

GitPro

---

# Predictive Modeling in Health Care

Taposh Dutta-Roy



---

## **Introduction**

Health Care industry

Analytics in Health Care

Examples of Predictive Models in Health Care

Predictive Modeling Process & Tools

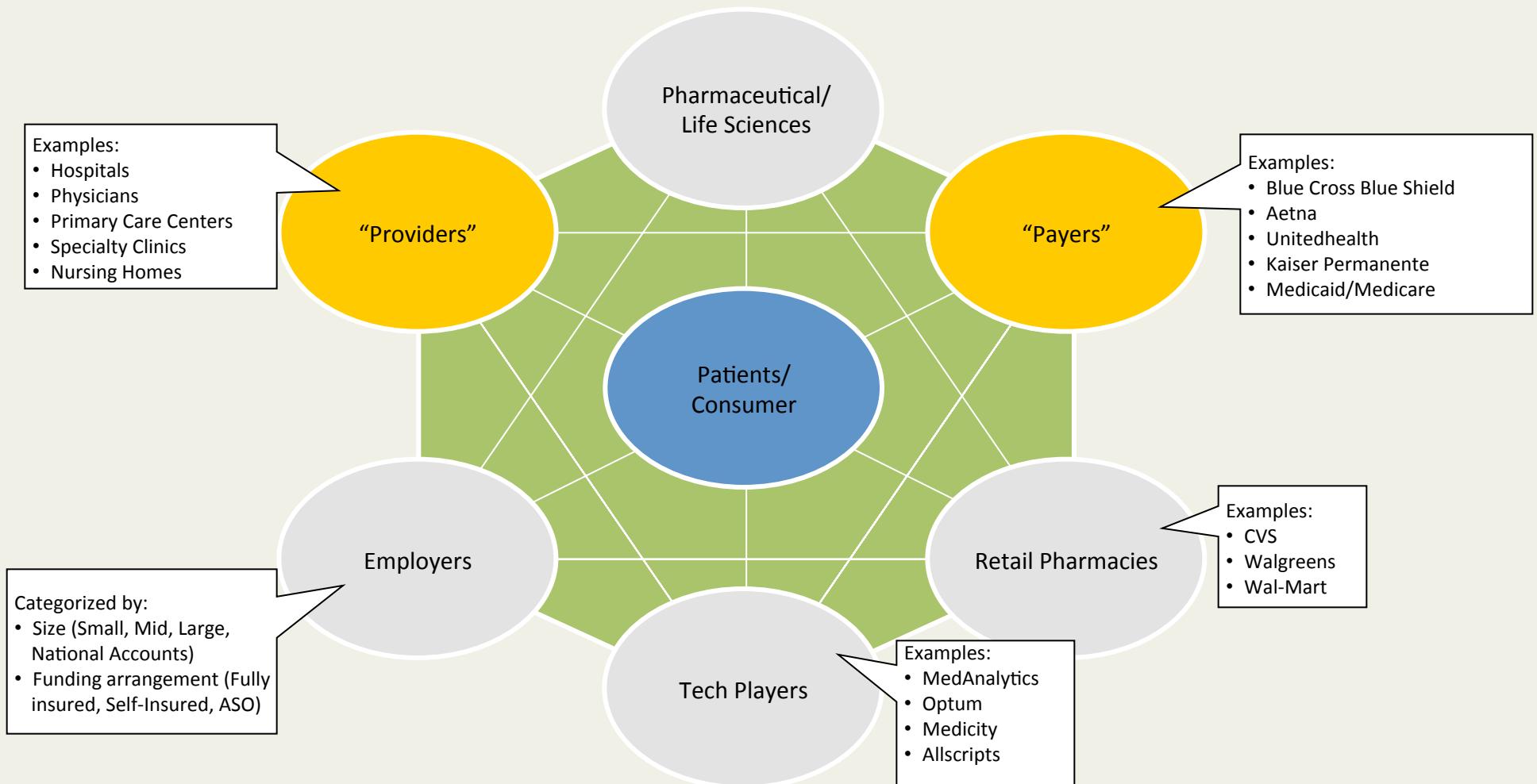
Kaggle - Bike Sharing Predictive Modeling

---

# Health Care industry

---

# Healthcare Ecosystem



---

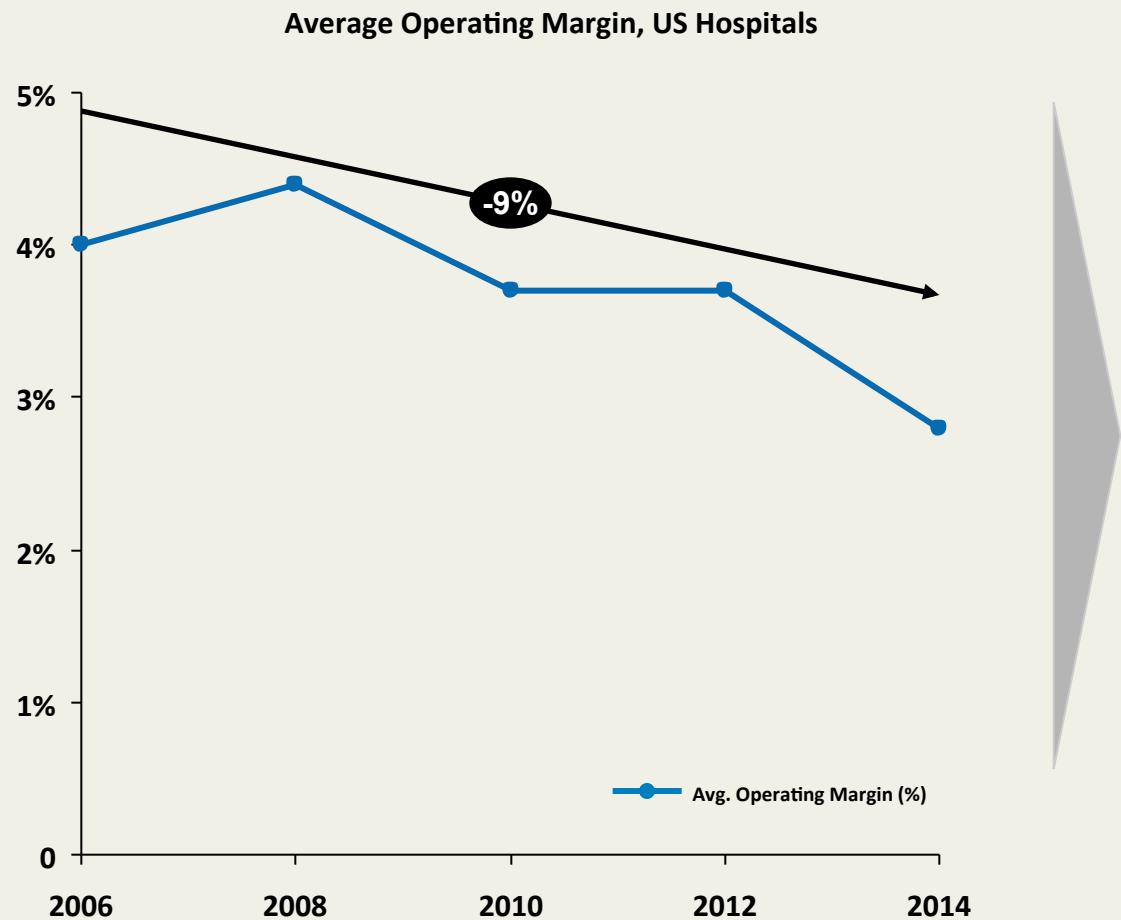
**“The healthcare ecosystem is made up of many players – some misaligned incentives exist across stakeholders”**

