

# Syllabus and Introduction

Petchara Pattarakijwanich

Introduction to Data Science, 19 August 2022

# Outline of This Week

- Syllabus of This Course
- What are Data Science and Machine Learning?
- Simple Examples of DS&ML
- Types of Problems in DS&ML  
(Regression vs Classification, Supervised vs Unsupervised)
- Example of Problems in Real-Life
- Methodology (Training → Validation → Application)
- Under/Over-fitting, Bias-Variance Tradeoff

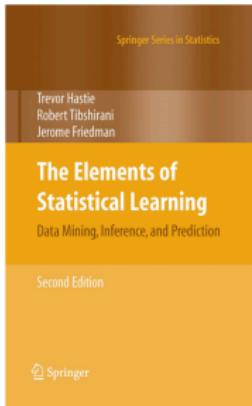
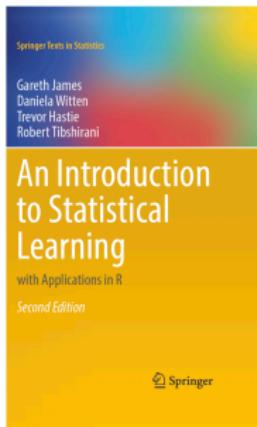
# Syllabus of This Course

Topics:

- Basic Programming in Python
  - Regression and Classification
  - Resampling Methods
  - Subset Selection
  - Shrinkage Methods (Ridge & Lasso)
  - Dimensionality Reduction Method
  - Decision Tree, Bagging, Boosting, Random Forest
  - Support Vector Machine
  - Neural Network
  - Unsupervised Learning
-

# Syllabus of This Course

- **An Introduction to Statistical Learning** (Gareth James et al.)  
<https://www.statlearning.com>
- **The Elements of Statistical Learning** (Trevor Hastie et al.)  
<https://hastie.su.domains/ElemStatLearn/>
- **Python Data Science Handbook** (Jake VanderPlas)  
<https://jakevdp.github.io/PythonDataScienceHandbook/>



# Syllabus of This Course

Grading:

- Homework and Exercise 50%
- Final Project 50%

# Syllabus of This Course

Calendar:

- Weeks with classes:
  - 19/26 Aug
  - 2/9/16/23/30 Sep
  - None in Oct
  - 4/11/18/25 Nov
  - 2 Dec
- Weeks with no class: 14/21/28 Oct
- Midterm exam Week: 3-7 Oct
- Final exam Week: 6-16 Dec

ฉบับที่ 1

# What is Data Science?

## What is a Data Scientist?

Data scientists use technology to glean insights from large amounts of data they collect. It's a field that requires statistics, quantitative reasoning and computer programming skills. On top of all that, you need to be a good communicator so you can report your research findings and explain how they address a larger question you're trying to answer.



"A data scientist really is a scientist at heart," says Scott Beliveau, chief of the enterprise advanced analytics branch within the U.S. Patent and Trademark Office's Office of the Chief Technology Officer. "But rather than using chemicals or other things, a data scientist uses data — numbers, zeros, sometimes it's textual information — to try and solve and answer problems."

MEDIAN SALARY  
**\$98,230**

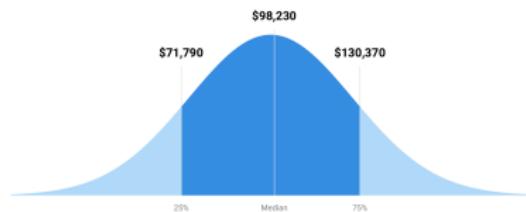
UNEMPLOYMENT RATE  
**2%**

NUMBER OF JOBS  
**19,800**

While data science is still a new career field, employers are increasingly recognizing the value of professionals with this expertise. Today, you'll find data scientists working at a range of organizations, including tech startups, government agencies, large companies and research institutions.

