

HONORARY PATRONAGE  
RECTOR OF



# IMPROVING PREDICTIVE MODELING THE WAY FOR PERSONAL CARRIER

ADAM KARWAN, PHD

7 NOVEMBER 2017, WARSAW

Partner



Sponsor



Media Patronage



# Kaggle

The Home of Data Science and Machine Learning

- **WHAT IS KAGGLE** (Data Science as a Competition)



# kaggle

What could the world's best analysts find in your data?

- Data Sets
- Contests

# What is Kaggle

Competitions Datasets Kernels Discussion

Active

All


Entered

Sort by Prize

19 active competitions


All Categories

Search competitions




**Passenger Screening Algorithm Challenge**  
Improve the accuracy of the Department of Homeland Security's threat recognition algorithms  
**Featured** · a month to go · terrorism, image, object detection

**\$1,500,000**  
364 teams



**Zillow Prize: Zillow's Home Value Prediction (Zestimate)**  
Can you improve the algorithm that changed the world of real estate?  
**Featured** · 2 months to go · housing, real estate

**\$1,200,000**  
3,780 teams



**Statoil/C-CORE Iceberg Classifier Challenge**  
Ship or iceberg, can you decide from space?  
**Featured** · 3 months to go · weather, shipping, binary classification

**\$50,000**  
719 teams

- Self-Learning
- Recruitment
- Algorithm for Companies

# Interview Task for Boeing Digital Aviation Research



Boeing Global Services  
Digital Aviation & Analytics Lab  
Gdansk

## BIKE SHARING CHALLENGE

<https://www.kaggle.com/c/bike-sharing-demand>

## GOAL

combine **historical usage patterns** with **weather data** in order to **forecast bike rental demand** in the Capital Bike Share Program in Washington, D.C.



