

# 應用機器學習

Brian Chan 陳醒凡

# 課程目標

1. 了解基本的數據分析
2. 了解基本的機器學習(Machine Learning)方法
3. 掌握Python的基本操作和一些有用的package
4. 處理及從網上下載數據
5. 在Python上應用機器學習

# 今天課堂 概要

## 1. Optimization algorithm in ML

- Optimization framework
- Application: Least square method

## 2. Statistical foundation of machine learning methods

- Probability distribution
- Estimation framework
- Basics of statistics
- Bayesian probability

# MACHINE LEARNING

- I. 掌握各種ML 方法
- II. 了解各種ML方法的原理
- III. 認識不同ML方法的優缺點



# 1. Optimization algorithm in ML

- Optimization framework
- Application: Least square method

I. 認識Machine Learning 背後操作及原理

II. 了解Optimization跟regression model的關係

# OPTIMIZATION ALGORITHM (BASIC)

$$\min_{\theta} J(\theta)$$

such that  $\theta \in \Theta$

1. Objective function,  $J(\theta)$
2. Control variables,  $\theta$
3. Constraints,  $\theta \in \Theta$

Optimization is a very useful issue in ML.

Many ML methods reply optimization algorithms to solve the parameters in the model.

The example includes regression methods and the class of NN models.

# EXAMPLES

Constrained optimization

$$\min_x x^2 + 2x + 1$$

$$\text{s.t. } x \in [1, 2]$$

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

Portfolio optimization

$$\min_x x^T \Sigma x$$

$$\text{s.t. } x^T \mathbb{1} = 1$$

$$x \geq 0$$

$$x \in \mathbb{R}^n, \Sigma \in S_+^n$$

Argmax (or argmin) are the points, or elements, of the domain of some function at which the function values are maximized (or minimized).

# GRADIENT DESCENT

Gradient descent is one of the most basic but commonly used optimization algorithms. Many advanced optimization algorithm is developed based on gradient descent.

Under some conditions, the optimization can be solved by gradient descent.

$$\theta \leftarrow \theta - \eta \cdot \frac{\partial J(\theta)}{\partial \theta}$$

[https://en.wikipedia.org/wiki/Gradient\\_descent](https://en.wikipedia.org/wiki/Gradient_descent)

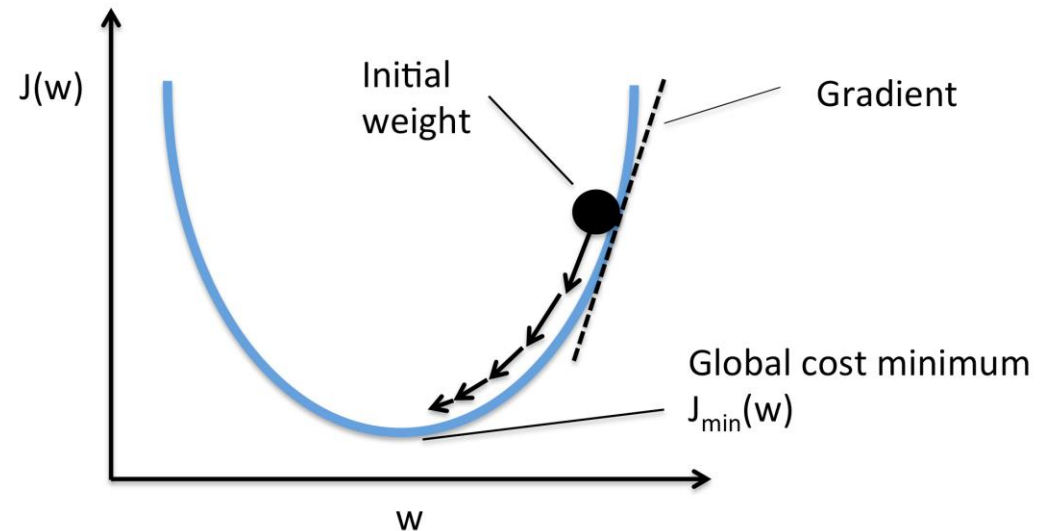


# GRADIENT DESCENT

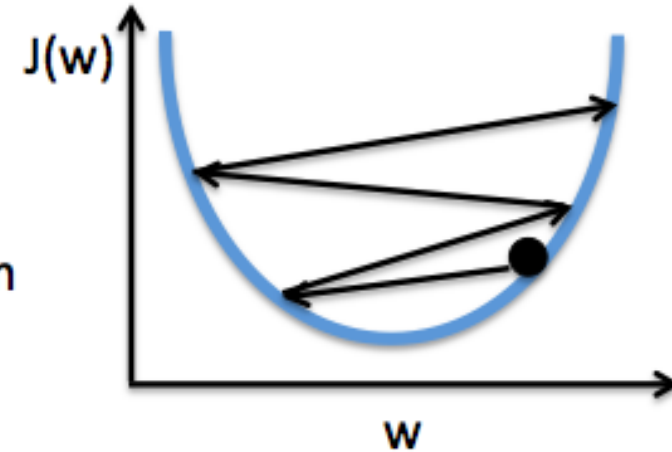
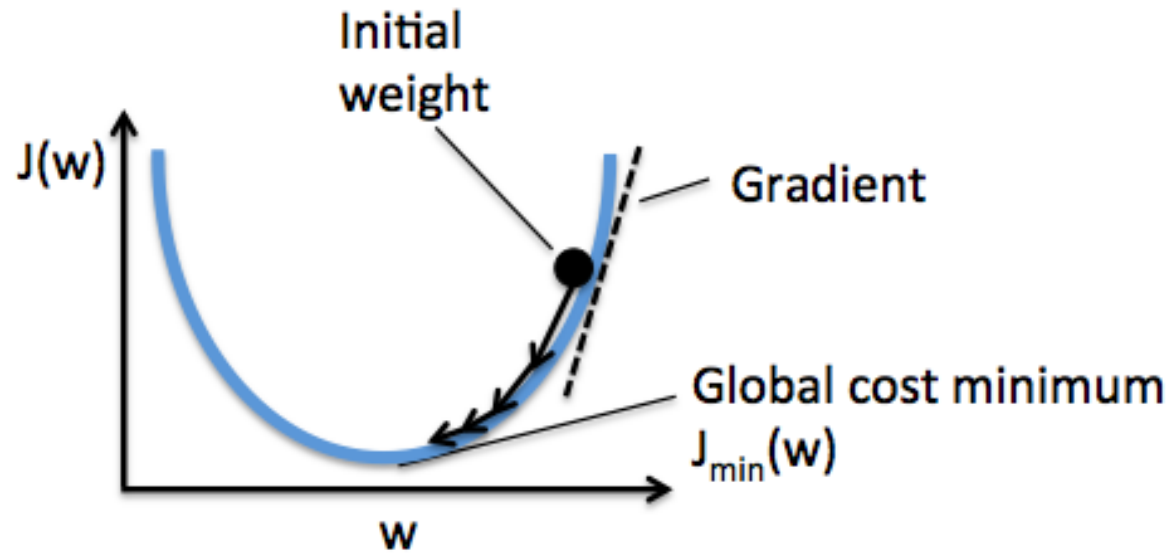
At initial step  $i = 1$ , the control variable is set to be  $\theta^1 = \bar{\theta}$ .