---

# Recent Trends in Health Care industry

---

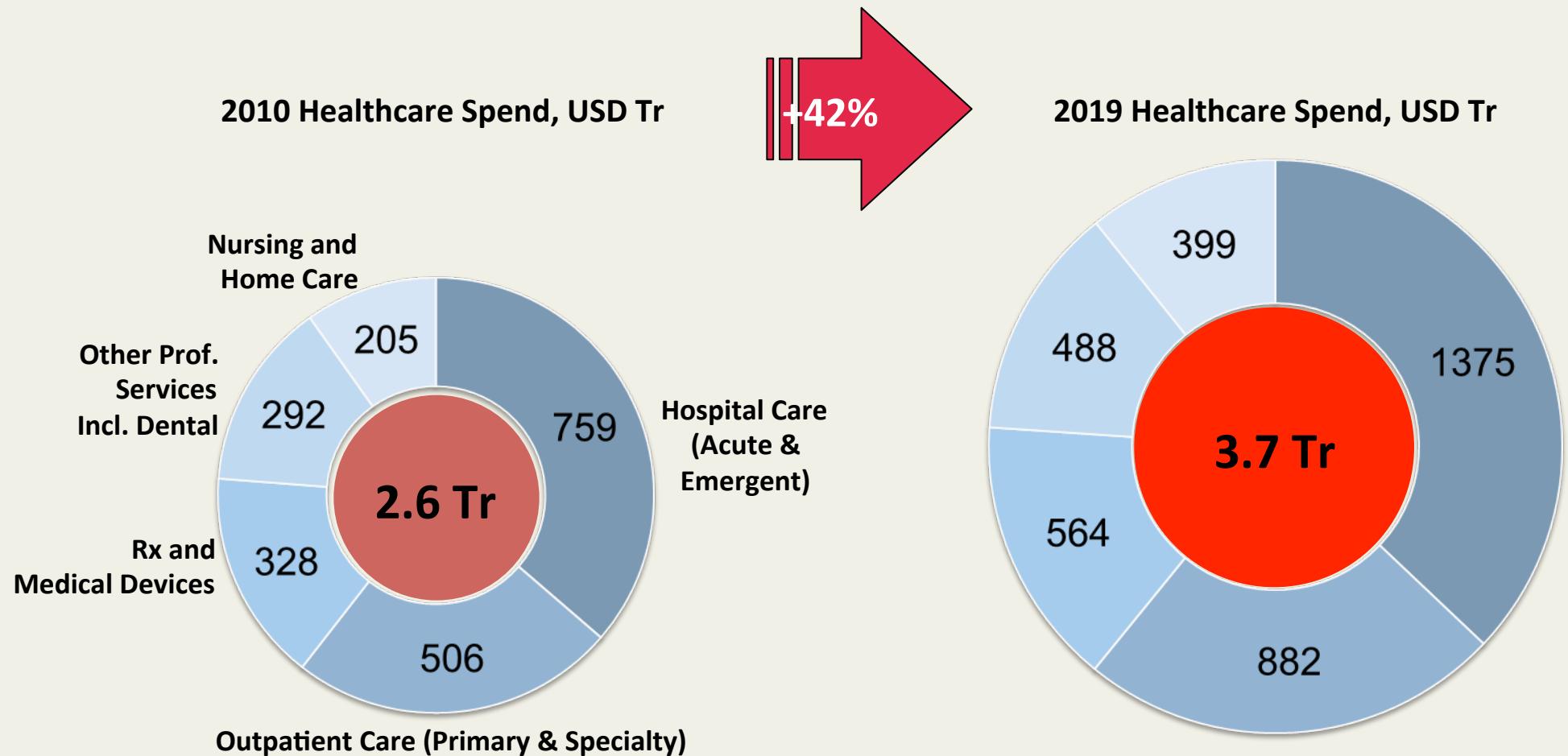
# U.S. hospitals face increasing financial pressures in the marketplace, threatening operating margins



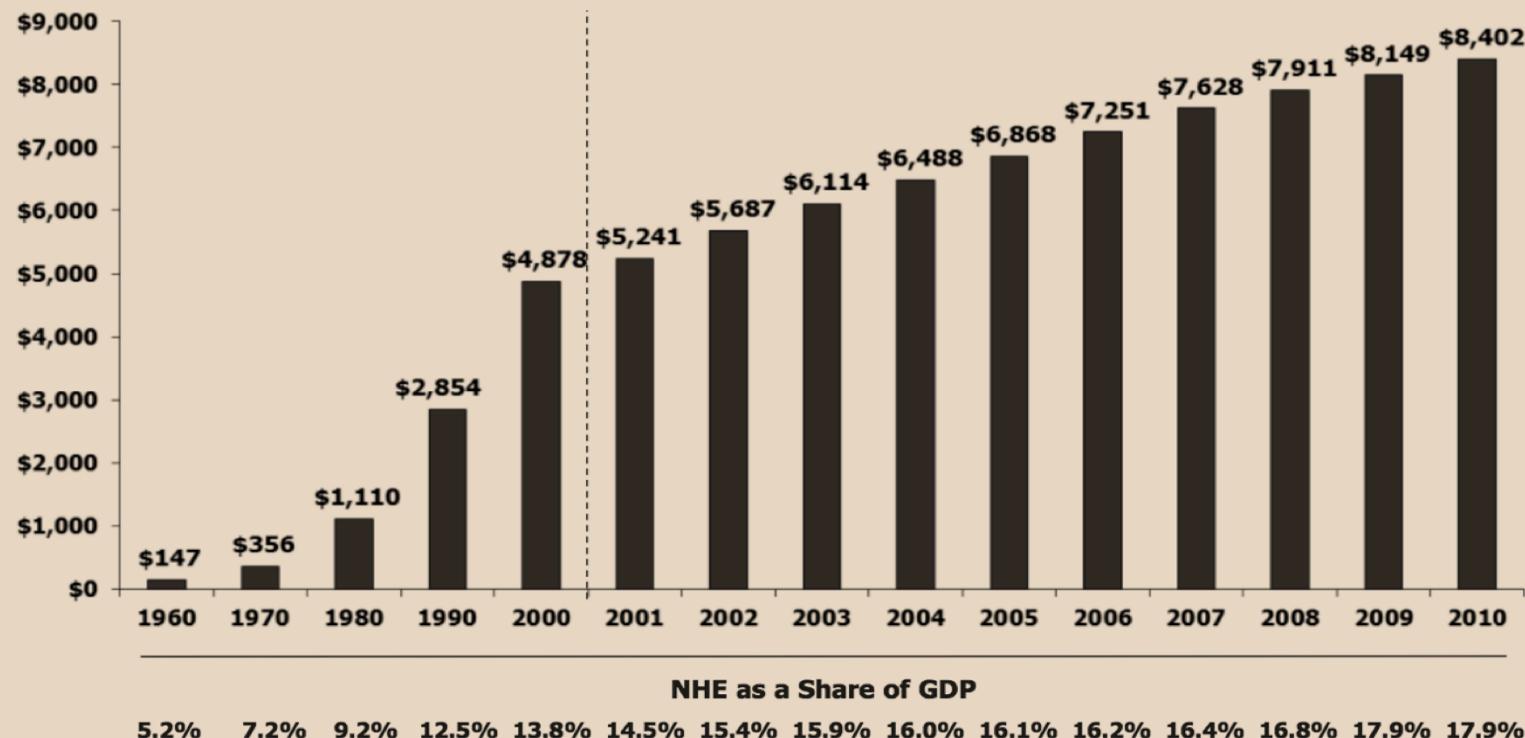
## Contributing Factors

- ❖ Uncompensated (i.e. charity) and low reimbursement (i.e. Medicaid) care levels growing faster than increase in net patient revenue
- ❖ Misaligned financial incentives between physicians and hospitals contribute to overutilization and excessive labor/non-labor costs

# Healthcare spend is set to grow almost 50% by 2019 – hospital spend and nursing home spend are fastest growing



# Figure 1: National Health Expenditures per Capita, 1960-2010



Notes: According to CMS, population is the U.S. Bureau of the Census resident-based population, less armed forces overseas.