# BIKE SHARING CHALLENGE

Final Result & Presentation Hints



Algorithm



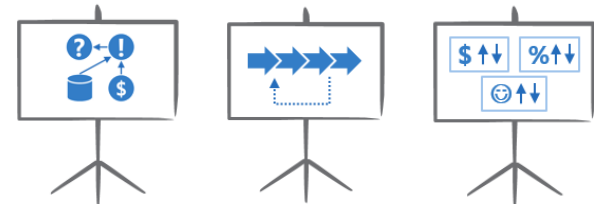
Presentation

- **Personal Goal** – create model in **TOP100** best solutions in 7 days

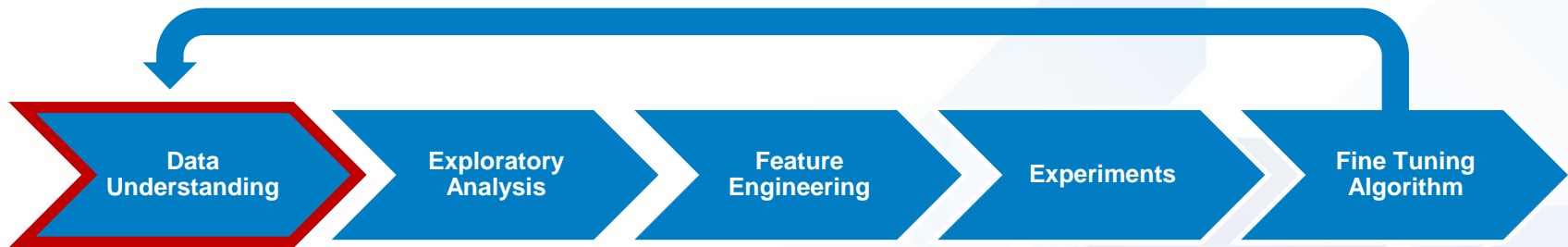


- **How to properly present a Data Mining project?**

- Start with big picture
- Overview of process
- Show the main outcome

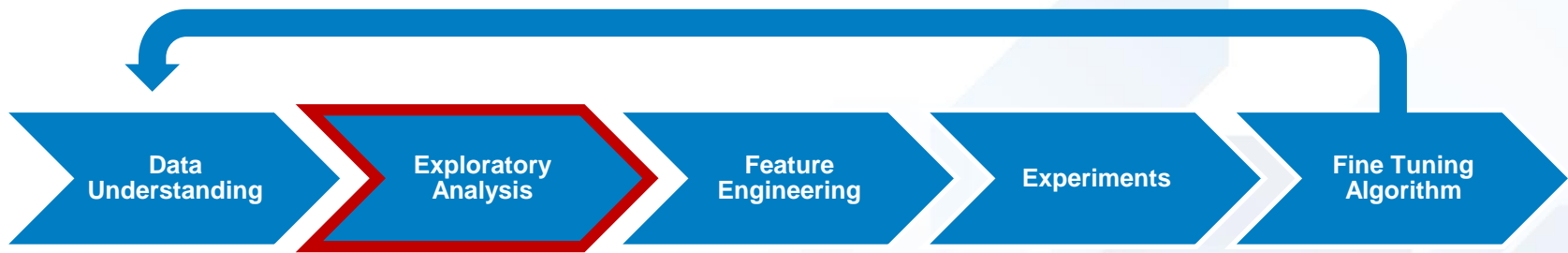


**Source** <http://algolytics.com/data-visualization-essentials-for-data-scientists>

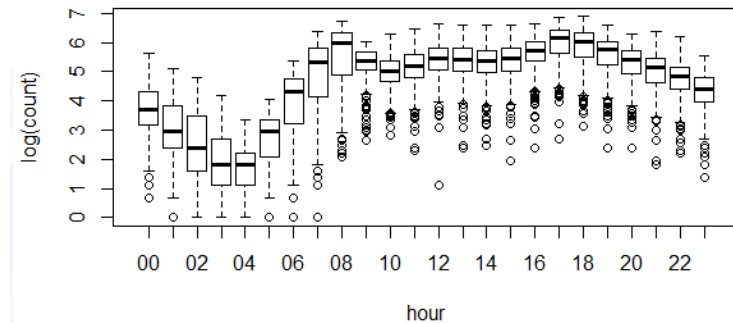
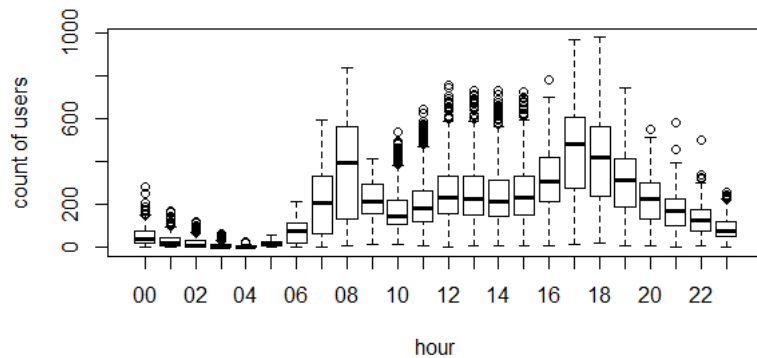


No	Attribute	Description
1	<b><i>datetime</i></b>	hourly date +timestamp
2	<b><i>season</i></b>	1 = spring 2 = summer 3 = fall 4 = winter
3	<b><i>holiday</i></b>	whether the day is considered a holiday (0; 1)
4	<b><i>workingday</i></b>	whether the day is neither a weekend nor holiday (0; 1)
5	<b><i>weather</i></b>	1: Clear, Few clouds, Partly cloudy, Partly cloudy 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds 4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
6	<b><i>temp</i></b>	temperature in Celsius
7	<b><i>atemp</i></b>	"feels like" temperature in Celsius
8	<b><i>humidity</i></b>	relative humidity
9	<b><i>windspeed</i></b>	wind speed
10	<b><i>casual</i></b>	number of non-registered user rentals initiated
11	<b><i>registered</i></b>	number of registered user rentals initiated
12	<b><i>count</i></b>	number of total rentals

