

Machine Learning in Business

John C. Hull

Chapter 1

Introduction



What is Machine Learning

- ✚ Machine learning is a branch of AI
- ✚ The idea underlying machine learning is that we give a computer program access to lots of data and let it learn about relationships between variables and make predictions
- ✚ Some of the techniques of machine learning date back to the 1950s but improvements in computer speeds and data storage costs have now made machine learning a practical tool



Software

- ⊕ There are several alternatives such as Python, R, MatLab, Spark, and Julia
- ⊕ Need ability to handle very large data sets and availability of packages that implement the algorithms.
- ⊕ Python seems to be winning at the moment
- ⊕ Scikit-Learn has freely available packages for many ML tasks



Traditional statistics

- ⊕ Means, SDs
- ⊕ Probability distributions
- ⊕ Significance tests
- ⊕ Confidence intervals
- ⊕ Linear regression
- ⊕ etc



The new world of statistics

- ⊕ Huge data sets
- ⊕ Fantastic improvements in computer processing speeds and data storage costs
- ⊕ Machine learning tools are now feasible
- ⊕ Can now develop non-linear prediction models, find patterns in data in ways that were not possible before, and develop multi-stage decision strategies
- ⊕ New terminology: features, labels, activation functions, target, bias, supervised/unsupervised learning.....



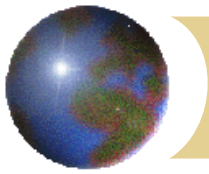
Types of Machine Learning

- ⊕ Unsupervised learning (find patterns)
- ⊕ Supervised learning (predict numerical value or classification)
- ⊕ Semi-supervised learning (only part of data has values for, or classification of, target)
- ⊕ Reinforcement learning (multi-stage decision making)



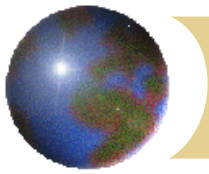
Applications of ML

- ✚ Credit decisions
- ✚ Classifying and understanding customers better
- ✚ Portfolio management
- ✚ Private equity
- ✚ Language translation
- ✚ Voice recognition
- ✚ Biometrics
- ✚ etc

