# Foundation of Data Science

## UCS548

# Course Information

- [Syllabus Copy](#)

- [Academic Calendar](#)

- Attendance Requirement: **>75%**

# Tentative Evaluation Plana

The Marks breakup of the course is as follows:

- LAB Evaluation : 15 Marks (Week 16, Nov 25-29, 2023 in scheduled LABs)
- Sessional Quiz : 10 Marks (Week 7, Sept 9-13, 2023 in scheduled LABs)
- Assignments from Course Era : 5 Marks (Before EST)[Grade will be weighted with number of attempts and attainment level]
- MST : 35 Marks
- EST :  35 Marks

# What is data science?

*data science is using data to answer questions*

Data science can involve:
- Statistics, computer science, mathematics
- Data cleaning and formatting
- Data visualization

"who combines the skills of software programmer, statistician and storyteller/artist to extract the nuggets of gold hidden under mountains of data"

- **Data Science is an interdisciplinary field that focuses on extracting knowledge from data sets which are typically huge in amount.**

- **The field encompasses analysis, preparing data for analysis, and presenting findings to inform high-level decisions in an organization.**

- **As such, it incorporates skills from computer science, mathematics, statics, inform.**

- **Some of the techniques utilized in Data Science encompass machine learning, visualization, pattern recognition, probability modeling data, data engineering, signal processing, etc.**

# Why do we need data science?

| Value | Prefix |
|-------|--------|
| $10^{24}$ | Yotta |
| $10^{21}$ | Zetta |
| $10^{18}$ | Exa |
| $10^{15}$ | Peta |
| $10^{12}$ | Tera |
| $10^{9}$ | Giga |
| $10^{6}$ | Mega |

There is an estimated 1.2 zettabytes worth of information currently available - and this number is growing exponentially.

# What is big data?

**Volume**
More and more data is becoming increasingly available

**Velocity**
Data is being generated at an astonishing rate

**Variety**
The data we can analyse comes in many forms

# What is a data scientist?

*Somebody who uses data to answer questions.*

Hacking Skills

Math & Statistics Knowledge

Machine Learning

Data Science

Danger Zone!

Traditional Research

**What skills does a data scientist embody?**

Substantive Expertise

# Why do data science?



Top 10 Emerging Jobs, 2017

| Job | Rate of Growth |
|---|---|
| Machine Learning Engineer | 9.8x |
| Data Scientist | 6.5x |
| Sales Development Representative | 5.7x |
| Customer Success Manager | 5.6x |
| Big Data Developer | 5.5x |
| Full Stack Engineer | 5.5x |
| Unity Developer | 5.1x |
| Director of Data Science | 4.9x |
| Brand Partner | 4.5x |
| Full Stack Developer | 4.5x |

Rate of Growth (2012 – 2017)

*According to Glassdoor, in which they ranked the top 50 best jobs in America, Data Scientist is THE top job in the US in 2017, based on job satisfaction, salary, and demand.*

# Examples of data scientists

# Data science in action!

## LETTERS

## Detecting influenza epidemics using search engine query data

Jeremy Ginsberg[1], Matthew H. Mohebbi[1], Rajan S. Patel[1], Lynnette Brammer[2], Mark S. Smolinski[1] & Larry Brilliant[1]

# Prerequisites for Data Science

The following are the three essential traits of Data Scientist:

**Curiosity**





**Communication Skills**



shutterstock.com · 1167081736

**Curiosity:** Only when you ask questions, you will have a better understanding of the business problem.

**Common Sense:** To identify new *ways to solve a business problems* and to detect priority problems.

**Communication Skills**: A Data Scientist needs to *communicate their findings* to business teams to act upon the insights

Skills required for Data Scientist

# Data Science Three skill tracks:

## Engineering , Analysis, Modeling

**Engineering**

- **Involves in building the data pipeline infrastructure.**

- **It involves the software and the hardware used to store the data and perform data ETL (i.e., extract, transform, and load).**

- **Store and compute data on the cloud.**

- **The fundamental building block for automation is maintaining the data pipeline through modular, well-commented code and version control.**

