Agenda

- Tooling
- Data Science Steps
- Demo

# Tooling

# Analytics!





# Oracle Analytics Cloud

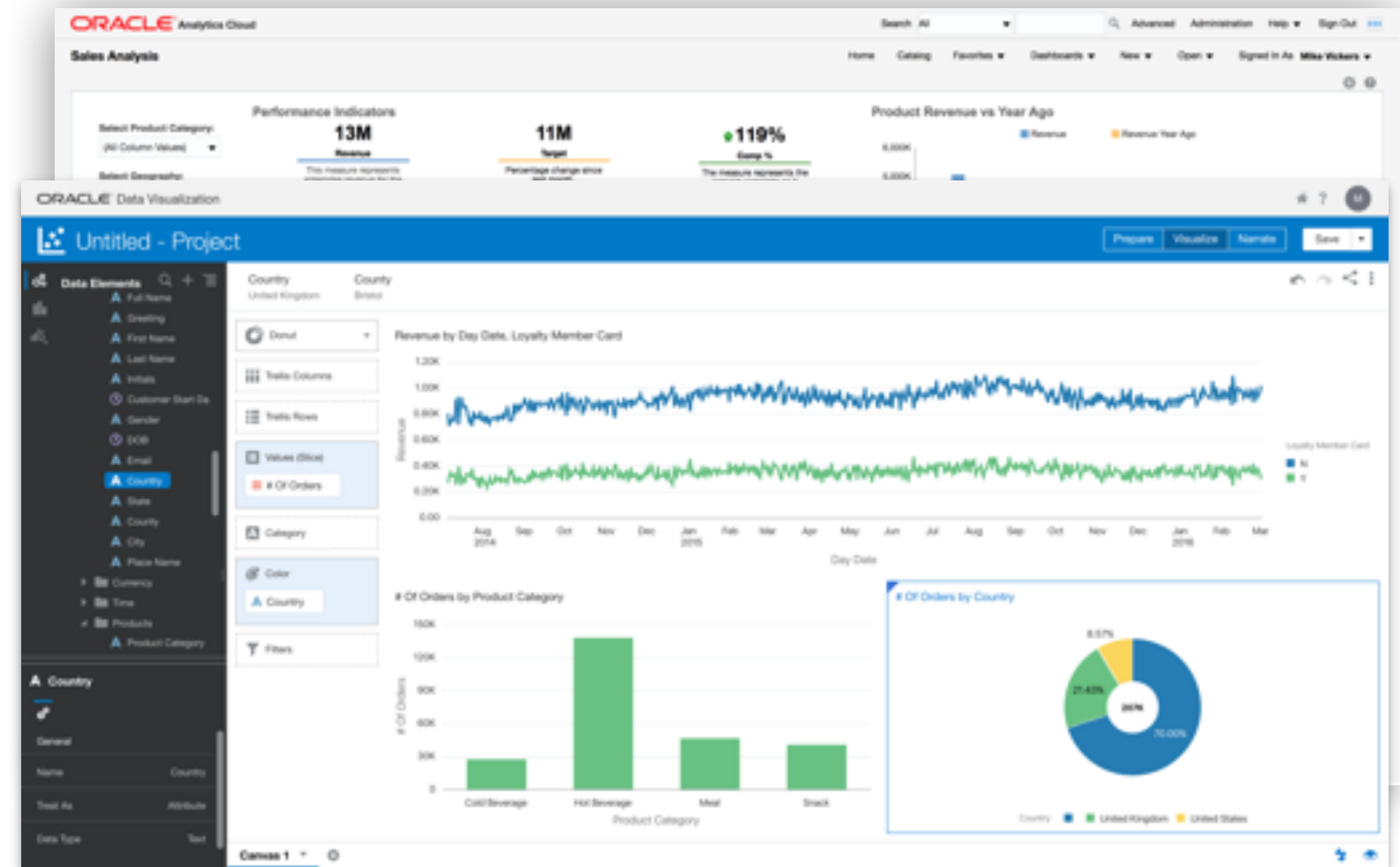
- Oracle's complete suite of Platform Services (PaaS) for unified analytics in the cloud
- Delivered entirely in the cloud:
  - ▶ No infrastructure footprint
  - ▶ Flexibility to scale up or down based on your immediate needs
  - ▶ Simplified, metered licensing
- Several options to suit your needs:
  - ▶ Oracle or customer/partner managed services
  - ▶ Functionality bundled into 3 editions





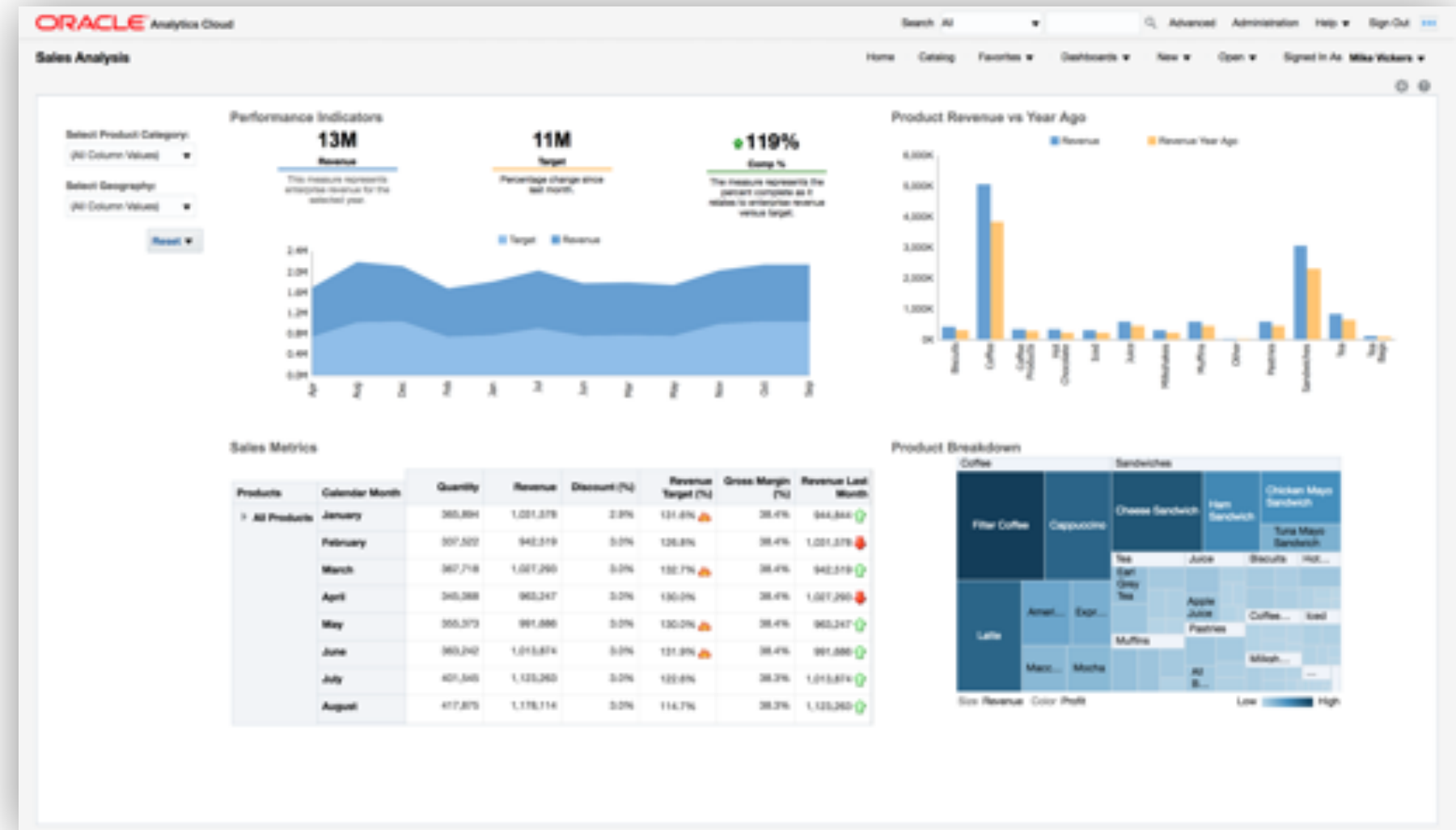
# Functions

- OAC supports **every** type of analytics workload across your organisation
- *Classic* enterprise BI:
  - ▶ Analysis & dashboarding
  - ▶ Published reporting
  - ▶ Enterprise Performance Management
- *Modern* departmental/personal discovery:
  - ▶ Extended data mashup & modelling
  - ▶ Data preparation, exploration & visualisation
  - ▶ Data science & machine learning