Source: Centers for Medicare and Medicaid Services, Office of the Actuary, National Health Statistics Group, at <http://www.cms.hhs.gov/NationalHealthExpendData/> (see Historical; NHE summary including share of GDP, CY 1960-2010; file nhegdp10.zip).



---

“The share of economic activity (gross domestic product, or GDP) devoted to health care has increased from 7.2% in 1970 to 17.9% in 2009 and 2010”

---

# Key Facts

- In 2010, the U.S. spent \$2.6 trillion on health care, an average of \$8,402 per person.
- Half of health care spending is used to treat just 5% of population.
- Although only 10% of total health expenditures, spending on prescription drugs has received considerable attention because of its rapid growth (114% from 2000 to 2010).
- Many policy experts believe new technologies and the spread of existing ones account for a large portion of medical spending & its growth.

# As Health Exchanges are established, the market will increasingly move toward a more competitive, transparent, retail environment

The image shows two screenshots of a Health Exchange website. The left screenshot displays a search interface for finding insurance plans, with tabs for 'Find Insurance', 'Health Care Reform', and 'About Us'. Below this is a section titled 'Find Insurance: Individuals & Families' with a sub-section 'STEP 4 of 6 - COMPARE PLANS (OVERVIEW)'. It includes a note: 'Click "View Plan" for more details. You can also compare up to 5 plans at a time. Check the box next to each plan to select it.' A table lists five plans with columns for Tier, Plan, Premium, Deductible, and Copayments.

Tier	Plan	Premium*	Deductible	Copayments*
S	<input checked="" type="checkbox"/> Tufts Health Plan HMO Select 20	\$645.16	\$1,000/2,000	\$30 / \$100 / 10% of Rx \$500/1,000/10%
S	<input checked="" type="checkbox"/> Fallon Community Health Plan FCHP Select Care	\$668.00	\$1000/1,000	\$30 / \$100 / 10% of Rx \$500/1,000/10%
S	<input type="checkbox"/> Neighborhood Health Plan NHPTwo Select	\$724.08	None/None	\$30 / \$100 / 10%
S	<input type="checkbox"/> Harvard Pilgrim Health Care Harvard Pilgrim Best Buy and 100% HMO	\$730.13	\$1,000/2,000	\$30 / \$100 / 10%
S	<input type="checkbox"/> Fallon Community Health Plan FCHP Select Care	\$762.00	\$1000/1,000	\$30 / \$100 / 10%

The right screenshot shows a 'COMPARE PLANS' feature. It includes a note: 'Here are the plans that are in the tier(s) you selected. Select up to 3 plans at a time to get more detail. Use the "Cost" or "Benefits" tabs to compare further. Adjust the contribution levels to change the employer or employee costs for your client.' It features sliders for 'The employer's contribution for employees' (set at 50%) and 'The employer's contribution for the spouses and dependents of employees' (set at 25%). A table compares the selected plans based on 'Employee Only', 'Employee + Spouse', and 'Family' contributions. The table has columns for Carrier & Plan, Employee Access, Total Monthly Cost, and Employer Monthly Cost.

Carrier & Plan	Employee Access	Total Monthly Cost	Employer Monthly Cost
Tufts Health Plan HMO Select 20	7/8	\$2,337.76	\$1,186.88
Neighborhood Health Plan NHPTwo Select	8/8	\$1,820.01	\$910.00
Blue Cross Blue Shield of Massachusetts HMO Blue Value with Basic Rx	8/8	\$2,448.43	\$1,224.21
Harvard Pilgrim Health Care Harvard Pilgrim Tiered Copayment HMO 30	8/8	\$2,579.52	\$1,289.76
Fallon Community Health Plan FCHP Select Care	8/8	\$2,064.00	\$1,032.00

---

# Analytics in Health Care

Trends in analytics used in Health Care



---

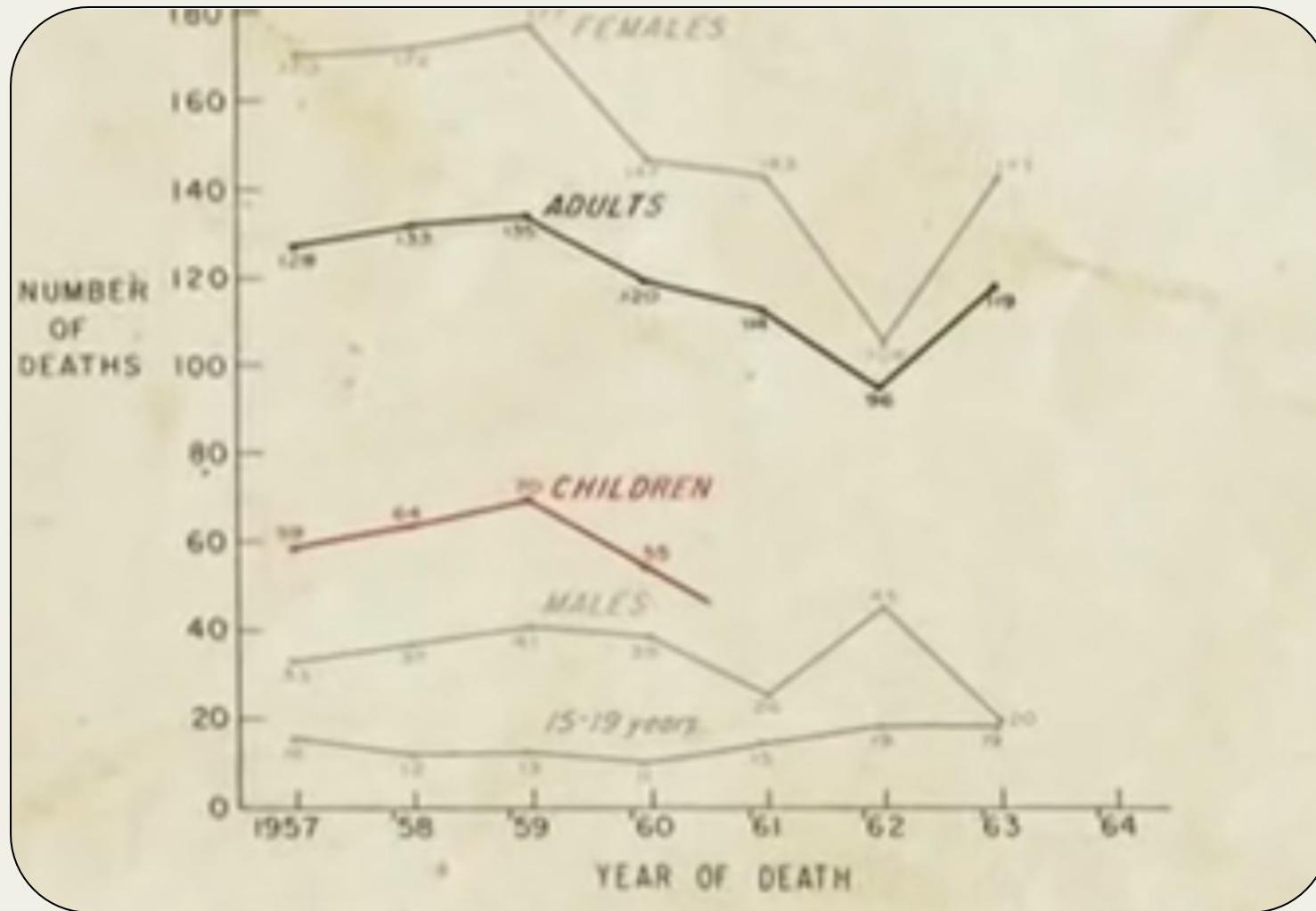
**Analytics in HealthCare has been used for a very long time.**

## (Circa 1950) Basic Statistics Uncovered Kuru : A medical detective Story



Source : <https://www.youtube.com/watch?v=ZH0y6hiUcyA>

(1950-1964) Dr. Michael Alpers saved data for decades to get to a final realization, what causes Kuru using Data Analysis.





