

Introduction (Cont.)

Petchara Pattarakijwanich

Introduction to Data Science, 26 August 2022

Types of Data Science and Machine Learning

- Supervised Learning

คือ machine learning ที่เราสอน (รู้เป้าหมาย)

(training data มาสอนแล้ว \Rightarrow ถ้า test data มาแล้ว \Rightarrow ให้รู้)

- Regression รู้ว่าคือ machine learning ที่เราสอนให้ทำนายค่าตัวเลข (เชิงตัวเลข) \rightarrow Quantitative

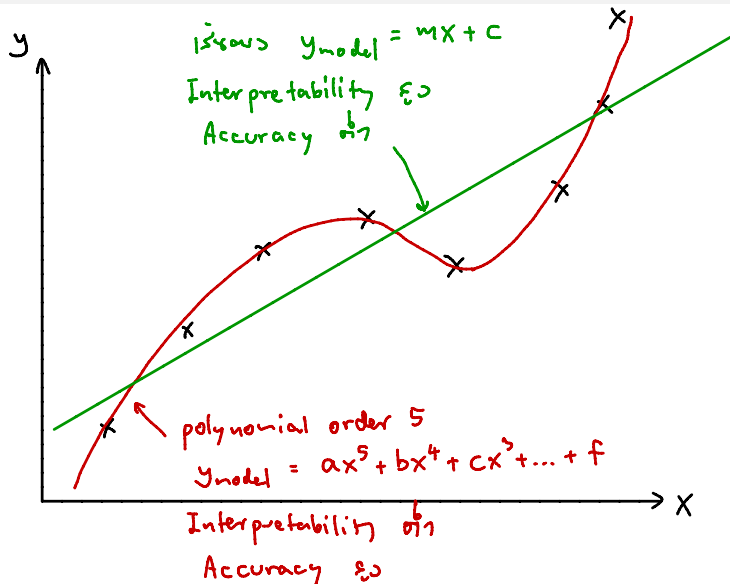
- Classification รู้ว่าคือ machine learning ที่เราสอนให้ทำนายค่าที่เป็นประเภท (เชิงคุณภาพ) \rightarrow Qualitative

- Unsupervised Learning

- Visualization
- Correlation
- Clustering

[Data Exploration]

Interpretability vs Accuracy



Methodology of Data Science

- Training Data

รู้ input รู้ output ฝึกฝนโมเดล

- Test (Validation) Data

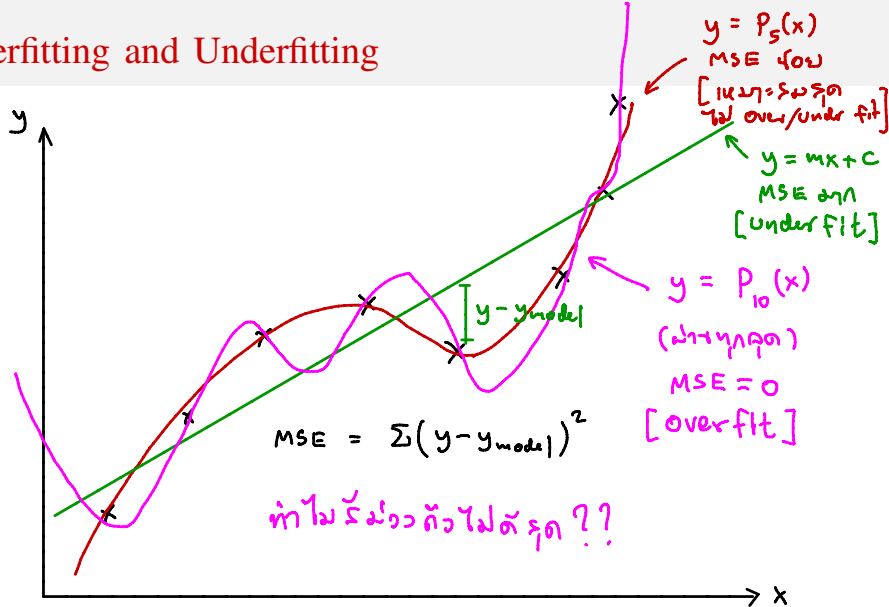
รู้ input รู้ output ฝึกทดสอบโมเดล

[ใช้โมเดลหา output → เปรียบกับค่าจริง]