At each iteration,  $\theta^i$  is updated with  $\theta^i \leftarrow \theta^{i-1} - \eta \cdot \frac{\partial J(\theta)}{\partial \theta} \big|_{\theta=\theta^{i-1}}$ .

Repeat this until  $|\theta^i - \theta^{i-1}| < \delta$ .



# GRADIENT DESCENT



# SOME MORE OPTIMIZATION ALGORITHM IN ML

**Stochastic gradient descent**

**Momentum**

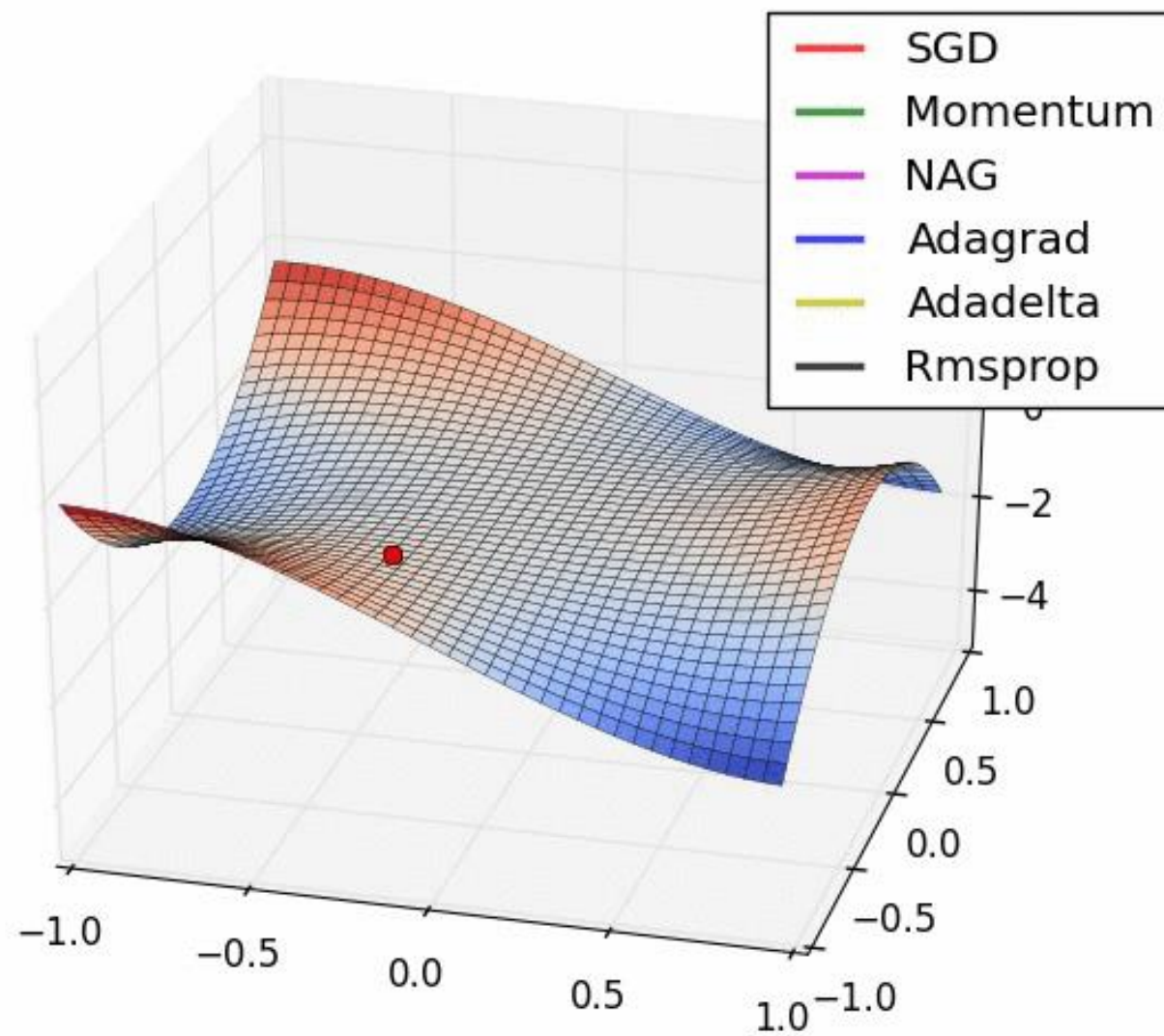