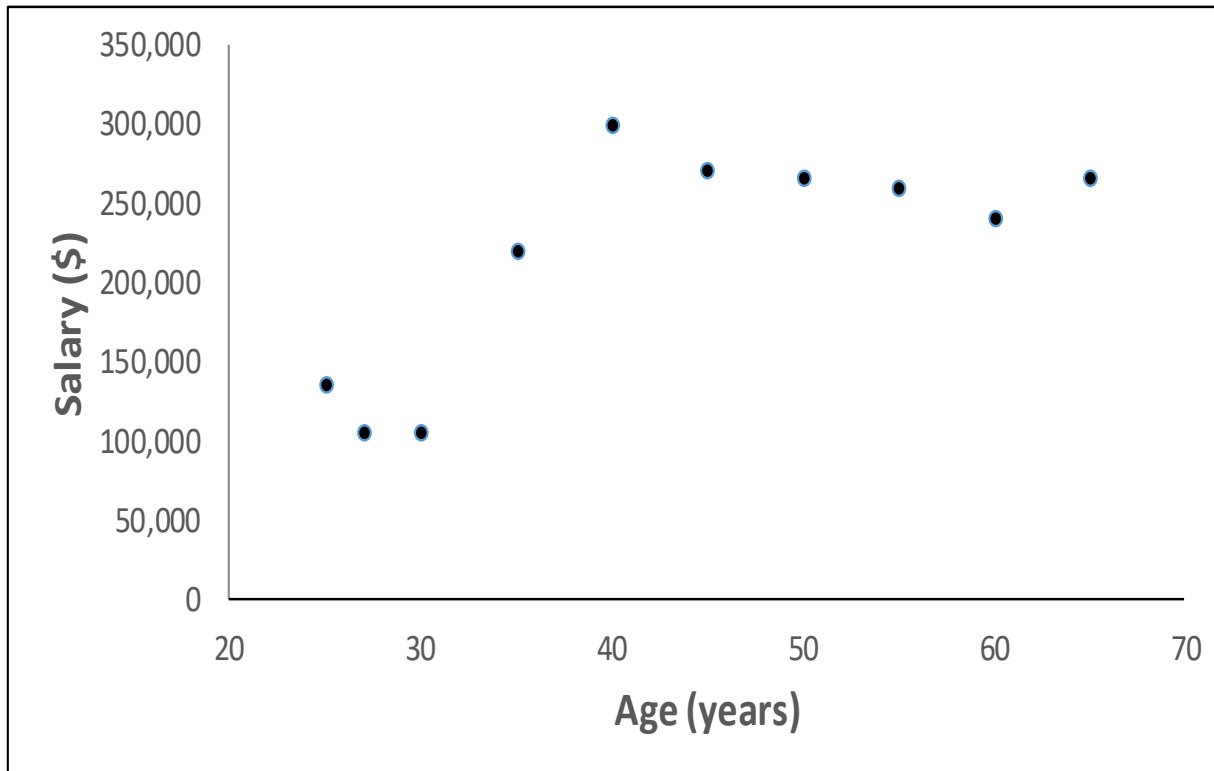


A Baby Data Training Set (Salary as a function of age for a certain profession in a certain area) Table 1.1

Age (years)	Salary (\$)
25	135,000
55	260,000
27	105,000
35	220,000
60	240,000
65	265,000
45	270,000
40	300,000
50	265,000
30	105,000



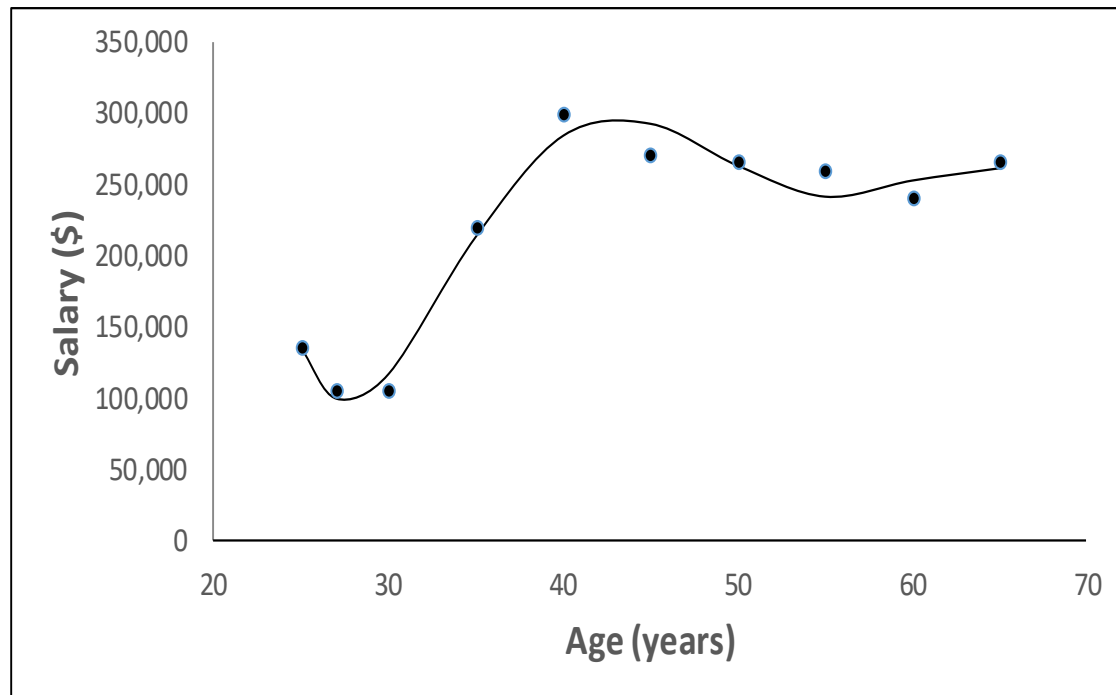
Scatter plot (Figure 1.1)