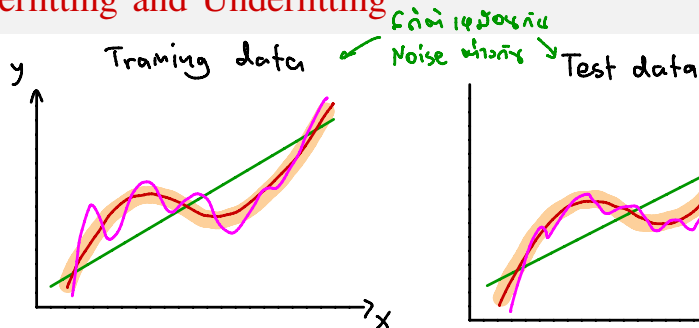
- Real-world use

รู้ input ไม่รู้ output ฝึกโมเดลหา output ⇒ ใช้งานจริง

Overfitting and Underfitting



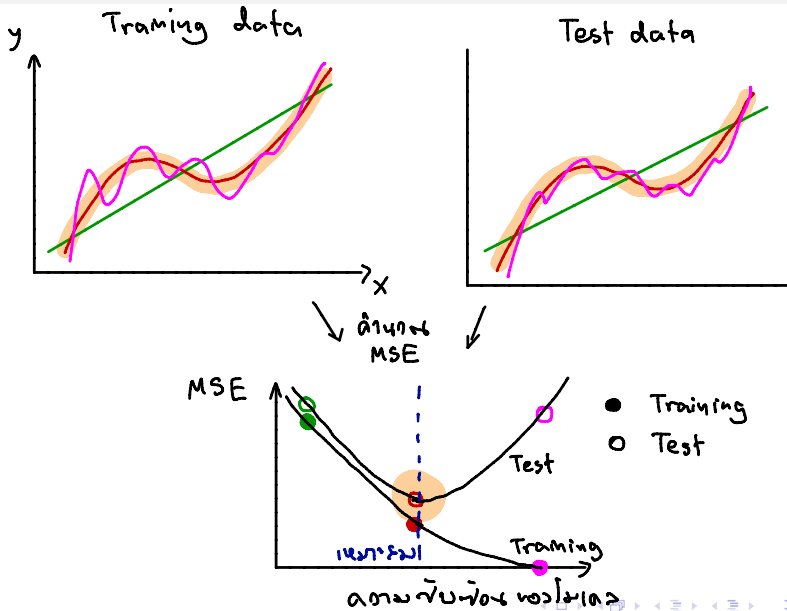
Overfitting and Underfitting



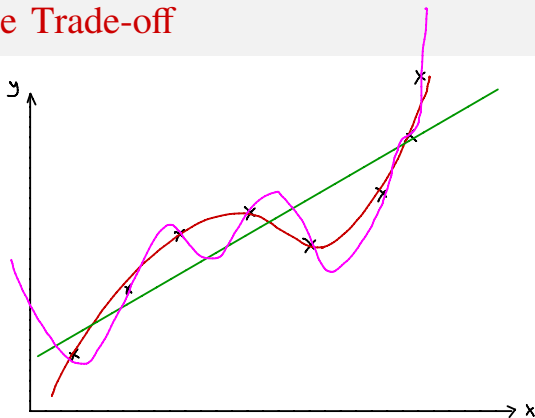
- model training data ၁၀၀%
- model training data ၁၀၀%
- model training data ၁၀၀%

- model test data ၇၀%
- model test data ၁၀၀%
- model test data ၁၀၀%

Overfitting and Underfitting



Bias-Variance Trade-off



Bias ไม่ตรงตามตัวจริง
 ไม่ละเอียดถี่ถ้วน

ง่าย

น้อย

น้อย ?????

Variance ไม่ sensitive
 ต่อการเปลี่ยนแปลง

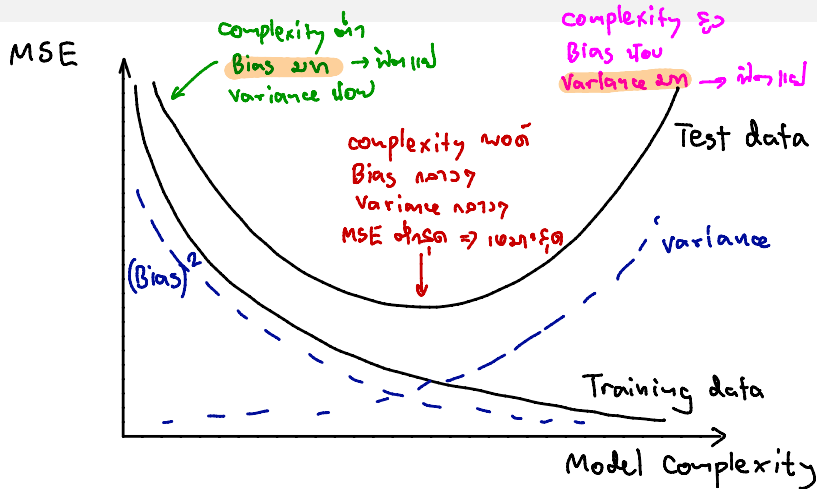
ง่าย

ง่าย

ง่าย ?????

↔
Bias-Variance
Trade off

Bias-Variance Trade-off



$$MSE_{test} = \text{Noise} + (Bias)^2 + \text{Variance}$$

Bias-Variance Trade-off

- The *bias* error is an error from erroneous assumptions in the learning *algorithm*. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting).
- The *variance* is an error from sensitivity to small fluctuations in the training set. High variance may result from an algorithm modeling the random *noise* in the training data (*overfitting*).