**“Traditional Reports”**

**“Correlation”**

**“Correlation & Causation”**

---

**“Predictive analytics has been part of healthcare for several decades. Its applying what doctors have been doing on a bigger scale. What has changed with big data is our ability to process more data(quantity), combine variety of data and measure faster.”**

---

# **Common Use Cases of Predictive Modeling:**

- 1. Clinical Decision Support**
- 2. Readmissions**
- 3. Chronic Disease management**
- 4. Patient Matching**

---

# Clinical Decision Support - Sepsis Alert

## Use Case:

Very common use case. Implementation of a cloud based alerting system for variety of inpatient activities – care, Labs.

**KPI:** Better quality of care and sending less patients to emergencies. Cost reduction for patients and hospitals.

# Clinical Decision Support Reduces Sepsis

Implementation of a cloud-based alerting system and change management were associated with a 53% drop in sepsis mortality in a single-site study.

The screenshot shows a news article from Medscape. The title is "Clinical Decision Support May Help Reduce Sepsis Mortality" by Ken Terry, published on May 20, 2015. The article discusses a study showing a 53% drop in sepsis mortality after implementing a cloud-based alerting system and change management. It includes three editor's recommendations: "Sepsis Screening Tool Spots Subtle Signs, Saves Lives", "Quality Improvement Methods Improve Sepsis Care in Children", and "EHR Data May Predict Sepsis Mortality, Study Shows". There is also a "Topic Alert" section and a "RELATED DRUGS & DISEASES" section for "Streptococcus Group D Infections". The URL of the article is http://www.medscape.com/viewarticle/845051.

<http://www.medscape.com/viewarticle/845051>

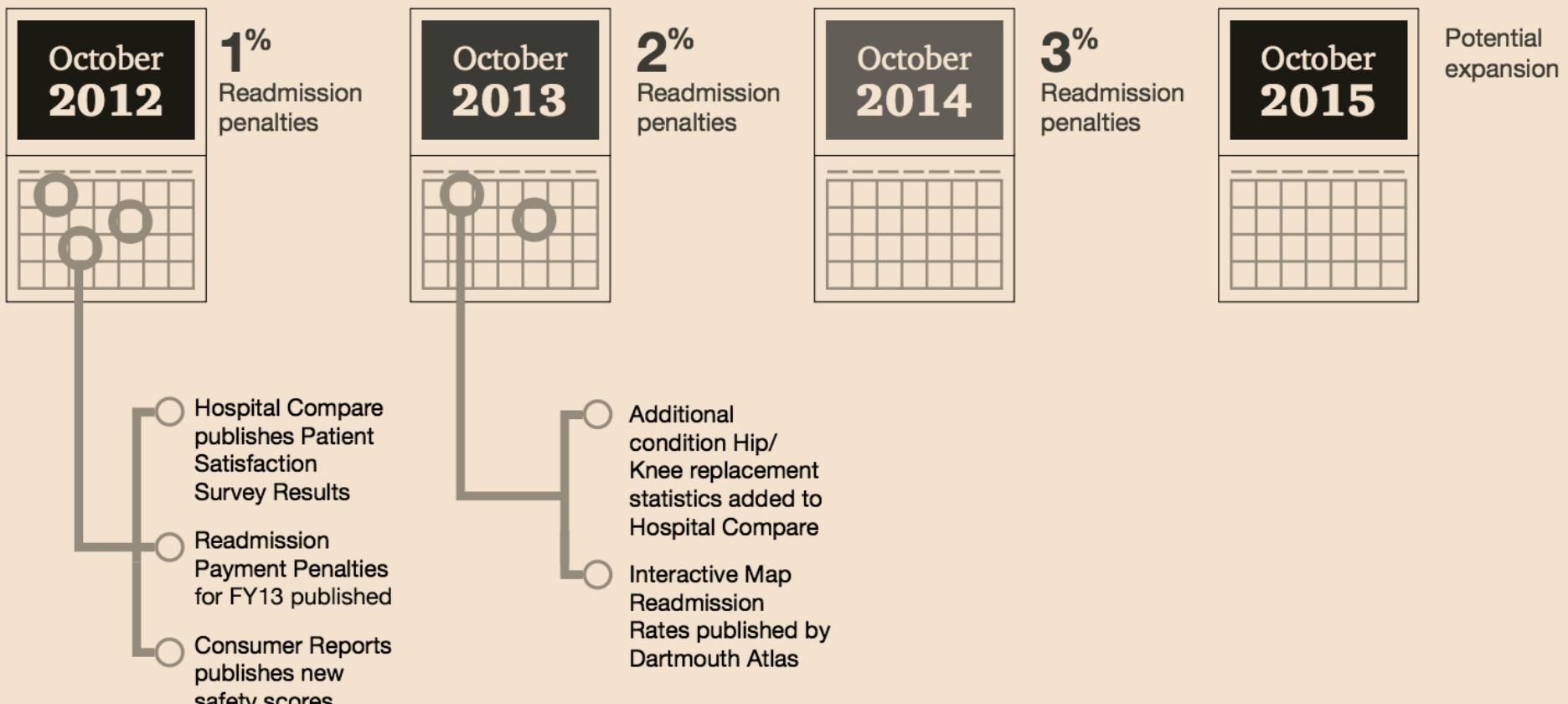
# Readmissions

**Use Case:** Recent changes to Centers for Medicare and Medicaid Services (CMS) compensation will no longer cover hospital expenses for congestive heart failure (CHF) patients who are readmitted within 30 days of discharge. In addition to the risk of degraded health outcomes, the hospital faces financial loss for preventable CHF readmits.

**KPI:** An improved understanding of the prime drivers affecting readmission will benefit both the quality of care and business revenues.