TRAIN	10886
TEST	6493

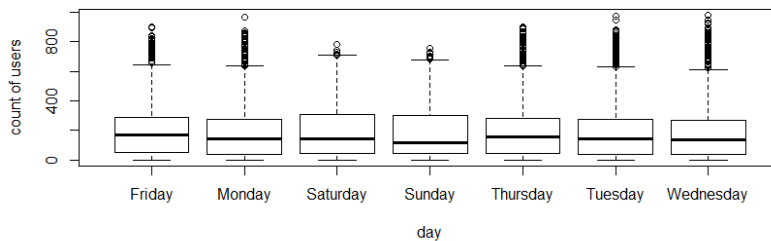


## ■ Hourly

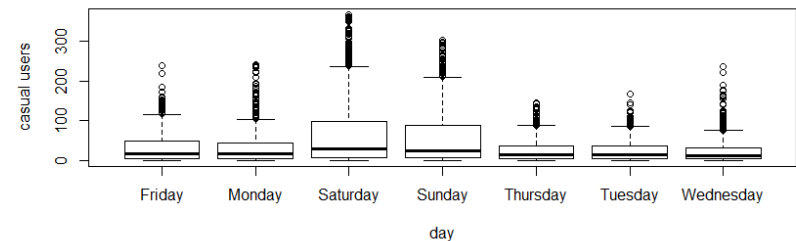


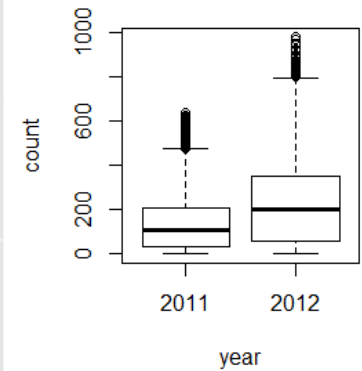
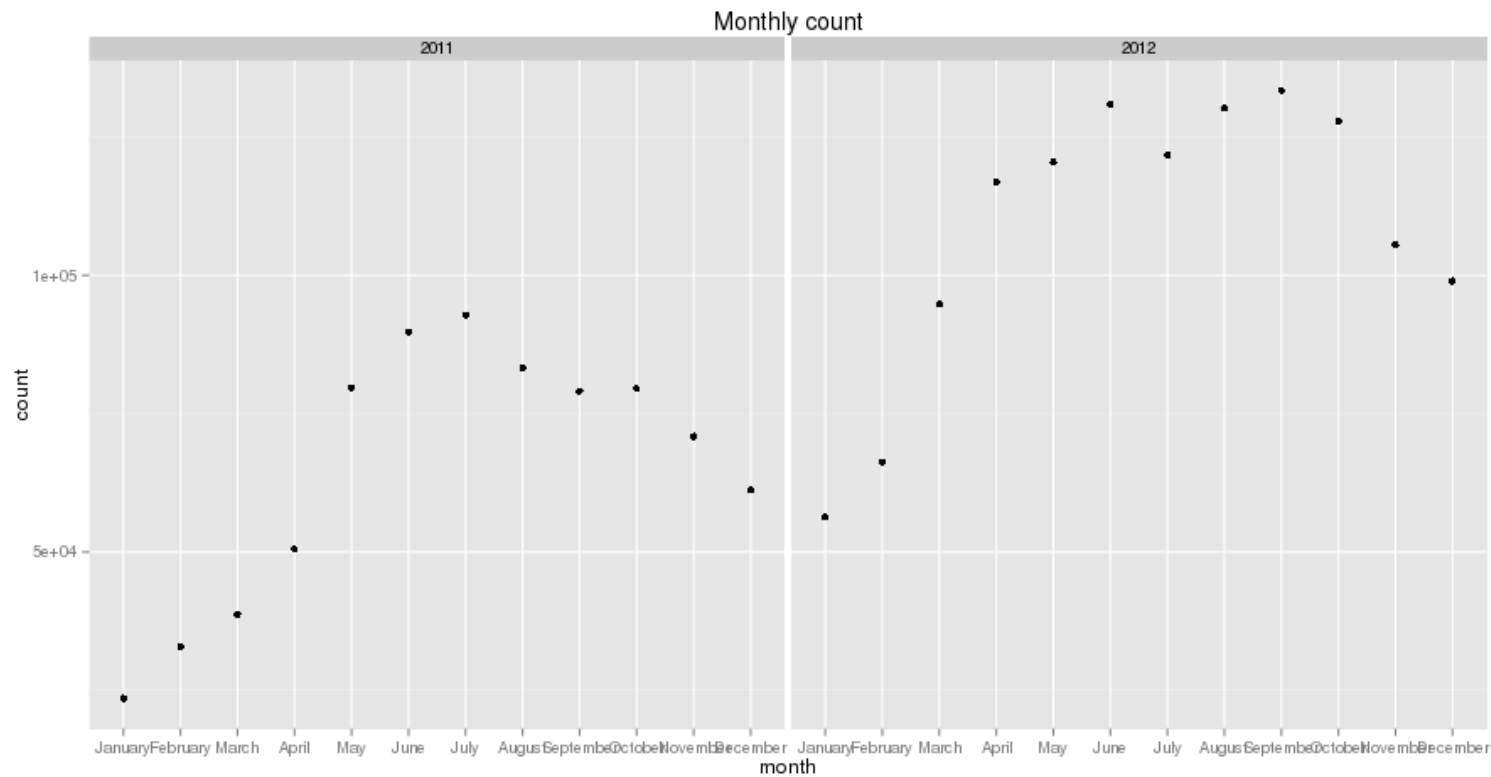
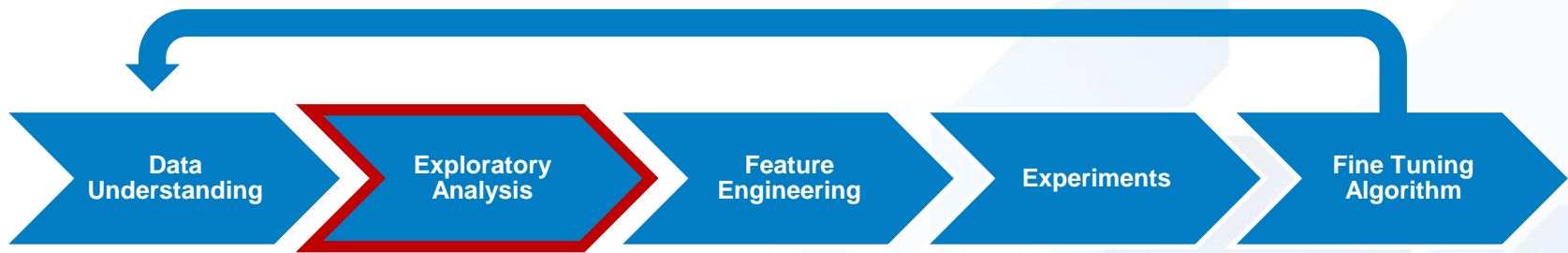
## ■ Daily

