# CDS501: PRINCIPLES & PRACTICES OF DATA SCIENCE & ANALYTICS

Chapter 1: Data Science Process

# Outline

- Introduction

- What is Data Science?

- Applications of Data Science

- Stages of Data Science Project

# How Much Data Do We Create Every Day?

# Introduction

**306.4 billion emails** are sent everyday.

**500 million Tweets** are made everyday.

**95 million photos** and **videos** are shared every day on Instagram.

**1.7MB of data** is created every second by every person during 2020.

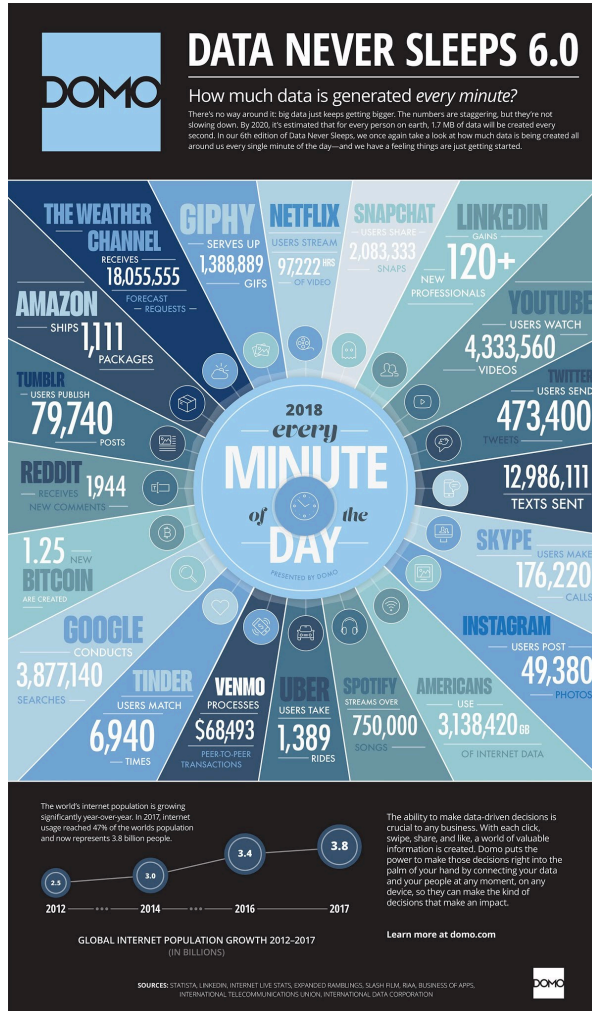**2.5 quintillion bytes of data** are produced by humans every day.