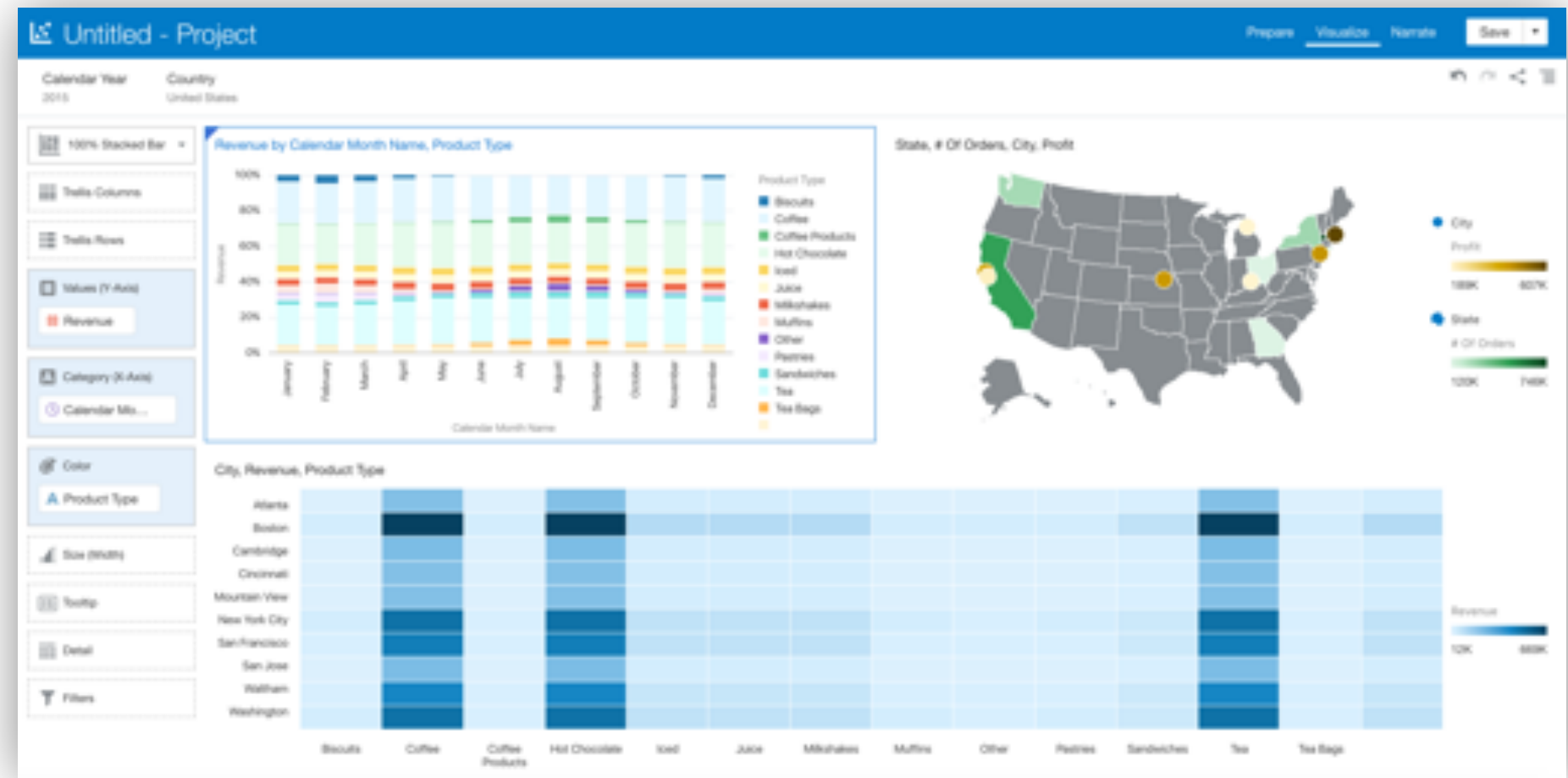
# Classic Enterprise BI

- Similar User Experience to OBIEE 12c
  - ▶ Centrally maintained & governed
  - ▶ Semantic model remains key
- Interactive Dashboards
  - ▶ Ideal for KPI measurement & monitoring
  - ▶ Guided navigation paths
- BI Publisher
  - ▶ Highly formatted, burst outputs
- Action Framework
  - ▶ Navigation actions
  - ▶ Scheduled agents



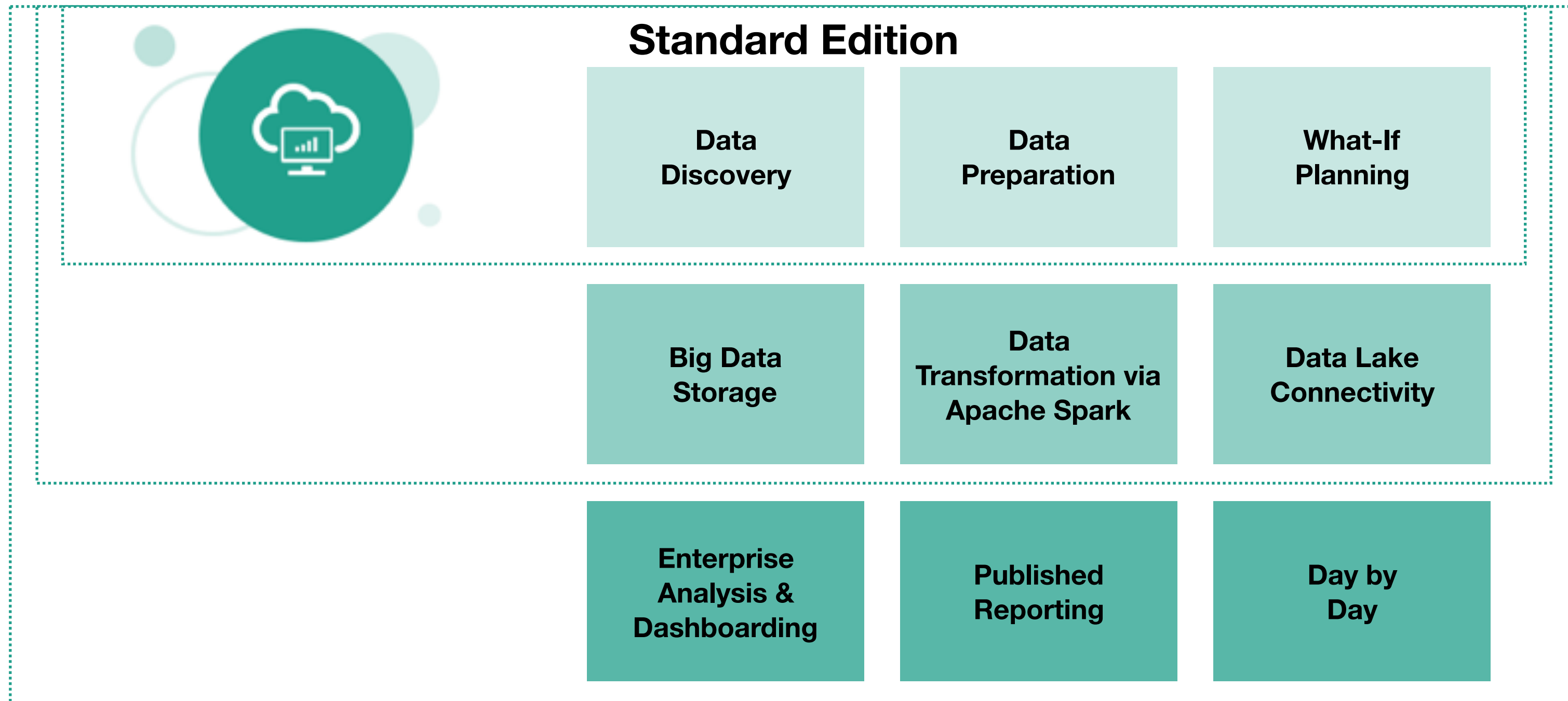
# Modern Data Discovery

- Data Preparation
  - ▶ Acquire data from multiple connections
  - ▶ Apply enrichments data prior to analysis
  - ▶ Define repeatable preparation flows
- Data Visualisation
  - ▶ Create visual insights rapidly
  - ▶ Construct narrated storyboards
  - ▶ Share findings
- Machine Learning
  - ▶ Build & train ML models
  - ▶ Apply model to new data sets





# Three Edition Options



# Two Purchasing Options



## Monthly Flex

Based on **Universal Credits** model

**12 month** minimum tenure

Payments made in **advance**

Unused credits are **forfeited**

### Suitable for:

Predictable, production workloads

Long running platforms



## Pay As You Go

Based on **Universal Credits** model

**No** minimum tenure

Payments made in **arrears**

Based on **consumption**

### Suitable for:

Rapid Prototyping

Testing & Sampling

Elastic Scalable



# OAC And Data Science



# Basic Operations



What are the  
Drivers for My  
Sales?



Based on my Experience  
I can Guess....



Statistically Significant  
Drivers for Sales Are ...

Augmented  
Analytics

# Basic Operations



Is this Client  
going to accept  
the Offer?

YES/NO  
**50%**

Basic ML  
Model

**70%**

**Before Starting.... Define the Problem!**





# Problem Definition: **Predicting Wine Quality**



# Rule Based

Italy or France -> Good

Rest of the World -> Bad

Price  $\geq$  10 Euros -> Good

Price < 10 Euros -> Bad

Price > 30 & Production Zone = Veneto & .... -> 6.5

**TEP**

Task



Estimate Wine  
Good/Bad

Experience



Corpus of Wines  
Descriptions with Ratings




Performance



Accuracy




# Accuracy

		Predicted Value	
		Good	Bad
Real Value	Good		
	Bad		

$$\text{Accuracy} = \text{happy face emoji} / (\text{happy face emoji} + \text{sad face emoji with sweat drop})$$

# Dataset



# Wine Reviews

130k wine reviews with variety, location, winery, price, and description

zackthoutt • updated a year ago (Version 4)

[Data](#)
[Overview](#)
[Kernels \(1,686\)](#)
[Discussion \(19\)](#)
[Activity](#)

[Download \(51 MB\)](#)
[New Kernel](#)

Data (51 MB)

### Data Sources

- winemag-data-130k... 130k x 14
- winemag-data\_first15... 151k x 11
- winemag-data-130k-v2.json

### About this file

Here is a CSV version of the data I scraped. This dataset has three new fields --Title (which you can parse the vintage from), Taster Name, and Taster Twitter Handle. This should also fix the duplicate entries problem in the first version of the dataset and add ~25k unique reviews to play with.