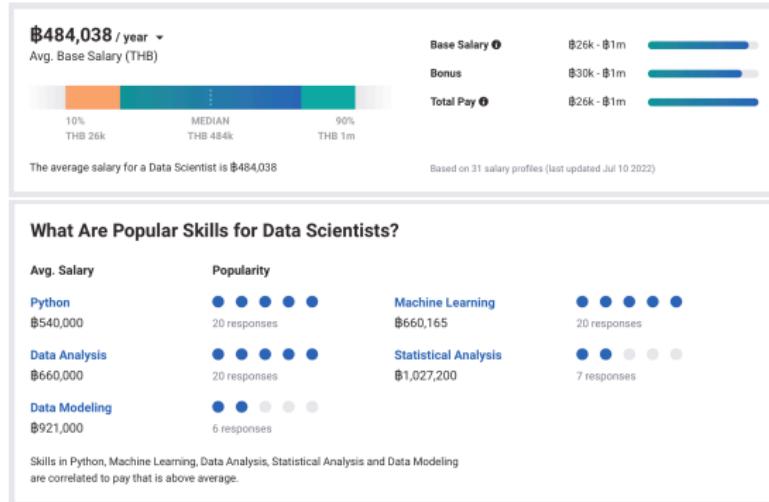
## How Much Does a Data Scientist Make?

Data Scientists made a median salary of \$98,230 in 2020. The best-paid 25 percent made \$130,370 that year, while the lowest-paid 25 percent made \$71,790.



<https://money.usnews.com/careers/best-jobs/data-scientist>

# What is Data Science?

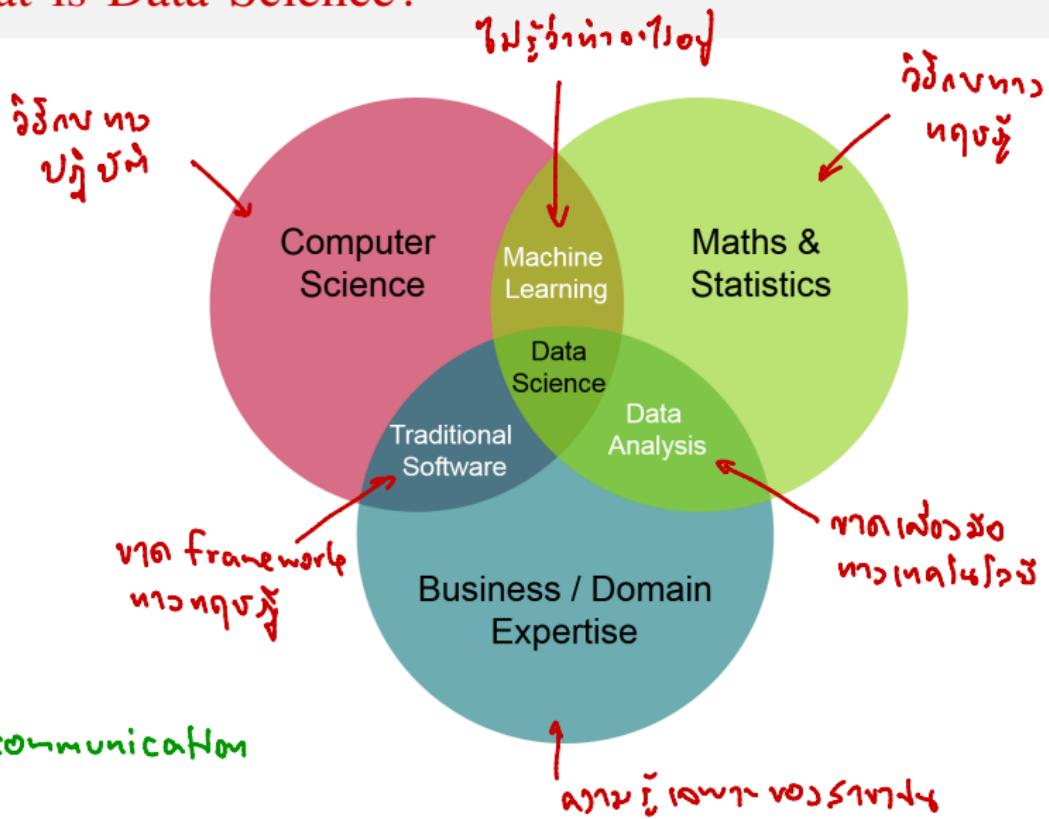


[https://www.payscale.com/research/TH/Job=Data\\_Scientist/Salary](https://www.payscale.com/research/TH/Job=Data_Scientist/Salary)