# Introduction

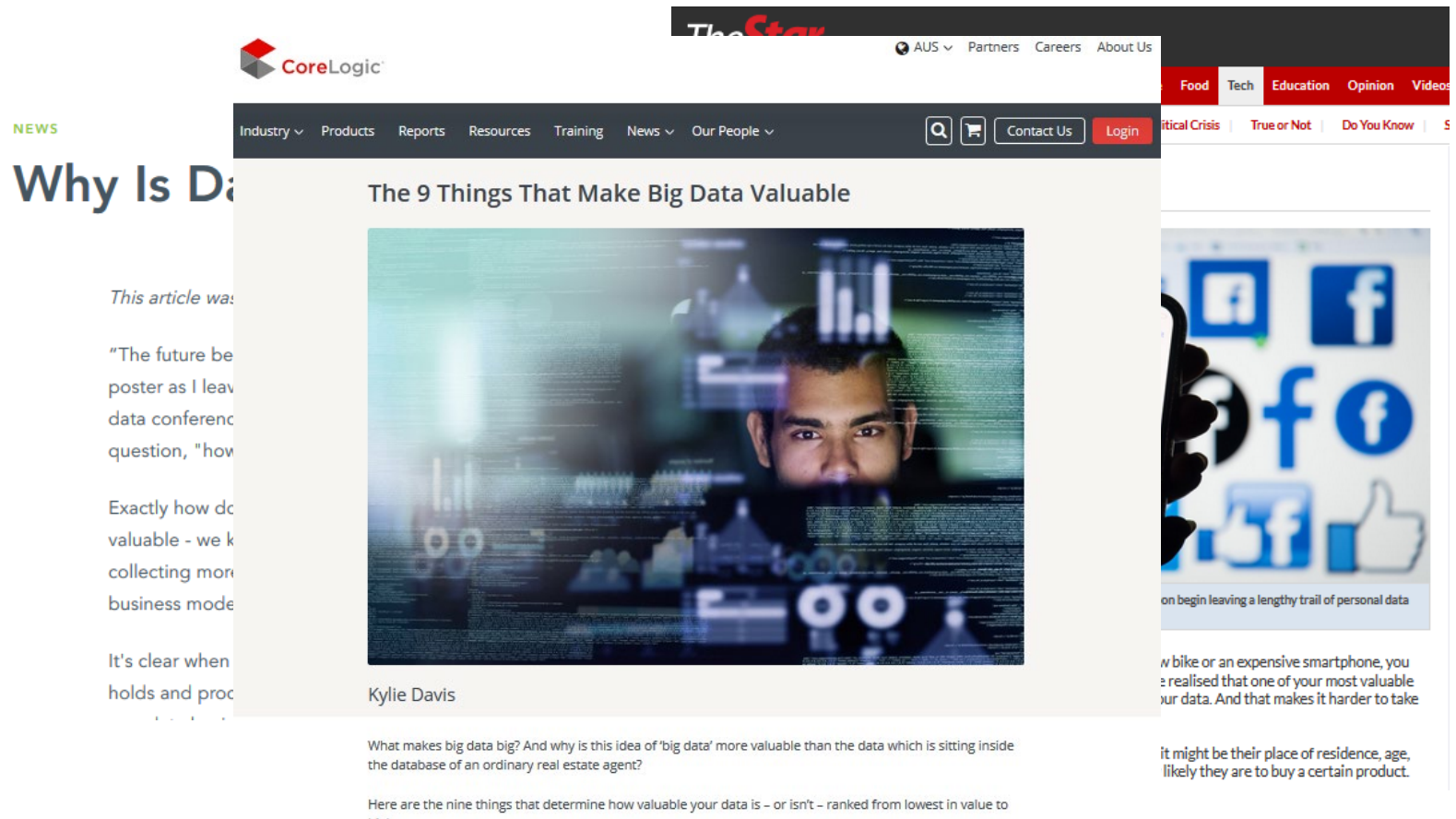1 quintillion = 1 000 000 000 000 000 000

A million million million

2.5 quintillion = 2 500 000 000 000 000 000

# Introduction



- Social media
- Online video sharing platform
- Media service provider
- Online retailer
- Mobile payment service

# Introduction



Data is valuable

# Introduction



Data is valuable

# Introduction



BLOG
## Oil vs. Data – Which is more Valuable?

**PETER SILVA**
PUBLISHED APRIL 09, 2019

It depends who you ask.

In recent years there's been a volley of sorts about data replacing oil as the world's most valuable resource. The basic premise is that in this new digital economy, data and what you extract from that data is similar to oil a century ago. An untapped, massive asset that—depending on how you extract and use it—can have enormous rewards. The raw material's value comes from the refinement into a commodity. For oil, it's the energy extracted; for data, it's in the knowledge extracted.

Economists, professors and even CEOs are touting that data is the new oil in today's economy while others are saying, "no way!" (See the [many] references below for examples.)

