# STT481 Capstone in Statistics

## Lecture 0: Introduction

Chih-Li Sung

09/02/2020

# Statistics in the news

- ▶ Quote of the Day, New York Times, August 5, 2009

  *"I keep saying that the sexy job in the next 10 years will be statisticians. And I'm not kidding."*

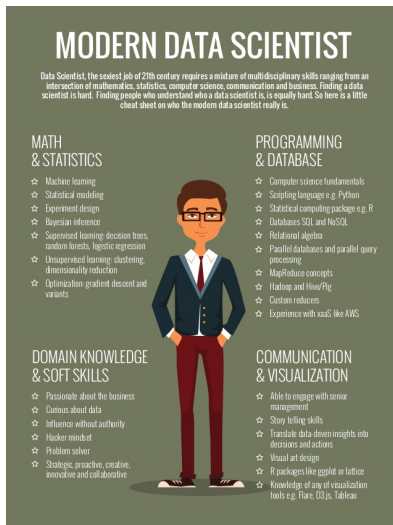– Hal Varian, chief economist at Google

- ▶ Quote of [Forbes](), Dec 11, 2017

  *"LinkedIn's Fastest-Growing Jobs Today Are In Data Science And Machine Learning"*

– Louis Columbus

# Are you ready to be a statistician/data scientist

▶ Data scientist skill-set ([source](#))

# What we have learned so far

- Basic Statistics
- Probability
- Hypothesis Testing
- Linear Regression (maybe?)
- . . .

# What this course is about

- ▶ Statistical capstone experiences are essential for statisticians/data scientists to perform an in-depth analysis of real-world data.
- ▶ Capstone experiences can develop statistical thinking by engaging in a consulting-like experience that requires skills outside the scope of traditional courses:
  - ▶ defining a complex problem,
  - ▶ analyzing data,
  - ▶ building a strong team,
  - ▶ programming techniques,
  - ▶ and communicating effectively.

# What you should expect to learn in this course

- ▶ Problem formulation
- ▶ Data collection
- ▶ Advanced statistical modeling, preliminary data analysis, and machine learning
- ▶ Statistical software (R)
- ▶ Thorough and elaborate statistical analyses of data
- ▶ Presentation and data visualization

# A standard procedure of statistical analysis

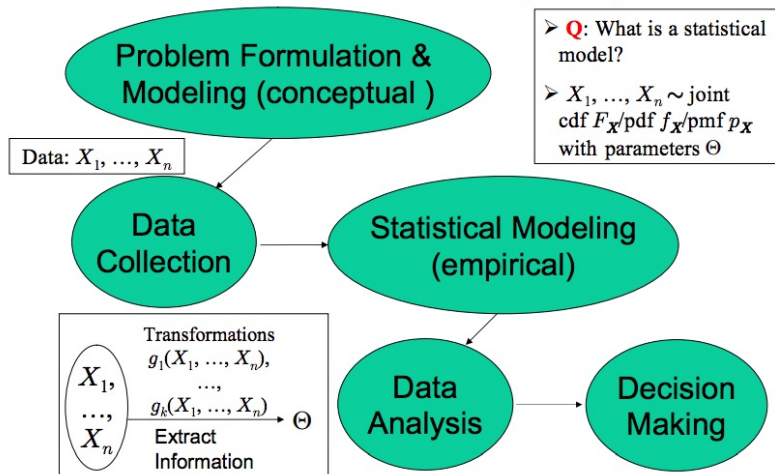• Statistics divides the study of data into *five* steps:



Figure 2:

# Data collection

*"In God we trust; all others bring data.""* – Edwards Deming

- ▶ Design of experiment (not covered)
- ▶ Survey sampling (not covered)
- ▶ Web scraping
- ▶ Text mining (not covered)

# Statistical modeling

*"Essentially, all models are wrong, but some are useful." – George Box*

► Supervised learning
  ► Linear regression model
  ► Logistic Regression
  ► Linear Discriminant Analysis
  ► KNN
  ► Linear Model Selection and Regularization
  ► Nonparametric Regression
  ► Tree-Based Methods
  ► Support Vector Machines

► Unsupervised learning
  ► Principal Components Analysis
  ► Clustering Methods

# Decision making/visualization

- ggplot2
- plotly
- shiny (myapp)
- D3.js

# Grading

- ▶ Final grades will be based on
  - ▶ one course project (20%)
  - ▶ five homework assignments ($50\% = 5 \times 10\%$)
  - ▶ two midterms ($30\% = 10\% + 20\%$)
- ▶ The final grade would be based on your total grade percentage and will be determined roughly as:

```
## Warning: package 'knitr' was built under R version 3.5.2

##
## %      90-100 85-89 80-84 75-79 70-74 65-69 60-64 0-59
## Grade      4   3.5     3   2.5     2   1.5     1    0
```

# Homework

- Homework assignments include conceptual and applied exercises. Typesetting your reports/solutions in *Latex* or *R markdown* (link1,link2) is strongly encouraged. Unreadable handwriting is subject to zero credit. Unreadable handwriting is subject to zero credit.
- Homework Policy: link

# Project

▶ Kaggle is a open platform for predictive modeling and analytics competitions, where companies and researchers can post their data and problems for users to solve.

▶ In this project, you are given the house dataset on Kaggle and the goal is to predict the final price of each home in Ames, Iowa. See link.

▶ Grading:

  ▶ Accurate prediction (10%) -

  min{15%, 0.12/(your prediction score)×10%}

  ▶ Final Report (10%)

▶ Final Report: See link.

# Important Dates and Policy

See syllabus

Any questions or feedback?

Homework 1, Q1: Create an account on Kaggle (take a screenshot like the image below). (Homework questions can be found on D2L/Assignments)
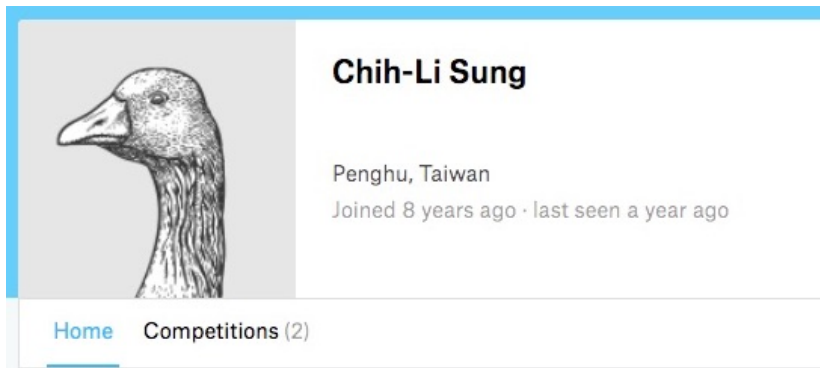


**Chih-Li Sung**

Penghu, Taiwan
Joined 8 years ago · last seen a year ago

Home   Competitions (2)

Figure 4:

# Introduce yourself

- Your name
- Major
- Senior or Junior
- What time is it now in your place