**Adagrad**

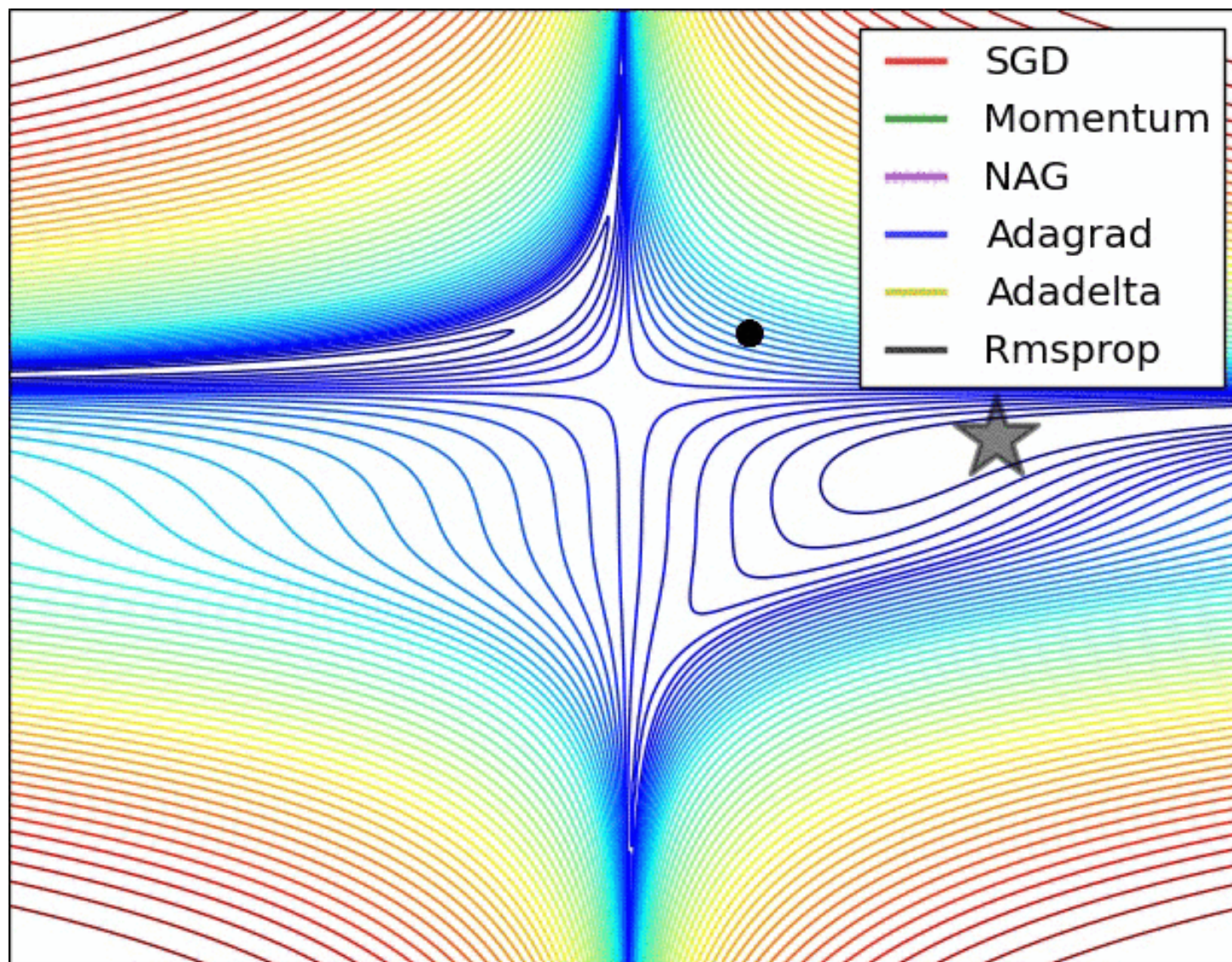
**RMSProp**

**Adam**

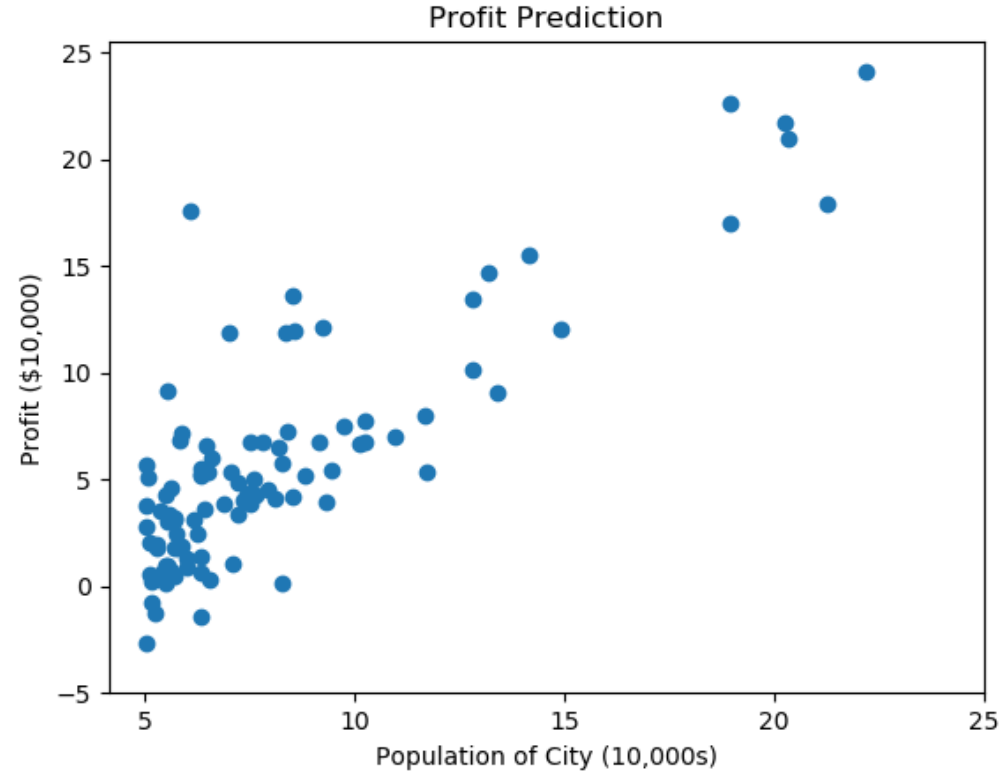
and some more ...

<http://ruder.io/optimizing-gradient-descent/index.html#adam>





# EXAMPLE: LEAST SQUARE METHOD





# LEAST SQUARE METHOD

Let  $y_i$  be observed data for  $i = 1, \dots, N$ .

$$J(\theta_0, \theta_1) = \frac{1}{N} \sum_{i=1}^N [y_i - (\theta_0 + \theta_1 x_i)]^2$$