The earliest mention of this notion is from 2006. UK Mathemetician and architect of Tesco's Clubcard, Clive Humby said, *"Data is the new oil. It's valuable, but if unrefined it cannot really be used. It has to be changed into gas, plastic, chemicals, etc. to create a valuable entity that drives profitable activity; so must data be broken down, analyzed for it to have value."*

Data is unrefined and cannot really be used

Data needs to be changed into useful form

Manage, Analyze and Model

# What is Data Science?

- Data science is a process of extracting knowledge and insight, and transforming hypothesis to actionable predictions from a huge and diverse set of data through managing, analyzing and modelling using various statistical and computing methods



towardsdatascience.com

**Visualization**

Graphical representation of data to understand and communicate information

**Machine Learning**

Subset of AI to learn and identify patterns in data without human assistance

**Data Science**

**Statistics**

$$s^2 = \frac{\sum (X - \bar{X})^2}{N - 1}$$

Methods of analyzing and interpreting data

# Types of Data Analytic



**What is likely to happen?**
Provide foresight by identifying patterns in data

**Prescriptive analytics**

**What happened?**
Provide insight based on past data

**Predictive analytics**

**What is the best course of action?**
Provide the best options to choose to achieve desired outcome

**Diagnostic analytics**

**Descriptive analytics**

**Why did it happen?**
Examines the cause of past results

ADDED-VALUE CONTRIBUTION

COMPLEXITY

www.scnsoft.com

# Applications of Data Science

## #TECH: Towards safer highways

By Hanum Afandi - January 18, 2021 @ 3:18pm

**Better monitoring**

Through the implementation of S3, the level of highway efficiency will be upgraded and the safety of drivers improved. S3 enables the monitoring and detection of accidents, foreign objects, wild animals, potholes, surface cracks and ponding. The system also covers problems such as water spots, guard rail and slope failure, liquid spillage and road signage damage.

It combines technologies like AI and machine learning to provide notification to the relevant parties for further action. Since the launch of S3 on 19 August 2020, 1,303 incidents were detected in the first month alone. So far, the S3 has helped operations in carrying out immediate rectification with the real-time notifications. Fifty per cent of surface damage and highway asset damage were detected by the system and repairs were made immediately.

**Improving Safety**

To improve security and safety, the company uses the Artificial Intelligence System Analytics (Aisya). By leveraging dashboard cameras and computers installed in every highway patrol car, it is able to obtain images of damage and accidents immediately.

## THE VERGE
TECH  REVIEWS  SCIENCE  CREATORS  ENTERTAINMENT  VIDEO  MORE

POLICY \ TECH \ PRIVACY

## Most US government agencies are using facial recognition

*A new GAO report finds 19 agencies are using some form of the technology*

By Russell Brandom | Aug 25, 2021, 1:23pm EDT

**verge deals**

Subscribe to get the best Verge-approved tech deals of the week.

## Computer Vision

# Applications of Data Science

**Google AI Blog**

The latest from Google Research

## Applying Deep Learning to Metastatic Breast Cancer Detection

Friday, October 12, 2018
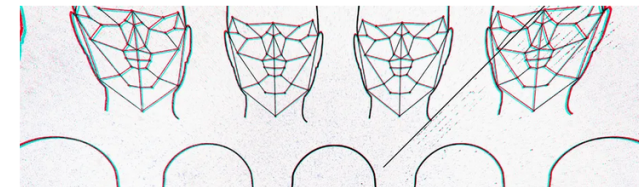
Posted by Martin Stumpe, Technical Lead and Craig Mermel, Product Manager, Healthcare, Google AI

A pathologist's microscopic examination of a tumor in patients is considered the gold standard for cancer diagnosis, and has a profound impact on prognosis and treatment decisions. One important but laborious aspect of the pathologic review involves detecting cancer that has spread (metastasized) from the primary site to nearby lymph nodes. Detection of nodal metastasis is relevant for most cancers, and forms one of the foundations of the widely-used TNM cancer staging.

In breast cancer in particular, nodal metastasis influences treatment decisions regarding radiation therapy, chemotherapy, and the potential surgical removal of additional lymph nodes. As such, the accuracy and timeliness of identifying nodal metastases has a significant impact on clinical care. However, studies have shown that about 1 in 4 metastatic lymph node staging classifications would be changed upon second pathologic review, and detection sensitivity of small metastases on individual slides can be as low as 38% when reviewed under time constraints.