**Figure 4. Hospital readmission penalties increase along with publicly reported results**

## **Hospital readmissions timeline and highlights of consumer ratings**



Source: PwC Health Research Institute<sup>26,27,28,29,30</sup>

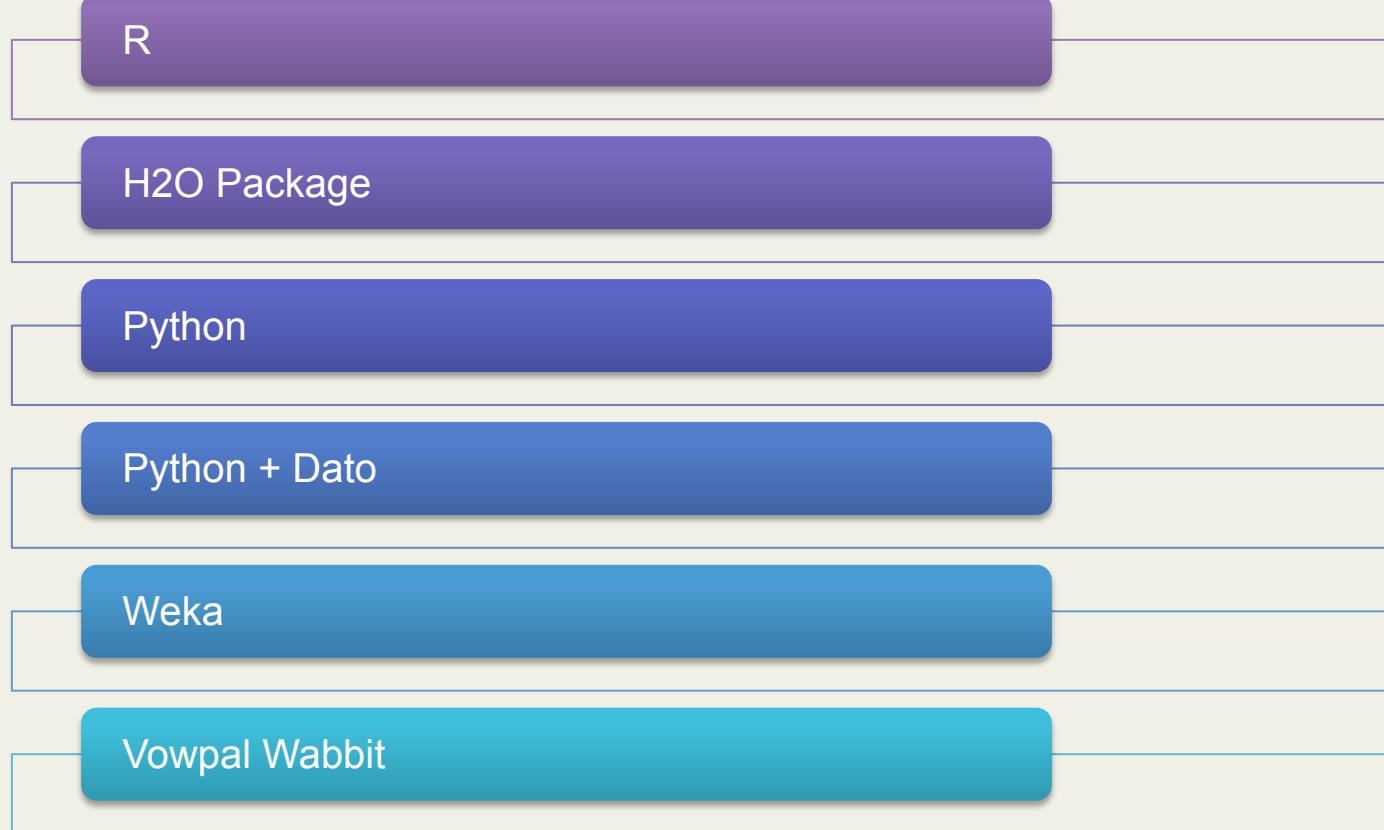
---

# Predictive Analytics Process & Tools

# Predictive Modeling Process

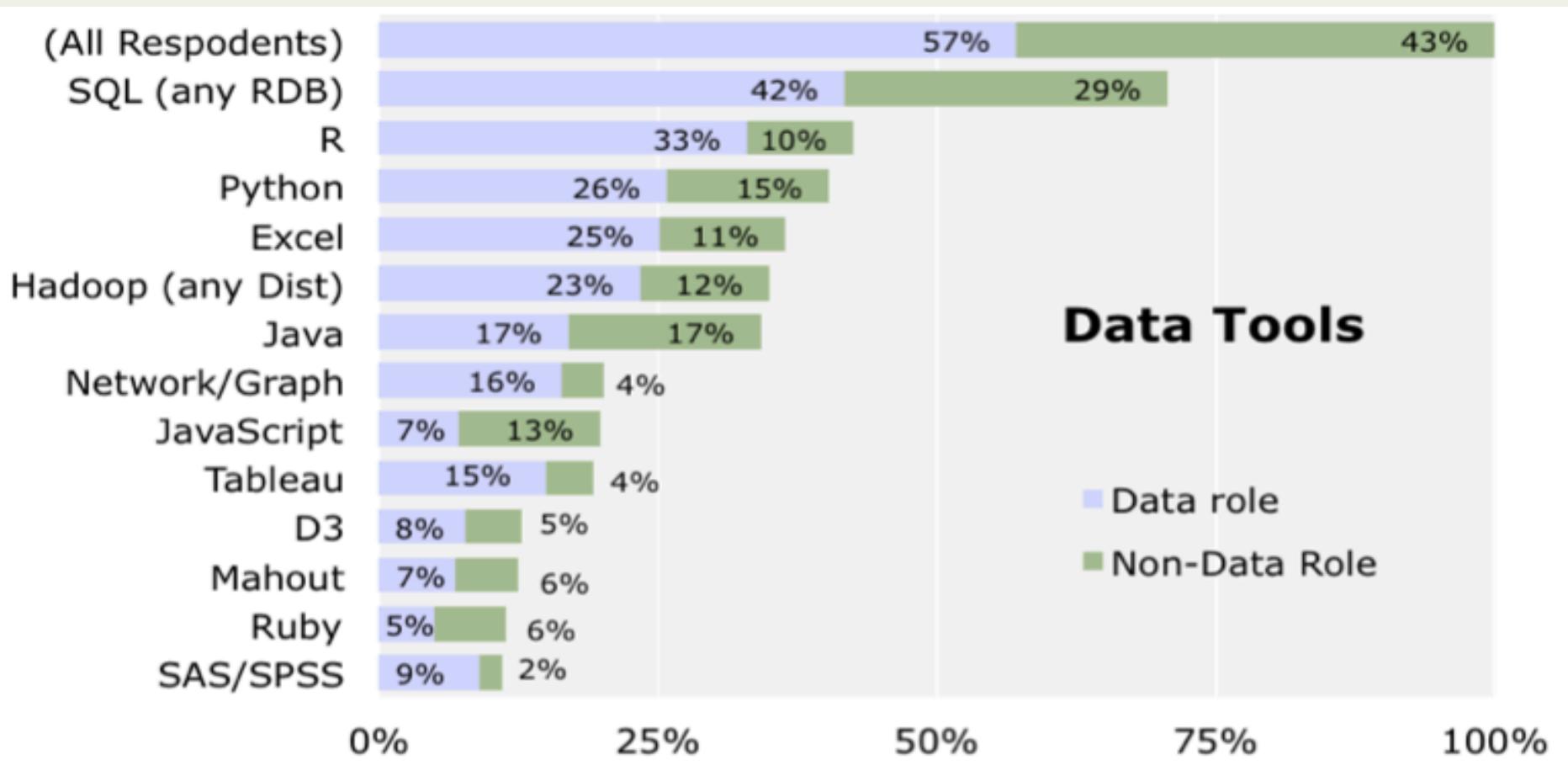
- Collect Data
- Visualizing & Understanding Data
- Determine Metrics – ROC, AUC (C-stat), Sensitivity/Specificity
- Dividing Data into Training, Validation, & Test
- Factor Engineering
- Applying Algorithm(s)
- Reviewing Training Results.
- Scoring
- Feedback new data into collection set

# Tools



# Survey on top tools used by Data Scientist

<http://blog.revolutionanalytics.com/2014/01/in-data-scientist-survey-r-is-the-most-used-tool-other-than-databases.html>



---

# Example - Kaggle Bike Sharing

# About Bike Share



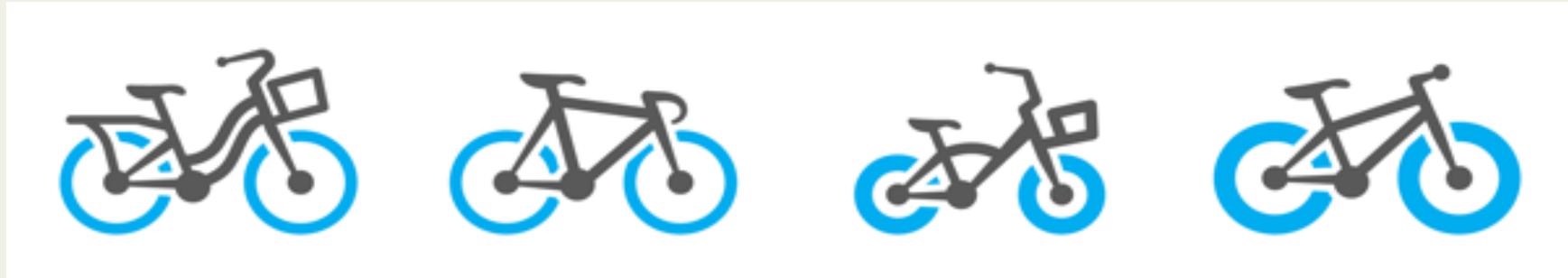
Competition:

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

Challenge:

Forecast use of a city's bike share system

# Source : Bike-Share Data



Publication :

Fanaee-T, Hadi, and Gama, Joao, Event labeling combining ensemble detectors and background knowledge, Progress in Artificial Intelligence (2013): pp. 1-15, Springer Berlin Heidelberg.