registered



casual

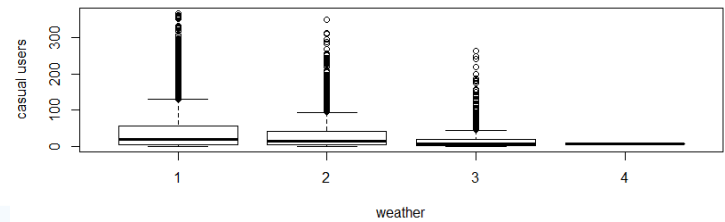
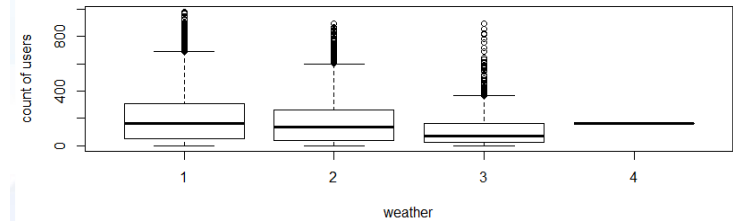
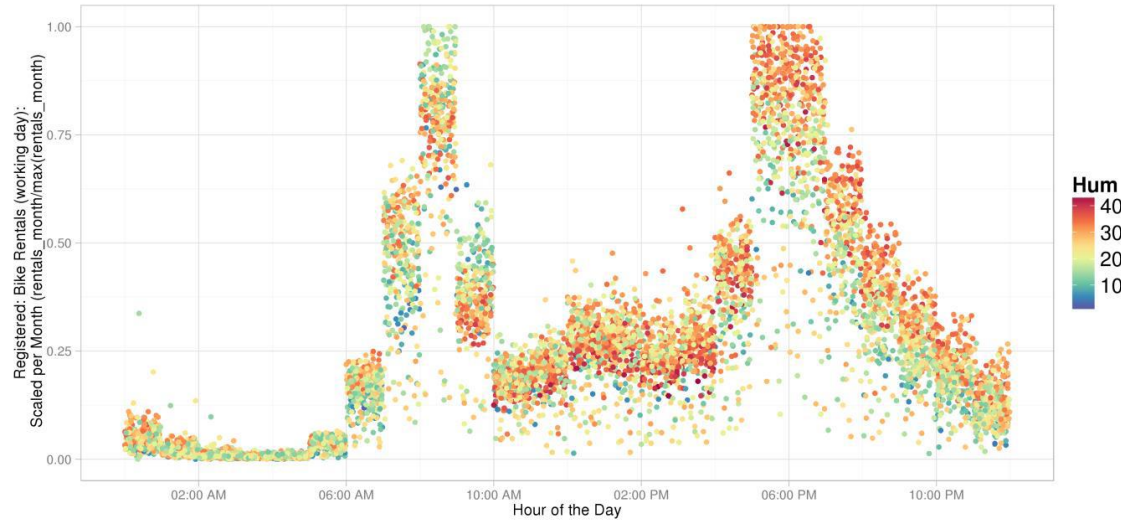




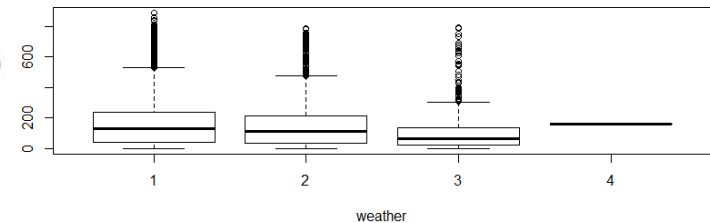
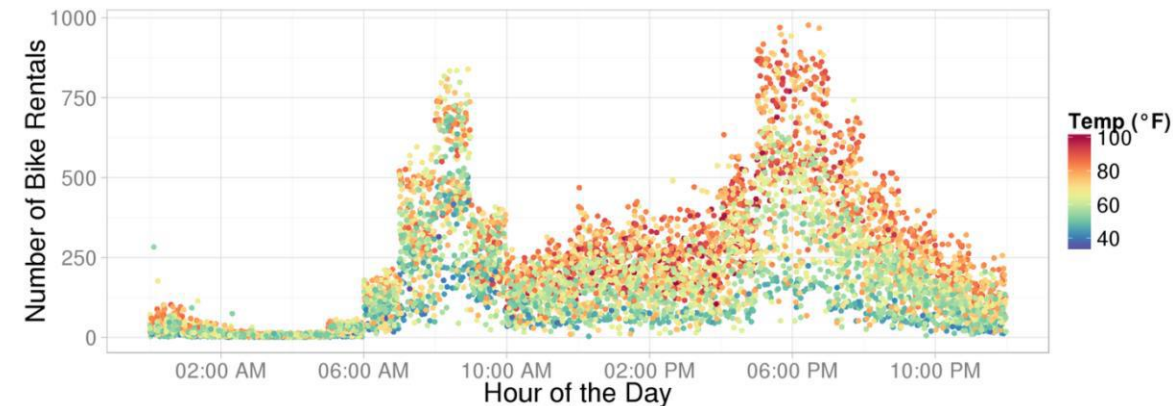


# Weather, Temperature, Humidity

On workingdays, any deducible effect of humidity, by any chance...?  
(taking the bike to work no matter what, but dry if back home...?)

