# Advanced Analytics with R

## Course Introduction

Forbes Earning Preview: H.J. Heinz

A quality first quarter earnings announcement could push shares of H.J. Heinz (HNZ) to a new 52-week high as the price is just 49 cents off the milestone heading into the company's earnings release on Wednesday, August 29, 2012.

The Wall Street consensus is 80 cents per share, up 2.6 percent from a year ago when H.J reported earnings of 78 cents per share. The consensus estimate remains unchanged over the past month, but it has decreased from three months ago

when it was 82 cents. Analysts are expecting earnings of $3.52 per share for the fiscal year. Analysts project revenue to fall 0.3 percent year-over-year to $2.84 billion for the quarter, after being $2.85 billion a year ago . For the year, revenue is projected to roll in at $11.82 billion.

Paging DR. WATSON

www.cnbc.com/2016/10/25/driverless-beer-run-bud-makes-shipment-with-self-driving-truck.html

**BEHIND THE WHEEL**    edited by PHIL LeBEAU

CONSUMER | RETAIL | MEDIA | **AUTOS** | FOOD AND BEVERAGE | RESTAURANTS | FASHION | GOODS

# 'Driverless' beer run; Bud makes shipment with self-driving truck

19K SHARES

177 COMMENTS    Join the Discussion

Phil LeBeau | @Lebeaucarnews
Tuesday, 25 Oct 2016 | 12:03 AM ET

CNBC

As beer runs go, this one stands out.

Anheuser-Busch hauled a trailer loaded with beer 120 miles in an autonomous-drive truck, completing what's believed to be the first commercial shipment by a self-driving vehicle.

The trip happened last week in Colorado as Anheuser-Busch, collaborated with Otto, a subsidiary of Uber that is developing self-driving truck technology. The semi drove autonomously on the highway between Fort Collins, Colorado and Colorado Springs, Colorado.

"The incredible success of this pilot shipment is an example of what is

# Our workforce and industries have changed dramatically over time.

**1940**

**2010**

## 45,166,083
Employees

## 139,033,928
Employees

**23.4%** Manufacturing

**18.5%** Agriculture

**14.0%** Retail trade

**8.9%** Personal services

**7.4%** Professional & related services

**23.2%** Educational services, health care & social assistance

**11.7%** Retail trade

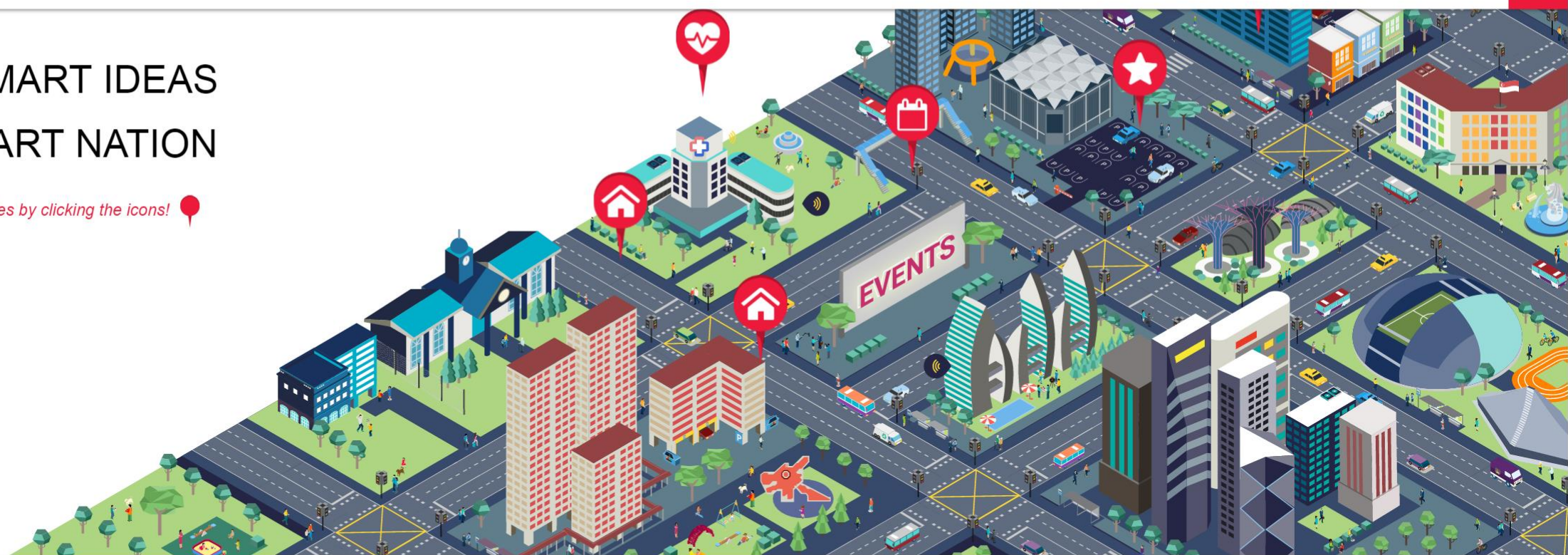**10.6%** Professional, scientific, management & administrative services, waste management services