## What is Data Science?

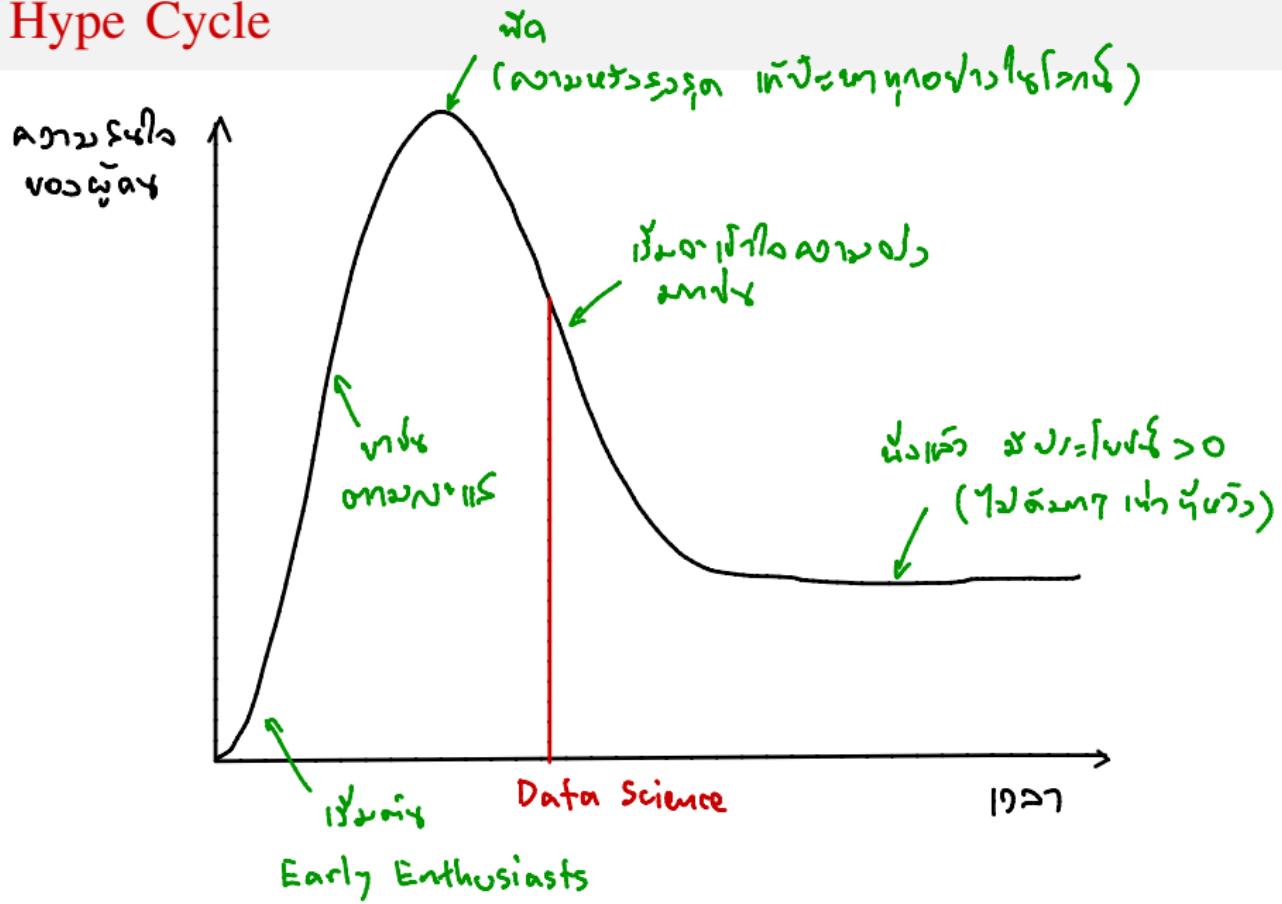
<https://th.linkedin.com/>

# What is Data Science?



[https://thedataScientist.com/data-science-without-programming/data\\_science\\_venn\\_diagram/](https://thedataScientist.com/data-science-without-programming/data_science_venn_diagram/)