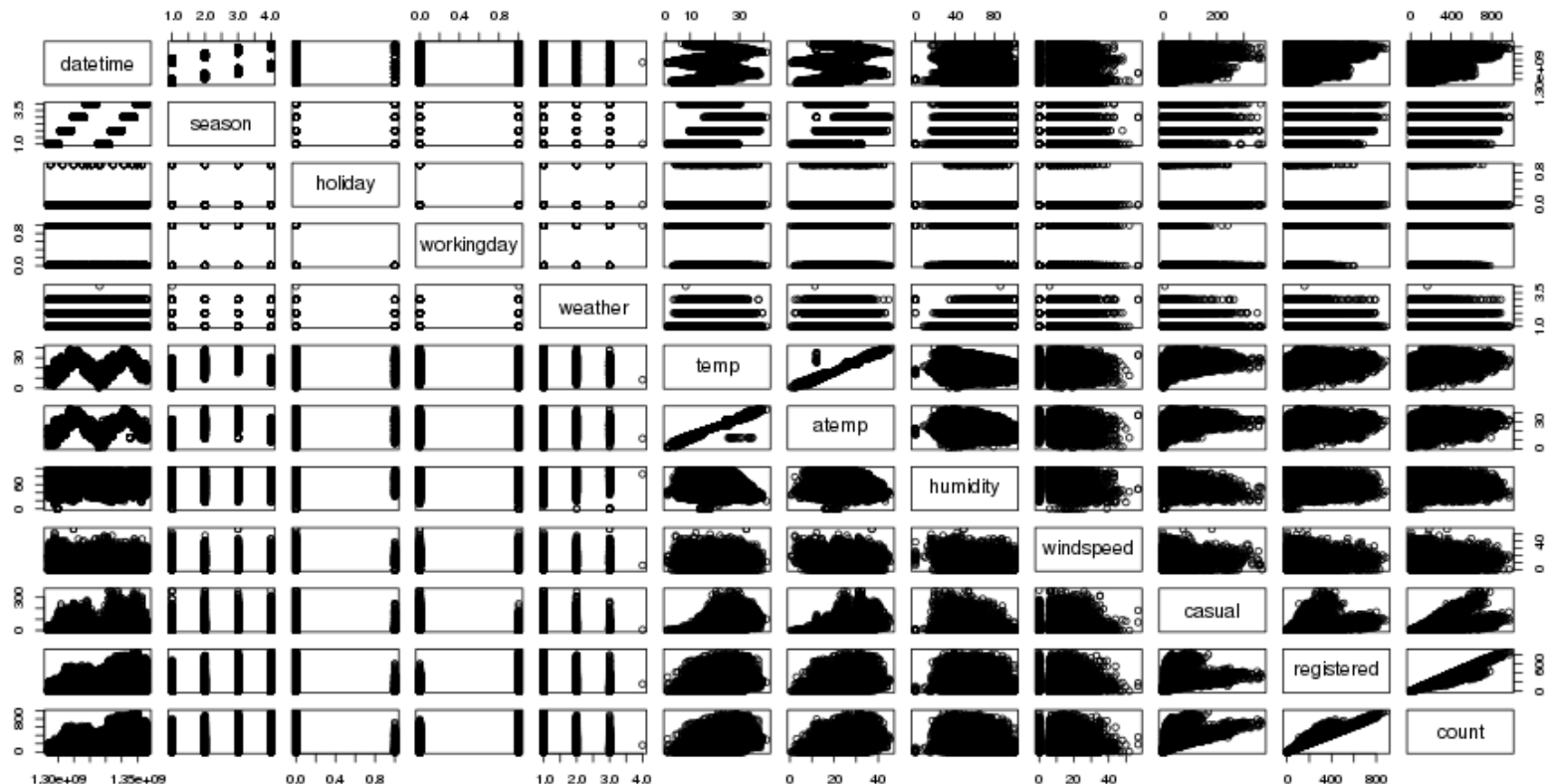


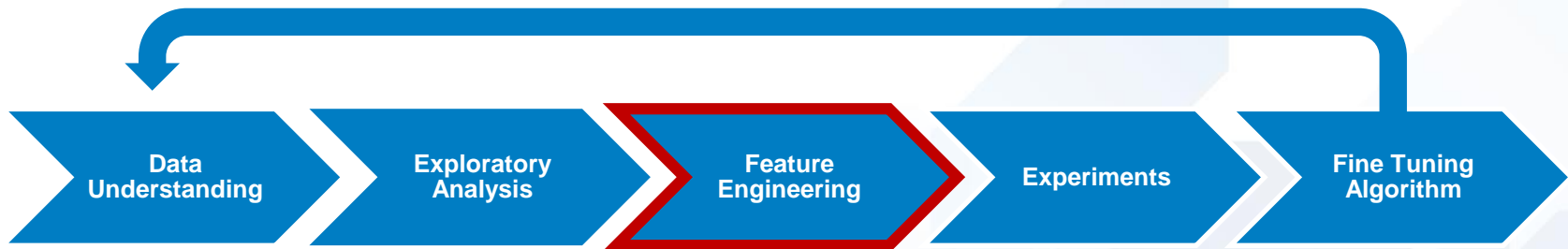
On workdays, most bikes are rented on warm mornings and evenings