### Columns

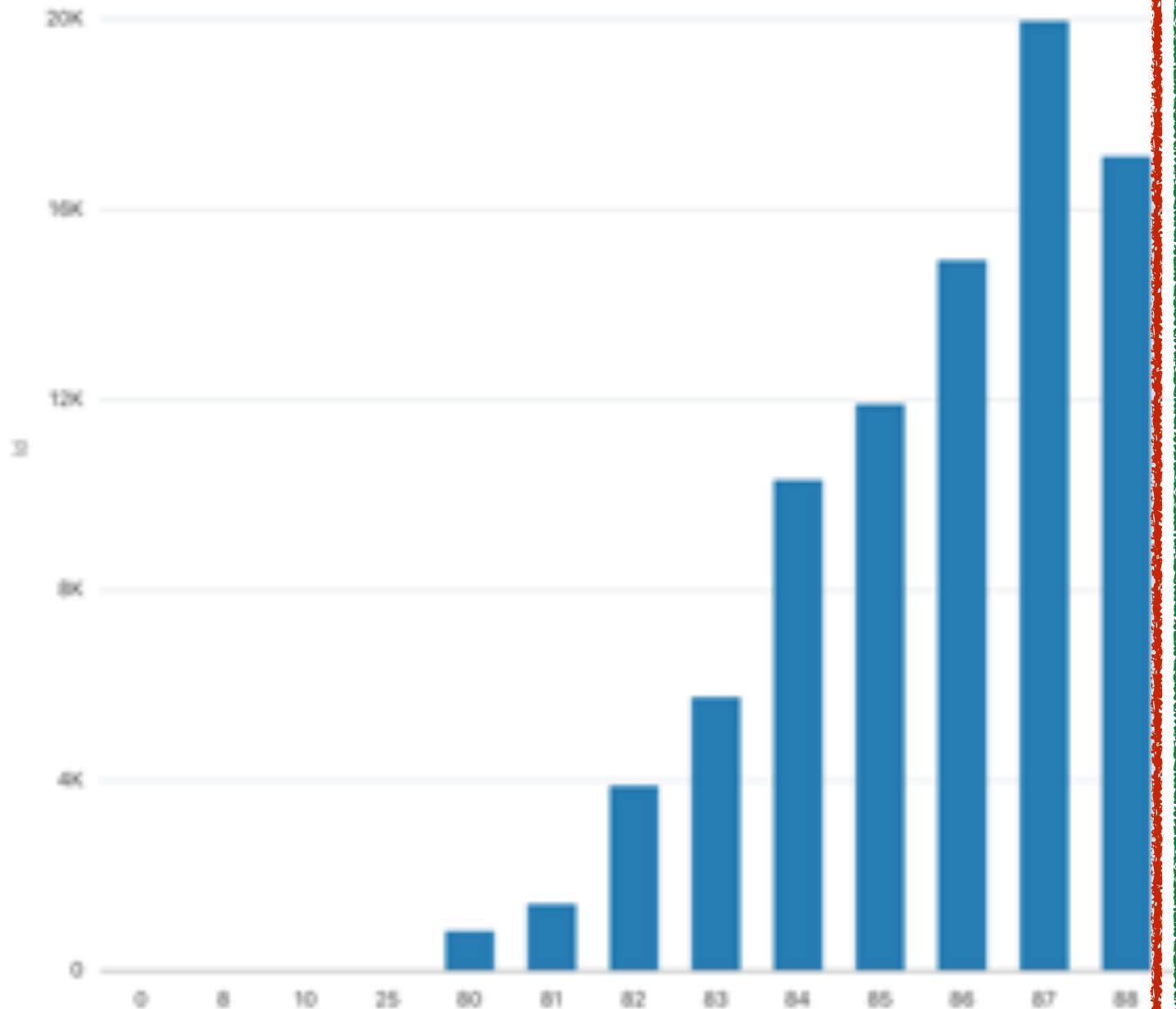
- #
- A **country** The country that the wine is from
- A **description**
- A **designation** The vineyard within the winery where the grapes that made the wine are from
- # **points** The number of points WineEnthusiast rated the wine on a scale of 1-100 (though they say they

# The Data

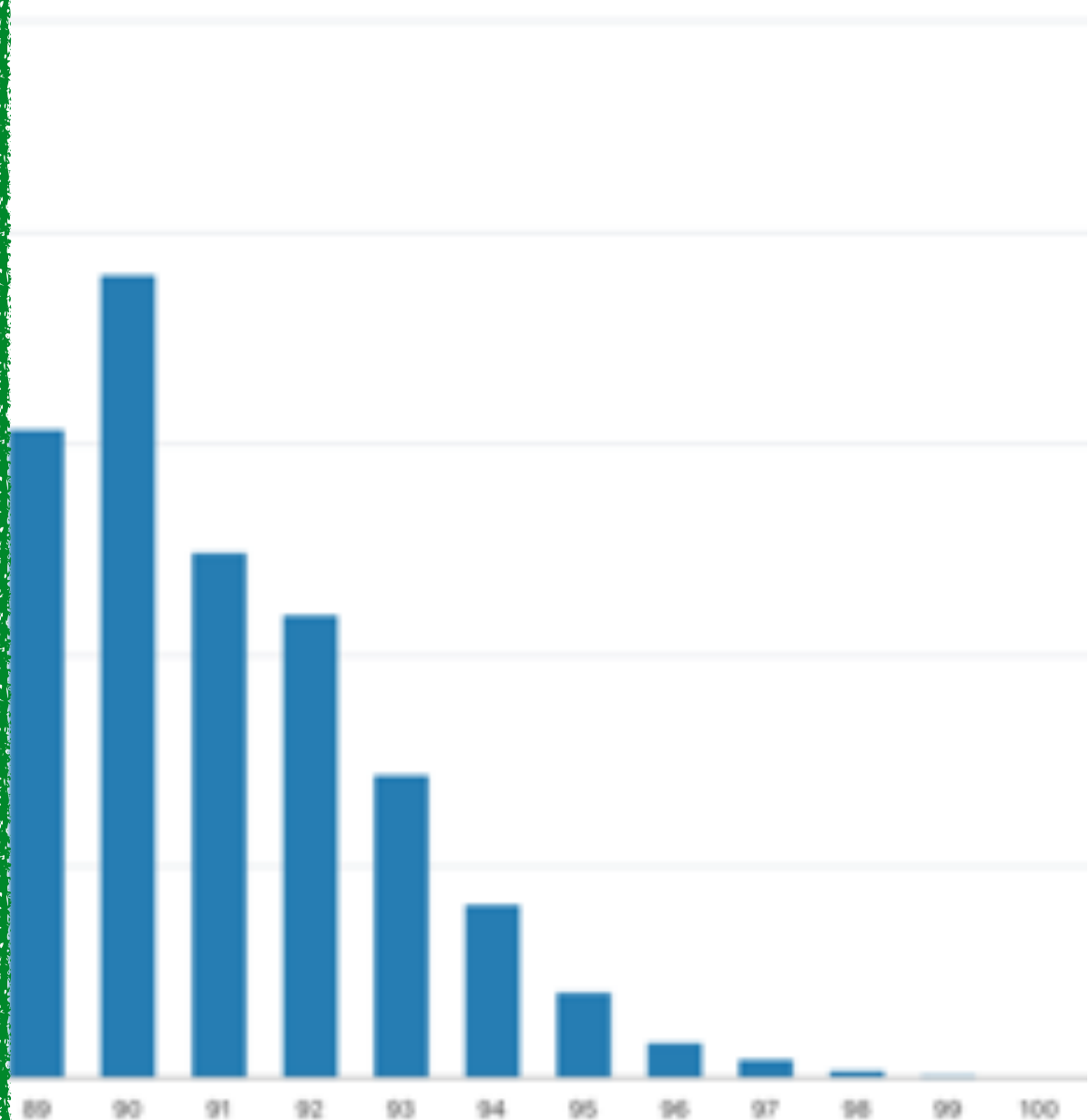
A	B	C	D	E	F	G	H	I	J	K
id	country	description	designation	points	price	province	region_1	region_2	variety	winery
0	US	This tremendous 100% Martha's Vineyard		96		235 California	Napa Valley	Napa	Cabernet Sauvignon	Heitz
1	Spain	Ripe aromas of fig, bla	Canodorum Selección	96		110 Northern Spain	Toro		Tinta de Toro	Bodega Carmen Rodríguez
2	US	Mac Watson honors ti	Special Selected Late	96		90 California	Knights Valley	Sonoma	Sauvignon Blanc	Macaulay
3	US	This spent 30 months	Reserve	96		65 Oregon	Willamette Valley	Willamette Valley	Pinot Noir	Ponzi
4	France	This is the top wine fr	La Brûlée	95		66 Provence	Bandol		Provence red blend	Domaine de la Brûlée
5	Spain	Deep, dense and pure	Numantia	95		73 Northern Spain	Toro		Tinta de Toro	Numantia
6	Spain	Slightly gritty black-fr	San Román	95		65 Northern Spain	Toro		Tinta de Toro	Maurinos
7	Spain	Lush cedary black-fru	Canodorum Vínico D	95		110 Northern Spain	Toro		Tinta de Toro	Bodega Carmen Rodríguez
8	US	This re-named vineyar	Slice	95		65 Oregon	Chehalem Mountains	Willamette Valley	Pinot Noir	Bergström
9	US	The producer sources	Gap's Crown Vineyard	95		60 California	Sonoma Coast	Sonoma	Pinot Noir	Blue Farm
10	Italy	Elegance, complexity	Ronco della Chiesa	95		80 Northeastern Italy	Collio		Friulano	Borgo del Tiglio
11	US	From 18-year-old vine	Estate Vineyard Waco	95		48 Oregon	Ribbon Ridge	Willamette Valley	Pinot Noir	Patricia Green Cellars
12	US	A standout even in th	Weber Vineyard	95		48 Oregon	Dundee Hills	Willamette Valley	Pinot Noir	Patricia Green Cellars
13	France	This wine is in peak c	Château Montus Pre	95		90 Southwest France	Madiran		Tannat	Vignobles Brumont
14	US	With its sophisticated	Grace Vineyard	95		185 Oregon	Dundee Hills	Willamette Valley	Pinot Noir	Domaine Serene
15	US	First made in 2006, th	Sigrid	95		90 Oregon	Willamette Valley	Willamette Valley	Chardonnay	Bergström
16	US	This blockbuster, pow	Rainin Vineyard	95		325 California	Diamond Mountain D	Napa	Cabernet Sauvignon	Hall
17	Spain	Nicely oaked blackber	6 Añitos Reserva Pre	95		80 Northern Spain	Ribera del Duero		Tempranillo	Valduero
18	France	Coming from a seven-	Le Pigeonnier	95		290 Southwest France	Cahors		Malbec	Château Lagry/Dette
19	US	This fresh and lively	Gap's Crown Vineyard	95		75 California	Sonoma Coast	Sonoma	Pinot Noir	Gary Farrell