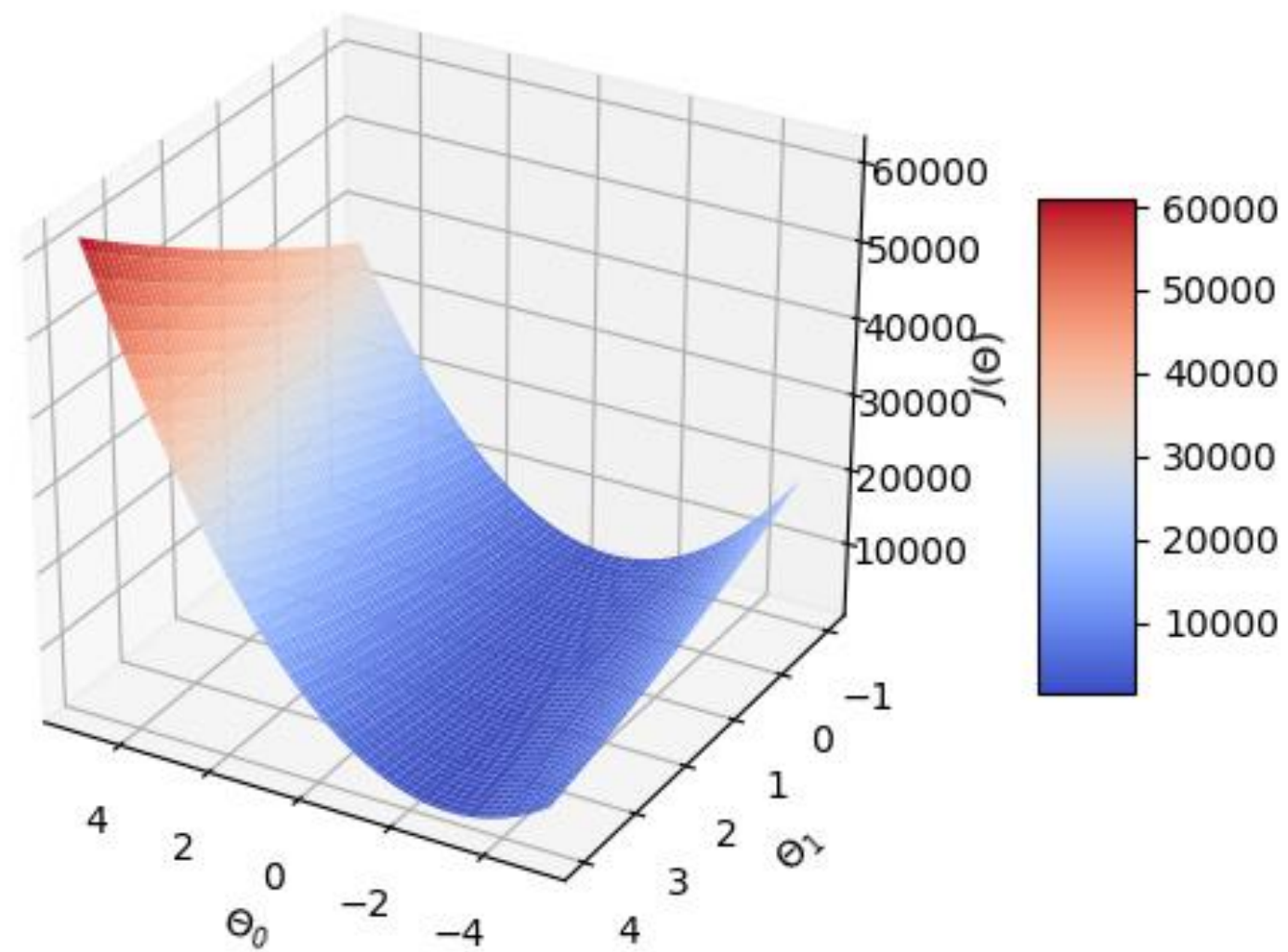


# LEAST SQUARE METHOD (COMPACT FORM)

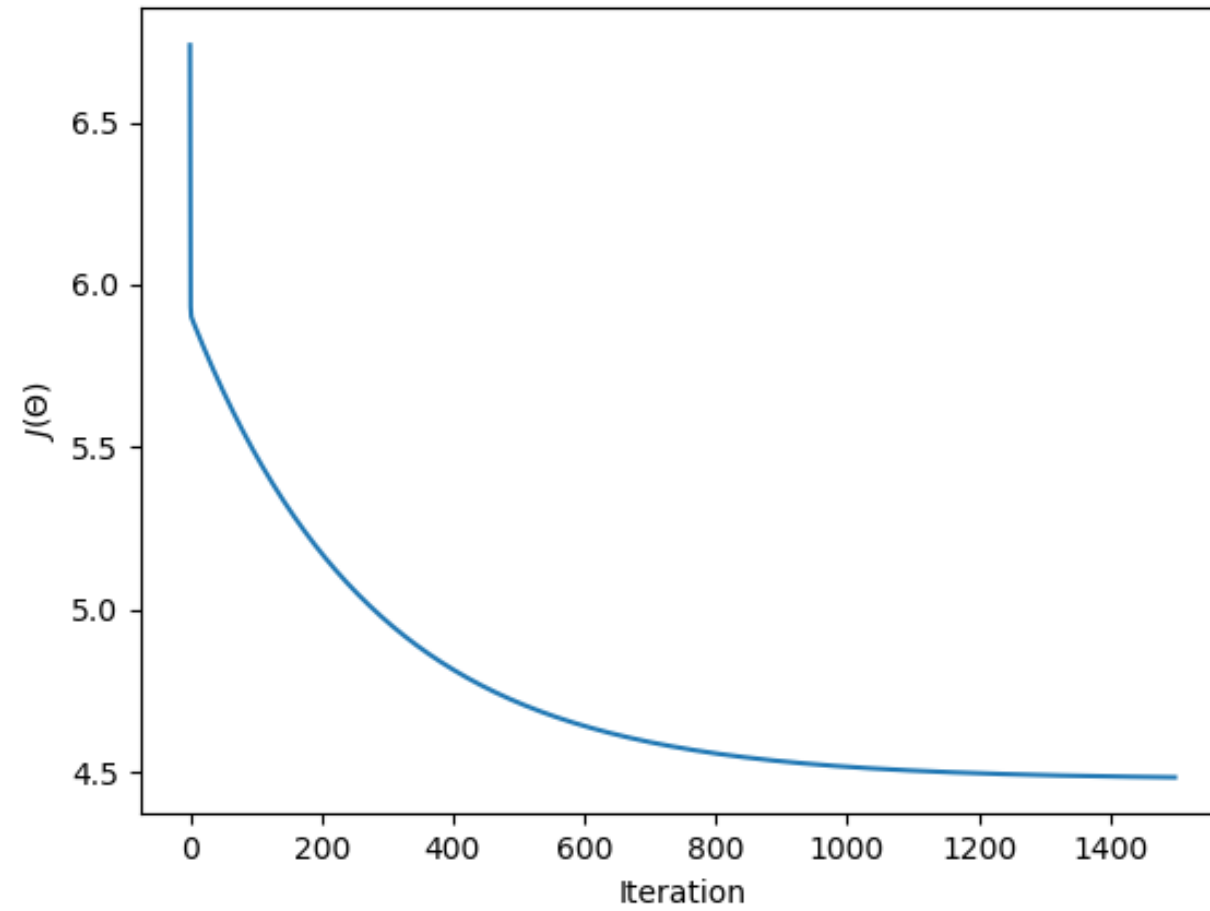
Let  $y$  be observed data such that  $y = (y_1, y_2, \dots, y_N)^T$  and  $Z = (1, X)^T$ .

