



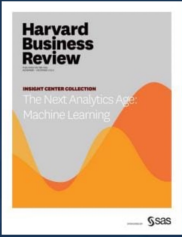
Machine Learning

Prof. Adil Khan

Who Am I?

- Adil Khan
- Professor in Innopolis University (IU), Russia
- Expertise: Machine Learning, Deep Learning
- Head of Machine Learning Lab at IU
- Director of the Institute of AI at IU





WHITE PAPER

The Next Analytics Age: Machine Learning

A Harvard Business Review Insight Center Collection

The top most emerging trend in Computing and Information Technology

Machine Learning

Large Scale Academic Research

One of the largest number of start-ups

Has changed business in almost every industry

Active Recruitment of ML Engineers

To Effectively use Machine Learning

- One must understand:
 - What is machine learning?
 - The available machine learning methods at your disposal
 - Characteristics of those methods
 - Circumstances under which a method would be most effective
 - Their theoretical underpinnings

Course Objectives

1. Teach you how to **implement** important machine learning methods and **apply** them for solving real-world problems
2. Provide you with both the **theoretical** and **practical** knowledge of machine learning methods
3. **Supervised Methods** (Linear, Polynomial and Logistic Regression, Decision Trees, ANNs, CNNs, etc.)
4. **Unsupervised Methods** (k-means, kmeans++, Hierarchical Clusternig, etc.)
5. **Working with ML Models** (Regularization, Dimensionality Reduction, Ensemble Learning)
6. **Generative Models** (GANs)

Prerequisites

- Basics of
 - Python programming
 - Linear algebra
 - Calculus
 - Probability

Plan for Today

- Foundations of Machine Learning
 - What does it mean for machines to learn?
 - ❖ Task, Experience and Performance
 - Regression
 - Classification
 - Train and Test Error
 - Overfitting and Underfitting
 - Bias-Variance Tradeoff

What Does it Mean for a Machine to Learn?

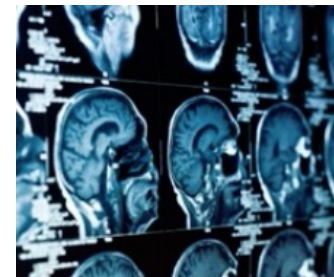
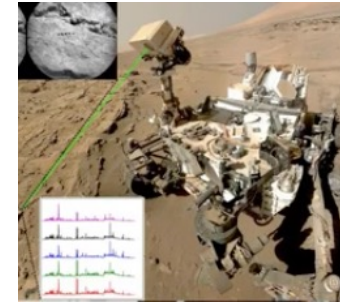
- If you downloaded a copy of Wikipedia, has your computer really learned something?
- Did it make your computer any smarter?

What is Machine Learning?

Machine Learning

- Computer programs that improve their performance at some task through experience

Usual Examples of **Tasks** where ML is being Used



Experience

Targets: For example, malware or Not

Data

$$D = \{(x_i, y_i)\}_{i=1}^N$$

Number of available examples

The diagram shows the equation $D = \{(x_i, y_i)\}_{i=1}^N$ in blue. A blue arrow points from the text 'Targets: For example, malware or Not' to the y_i term. Another blue arrow points from the text 'Predictors: For example, execution behavior of a program' to the x_i term. A third blue arrow points from the text 'Number of available examples' to the superscript N .

Predictors: For example, execution behavior of a program

$$x \in \mathbb{R}^d$$

Examples of Predictors and Response

Input (x)	Output (y)	Application
Home Features	Price	Real Estate
Ad, User info	Click ad? (0/1)	Online Advertising
Image	Object (1, ..., 1000)	Photo Tagging
Audio	Text Transcript	Speech Recognition
English	Chinese	Machine Translation
Image, Radar Info	Position of other cars	Autonomous Driving

Performance

- Needs to be defined according to the given task
- For example: the ratio of correctly identified malwares

Goal of Learning

- Learning or inferring a “functional” relationship between predictors and target


$$D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$$

$$\boldsymbol{x} \in \mathbb{R}^d$$

$$y = f(\boldsymbol{x})$$

$$\hat{f} \approx f \quad \textit{Goal of learning}$$

Putting It All Together (1)

$$D = \{(\mathbf{x}_i, y_i) \in X \times Y : 1 \leq i \leq N\}$$


Learning
Algorithm



$$\hat{f}: X \rightarrow Y$$

Putting It All Together (2)

$$D = \{(\mathbf{x}_i, y_i) \in X \times Y : 1 \leq i \leq N\}$$

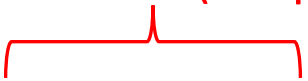
The diagram illustrates the supervised learning process. It starts with a dataset D defined as a set of pairs (\mathbf{x}_i, y_i) from $X \times Y$ for i from 1 to N . A blue arrow points from the dataset D to a gray box labeled "Supervised Learning Algorithm". Another blue arrow points from the algorithm box to the learned function $\hat{f}: X \rightarrow Y$.