A Good Fit, Figure 1.2 ($Y = \text{Salary}$, $X = \text{Age}$)

$$Y = a + b_1X + b_2X^2 + b_3X^3 + b_4X^4 + b_5X^5$$



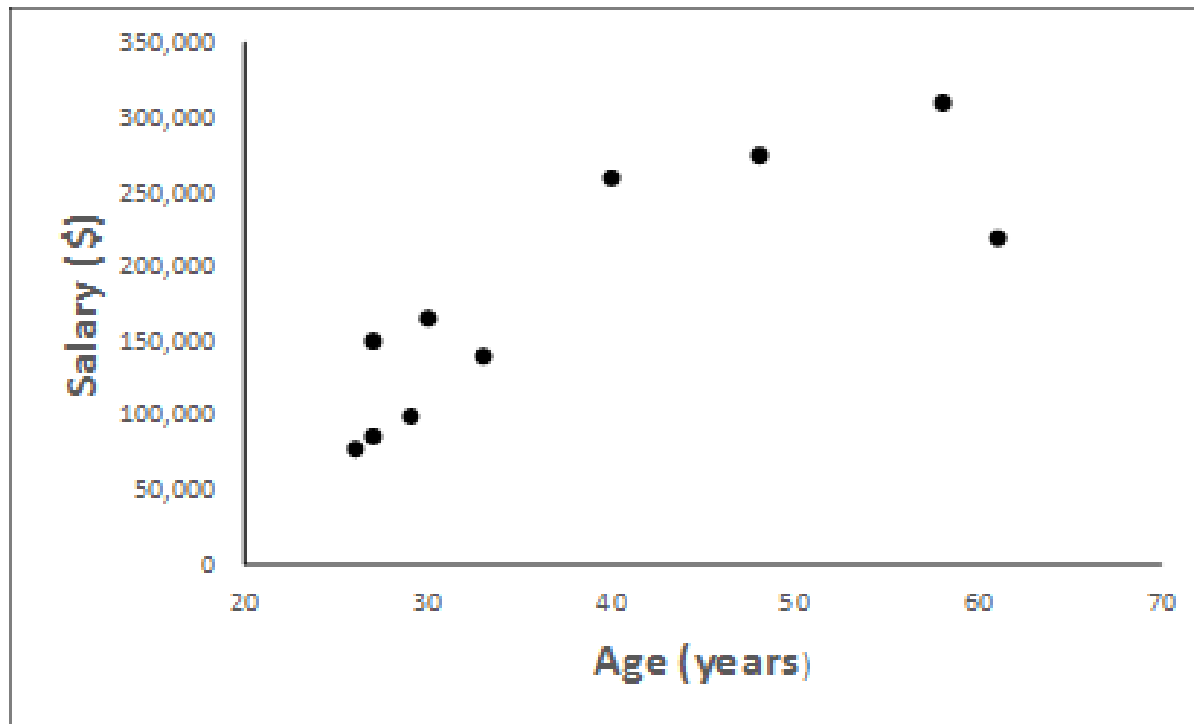


An Out-of-Sample Validation Set (Table 1.2)

Age (years)	Salary (\$)
30	166,000
26	78,000
58	310,000
29	100,000
40	260,000
27	150,000
33	140,000
61	220,000
27	86,000
48	276,000



Scatter Plot for Validation Set (Figure 1.3)





The Fifth Order Polynomial Model Does Not Generalize Well

- ✚ The root mean squared error (rmse) for the training data set is \$12,902
- ✚ The rmse for the test data set is \$38,794
- ✚ We conclude that the model overfits the data