# Hype Cycle



# What is Data Science?

- Extract insights from data → Actionable Plan
- Use scientific methods
- Techniques from mathematics and statistics
- Algorithms and data management from computer science

# Simple Example of Data Science

ជំនាញតែងចូលរួម្រោះ 2 ពាន់

① តើកីឡា និងវគ្គិស្ស 5 ខ្លួនអាមេរិក ( $\approx 1500$  នាក់)

- GPA > 2.5
- នាមឈូណី ONET នឹង 1, 2, 3, ...
- នាមឈូណី
- ទំនាក់ទំនង
- ទិន្នន័យអ៊ូរុប (តិច/លើម, ទិន្នន័យ/ទិន្នន័យ/ទិន្នន័យ)
- នៅក្នុងសម្រាប់ 1
- ជាកីឡា ដែល?

② តើកីឡា និងវគ្គិស្ស ឱ្យបាន ( $\approx 300$  នាក់)

- តិច
- មិនមែនជាកីឡា

រូបរាង  
Pattern

} ឱ្យ Pattern ដោយ  
លាងការពាណិជ្ជកម្ម

រូបរាង  
 $\Rightarrow$  លាងការពាណិជ្ជកម្ម

# Simple Example of Data Science

ជំនាញតាមលក្ខណៈ 2 រយៈ  
នៅពេលវិភាគ នឹង ឯកសារ  $\vec{y} = f(\vec{x})$

① តើកីឡា មានប័ណ្ណីស ៥ ខែកាត់ខ្លួន ( $\approx 1500$  នាព.)

- GPA ឬរាយរោគ ( $x_1$ )

- នាមឈាម ONET និង 1, 2, 3, .. ( $x_2, x_3, x_4$ )

- ការងារអនុវត្ត ( $x_5$ )

- នាយករដ្ឋមន្ត្រី ( $x_6$ )

- ទីតាំងនៃបុរាណ (ខេត្ត/ឈូរ, ស្វែន/ឈូរ/សង្កាត់)

- នាមឈាមលើក 1 ( $y_1$ ) } Dependent variables

- តម្លៃតុលាទី? ( $y_2$ ) } Targets output

② តើកីឡា មានប័ណ្ណីស ៥៩៩ ( $\approx 300$  នាព.)

- តម្លៃតុលាទី  $\vec{x}$

- តម្លៃតុលាទី ធនធាន  $\vec{y}_{predicted} = f(\vec{x})$

Independent Variables  
Features  
Inputs

# Simple Example of Data Science

①

x <sub>1</sub>	x <sub>2</sub>	..	..	y <sub>1</sub>	y <sub>2</sub>

②

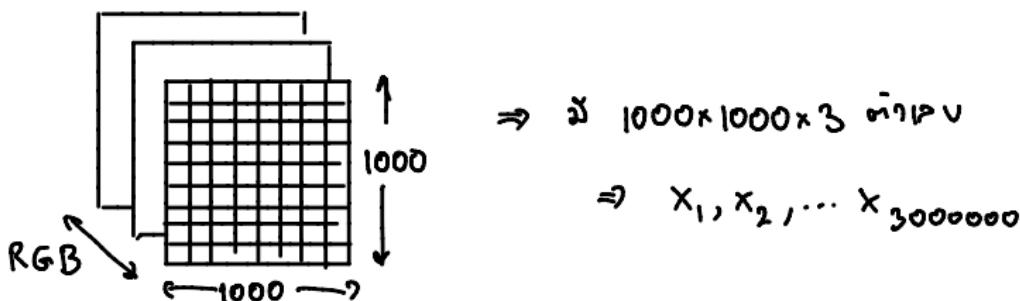
x <sub>1</sub>	x <sub>2</sub>	..	..	y <sub>1</sub>	y <sub>2</sub>
				y <sub>1</sub>	y <sub>2</sub>

ນັ້ນ  $f$  ສັບອຸປະກອດ  $y = f(x)$

ອີງ  $f$  ສັບອຸປະກອດ  $y = f(x)$

# Simple Example of Data Science

ຫຍຸ້ງ : Image classification ລາວພິບເຂດ ແກ້ໄຂແລ້ວ ?