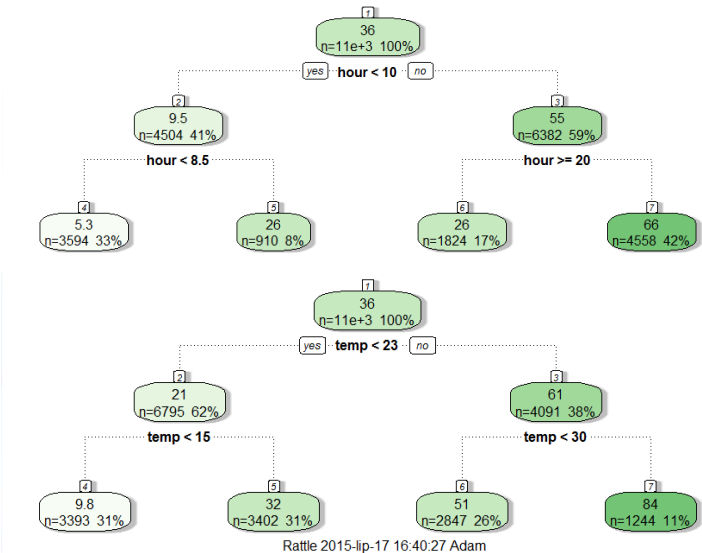
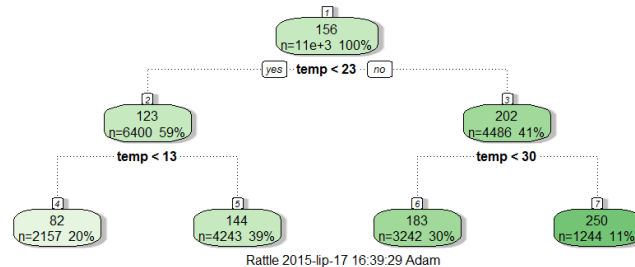
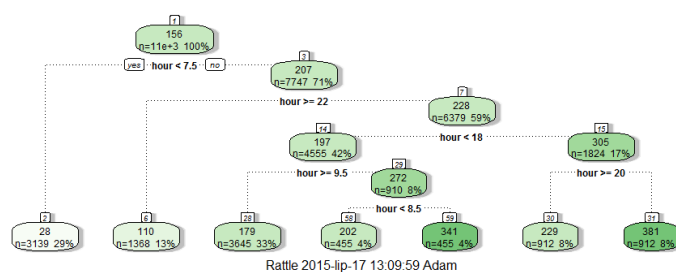
# Correlations

	train.registered	train.casual	train.count	train.temp	train.humidity	train.atemp	train.windspeed
train.registered	1.00	0.50	0.97	0.32	-0.27	0.31	0.09
train.casual	0.50	1.00	0.69	0.47	-0.35	0.46	0.09
train.count	0.97	0.69	1.00	0.39	-0.32	0.39	0.10
train.temp	0.32	0.47	0.39	1.00	-0.06	0.98	-0.02
train.humidity	-0.27	-0.35	-0.32	-0.06	1.00	-0.04	-0.32
train.atemp	0.31	0.46	0.39	0.98	-0.04	1.00	-0.06
train.windspeed	0.09	0.09	0.10	-0.02	-0.32	-0.06	1.00