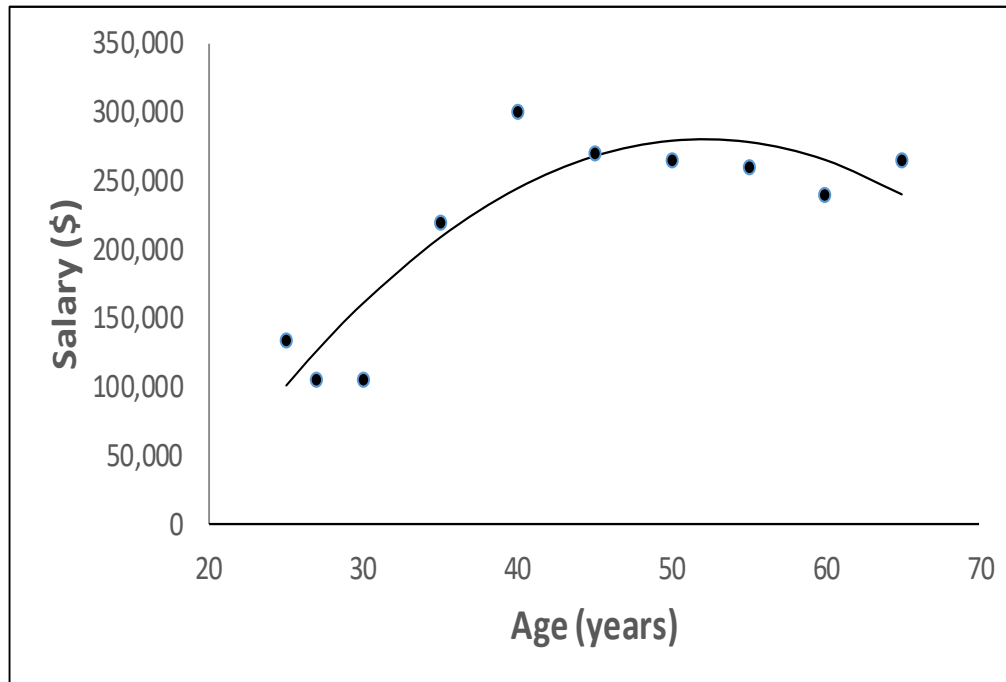
ML Good Practice

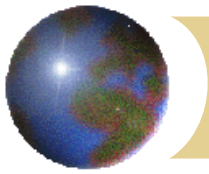
- ⊕ Divide data into three sets
 - ⊠ Training set
 - ⊠ Validation set
 - ⊠ Test set
- ⊕ Develop different models using the training set and compare them using the validation set
- ⊕ Rule of thumb: increase model complexity until model no longer generalizes well to the validation set
- ⊕ The test set is used to provide a final out-of-sample indication of how well the chosen model works



Quadratic Model for Baby Data Set (Figure 1.4)

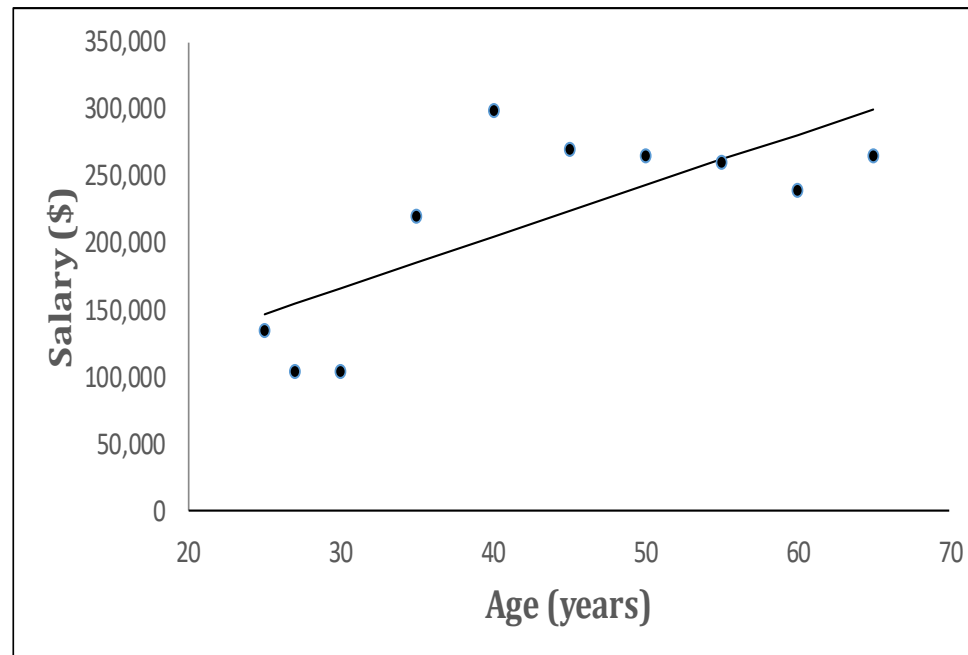
✚ $Y = a + b_1X + b_2X^2$

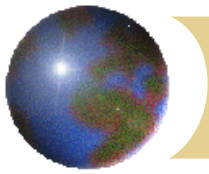




Linear Model for Baby Data Set (Figure 1.5)

$$Y = a + b_1X$$





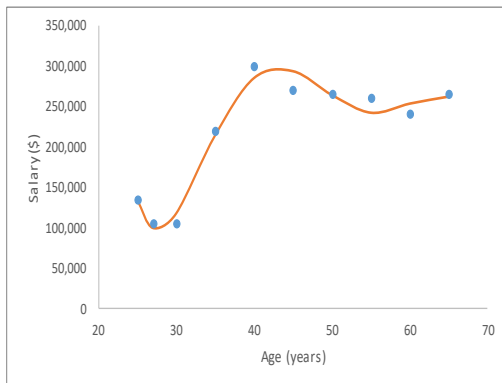
Summary of Results: The linear model under-fits while the 5th degree polynomial over-fits (Table 1.3)