**Imperial College London**

News

Home | College and Campus | Science | Engineering | Health | Business | Search here... | Go ▸

## AI breast cancer screening project wins government funding for NHS trial

*by Ryan O'Hare*
*16 June 2021*

Be the first to comment
Share this
Tweet this
Share on reddit
Share on LinkedIn
Print this story

The partnership, which includes Imperial College London, Google Health, Imperial College Healthcare NHS Trust, St George's Hospitals NHS Foundation Trust, and the Royal Surrey NHS Foundation Trust builds on previous work, in which the researchers trained the algorithm on depersonalised patient data and mammograms from patients in the UK and US.

The findings, published in Nature in January 2020, showed the AI system was able to correctly identify cancers from the images with a similar degree of accuracy to expert radiologists, and demonstrated potential to assist clinical staff in practice.
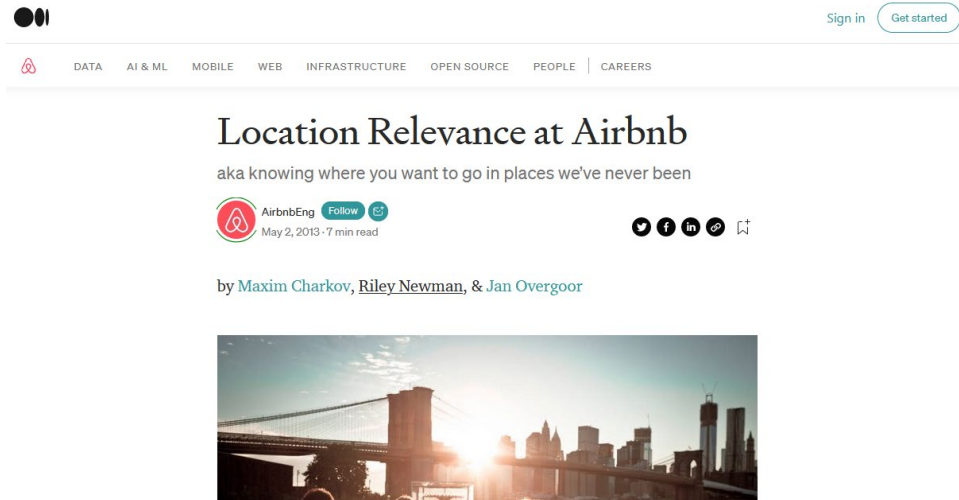
**LATEST NEWS**

**THE KIDS AREN'T ALRIGHT**
Children will face huge increases in extreme climate events in their lifetimes

**IMPACT GRADUATION**
Imperial's flagship BAME talent programme celebrates its seventh cohort

## Healthcare Services

# Applications of Data Science



Recommendation Systems

# Applications of Data Science

- Is there any relationship between weather and sales?

- Is weather is influencing the sales of your shop?