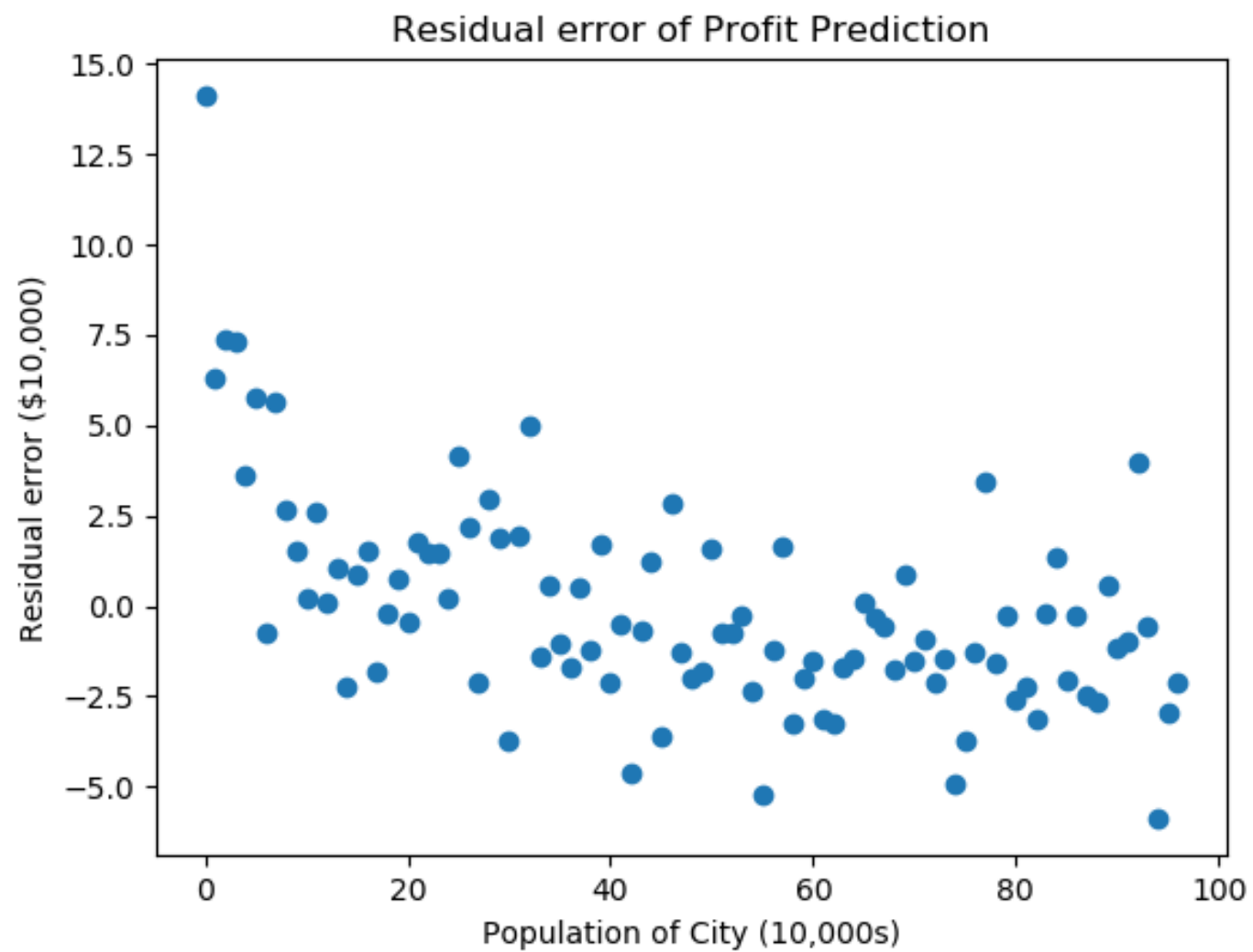
$$J(\theta) = \frac{1}{N} ||y - Z \cdot \theta||^2,$$

where  $\theta = (\theta_1, \theta_2)^T$ .



Cost function using Gradient Descent







## 2. Statistical foundation of machine learning methods

- Probability distribution
- Estimation framework
- Basics of statistics
- Bayesian probability

# PROBABILITY (DISCRETE)

- Event
- Probability measure

## Discrete probability distribution:

$$P(\{\text{Head}\}) = P(\{\text{Tail}\}) = 0.5$$

$$P(\{1\}) = P(\{2\}) = P(\{3\}) = P(\{4\}) = P(\{5\}) = P(\{6\}) = 1/6$$

$$P(\{1\}, \{2\}) = P(\{1\}) + P(\{2\}) = 1/6 + 1/6 = 1/3$$



# CREDIT MIGRATION

	Moody's	S&P	Fitch	Meaning
Investment Grade	Aaa	AAA	AAA	Prime
	Aa1	AA+	AA+	High Grade
	Aa2	AA	AA	
	Aa3	AA-	AA-	
	A1	A+	A+	Upper Medium Grade
	A2	A	A	
	A3	A-	A-	
	Baa1	BBB+	BBB+	Lower Medium Grade
	Baa2	BBB	BBB	
	Baa3	BBB-	BBB-	
Junk	Ba1	BB+	BB+	Non Investment Grade Speculative
	Ba2	BB	BB	
	Ba3	BB-	BB-	
	B1	B+	B+	Highly Speculative
	B2	B	B	
	B3	B-	B-	
	Caa1	CCC+	CCC+	Substantial Risks
	Caa2	CCC	CCC	Extremely Speculative
	Caa3	CCC-	CCC-	In Default w/ Little Prospect for Recovery
	Ca	CC	CC+	
		C	CC	
			CC-	In Default
	D	D	DDD	

Table 1.8

One-year transition matrix (%)