	Polynomial of degree 5	Quadratic model	Linear model
Training set	12, 902	32,932	49,731
Validation set	38,794	33,554	49,990

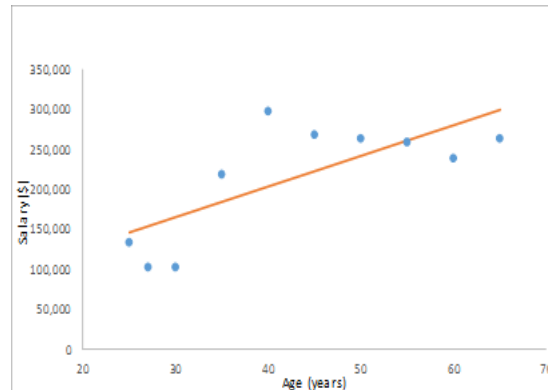


Overfitting/Underfitting;

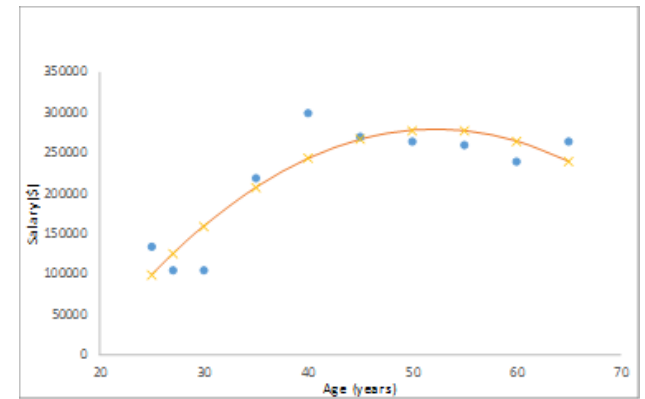
Example: predicting salaries for people in a certain profession in a certain area (only 10 observations)



Overfitting



Underfitting



Best model?