“Supervised”
Learning
Algorithm

$$\hat{f}: X \rightarrow Y$$

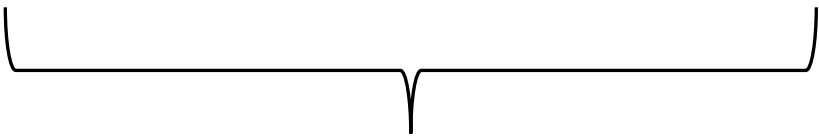
Example: Dataset 1

Dependent Variable
(Response)



Country	Age	Salary	Purchased
France	44	72000	No
Spain	27	48000	Yes
Germany	30	54000	No
Spain	38	61000	No
Germany	40		Yes
France	35	58000	Yes
Spain		52000	No
France	48	79000	Yes
Germany	50	83000	No
France	37	67000	Yes

Independent Variables
(Predictors)



Exampe: Dataset 2

Dependent Variable (Response)

YearsExperience	Salary
1.1	39343
1.3	46205
1.5	37731
2	43525
2.2	39891
2.9	56642
3	60150
3.2	54445
3.2	64445

Independent Variables (Predictors)

Compare the Response Variable

Country	Age	Salary	Purchased
France	44	72000	No
Spain	27	48000	Yes
Germany	30	54000	No
Spain	38	61000	No
Germany	40		Yes
France	35	58000	Yes
Spain		52000	No
France	48	79000	Yes
Germany	50	83000	No
France	37	67000	Yes

$$y = \{c_1, c_2, \dots, c_k\}$$

YearsExperience	Salary
1.1	39343
1.3	46205
1.5	37731
2	43525
2.2	39891
2.9	56642
3	60150
3.2	54445
3.2	64445

$$y \in \mathbb{R}$$

Classification and Regression

Country	Age	Salary	Purchased
France	44	72000	No
Spain	27	48000	Yes
Germany	30	54000	No
Spain	38	61000	No
Germany	40		Yes
France	35	58000	Yes
Spain		52000	No
France	48	79000	Yes
Germany	50	83000	No
France	37	67000	Yes

Classification

YearsExperience	Salary
1.1	39343
1.3	46205
1.5	37731
2	43525
2.2	39891
2.9	56642
3	60150
3.2	54445
3.2	64445

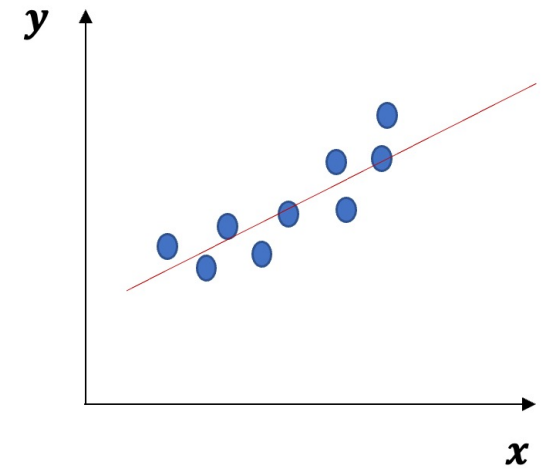
Regression

Estimating f

$$y = f(x; \text{parameters})$$


$$y = f(x; \mathbf{w})$$

$$y = f(x; w_0, w_1) = w_0 + w_1 x$$



$$\underbrace{y = w_0 + w_1 x}_f$$

Assessing the Quality of Learning

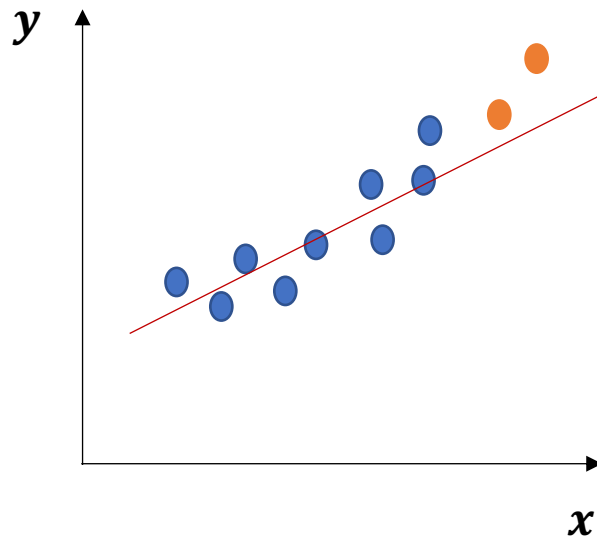
- Let $T_r = \{x_i, y_i\}_{i=1}^N$ be the training data we used to estimate $\hat{f}(x)$.
- To assess the quality of estimate, we can compute