Initial rating	Rating at year-end (%)							
	AAA	AA	A	BBB	BB	B	CCC	Default
AAA	90.81	8.33	0.68	0.06	0.12	0	0	0
AA	0.70	90.65	7.79	0.64	0.06	0.14	0.02	0
A	0.09	2.27	91.05	5.52	0.74	0.26	0.01	0.06
BBB	0.02	0.33	5.95	86.93	5.30	1.17	0.12	0.18
BB	0.03	0.14	0.67	7.73	80.53	8.84	1.00	1.06
B	0	0.11	0.24	0.43	6.48	83.46	4.07	5.20
CCC	0.22	0	0.22	1.30	2.38	11.24	64.86	19.79

Source: Standard & Poor's CreditWeek (15 April 96)

## Markov Chain

[https://www.moodys.com/sites/products/ProductAttachments/DRD/CTM\\_Methodology.pdf](https://www.moodys.com/sites/products/ProductAttachments/DRD/CTM_Methodology.pdf)

# MOMENTS

Expectation:

$$E[X] = \sum P(x_i)x_i$$

Variance:

$$Var[X] = \sum P(x_i)(x_i - \mu)^2$$

Covariance:

$$Cov[X, Y] = E[(X - \mu_x)(Y - \mu_y)]$$

Correlation:

$$Corr[X, Y] = \frac{Cov[X, Y]}{\sigma_x \sigma_y}$$

Expectation:

$$E[X] = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 4 + \frac{1}{6} \cdot 5 + \frac{1}{6} \cdot 6 = 3.5$$

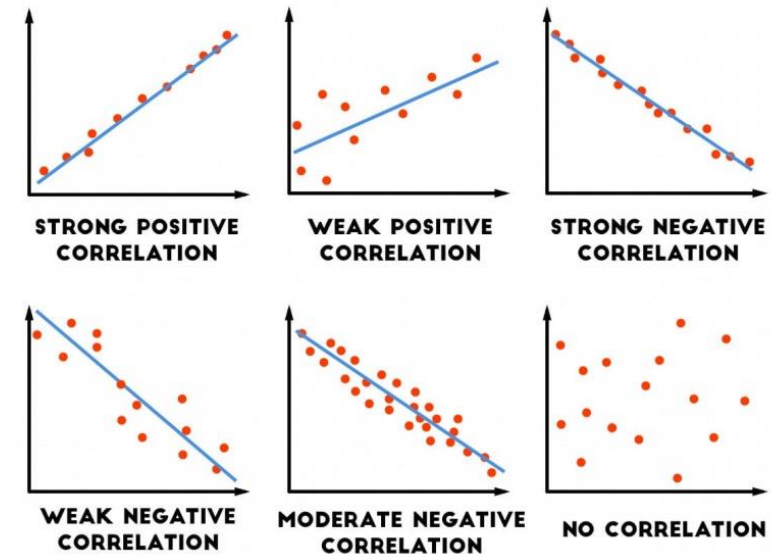
Variance:

$$Var[X] = \frac{1}{6} \cdot (1 - 3.5)^2 + \frac{1}{6} \cdot (2 - 3.5)^2 + \frac{1}{6} \cdot (3 - 3.5)^2 + \frac{1}{6} \cdot (4 - 3.5)^2 + \frac{1}{6} \cdot (5 - 3.5)^2 + \frac{1}{6} \cdot (6 - 3.5)^2$$



# COVARIANCE & CORRELATION

Dice1\Dice2	0	1
0	1/4	1/4
1	1/4	1/4



Covariance:

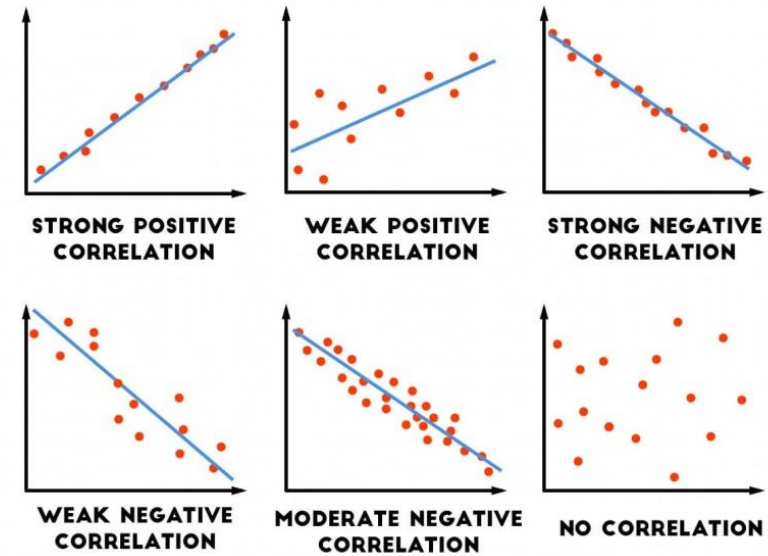
$$\text{Cov}[X, Y] = E[(X - \mu_x)(Y - \mu_y)] = \frac{1}{4} \cdot (0 - 0.5)(0 - 0.5) + \frac{1}{4} \cdot (1 - 0.5)(0 - 0.5) + \frac{1}{4} \cdot (0 - 0.5)(1 - 0.5) + \frac{1}{4} \cdot (1 - 0.5)(1 - 0.5) = 0$$

Correlation:

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sigma_x \sigma_y}$$

# COVARIANCE & CORRELATION

Dice1 \ Dice2	0	1
0	2/4	0
1	0	2/4

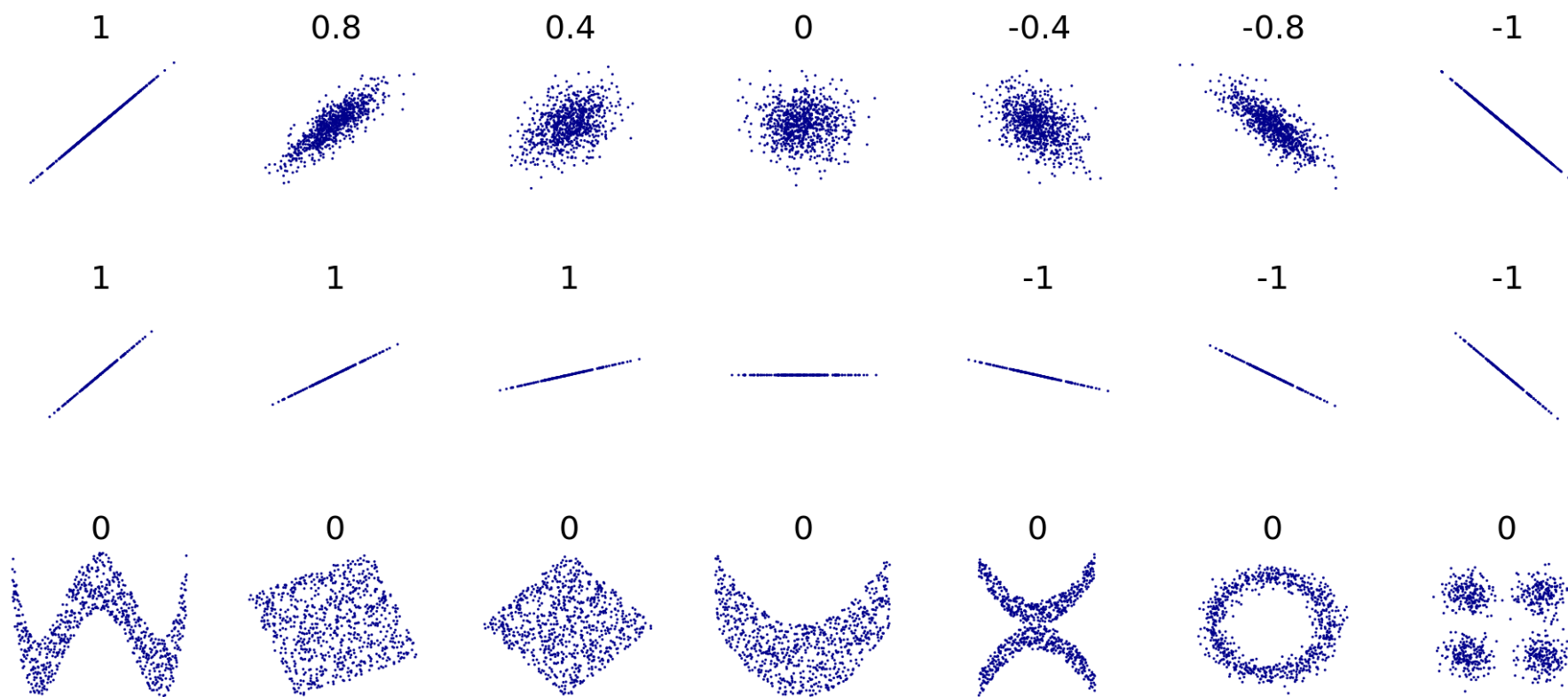


Covariance:

$$\text{Cov}[X, Y] = E[(X - \mu_x)(Y - \mu_y)] = \frac{2}{4} \cdot (0 - 0.5)(0 - 0.5) + \frac{2}{4} \cdot (1 - 0.5)(1 - 0.5) = 0.25$$

Correlation:

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sigma_x \sigma_y}$$



# REMINDER ABOUT CORRELATION

Relationship with regression

Interpretation of correlation:

When correlation = 1

When correlation = 0

Input attributes: prefer low correlation attributes (almost always)

# PROBABILITY (CONTINUOUS)

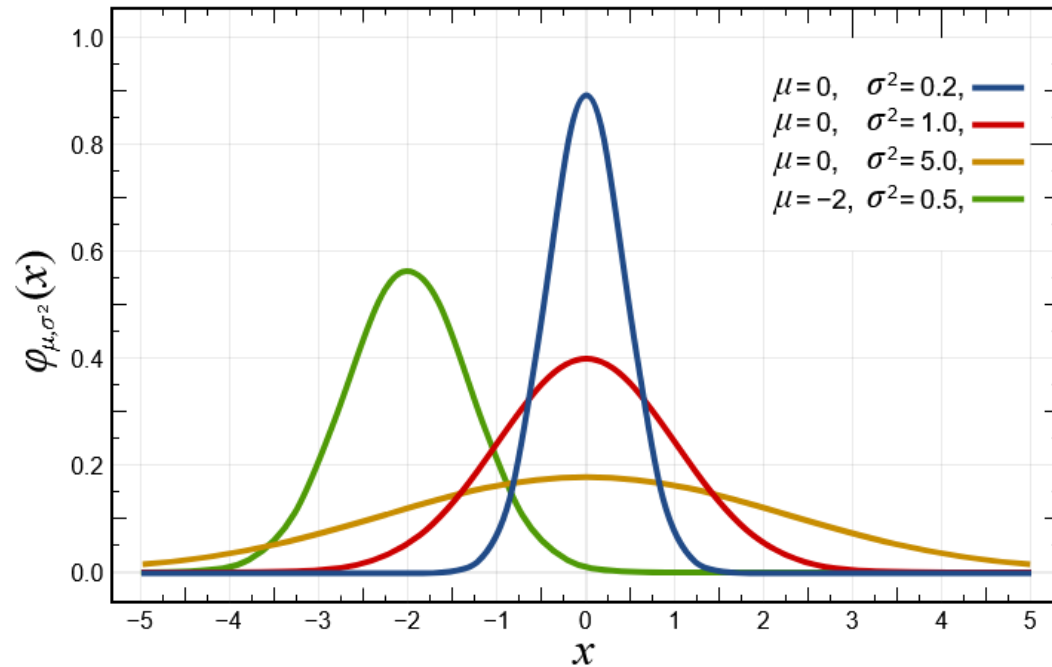
**Continuous probability distribution:**