Humidity  Temperature  Sales  Quantity Sold

# Applications of Data Science

- Can we detect a person with depression disorder?
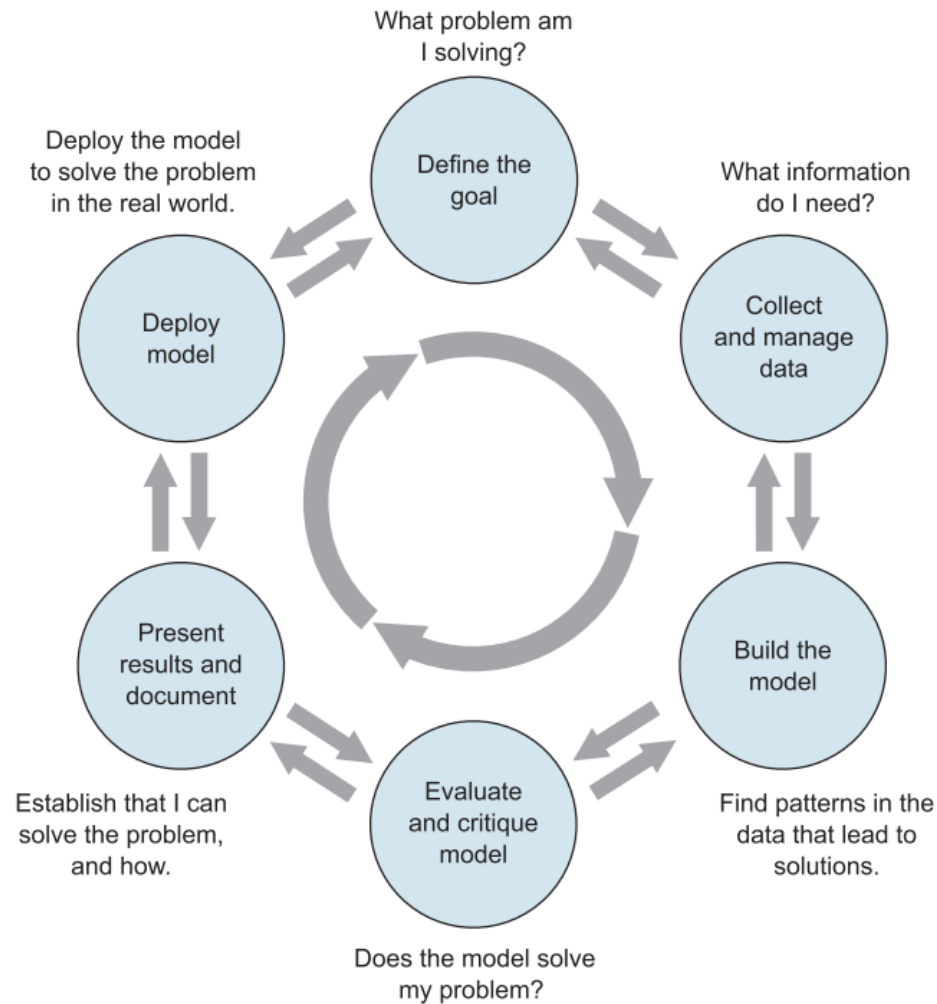
Social Media Activity

Medical Data

Person with and without Depression

# The Roles in Data Science Projects

| Project sponsor | Client | Data scientist | Data architect | Operations |
|---|---|---|---|---|
| • Person who wants the result<br>• Decide the project success/failure<br>• Keep them updated | • The end user<br>• Has interest in your model<br>• Who will be using your model | • Design the project steps<br>• Apply the process<br>• Pick the statistic and machine learning techniques that will be used | • Responsible for data and its storage<br>• Manage data warehouse<br>• e.g. database administrator | • Responsible for acquiring data and deliver the final results<br>• Responsible for deployment |

# Stages of a Data Science Project

# Defining the Goal

- Discuss and work with stakeholders/sponsor to understand and identify business problems

- An online shop is making losses
  - Products – which products are not selling well?
  - Customers – how to identify customers who are more likely to buy?
  - Fraudulent orders – how to identify fraudulent orders?

# Defining the Goal

- Define a specific, quantifiable and achievable goal

- "The detection accuracy rate must be at least 85%"

# Defining the Goal

- Specific goal allows stopping condition and acceptance criteria to be defined

- Otherwise the project will go unbounded

# Data Collection and Management

- The most time-consuming step in the process

- Most important and crucial step

# Data Collection and Management

- Identify data that is relevant to the question

- What are the attributes that are related to the target?

- Do you have the attributes that are related to the target?

# Data Collection and Management

- Number of products in the order – unusually large orders could be fraudulent

- Type of products in the order – wary of orders that are uncommonly purchased together.

- Billing/Shipping address of the order – not a real location or not a residential location

- Number of receiving orders in a timeframe – placing multiple orders at the same time

# Data Collection and Management

- Do you need additional attributes to address problem?

- Do you have sufficient examples or not?