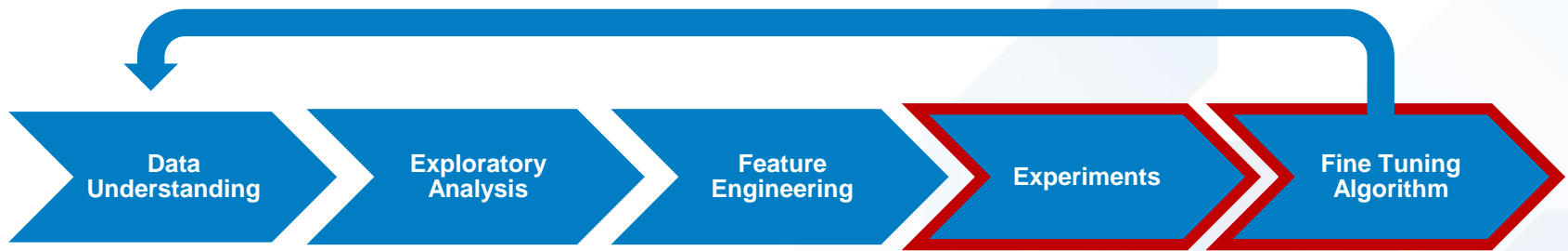




- *dp\_reg*
- *dp\_cas*
- *temp\_reg*
- *temp\_cas*

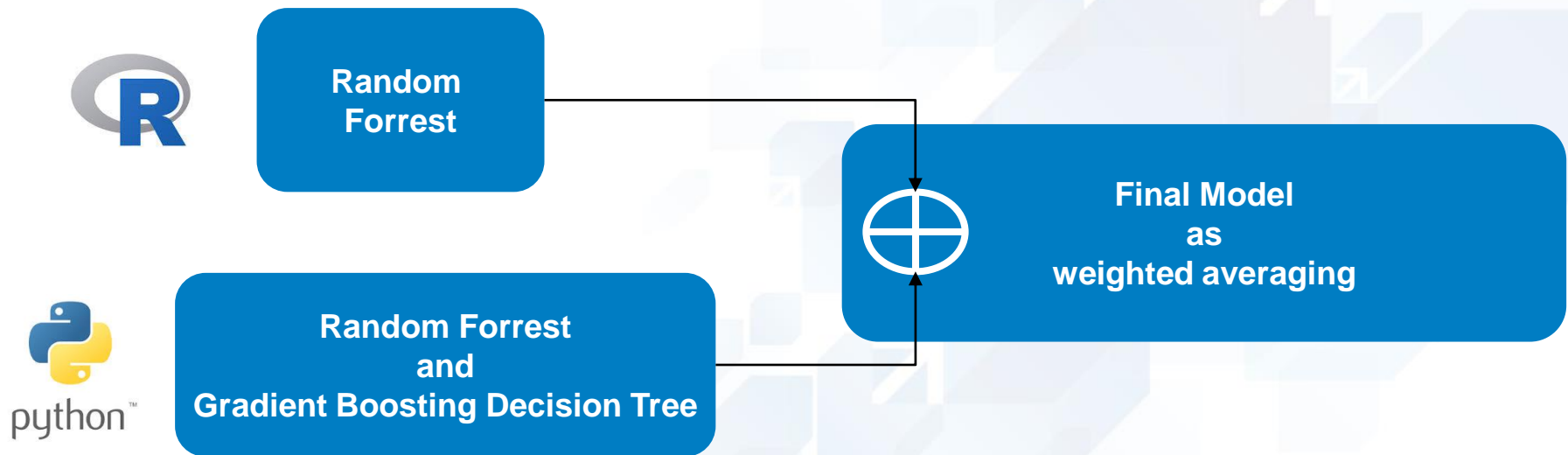


- *day\_type*
  - Holiday [holiday=0 and workingday=0]
  - Weekend [holiday=1]
  - Working Day [holiday=0 and workingday=1]
- *year\_part*
  - (1 – 8) from first quarter of 2011 till fourth of 2012



- **R model** bases on two separately computed **Random Forests** for *casual* and *registered* users. Final result is a sum of predicted values daily. Due to differences in *count* of users *hourly* we use logarithm to normalize values. Moreover each model was trained on 250 trees. In first model we estimate number of registered users therefore *dp\_reg* and *temp\_reg* are used, analogously *dp\_cas* and *temp\_cas* for second one. Other attributes used: *hour*, *day*, *day\_type*, *holiday*, *season*, *year*, *year\_part*, *weekend*, *workingday*, *atemp*, *humidity*, *weather*, *windspeed*.
- **Python model** bases on a combination of **RF (Random Forrest) and GBDT (Gradient Boosting Decision Trees)**. Twelve attributes are used: *hour*, *day*, *holiday*, *season*, *weekday*, *workingday*, *year*, *atemp*, *temp*, *humidity*, *weather*, *windspeed*. *Year* is normalized by substraction 2011. GBDT is computed on 100 and RF on 1000 trees. Before computing final result we compute average regression of two instances for each approach RF and GBDT. For estimated variables logarithm is used to normalize results.

## Ensamble Method



Final result **0.3704** that was **64** place out of **3252**

63	↓11	Dhruv Singal	0.37011	6	Sat, 11 Apr 2015 23:34:25 (-9.1d)
-		a.karwan	0.37040	-	Tue, 21 Jul 2015 20:24:47 Post-Deadline
<b>Post-Deadline Entry</b> If you would have submitted this entry during the competition, you would have been around here on the leaderboard.					
64	↓11	prashkr	0.37041	6	Thu, 09 Apr 2015 20:53:07 (-23.1h)



# JEPPESSEN®

A BOEING COMPANY



**Boeing Global Services**

Digital Aviation & Analytics Lab

Gdansk

## DATA SCIENCE in AVIATION INDUSTRY

- PREDICTIVE MX
- POST FLIGHT ANALYTICS
- REAL TIME MACHINE LEARNING





# Titanic Machine Learning from Disaster

<https://www.kaggle.com/c/titanic>

## More likely to survive

- Females
- Children
- 1<sup>st</sup> Class Passengers
- Traveling with Family

## More likely to perish

- Males
- Adults
- 2<sup>nd</sup> and 3<sup>rd</sup> Class Passengers
- Traveling alone



**Trevor Stephens**

Regular Data Scientist,  
Occasional Blogger.

📍 San Francisco, CA

🌐 Website

🐦 Twitter

🌐 LinkedIn

🐙 Github



43 1499 Total Recall

0.81340 8

### Your Best Entry

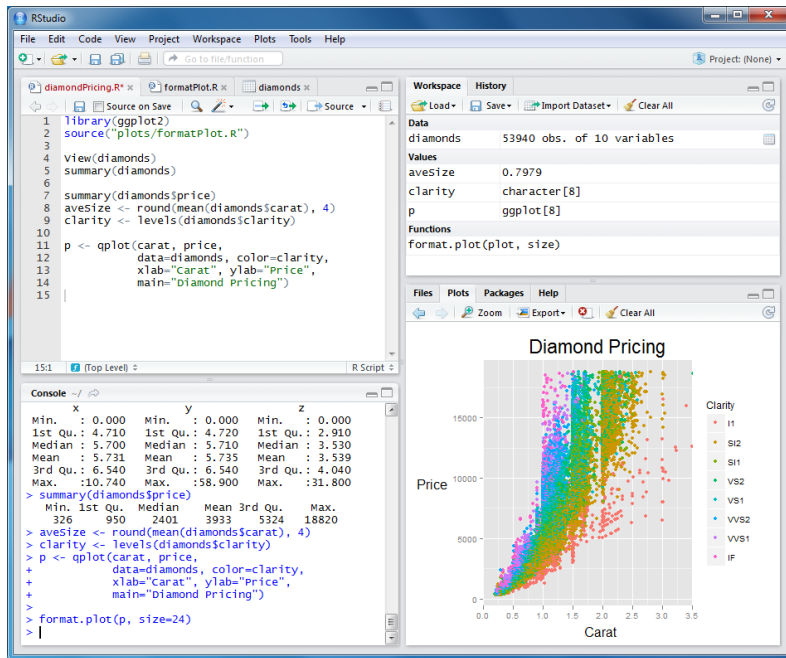
You improved on your best score by 0.01914.  
You just moved up 219 positions on the leaderboard.