**10.4%** Manufacturing

**6.2%** Construction

www.smartnation.sg

beta

About Smart Nation   Initiatives   Apps   Ideas   Resources   Happenings

# MANY SMART IDEAS
# ONE SMART NATION

*Uncover the possibilities by clicking the icons!*

EVENTS

# SMART NATION

A Smart Nation is one where people are empowered by technology to lead meaningful and fulfilled lives.

A Smart Nation harnesses the power of networks, data and info-comm technologies to improve living, create economic opportunity and build a closer community.

A Smart Nation is built not by Government, but by all of us - citizens, companies, agencies. This website chronicles some of our endeavours and future directions.

Learn more

BUILDING A FUTURE-READY
Singapore

# PROGRAMMING IS NEW O-LEVEL SUBJECT FROM 2017: WHAT PARENTS MUST KNOW

MARCH 7, 2016



Bukit View Secondary School students Tan Whee Lee, Nicholas Chong, Tio Hui Li, M. Charulatha, head of department for science Ng Wuay Boon, level head for science and Applied Learning Programme Reena Lloyd, Brandon Ng and (squatting) Neo Jun Wei, with solar cars they built at the school on 19 February 2016. The school will be among 19 secondary schools to offer programming as part of a new O-level subject called computing in 2017.

# Mining

# Data Mining

**Job Trends** from Indeed.com

— big-data — data-science

# Job Trends from Indeed.com

— R and "statistics"  — SAS and "statistics"  — Python and "statistics"



Percentage of Matching Job Postings

0.2

0.15

0.1

0.05

0

Jan '06   Jan '07   Jan '08   Jan '09   Jan '10   Jan '11   Jan '12   Jan '13

## Data Scientist

SAS
3%

Python
53%

R
44%

## Predictive Analytics

SAS
43%

R
41%

Python
16%

# Why R?

- R is free.
- R is a language. (You have the full flexibility and control)
- R is powerful.
- R has an amazing ecosystem of developers.

# Learning Objectives

- Master basic R programming techniques to model and analyse data.

- Self-explore R packages for analytics needs.

- Understand some common analytics methodologies used in current business environment.

- Appreciate what data science is.

# I am training what industry wants ...

# Right expectations

- Learning programming can be very painful
  - But, believe me, it will be very rewarding
- Be prepared to solve problems yourself
- Be prepared to be better than the professor!!!

# Readings

# Topics

**Fundamentals**

Basic components of R

Manipulating data

Exploring data

Data Visualization

Programming structure

**Applications**

Simulation

Regression Modelling

Generalized Regression Modelling

Classification

Sentiment Analytics

# Assessment

- Class Participation                        20%
- Group Project                                 40%
- Test 1                                            20%
- Test 2                                            20%

# Public datasets

- http://hadoopilluminated.com/hadoop_illuminated/Public_Bigdata_Sets.html

- https://www.quora.com/Data/Where-can-I-find-large-datasets-open-to-the-public

- http://www.kdnuggets.com/datasets/index.html

- http://www.researchpipeline.com/mediawiki/index.php?title=Main_Page

# Places to get help

- IVLE
- https://www.r-project.org/
- Stackoverflow.com
- Cross Validated

# R Communities

- http://www.r-bloggers.com/
- https://twitter.com/search?q=%23rstats&src=typd
- http://www.revolutionanalytics.com/r-community

# Good about R

- What you can think of, someone may have done it for you.