# Data Ingest



	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
1	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0000	3	13	16
2	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0000	8	32	40
3	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0000	5	27	32
4	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0000	3	10	13
5	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0000	0	1	1
6	2011-01-01 05:00:00	1	0	0	2	9.84	12.880	75	6.0032	0	1	1

The goal is to predict counts either based on sum of casual & registered or directly

# Data Fields

Datetime : hourly date + timestamp

Season : 1 = spring, 2 = summer, 3 = fall, 4 = winter

Holiday : whether the day is considered a holiday

Workingday : whether the day is neither a weekend nor holiday

Weather :

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

Temp : temperature in Celsius

Atemp : "feels like" temperature in Celsius

Humidity : relative humidity

Windspeed : wind speed

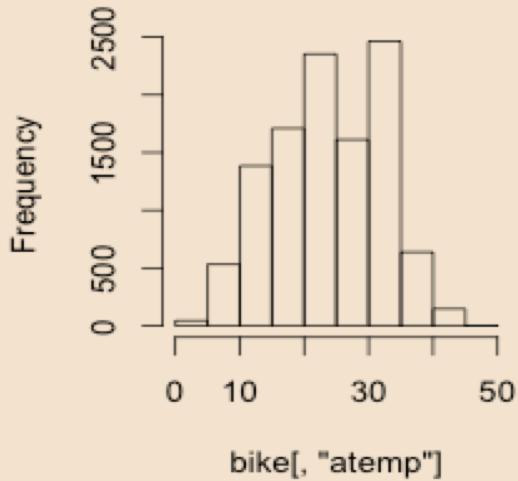
Casual : number of non-registered user rentals initiated