- **Key task involved are:-**

## Engineering

- **Key task involved are:-**

  1. *Data Environment*: **Designing and setting up the entire environment to support data science workflow is the prerequisite for data science projects. It may include setting up storage in the cloud, Kafka platform, Hadoop and Spark cluster, etc**

  2. *Data Management*: **Automated data collection, that includes parsing the logs (depending on the stage of the company and the type of industry you are in), web scraping, API queries, and interrogating data streams. Determine and construct data schema to support analytical and modeling needs. Use tools, processes, guidelines to ensure data is correct, standardized, and documented.**

  3. *Production*: **Involves the whole pipeline from data access, preprocessing, modeling to final deployment. It is necessary to make the system work smoothly with all existing software stacks.**

# Data Science Three skill tracks:

## Analysis

- **Analysis turns raw information into insights in a fast and often exploratory way.**
- **In general, an analyst needs to have decent domain knowledge, do exploratory analysis efficiently, and present the results using storytelling.**
- **Key point includes are:-**

1. *Domain Knowledge*: **understanding of the organization or industry where you apply data science. You can't make sense of data without context.**

2. *Exploratory Analysis*: **team look at as much data as possible so that the decision-makers can get a sense of what's worth further pursuing. It often involves different ways to slice and aggregate data.**

3. *Storytelling*: **It is the art of telling people what the numbers signify. It usually requires data summarization, aggregation, and visualization. It is crucial to answering the following questions before you begin down the path of creating a data story.**
   - **Who is your audience?**
   - **What do you want your audience to know or do?**
   - **How can you use data to help make your point?**

# Data Science Three skill tracks:

**Modeling**

- **A process that dives deeper into the data to discover the pattern we don't readily see.**
- **A model only occupy a small part of a typical data scientist's day-to-day time.**
- **Some of the models are :-**

1. *Supervised Learning*: **In supervised learning, each sample corresponds to a response measurement. There are two flavors of supervised learning: regression and classification.**
   - ✓ **In regression, the response is a real number, such as the total net sales in 2017 for a company or the yield of corn next year for a state.**
   - ✓ **The goal for regression is to approximate the response measurement as much as possible.**
   - ✓ **In classification, the response is a class label, such as a dichotomous response of yes/no.**
   - ✓ **The response can also have more than two categories, such as four segments of customers**

1. *Un-supervised Learning*: In unsupervised learning, there is no response variable. Clustering approach is used for data analysis.

2. *Customized model development*: A data scientist may need to develop new models to accommodate the subtleties of the problem at hand. For example, people may use Bayesian models to include domain knowledge as the modeling process's prior distribution.

# Data Science Prerequisites

**What type of problem you are solving?**

*Description***:**

- The primary analytic problem is to summarize and explore a data set with descriptive statistics (mean, standard deviation, and so forth) and visualization methods.

- Data description is often used to check data, find the appropriate data preprocessing method, and demonstrate the model results.

*Comparison* **:**

- The first common modeling problem is to compare different groups. Is A better in some way than B? Or more comparisons: Is there any difference among A, B, and C in a particular aspect?

- The commonly used statistical tests are chisquare test, t-test, and ANOVA. There are also methods using Bayesian methods.

# Data Science Prerequisites

*Clustering* :

- Please note that clustering is unsupervised learning; there are no response variables. The most common clustering algorithms include K-Means and Hierarchical Clustering.

*Classification* :

- For classification problems, there are one or more label columns to define the ground truth of classes. We use other features of the training dataset as explanatory variables for model training. We can use the trained classifier to predict the labels of a new observation.

- The random forest algorithm is usually used as the baseline model to set model performance expectations.

# Data Science Prerequisites

*Regression* :

- Generally used for prediction and to answer the questions:-

- What will be the temperature tomorrow? What is the projected net income for the next season?  How much inventory should we have?

*Optimization* :

- It is an expansion of comparison problem and can solve problems such as:

- What is the best route to deliver the packages? What is the optimal advertisement strategy to promote a new product?.