①

$x_1, x_2, \dots$	$x_{3000000}$	$y$
	ເຊົາ ; ເຊົາ	ໜີ

ມີຕົວຢ່າງ  $f$  ລາຍກູດ ①

10,000

②

$x_1, x_2, \dots$	$x_{3000000}$	$y$
	ເຊົາ ;	ໜີ

ມີຕົວຢ່າງ  $f$  ມີຄຸນ ②

# Simple Example of Data Science

①

	$x_1$	$x_2$	...	...	$y_1$	$y_2$
1500			Training data			
			පුරුෂ අංක			

②

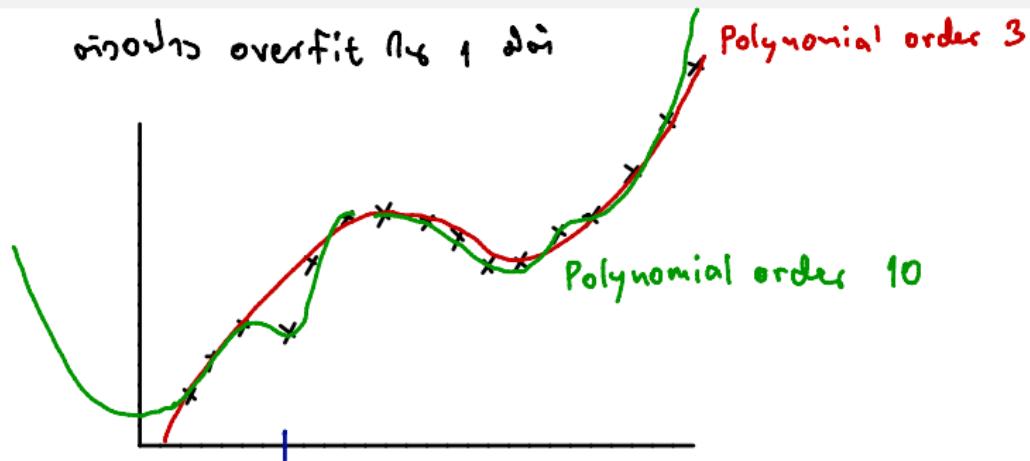
					$y_1$	$y_2$
300						
					යුතු	

මෙහි  $f$  සඳහා  $y = f(x)$

දැක්වාගැනීමේදී overfit?

මෙහි  $f$  විභාගය  $y = f(x)$

# Simple Example of Data Science



# Simple Example of Data Science

Input  $x_1, x_2, \dots, x_i \Rightarrow$  Output  $y_1, y_2$

Classification

- Quantitative [ GPA, ACT, SAT, ... ]  
(Real Number)
- Qualitative/Categorical [ Sex, Marital Status, ... ]  
(属性特征)
- Ordered Categorical [ Letter Grade, Exam/Pass/Fail ]  
A, B, C, ...

Output

- Quantitative  $\rightarrow$  Regression
  - Qualitative  $\rightarrow$  Classification
- } Supervised Learning

# Real-life Example of Data Science

- Identification of cells, tissues, bacteria, etc. from medical images
- Identification of dog breed from images
- Reading handwritten text from image (image to text)
- Personal identification from image (facial recognition)
- Automatic caption (speech to text)
- Predict property prices from location and other information
- Predict sales of businesses
- Predict whether a debt will be defaulted

# Real-life Example of Data Science

- Predict what diseases a person might have from medical record
- Predict a person's health from DNA (personalized medicine)
- Predict performance of employees from their CVs
- Predict health/activity from smart watch data
- Identify writers from writing styles
- Predict of credit card fraud
- Predict consumer buying behavior
- Predict election outcome from social media trend