# Data Science Project



What problem am I solving?

**Define the goal**

Deploy the model to solve the problem in the real world

**Deploy model**

What information do I need?

**Collect and management data**

Establish that I can solve the problem and how

**Present results and document**

Find patterns in the data that lead to solutions

**Build the model**

**Evaluate and critique model**

Does the model solve my problem?

Source: "Practical Data Science with R"

# Basic Analytics Steps

# Fundamental of R Programming

# Installation

- R: the core engine

  Download at http://www.r-project.org/ and install

- RStudio: the IDE

  Download at http://www.rstudio.com/ide/ and install

# Introduction to R Studio

# Making notes using R Markdown

# Working Environment

- getwd() - Get the current working directory
- setwd() - Set the working directory to a specific folder
- list.files()  - list all the files/directories under a specific folder
- objects() – list all the objects stored in current environment

# Create Objects

```
> x<-1
> y<-2
> y<-TRUE
> z<-c("A","B","C")
```

Create three new objects

```
> ls()|
```

Check that the new objects are in the working environment

# Remove Object

```
> rm(x)
> ls()
```

Remove x from the environment

```
> rm(list = ls())
```

Remove all objects from the environment

Warning: Don't do this unless necessary!!!

# Exercise

- Open "Lasagna Triers.csv" file.
- Discuss the various types of fields in the file.

# Data Type

- R has five basic or "atomic" classes of objects
  - Character
  - Numeric
  - Integer
  - Complex
  - logical



- On top of basic classes, you can build objects such as
  - Vector
  - List
  - Matrix
  - Factor
  - Data Frames

# Character

- A **_character_** object is used to represent string values in R. For example, "James", "China", etc.

```
> x <- "James"
> class(x)
```

# Commonly used functions on character

| Function | Description | Example |
|---|---|---|
| nchar() | Get the number of characters in the string | *a<-"Singapore"* <br> *nchar(a)* |
| regexpr() | Find the starting position of a small string in a large string | *regexpr("ex", "longtext")* |
| gregexpr() | find positions of every match of a small sting in a large string | *gregexpr("a","banana")* |
| grep() | find the positions of a regular expression in a vector of text strings | *txt<-c("arm","foot","lefroo", "bafoobar")* <br> *grep("foo", txt)* |
| substr() | extract part of a text string based on position in the text string | *substr("Singapore",2,4)* |
| sub() | replace the first match of a string with a new string | *sub("or","es","Singapore")* |
| gsub() | replace every match of a given sub string in a string with a new string | *gsub("a","o","banana")* |
| paste() | combine two strings into a new string | *paste("Liu", "Qizhang")* |

# Logical

- A *logical* value is normally produced by some logical operations.

| Operator | Description | Example |
|---|---|---|
| > | Check if a number is more than the other | 5 > 3 |
| < | Check if a number is less than the other | 3 < 5 |
| == | Check if two values are the same | 3 == 5 |
| <= | Check if a number is less than or equal to the other | 3 <= 5 |
| >= | Check if a number is greater than or equal to the other | 3 >= 5 |
| != | Check if a number is unequal to the other | 3 != 5 |
| & | (AND) It combines two logical values. It is TRUE only if both logical values are TRUE | (5>3) & (4>2) |
| \| | (OR) It combines two logical values. It is FALSE only if both logical values are FALSE | (5>3) \| (4>2) |
| ! | (NOT) Gives opposite value of the given logical value | !(5 > 3) |

# Numeric

- Decimal values are called numeric in R. It is the **default** computational data type.

```
> x <- 1.20
> class(x)
[1] "numeric"
```

```
> x <- 2
> class(x)
[1] "numeric"
```

# Integer

- Integer is a special data type for integer values.

```
x<-as.integer(2)
class(x)

x<-2L
class(x)
```

```
x<-as.integer(5>6)
x
```

# Coercion

- Use as(object) function to convert an object from one data type to another. This is called **_coercion_**.