Registered : number of registered user rentals initiated

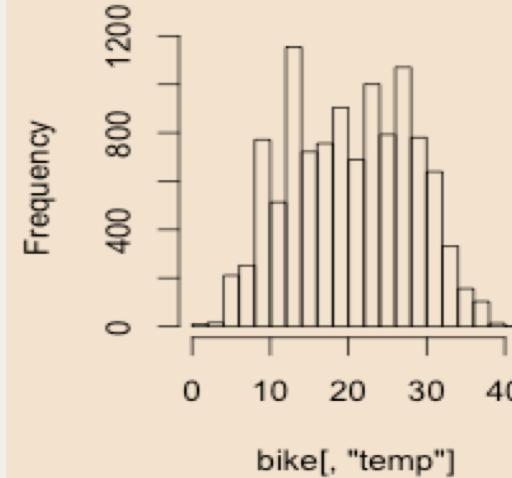
Count : number of total rental

# Data - Continuous

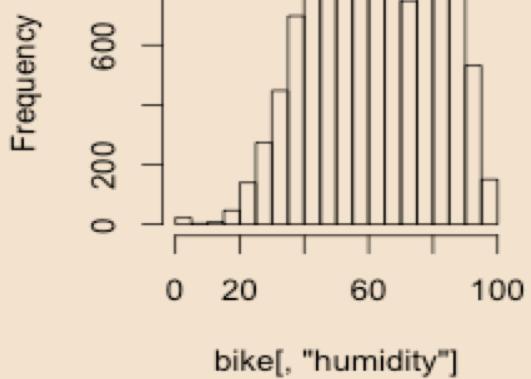
**Histogram of bike[, "atemp"]**



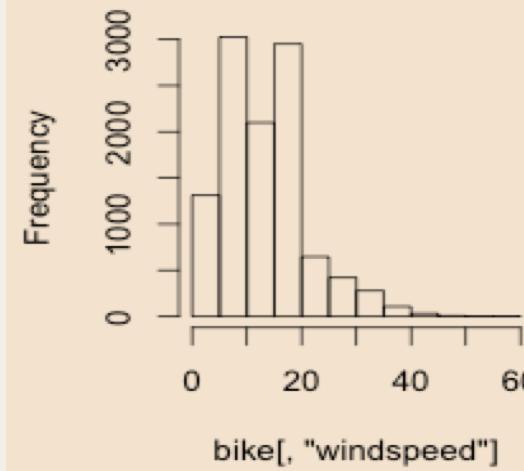
**Histogram of bike[, "temp"]**



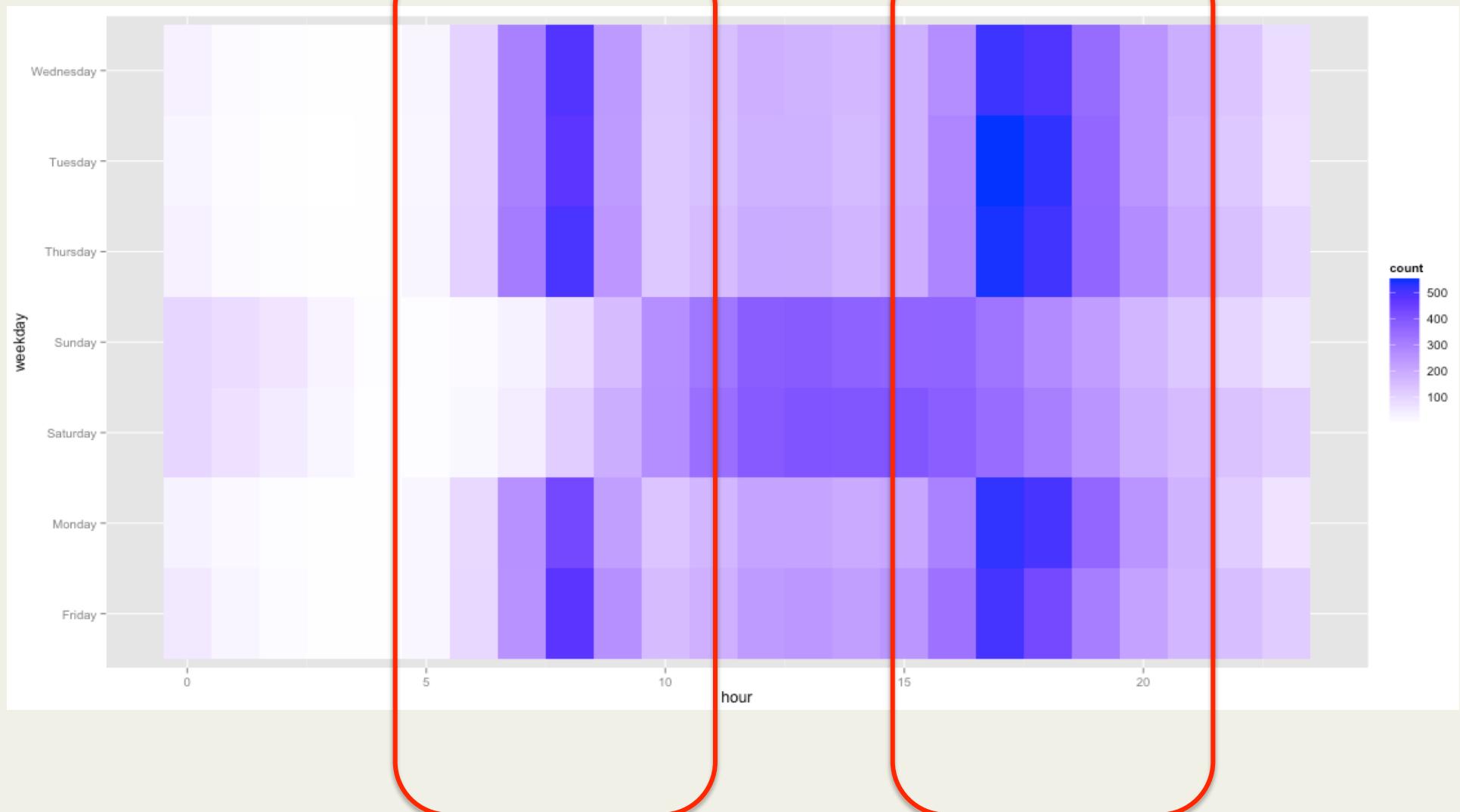
**Histogram of bike[, "humidity"]**



**Histogram of bike[, "windspeed"]**

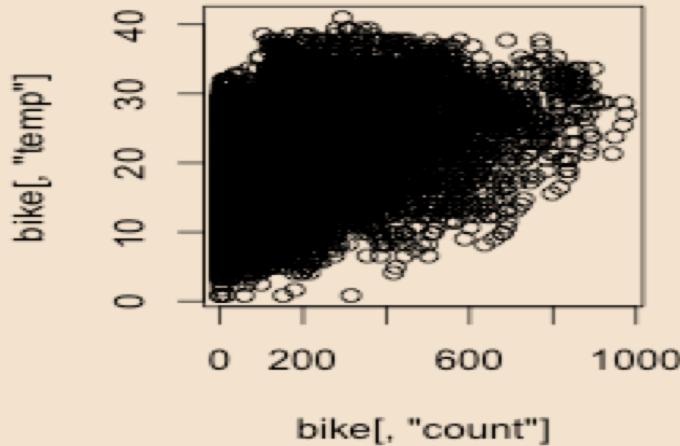


# Workday busy hours

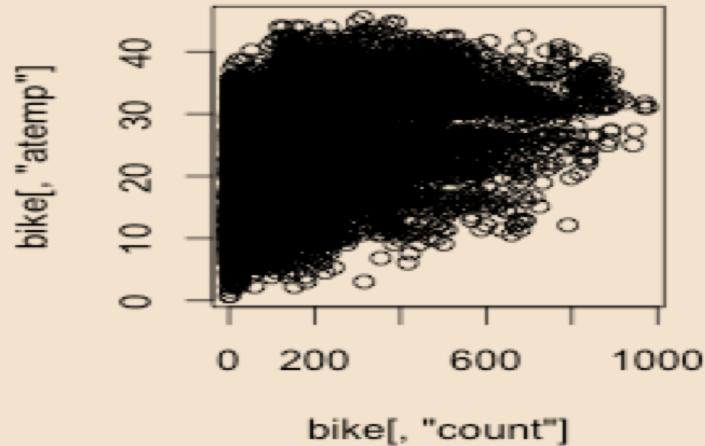


# Data

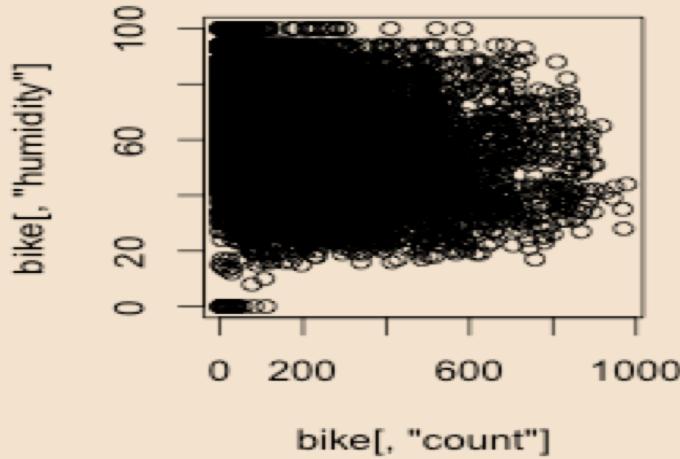
**Scatterplot of count vs. temp**



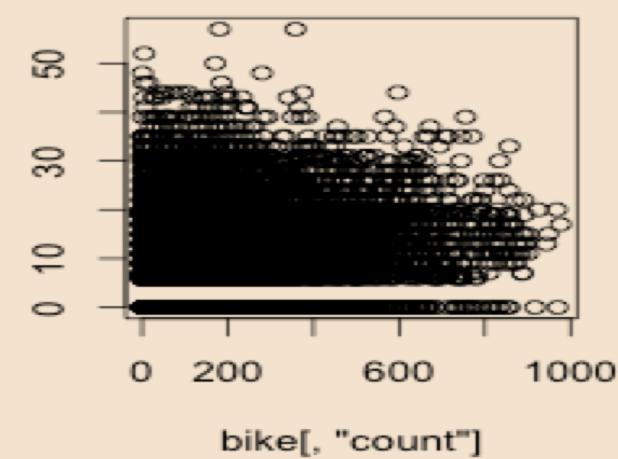
**Scatterplot of count vs. atemp**



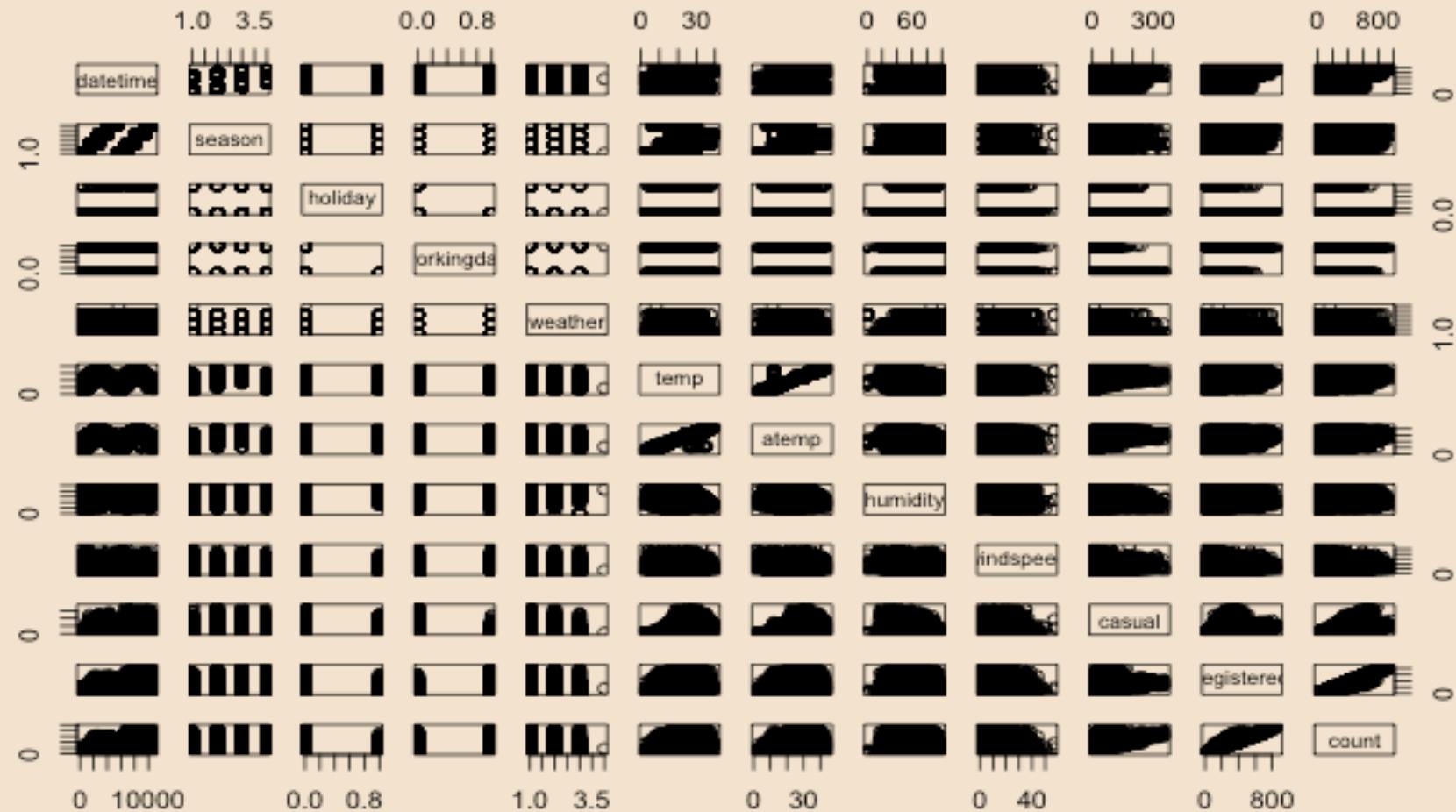
**Scatterplot of count vs. humidity**



**Scatterplot of count vs. windspeed**



# Pairs Command - R



# Using R

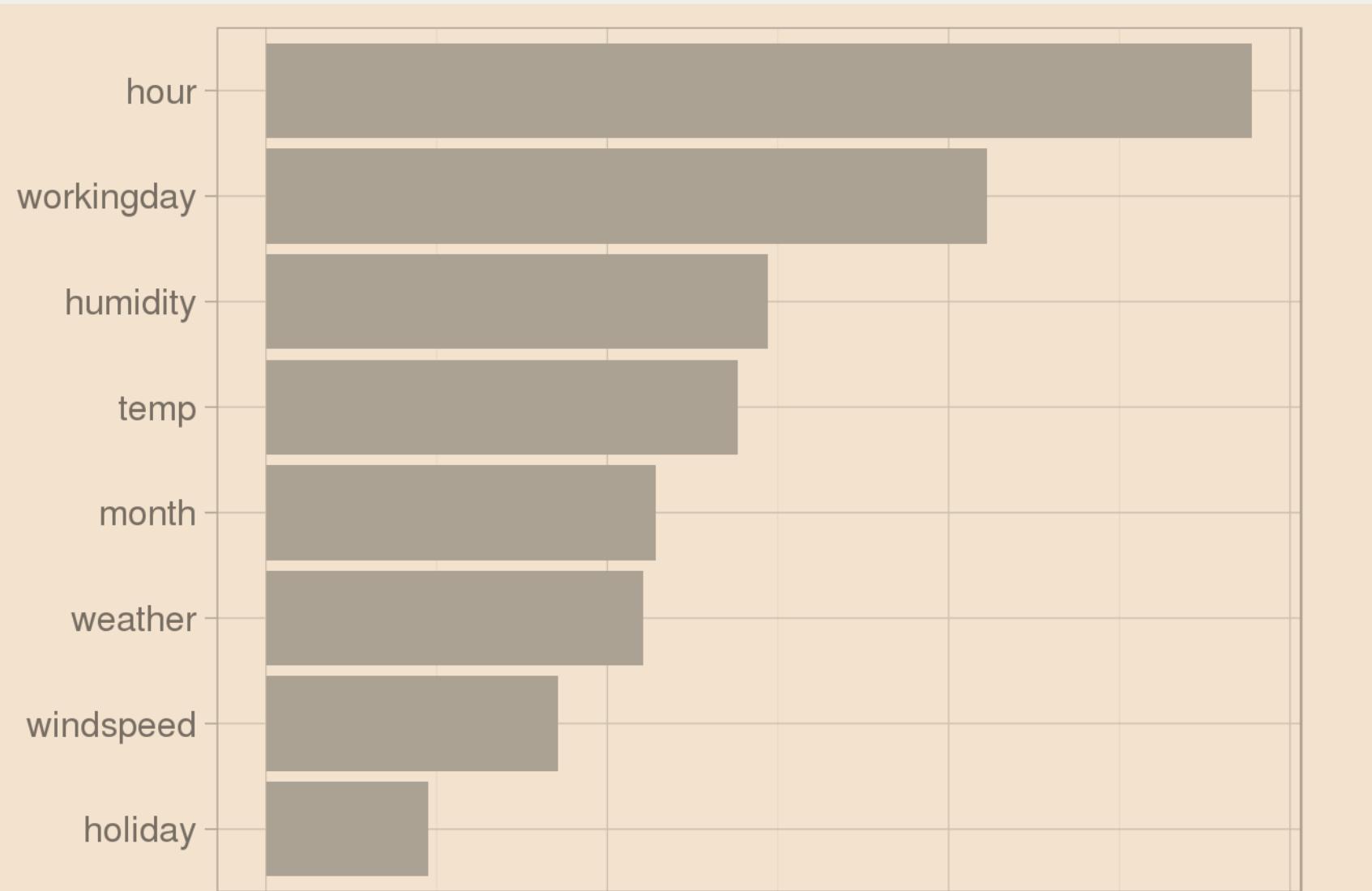
```
1 #-----  
2 # Main Program for kaggle-bike-sharing  
3 # Taposh Roy @taposh_dr  
4 #-----  
5 setwd("/Users/taposh/workspace/kaggle/bikeshare/")  
6 #sink the output  
7 sink("bikeshare.log", split = T)  
8 #source the libraries  
9 source("mylibraries.R")  
10 #input data  
11 source("inputdata.R")  
12 #View the data  
13 head(bike)  
14 #Visualize the data  
15 source("/Users/taposh/workspace/kaggle/bikeshare/visualize.R")  
16 #Utility Functions  
17 source("/Users/taposh/workspace/kaggle/bikeshare/utils.R")  
18  
19 #####  
20 ### Method #1 Matrix Computation  
21 #####  
22 #Factorengineering  