# Good/Bad

Bad



Good



# Become a Data Scientist with OAC

Connect

Clean

Transform  
&  
Enrich

Analyse

Train  
&  
Evaluate

Predict

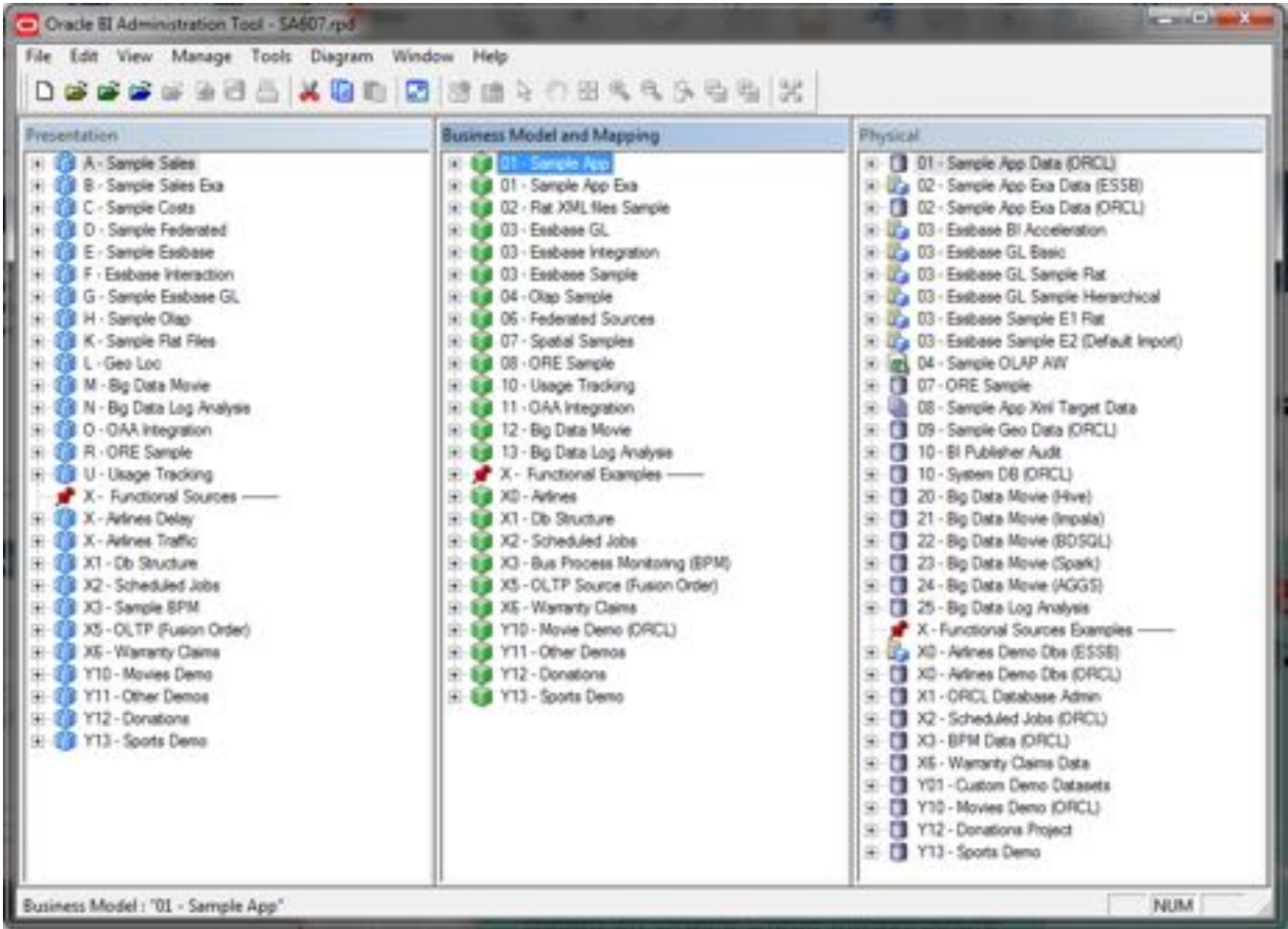


# Connect



## Data Sources

## Pre-Defined Data Models



# Clean

N/A

Missing Values

Mark <> MArk

Wrong Values

City  
“Rome”

Irrelevant Observations

Col1 -> Name

Labelling Columns

Role: CIO  
Salary:500 K\$

Handling Outliers

0-200k  
0-1

Feature Scaling

# Of Clicks

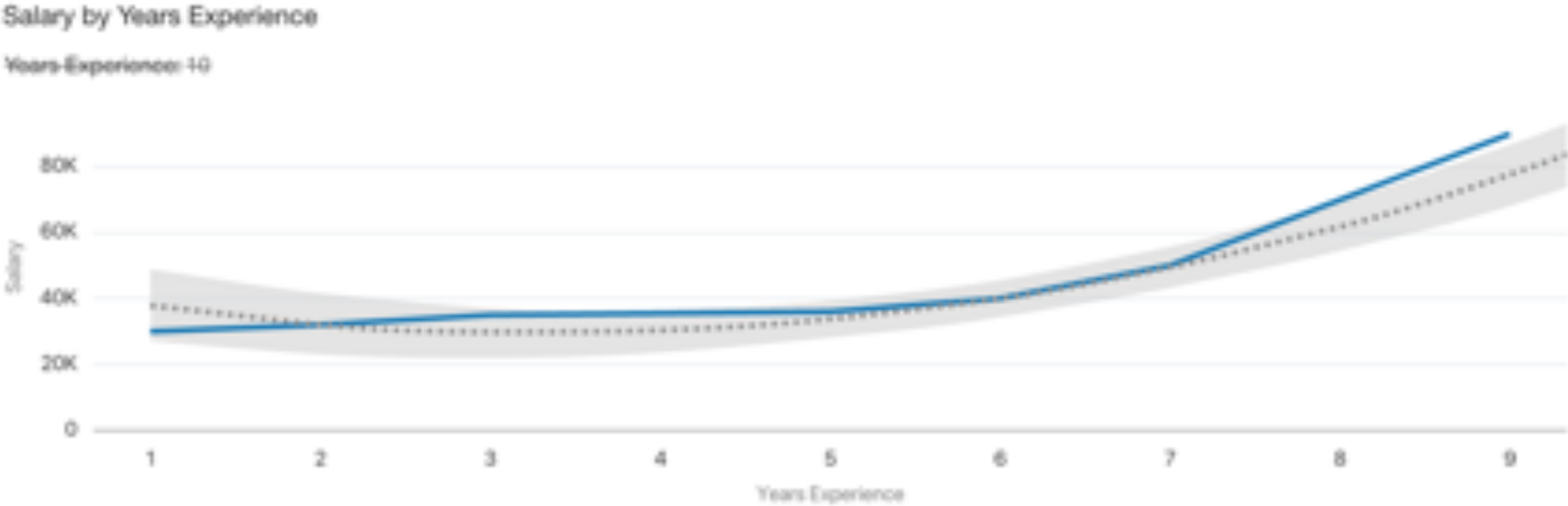
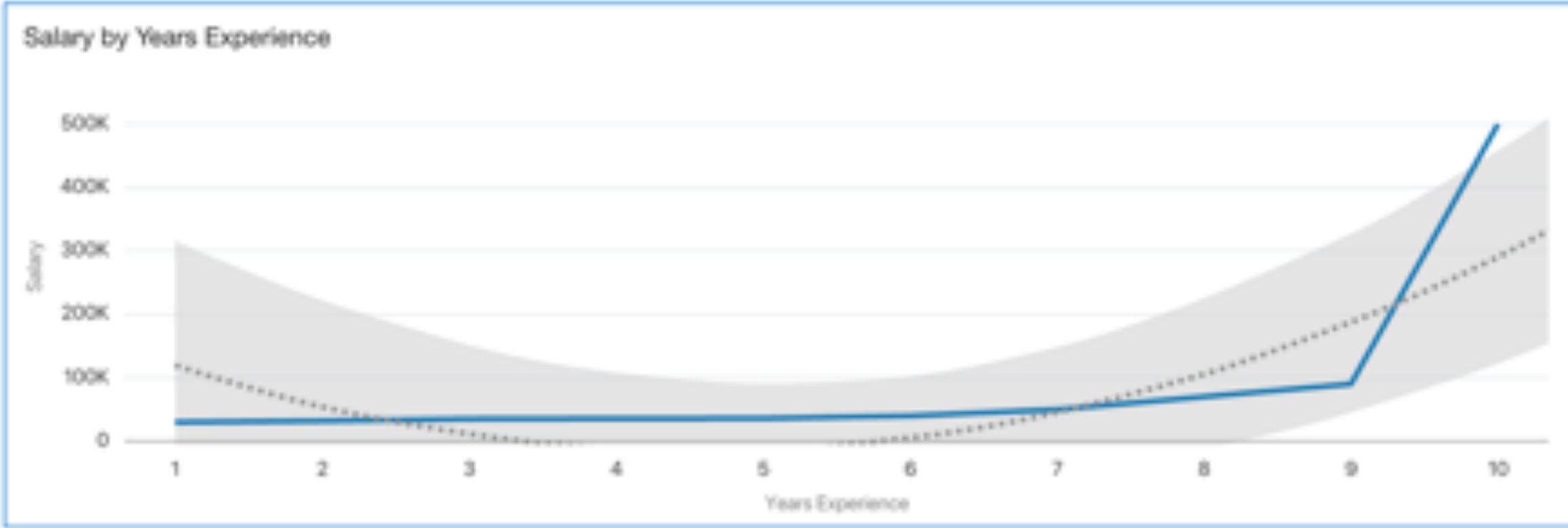
Aggregation