**Solution in Top 5%**

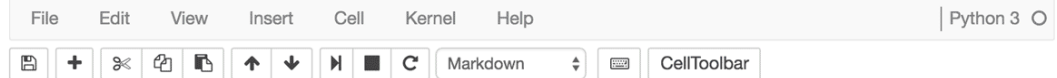
## Tutorial

<http://trevorstevens.com/kaggle-titanic-tutorial/getting-started-with-r>

# R Studio, Python Anaconda, Jupyter



jupyter spectrogram (autosaved)



## Simple spectral analysis

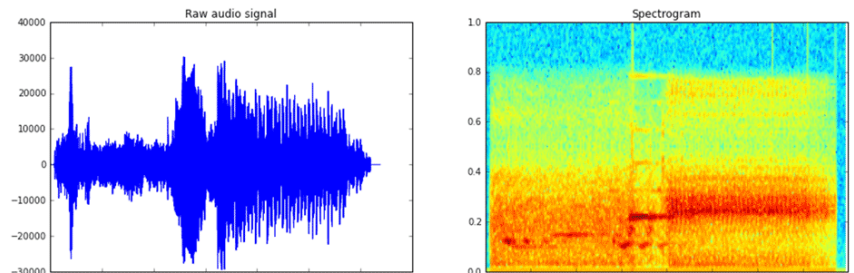
An illustration of the [Discrete Fourier Transform](#)

$$X_k = \sum_{n=0}^{N-1} x_n \exp\left(-\frac{2\pi i}{N} kn\right) \quad k = 0, \dots, N-1$$

```
In [2]: from scipy.io import wavfile
rate, x = wavfile.read('test_mono.wav')
```

And we can easily view it's spectral structure using matplotlib's builtin spectrogram routine:

```
In [5]: fig, (ax1, ax2) = plt.subplots(1,2,figsize=(16,5))
ax1.plot(x); ax1.set_title('Raw audio signal')
ax2.spectrogram(x); ax2.set_title('Spectrogram');
```

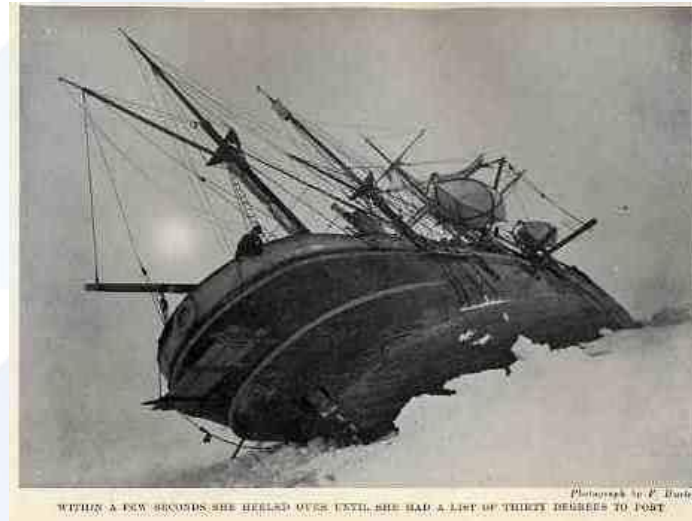




# Kaggle

## Tips and tricks

- Be patient
- Understand data
- Think more, code less
- Follow other approaches (Kaggle Forum, Github, Slideshare)
- Fit appropriate algorithm to data



355	.31	Luminous Logic	0.41089	16	Sat, 27 Dec 2014 18:37:49 (-24.7h)
-		a.karwan	0.41091	-	Mon, 20 Jul 2015 14:25:45 Post-Deadline
<b>Post-Deadline Entry</b> If you would have submitted this entry during the competition, you would have been around here on the leaderboard.					
356	.31	rachel wang	0.41092	15	Sun, 12 Apr 2015 01:57:48
1597	new	Chai Kothari	0.49523	19	Fri, 29 May 2015 03:36:57 (-0.5h)
-		a.karwan	0.49523	-	Mon, 20 Jul 2015 11:18:56 Post-Deadline
<b>Post-Deadline Entry</b> If you would have submitted this entry during the competition, you would have been around here on the leaderboard.					
1598	.103	batussi	0.49523	2	Mon, 24 Nov 2014 14:29:33 (-0.1h)