# Data Collection and Management

- Remove redundant and irrelevant examples

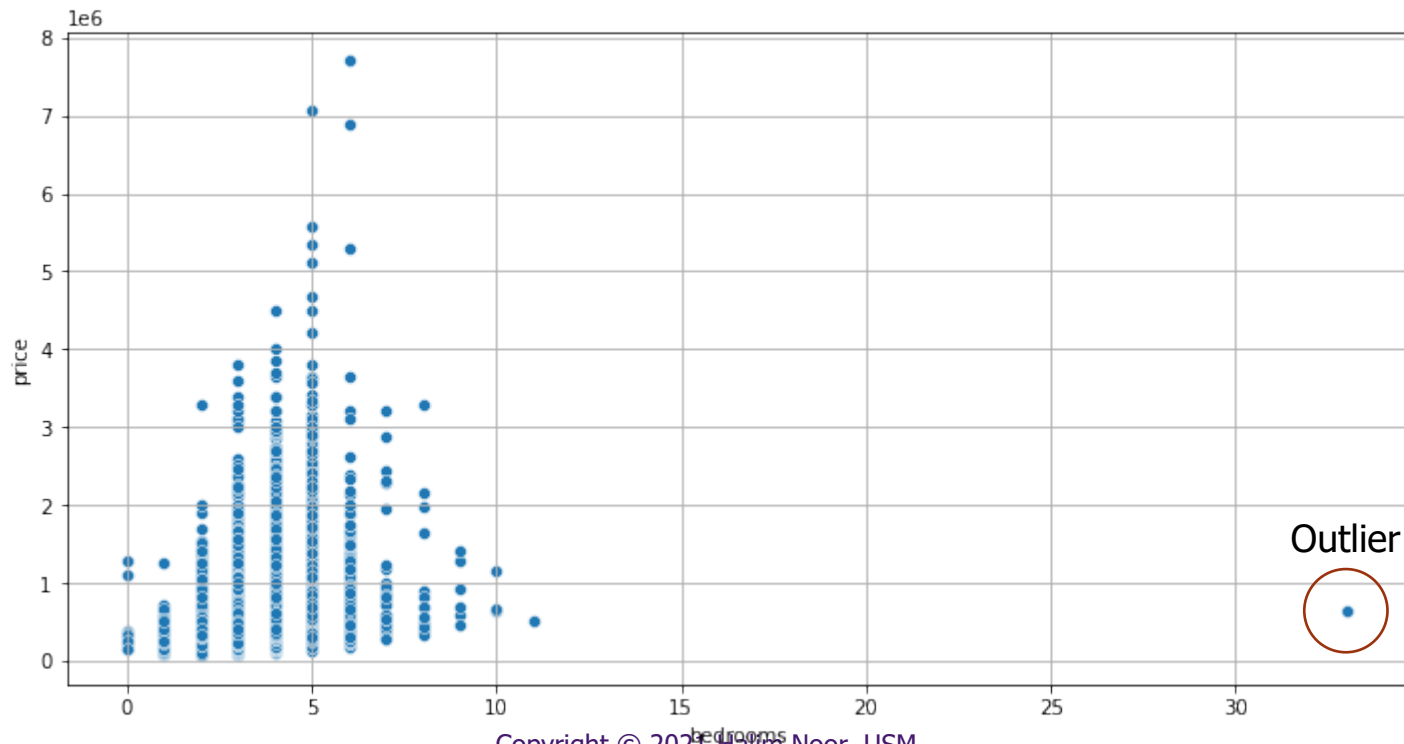- The collected data often needs to be cleaned from missing values and outliers,

# Data Collection and Management

- An attribute of an example (row) has no value

- Are there missing values? How do handle them?

# Data Collection and Management

- Values that are significantly differ from others

- Are there outliers? How do we handle them? Should we remove it or keep it?