Train: 80%  
Test: 20%

Train/Test Set Split

# Why Removing an Outlier?

Years Experience	Salary
1	30.000
2	32.000
3	35.000
4	35.500
5	36.000
6	40.000
7	50.000
8	70.000
9	90.000
10	500.000



# Enrich - Feature Engineering

Location -> ZIP Code

**Additional  
Data  
Sources?**

Name -> Sex

2 Locations -> Distance

**Data Flow**

Day/Month/Year -> Date

# Data Preparation Recommendations

[illegible]



# Analyse - Data Overview

Results

Metadata

Data Element	Data Type	Treat As	Aggregation	Sample Values
id	varchar(80)	Attribute	none	1470; 817; 1028; 632; 3689; 4148; 2576; 963; 4979; 281
country	varchar(137)	Attribute	none	US; France; Italy; Spain; Portugal; Germany; Argentina; Chile; Austria; Greece
country_continent	varchar(4000)	Attribute	none	NA; EU
country_fips	varchar(4000)	Attribute	none	US; IT; FR
country_iso3	varchar(4000)	Attribute	none	USA; FRA; ITA
country_iso_numeric	number	Measure	sum	840; 380; 250
country_iso2	varchar(4000)	Attribute	none	US; IT; FR
description	varchar(1247)	Attribute	none	This elegant wine combines subtle nutmeg and cardamom aromas with crisp app...
designation	varchar(122)	Attribute	none	Reserve; Estate; Reserva; Riserva; Estate Bottled; Vieilles Vignes; Crianza; Classic...
points	number	Measure	sum	90; 89; 88; 87; 91; 86; 92; 93; 85; 94
price	varchar(15)	Attribute	none	25; 20; 40; 18; 60; 30; 28; 35; 50; 15
province	varchar(53)	Attribute	none	California; Oregon; Bordeaux; Tuscany; Piedmont; Washington; Northern Spain; M...
region_1	varchar(75)	Attribute	none	Willamette Valley; Napa Valley; Barolo; Brunello di Montalcino; Russian River Valle...
region_2	varchar(35)	Attribute	none	Central Coast; Sonoma; Willamette Valley; Napa; Columbia Valley; Mendocino/La...
variety	varchar(53)	Attribute	none	Pinot Noir; Chardonnay; Bordeaux-style Red Blend; Cabernet Sauvignon; Red Ble...
winery	varchar(84)	Attribute	none	Tarara; Heron Hill; Byron; Bergström; Herdade do Rocim; Rusack; Sarah's Viney...

# Analyse - Explain

# points

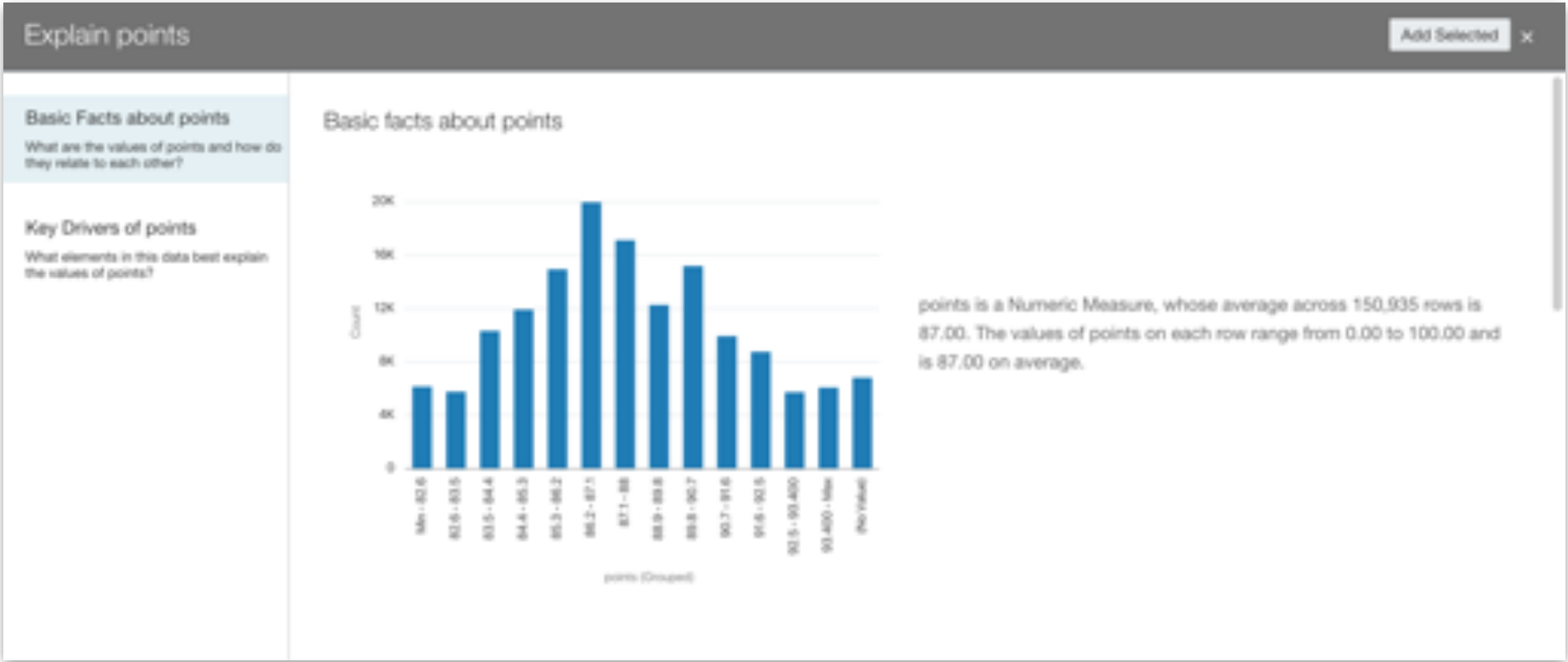
Add to Selected Visualization

Create Best Visualization

Pick Visualization...