# What is data?

"Information, especially facts or numbers, collected to be examined and considered and used to help decision making"

Cambridge Dictionary

WIKIPEDIA
The Free Encyclopedia

"A set of values of qualitative or quantitative variables"

https://dictionary.cambridge.org/dictionary/english/data; https://en.wikipedia.org/wiki/Data

# What is data?

"A set of values of qualitative or quantitative variables"

Set: In statistics, the population you are trying to discover something about

Variable: Measurements or characteristics of an item

Qualitative variable: Measurements or information about qualities

Quantitative variable: Measurements or information about quantities or numerical items

# What can data look like? (rarely)

| Name | Country of origin | Sex | Weight (kg) | Height (cm) |
|---|---|---|---|---|
| A. Bee | Canada | M | 75 | 163 |
| C. Dee | UAE | M | 80 | 180 |
| E. Eff | China | F | 72 | 175 |
| G. Haitch | South Africa | F | 68 | 172 |
| I. Jay | Poland | M | 77 | 168 |
| K. Elle | Japan | N/A | 76 | 173 |
| M. Enn | Chile | M | 80 | 190 |

# More common types of messy data

- Sequencing data
- Population census data
- Electronic medical records (EMR), other large databases
- Geographic information system (GIS) data (mapping)
- Image analysis and image extrapolation
- Language and translations
- Website traffic
- Personal/Ad data (eg: Facebook, Netflix predictions, etc)

# Messy data: Sequencing



A volcano plot is produced at the end of a long process to wrangle the raw FASTQ data into interpretable expression data

# Messy data: Census information



A population pyramid plot

# Messy data: Electronic medical records (EMR)

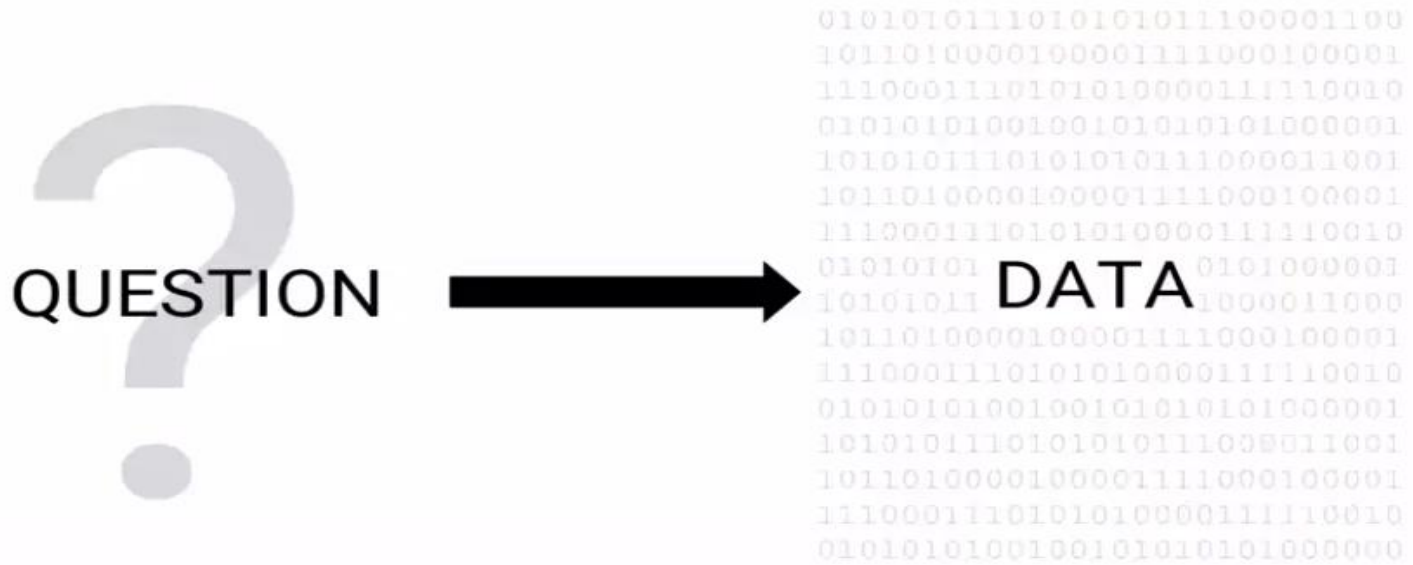https://deepdreamgenerator.com/#tools

The DeepDream software is trained on your image and a famous painting and your provided image is then rendered in the style of the famous painter

# Data is of secondary importance



A good data scientist asks questions first and seeks out relevant data second.