$$\text{MSE}_{\text{Tr}} = \text{Ave}_{i \in \text{Tr}} [y_i - \hat{f}(x_i)]^2$$

- But this is **not a reliable** approach

What can go wrong?

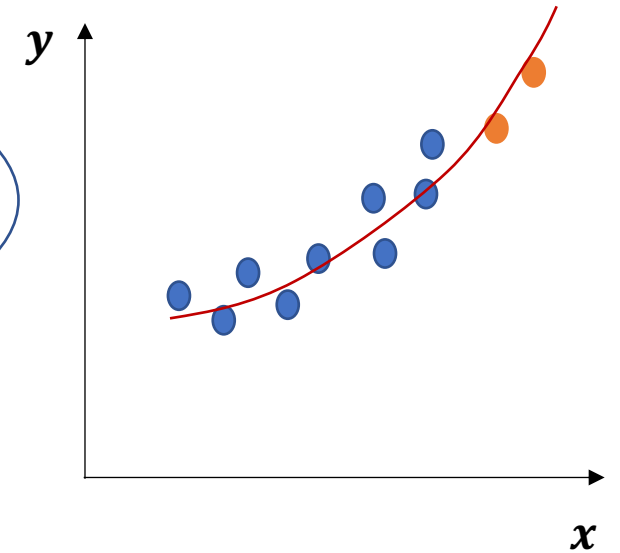
- Training Data
- Unseen Test Data



$$\frac{y = w_0 + w_1 x}{f}$$

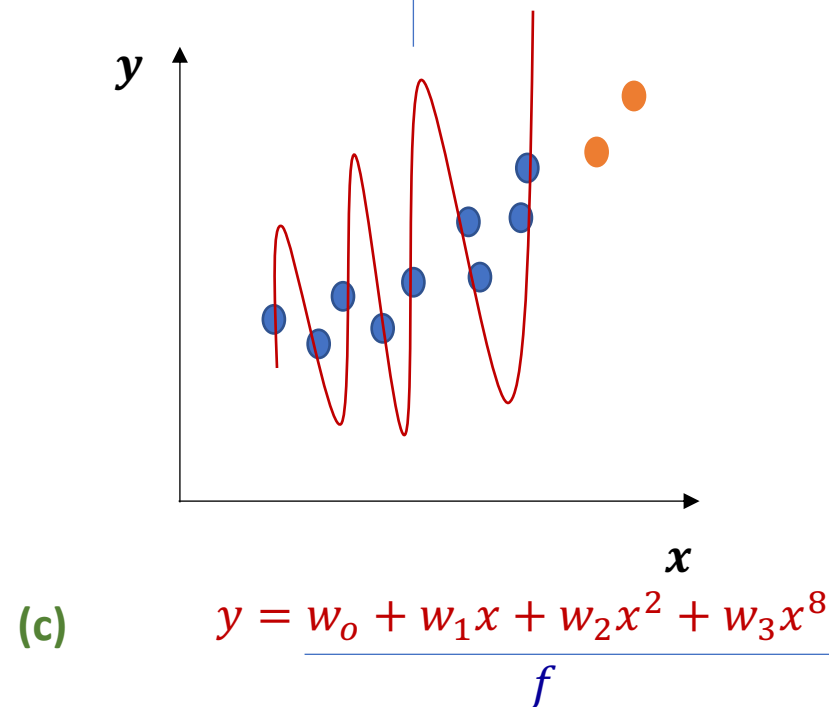
(a)

Smallest Training Error
Most Complex
Largest Test Error



$$\frac{y = w_0 + w_1 x + w_2 x^2}{f}$$

(b)