Create Filter

Explain points

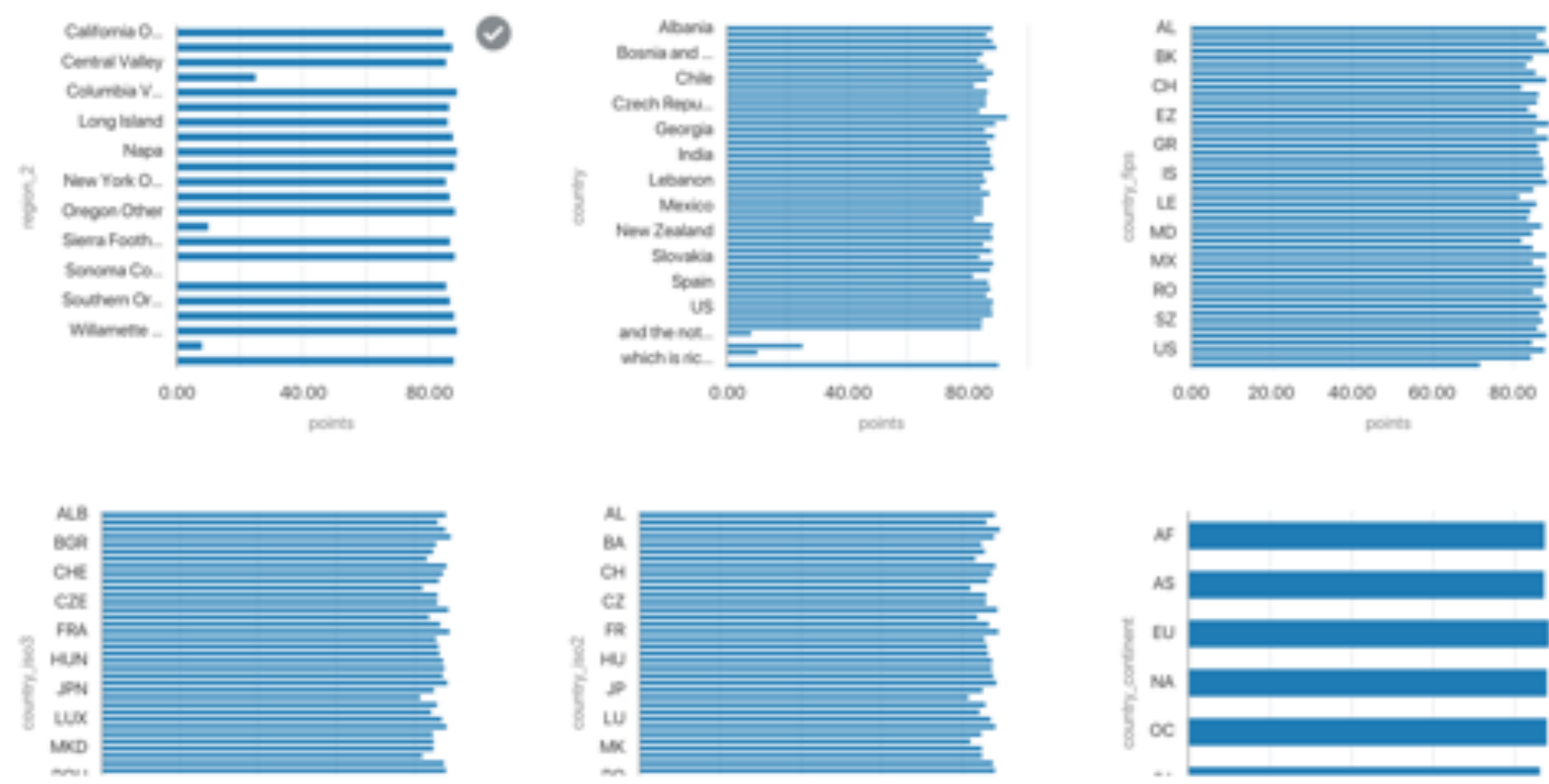


# Explain - Key Drivers

## Key Drivers of points

The 6 attributes most strongly correlated to outcomes for points are: **region\_2**, **country**, **country\_fips**, **country\_iso3**, **country\_iso2**, **country\_continent**

The charts below show the distribution of points values across each of the key drivers. Click the checkmarks above any of the visuals to add them to your project when done.

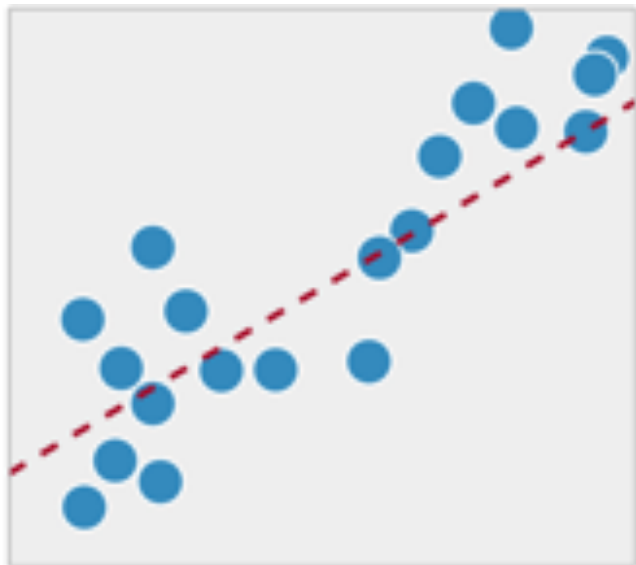


# Train - What Problem are we Trying to Solve?

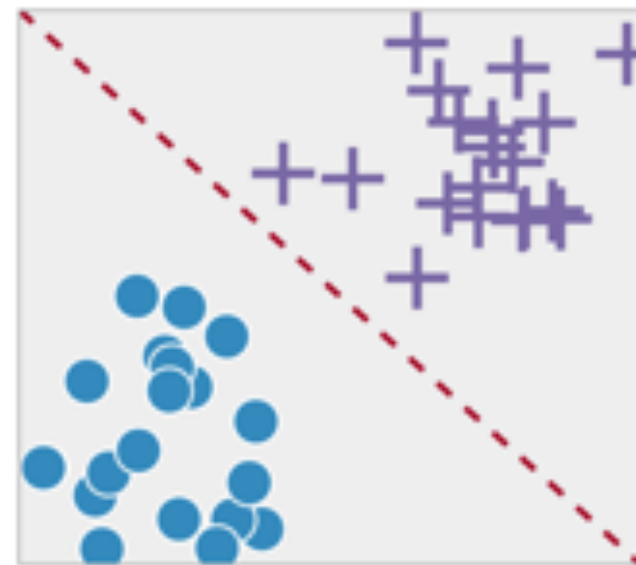
## Supervised

“I want to predict the value of Y,  
here are some examples”