- Coercion order: logical < integer < numeric < complex < character

# Vector

- Vector is the most basic data structure. It stores an ordered array of elements of the <span style="color:blue">same</span> data type.

```r
a<-c(3,1,5)
a
```

```r
b<-1:4
b
```

```r
c<-c(a,b)
c
```

```r
d<-c("Name",1)
d
```

# Special vectors

- seq() creates a vector of numbers in a special sequence.

```
seq(1,9,2)
```

- rep() creates a vector which is a replication of a number or a vector

```
rep(c(2,3,4), 3)
```

# Computing basic statistics

- mean(x)
- median(x)
- sum(x)
- sd(x)
- var(x)
- cor(x,y)
- cov(x,y)
- max(x)
- min(x)

# Selecting Vector Elements

- Use [] to select vector elements by their positions.

- Use negative indexes to exclude elements

- Use a vector of indexes to select multiple values

- Use a logical vector to select elements based on a condition

- Use names to access named elements

# List

- List is more flexible than vector. It allows multiple types of elements, including list itself.

- A list is constructed by *list()* command.

```
Mike<-list(Name="Mike",Salary=10000,Age=43,Children=c("Tom","Lily","Alice"))
Mike
```

# Matrix

- **_Matrix_** is a vector of elements arranged in two dimensions.

```
m1<-matrix(3:8,ncol=3,nrow=2)
m1
```

- It can be formed by introducing *dim()*.

```
m2<-3:8
```

```
dim(m2)<-c(3,2)
m2
```

```
dim(m2)<-c(2,3)
m2
```

- Matrix is the data structure in R corresponding to matrix in linear algebra.

# Naming a matrix

```r
m3<-matrix(3:8,ncol = 3,nrow = 2,byrow = TRUE)
```

```r
rownames(m3)<-c("Row1","Row2")
colnames(m3)<-c("Col.1","Col.2","Col.3")
m3
```

```r
m3["Row1","Col.2"]
```

# Data Frame

- The most commonly used data structure in R is *data frame,* which is an extension of matrix by allowing co-existence of data of various types.

- A column in a data frame normally represents a data field.

- A row in a data frame normally corresponds to a record of the data.

# Construct a data frame

```r
df<-data.frame(ID=c(1,2,3),Names=c("James","Jack","Tom"),Values=rnorm(3,mean=100,sd=10))
df
```

# Explore data set

- Explore data set *mtcars*

# Factor

- **Factors** are special variables used to store categorical variables.

- Factor is not an atomic data type and it can be either a numeric data or a character data.

# Advantage of using factor

- Factor variables are stored as a vector of integer values, thus it is a more efficient use of memory. Yet, the original set of character values will be used when the factors are displayed, which will make the data presentation more meaningful.

- Many statistical models will automatically handle factor variables properly.

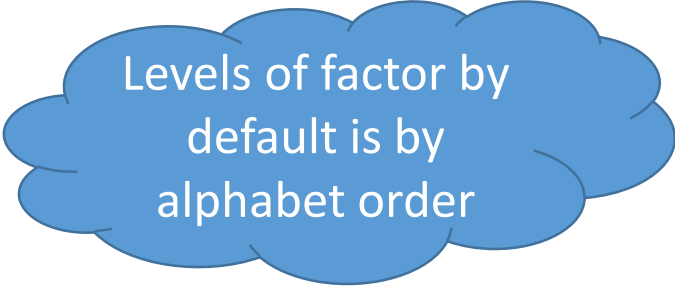- Factor variables are also very useful in many different types of graphics.
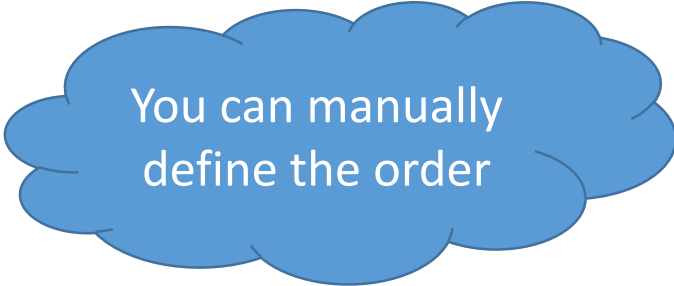
# Example

a<-c("Snake", "Dog", "Elephant", "Dog", "Cat" )

b <- factor(a)
b

as.integer(b)

> Levels of factor by default is by alphabet order

b <-factor(a, levels=c("Snake", "Cat", "Dog", "Elephant"))
b

as.integer(b)

> You can manually define the order