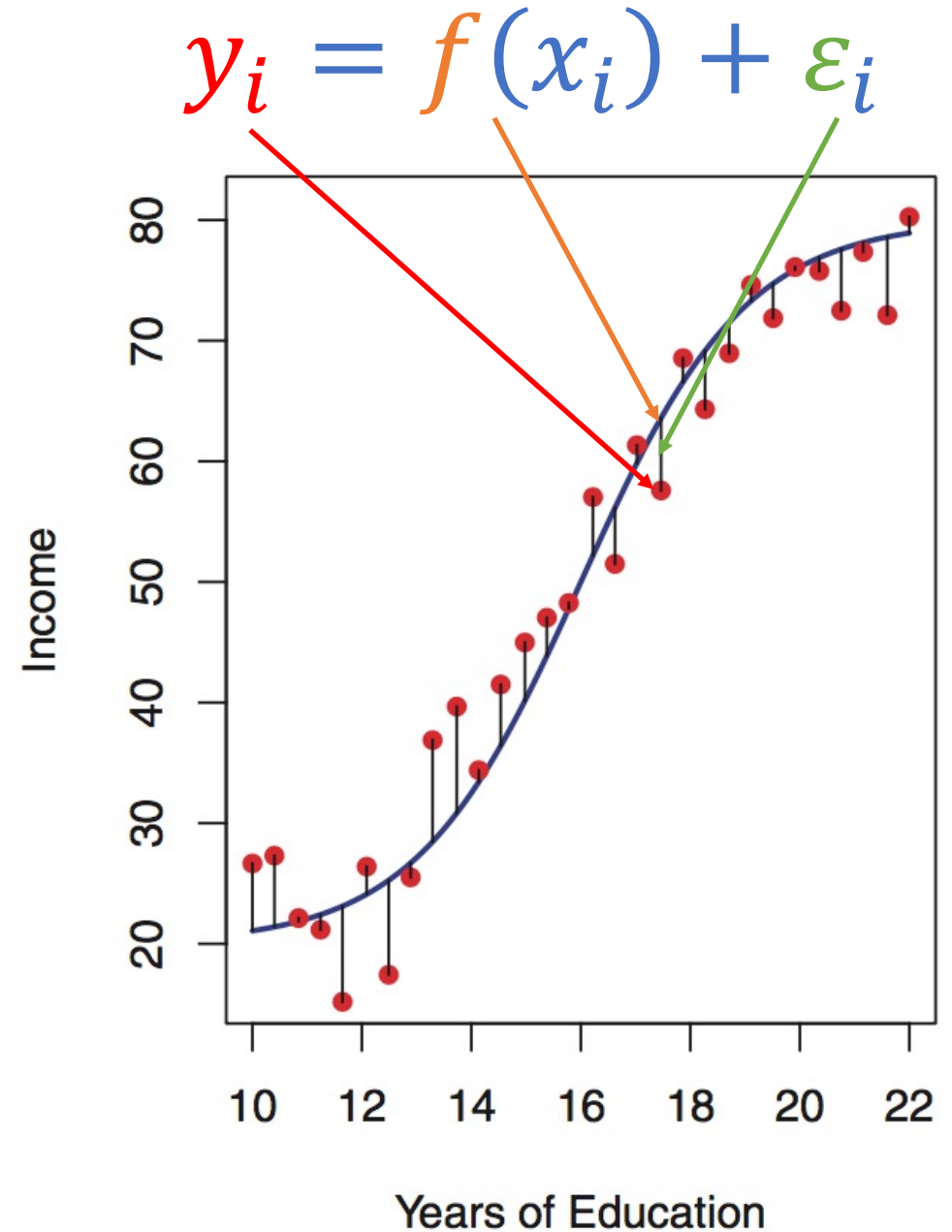


(c)

$$\frac{y = w_0 + w_1 x + w_2 x^2 + w_3 x^8}{f}$$

Why Does it Happen?

- Data is inherently **noisy**




Thus Use Separate Test Data

- Thus, if possible, we should try to use the test data $T_e = \{x_i, y_i\}_{i=1}^M$


$$\text{MSE}_{T_e} = \text{Ave}_{i \in T_e} [y_i - \hat{f}(x_i)]^2$$

Putting It All Together (3)

Training Phase

$$Tr = \{(\mathbf{x}_i, y_i) \in X \times Y : 1 \leq i \leq N\}$$


“Supervised”
Learning
Algorithm

$$\hat{f}: X \rightarrow Y$$


Evaluation Phase

$$Te = \{(\mathbf{x}_i, y_i) \in X \times Y : 1 \leq i \leq n\}$$


Learned
Model

Important Take-away (1)

Goal of Machine Learning

Learn how do features relate to labels?