### Regression



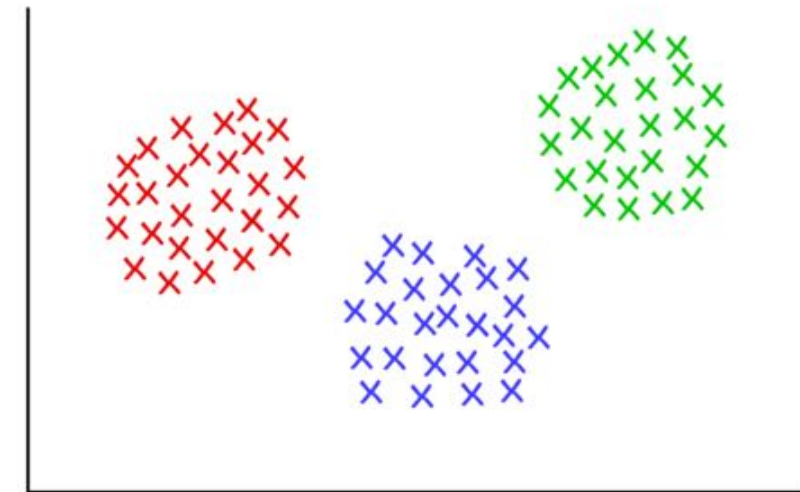
### Classification



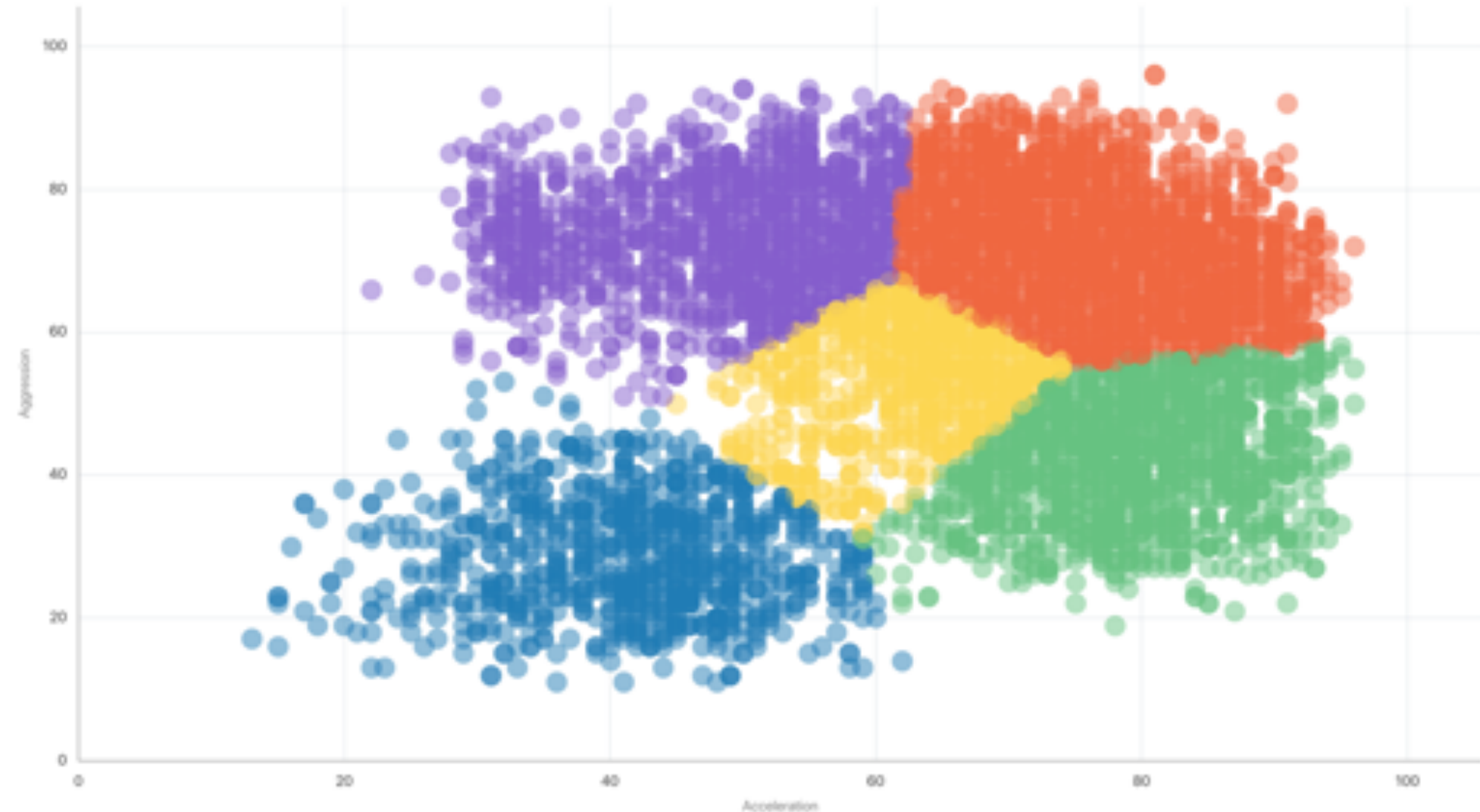
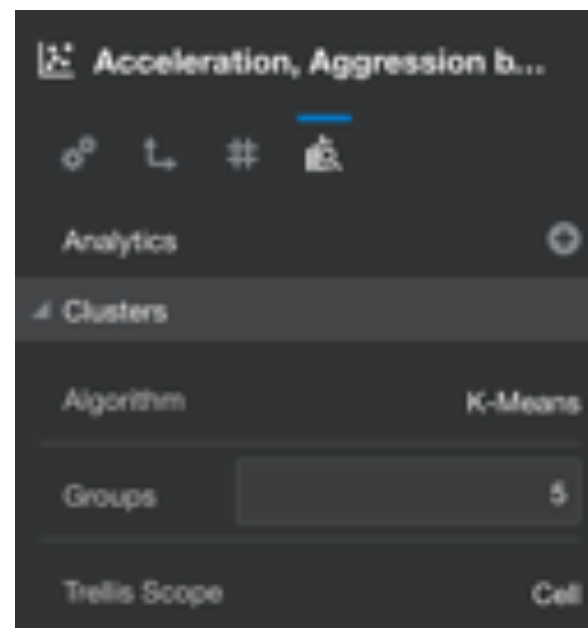
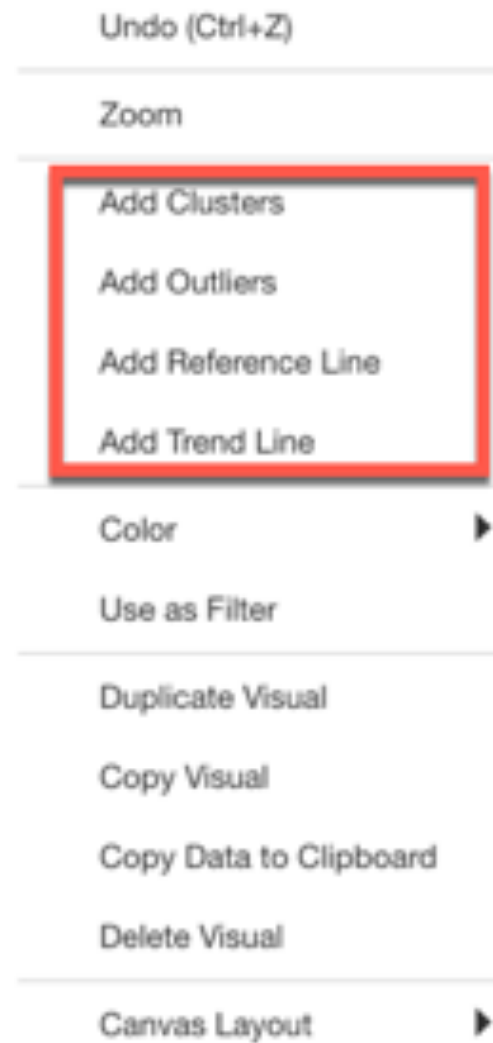
## Unsupervised

“Here is a dataset,  
make sense out of it!”

### Clustering



# Easy Models










# DataFlow Train Model





Select Train Numeric Prediction Model Script

 Linear Regression for model training




 Elastic Net Linear Regression for model training


 Random Forest for Numeric model training

 CART for Numeric Prediction training


# Which Model - Parameters To Pick?

Select Train Numeric Prediction Model Script


  




Linear Regression for model training



Elastic Net Linear Regression for model training



Random Forest for Numeric model training



CART for Numeric Prediction training

Train Numeric Prediction

Model Training Script

Linear Regression for model training

• Target

Select a column  
target, the target(label) to learn/predict

Regression Method

Lasso

Method for linear regression training.

Regularization Weight

1

Regularization Weight(L1 Ratio or L2 Ratio). Please enter 0 if it is Ordinary Least Squares linear regression.

Categorical Column Imputation

Most Frequent

The mode method for categorical features to fill NA. Two options: most frequent and least frequent. Default is most frequent.

Numerical Column Imputation

Mean

The mode method for numeric features to fill NA. Four options: mean, max, min, median. Default is mean.

Categorical Encoding Method

Indexer

Encoding method.

Maximum Null Value Percent

80

Maximum Null Value Percent

Train Partition Percent

80

# Select, Try, Save, Change, Try, Save .....



## Train Numeric Prediction

Model Training Script **Elastic Net Linear Regression for model training**

Target **points**  
target, the target(label) to learn/predict

L1 Ratio   
L1 Ratio

L2 Ratio   
L2 Ratio

Categorical Column Imputation **Most Frequent**  
The mode method for categorical features to fill NA  
Two options: most frequent and least frequent. Default is most frequent.