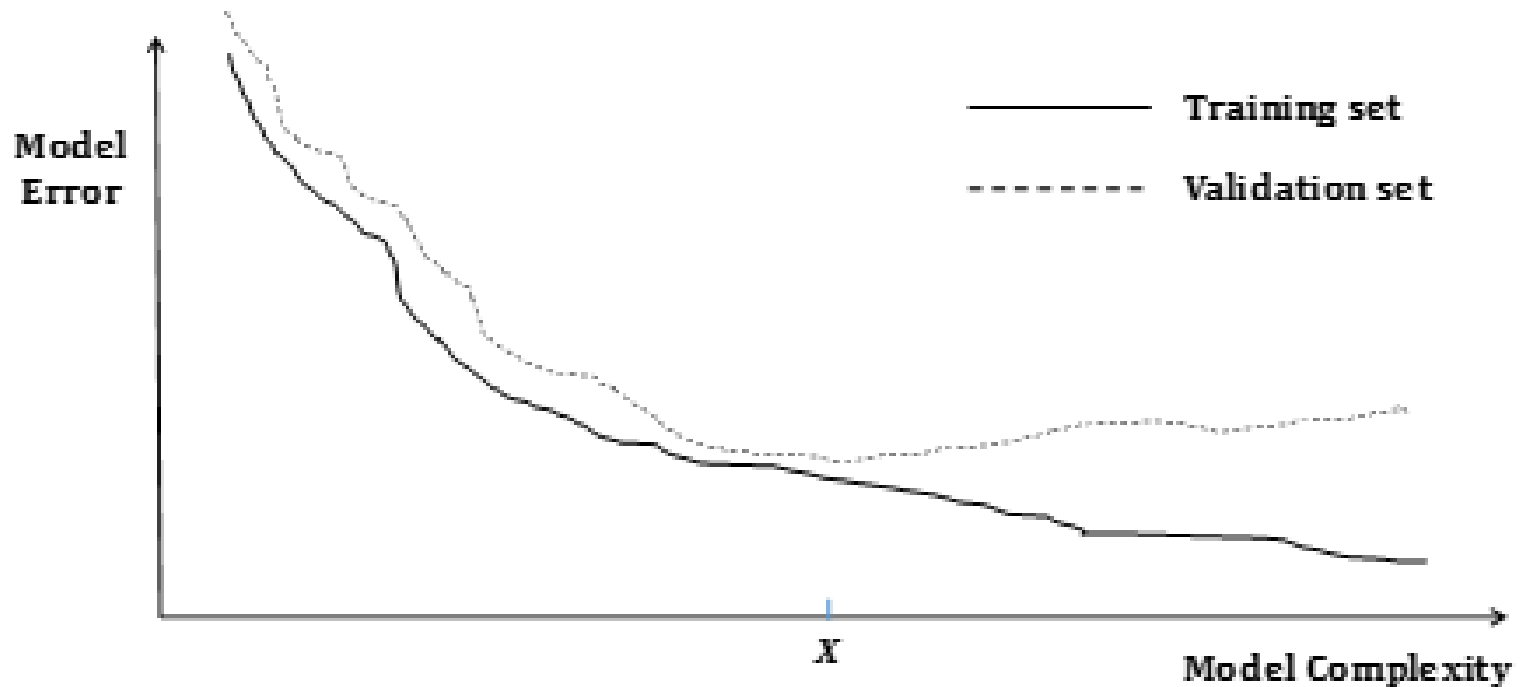
Test Set Results for Quadratic Model

Age (years)	Salary (\$)	Predicted salary (\$)	Error (\$)
26	110,000	113,172	-3,172
52	278,000	279,589	-1,589
38	314,000	232,852	+83,148
60	302,000	264,620	+37,380
64	261,000	245,457	+15,543
41	227,000	249,325	-22,325
34	200,000	199,411	+589
46	233,000	270,380	-37,380
57	311,000	273,883	-37,117
55	298,000	277,625	+20,375

SD of error is \$34,273



Typical Pattern of Errors for Training Set and Validation Set





Cleaning data (page 14-16)

- ✚ Dealing with inconsistent recording
- ✚ Removing unwanted observations
- ✚ Removing duplicates
- ✚ Investigating outliers
- ✚ Dealing with missing items

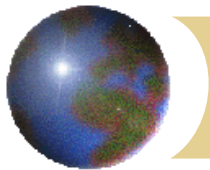


Bayes Theorem (useful when we want an uncertainty estimate as well as just a prediction)

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

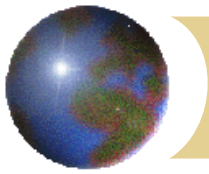
Example: We observe that 90% of fraudulent transactions are for large amounts late in the day. Also 3% of transactions are for large amounts late in the day and 1% of transactions are fraudulent

$$P(\text{fraud}|\text{large\&late}) = \frac{P(\text{large\&late}|\text{fraud})P(\text{fraud})}{P(\text{large\&late})} = \frac{0.9 \times 0.01}{0.03} = 0.3$$



Bayes can be counterintuitive

- ✚ One person in ten thousand has a certain disease
- ✚ A test is 99% accurate (i.e., if person has the disease the test gets this right 99% of the time; similarly when the person does not have the disease the test is right 99% of the time)
- ✚ You test positive
- ✚ What is the chance that you have the disease?
- ✚ X =test positive, Y =has disease, \bar{Y} = does not have disease
- ✚ $P(X|Y) = 0.99$; $P(Y) = 0.0001$
- ✚ $P(X) = P(X|Y)P(Y) + P(X|\bar{Y})P(\bar{Y}) = 0.99 \times 0.0001 + 0.01 \times 0.9999 = 0.0101$
- ✚ $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \frac{0.99 \times 0.0001}{0.0101} = 0.0098$



The Terminology

- ⊕ Features
- ⊕ Target
- ⊕ Labels
- ⊕ Supervised learning
- ⊕ Unsupervised learning
- ⊕ Semi-supervised learning
- ⊕ Reinforcement learning
- ⊕ And more to come..