$$f: x \rightarrow y$$

$$D = \{(x_i, y_i)\}_{i=1}^N$$

Features

Labels

For example: risk score of the patient

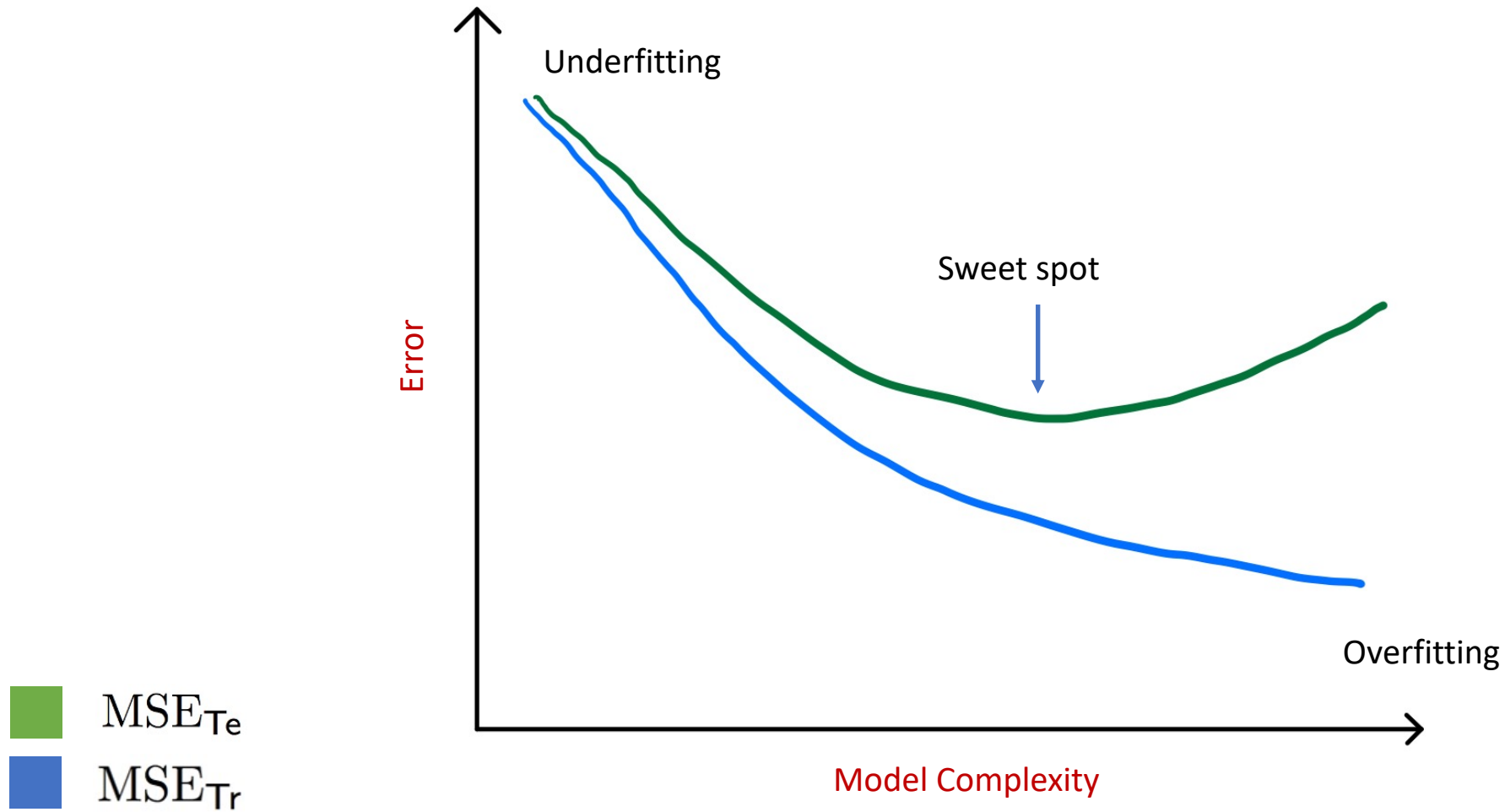
$$y \in \mathbb{R} \quad \text{Regression}$$

For example: COVID + OR COVID -

$$y \in \{c_i\}_{i=1}^k \quad \text{Classification}$$

For example: clinical Data of a patient

Important Take-away (2)



Bias Variance Tradeoff

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon).$$

Where (x_0, y_0) is a test observation

Typically, as the **flexibility or complexity** of \hat{f} increases, its variance increases, and its bias decreases. So choosing the flexibility based on average test error amounts to a **bias-variance trade-off**.

Summary

- Course Introduction
- Overview of Machine Learning
 - Three components: Task, Experience, Performance
 - Predictors and Response
 - Goal of Learning
 - Classification and Regression
 - Parametric Models
 - Model Evaluation
 - Bias-Variance Tradeoff