23 source("/Users/taposh/workspace/kaggle/bikeshare/factorengineering_bike.R")  
24 source("matrix.R")  
25 #####  
26 ### Method #2 Models  
27 #####  
28 #Factorengineering  
29 source("/Users/taposh/workspace/kaggle/bikeshare/factorengineering_v1.R")  
30 #colnames(bike)  
31 #colnames(test)  
32 source("models_nongrid.R")  
33 #####  
34 ### Method #2 Models  
35 #####  
36 #Factorengineering  
37 source("models.R")  
38
```

# Feature Engineering

```
1 #-----  
2 # Factor Engineering for kaggle-bike-sharing  
3 # Taposh Roy  
4 # @taposh_dr  
5 #-----  
6  
7 #Separating out the outputs from training  
8 actual<-c()  
9 countresult <-c()  
10 causal<-c()  
11 registered<-c()  
12 countresult<-cbind("count"=countresult,bike[, "count"] )  
13 actual<-cbind("count"=actual,bike[, "count"] )  
14 causal<-cbind(causal,bike[, "casual"] )  
15 registered<-cbind(registered,bike[, "registered"] )  
16  
17 causal<-log(causal)  
18 causal[causal[,1]<0,1]<-0  
19
```

# Models

```
1 #-----  
2 # Factor Engineering for kaggle-bike-sharing  
3 # Taposh Roy  
4 # @taposh_dr  
5 #-----  
6  
7 #Separating out the outputs from training  
8 actual<-c()  
9 countresult <-c()  
10 causal<-c()  
11 registered<-c()  
12 countresult<-cbind("count"=countresult,bike[, "count"] )  
13 actual<-cbind("count"=actual,bike[, "count"] )  
14 causal<-cbind(causal,bike[, "casual"] )  
15 registered<-cbind(registered,bike[, "registered"] )  
16  
17 causal<-log(causal)  
18 causal[causal[,1]<0,1]<-0  
19
```



Relative Importance

---

# Thank You !!