# Data Modeling

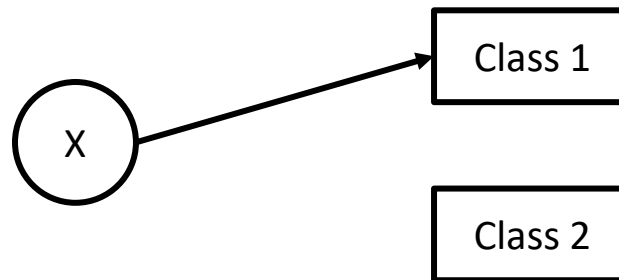- To confirm our insight and hypothesis about the data
  - Number of rooms and House prices – positive relationship

- Fit the data using linear model
  - If slope is positive = positive relationship
  - If slope is negative = negative relationship

# Data Modeling

- A model is an approximation of the data that describes the relationship between the attributes

- A model can be used to make predictions

# Data Modeling
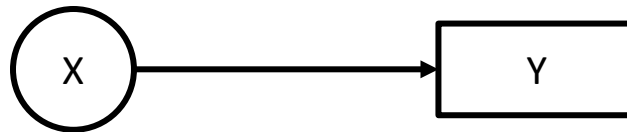
- Classification – deciding if something belongs to one category or another

- Placed order and Fraudulent or Not Fraudulent



X is an order, fraudulent order (Class 1) or valid order (Class 2)
X is social media activity, depression (Class 1) or not (Class 2)

# Data Modeling

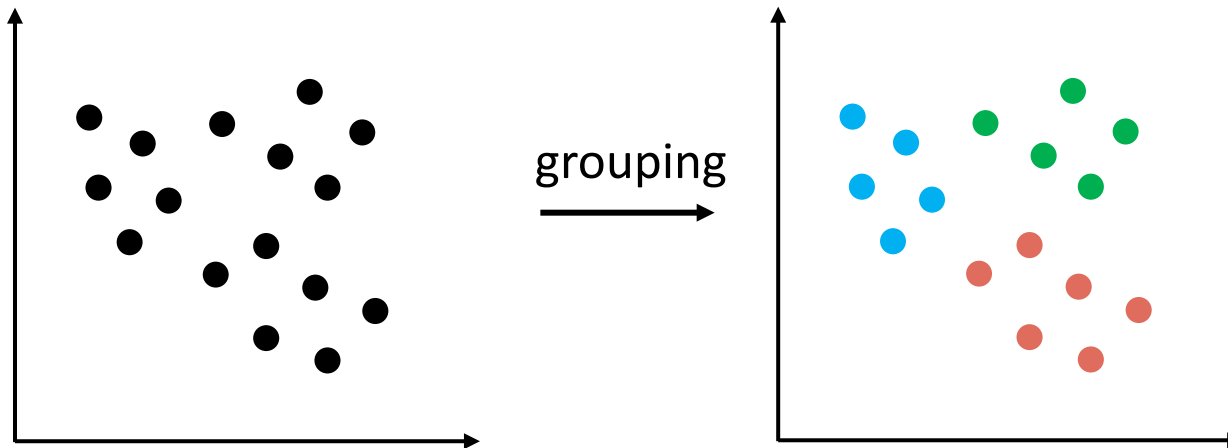- Regression (scoring) – estimating a numeric value



X is house attributes, house price (Y)
X is weather, sales (Y)

# Data Modeling

- Clustering – grouping items into most similar groups



grouping

# Model Evaluation

- Determine if the model meets the goals

- Is it accurate enough for your needs? Does it perform better than whatever estimate that is currently being used?

- Do the results of the model make sense in the context of the problem domain?

- If no, then repeat the Modeling step (or the steps before it)

# Presentation and Documentation

- Document the model for those who are responsible for using, running, and maintaining the model once it has been deployed

- Presentation must not be technical

- Presentation must make sense to the human brain and easy to understand – use visualization

- Highlight the most interesting findings or recommendation (if any)

# Model Deployment

- Ensure the model run smoothly

- Model can be updated when needed

- Monitor the performance of the model

- Why the model's decision is being overridden frequently?

- Is the model incomplete?

# Summary

- Data science is a process of extracting knowledge or insight from data

- Data science project involves many roles and skills – back-and-forth between data scientist and project stakeholders

- Project goal must be specific, measurable and quantifiable

End