### Select Train Numeric Prediction Model Script

Linear Regression  
for model training

Elastic Net Linear  
Regression for  
model training

Random Forest for  
Numeric model  
training

CART for Numeric  
Prediction training

Data SetsConnectionsData FlowsSequences

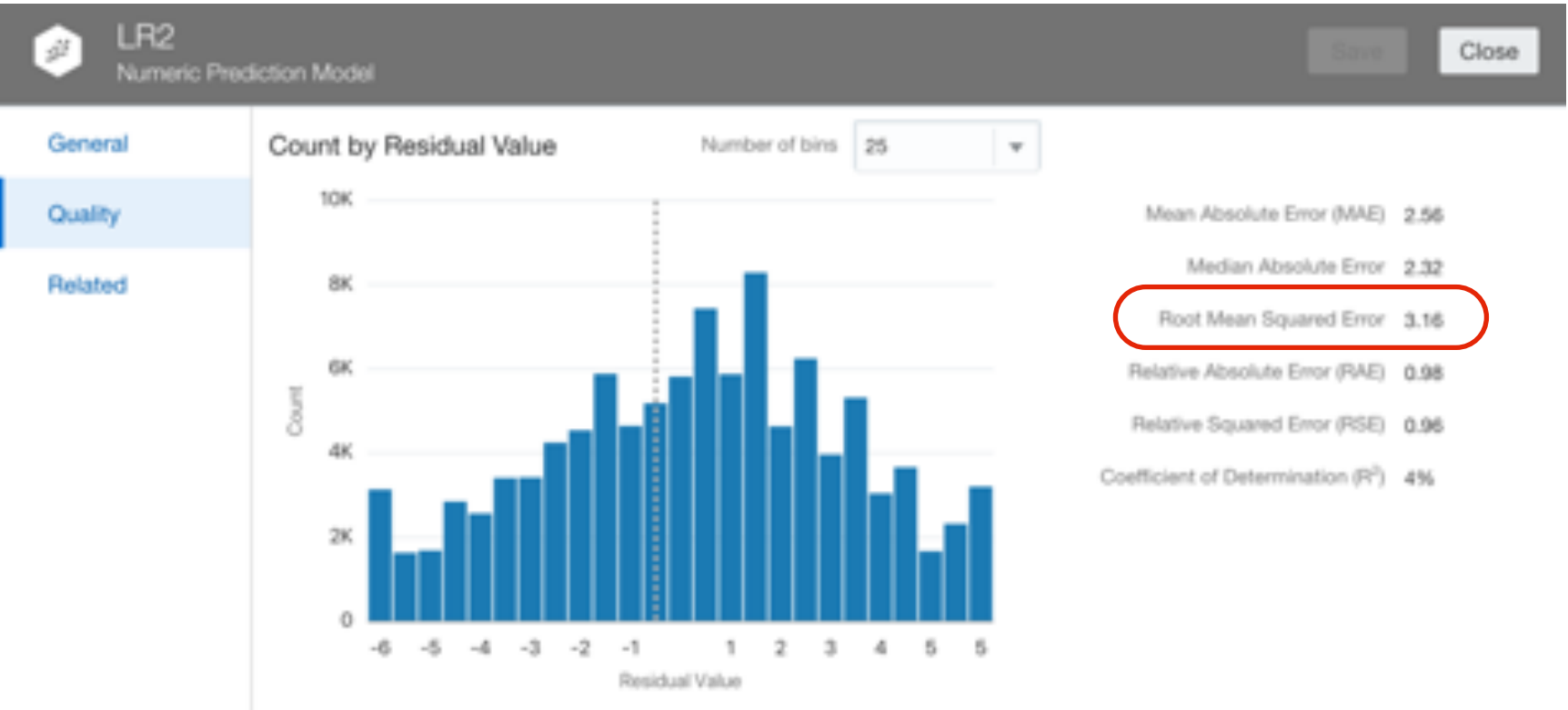
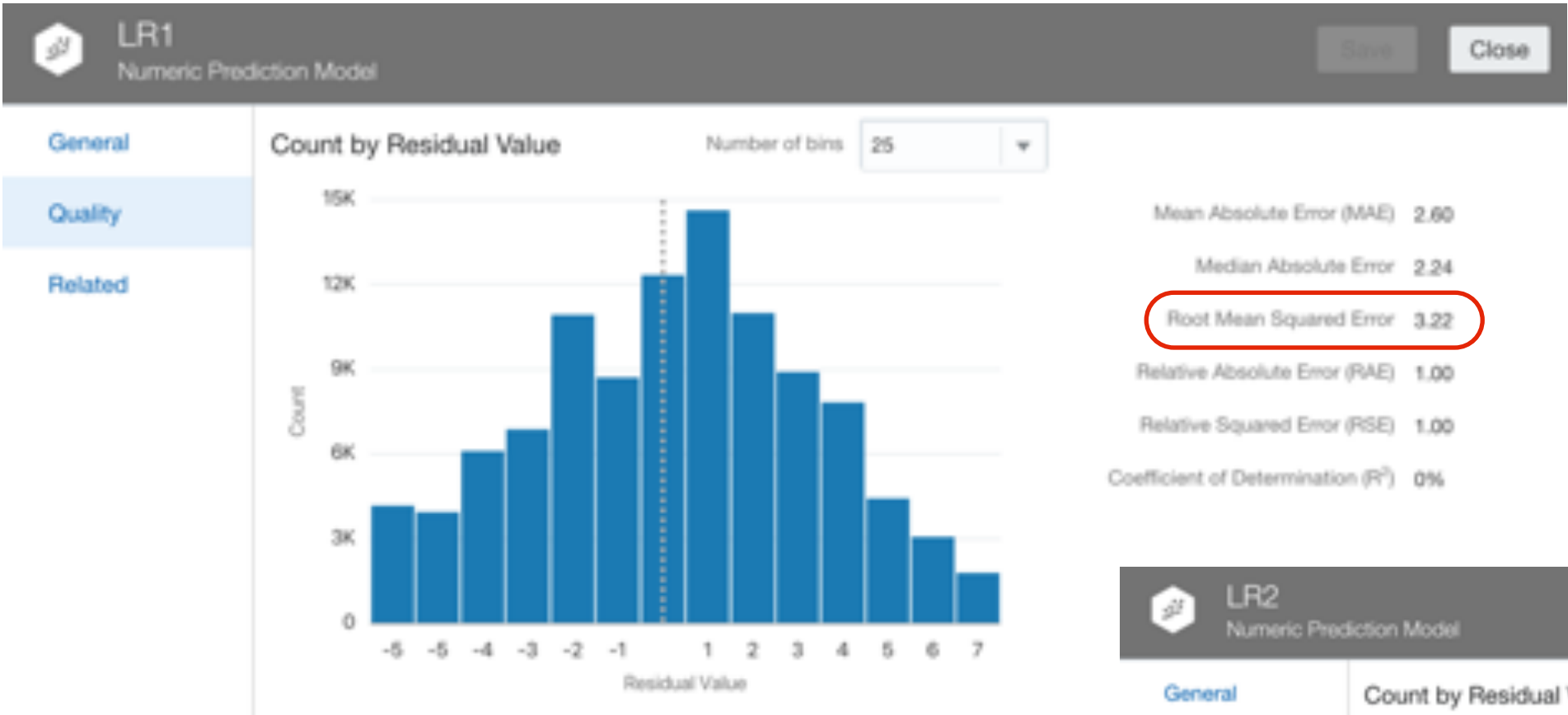
Type	Name
>>>	ELN1
>>>	LR2
>>>	LR1

Machine Learning

ScriptsModels

Type	Name
	ELN1
	LR2
	LR1

# Compare



# Compare

		Predicted Values		
Actual Values		0.0	1.0	Total
	0.0	40439	471	40910 (90%)
	1.0	3761	866	4627 (10%)
	Total	44200 (97%)	1337 (3%)	45537 (100%)



# There is No Single Truth...

Predicted Values				
	0.0	1.0	Total	
Actual Values	0.0	40408	502	40910 (90%)
	1.0	3731	896	4627 (10%)
	Total	44139 (97%)	1398 (3%)	45537 (100%)

$502/(502+896) = 64.09\%$

$471/(471+866)=64.77\%$

Predicted Values			
	0.0	1.0	Total
0.0	40439	471	40910 (90%)
1.0	3761	866	4627 (10%)
Total	44200 (97%)	1337 (3%)	45537 (100%)

# Predict - Use On the Fly







Add Data Set...

Create Scenario...

Add Calculation...

Create Scenario - Select Model

Search

Type	Name
	BinaryCart2
	BinaryCart1
	BinaryLogistic1
	ELN1
	LR2
	LR1

Edit Scenario - Map Your Data

Select which Data Set you want to use with the Model

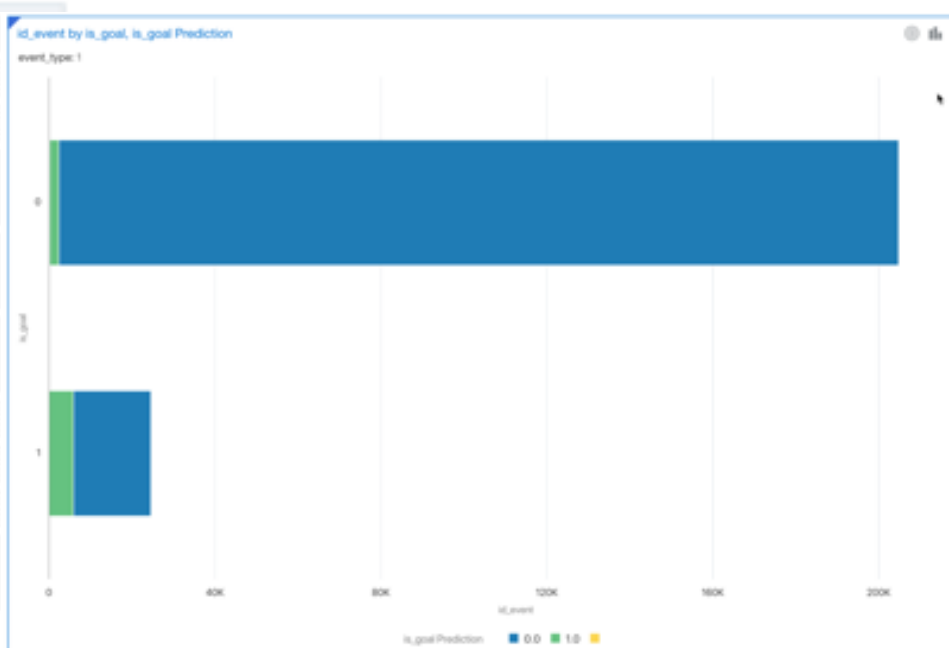
Data Set

FootballEvents

For each model input listed on the left, select a corresponding data element from your project

Model Input	Map To	
bodypart	bodypart	
location	location	
player	player	
situation	situation	
is_goal	is_goal	

Required Fields



# Predict - Step of a Data Flow



## Select Model

<div>Search</div>			
Type	Name	Outputs	Modified
	BinaryCart2	is_goal	54 minutes ago
	BinaryCart1	is_goal	51 minutes ago
	BinaryLogistic1	is_goal	12:25 PM
	ELN1	points	11:56 AM
	LR2	points	11:36 AM
	LR1	points	10:31 AM

Model **BinaryCart2**

### Outputs

Create	Output	Column Name
<input checked="" type="checkbox"/>	PredictedValue	<div>PredictedValue</div>
<input checked="" type="checkbox"/>	PredictionConfidencePercentage	<div>PredictionConfidencePercentage</div>
<input checked="" type="checkbox"/>	PredictionGroup	<div>PredictionGroup</div>

### Inputs

Model	Input
situation	<div>situation</div>
bodypart	<div>bodypart</div>
player	<div>player</div>

# Demo





# Conclusions

73% > 63% > 50%

Data Cleaning & Transformation

Model Creation & Evaluation

Trial Error

Visual -> UI Driven

Existing Skillset



# Example

## Wine Ratings Prediction using Machine Learning



Olivier Goutay [Follow](#)

Jun 14, 2018 · 5 min read



Image from [wall2born.com](http://wall2born.com)

# Custom ML Model

Available only in OAC Classic

A MEAL  
WITHOUT WINE  
IS CALLED  
BREAKFAST

# Is It Corked? Wine Machine Learning Predictions with OAC

Francesco Tisiot  
BI Tech Lead at Rittman Mead



ODTUG  
Kscope19   
SEATTLE, WASHINGTON • JUNE 23-27

PLEASE FILL OUT  
YOUR EVALUATIONS

SEATTLE

 Washington State  
Convention Center