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

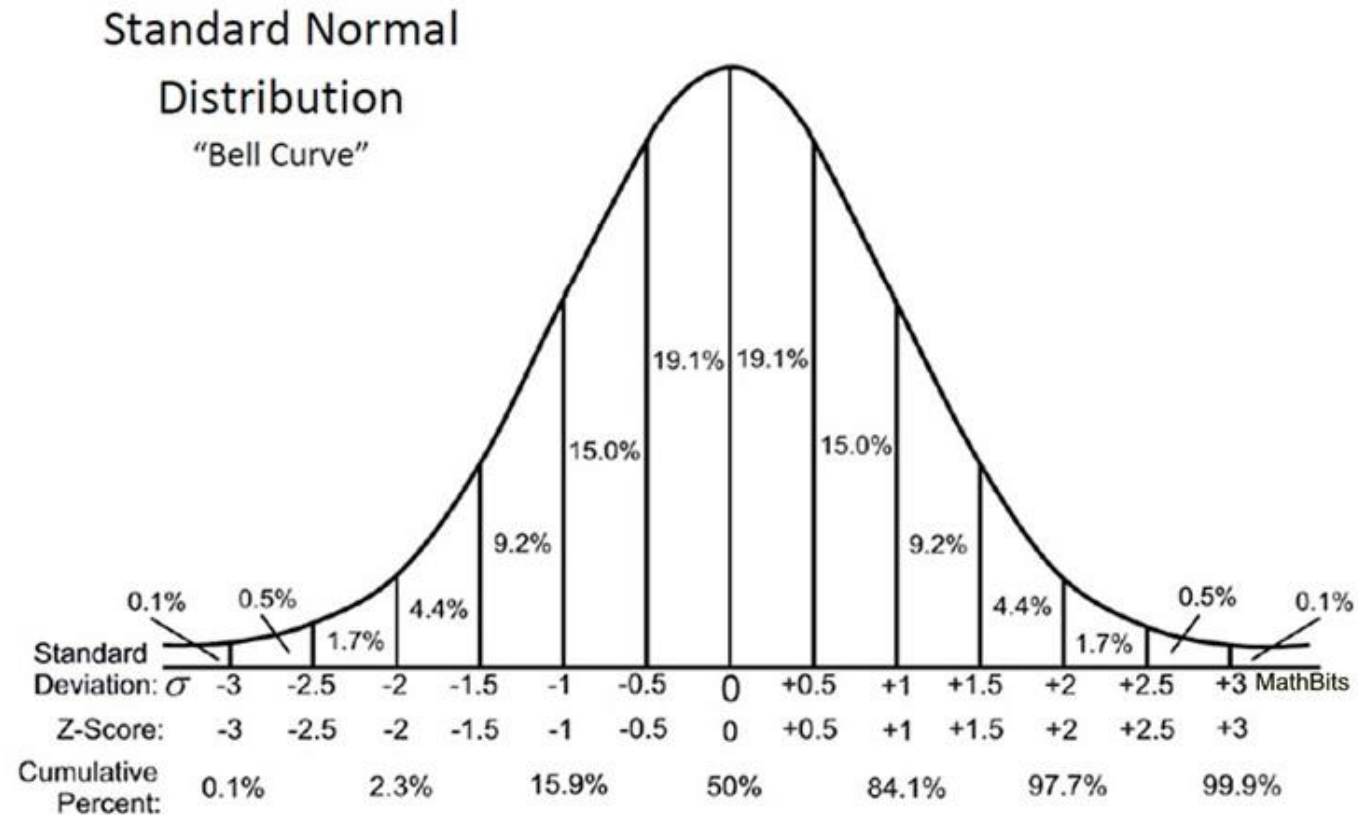
*Conditions:*

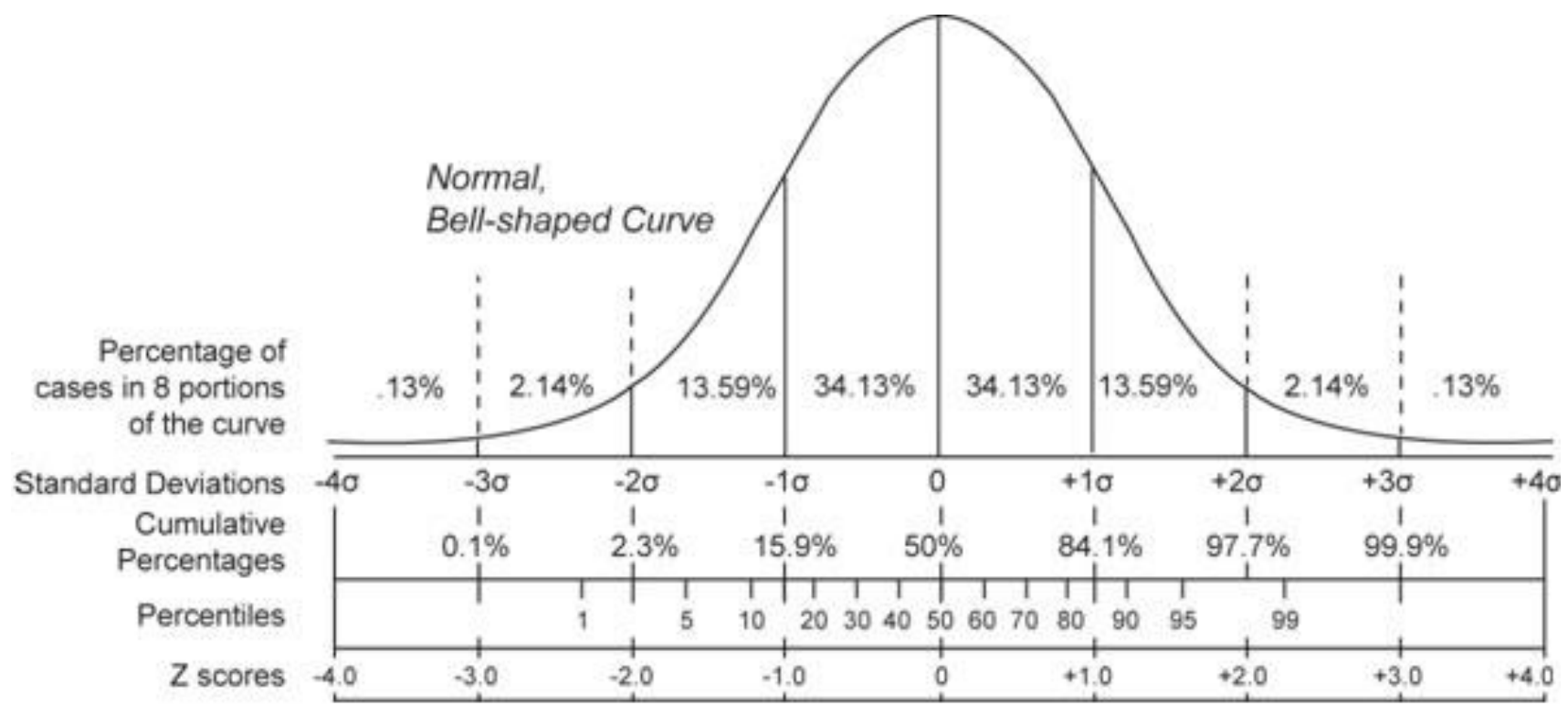
1.  $\int f(x) dx = 1$

2.  $f(x) \geq 0$

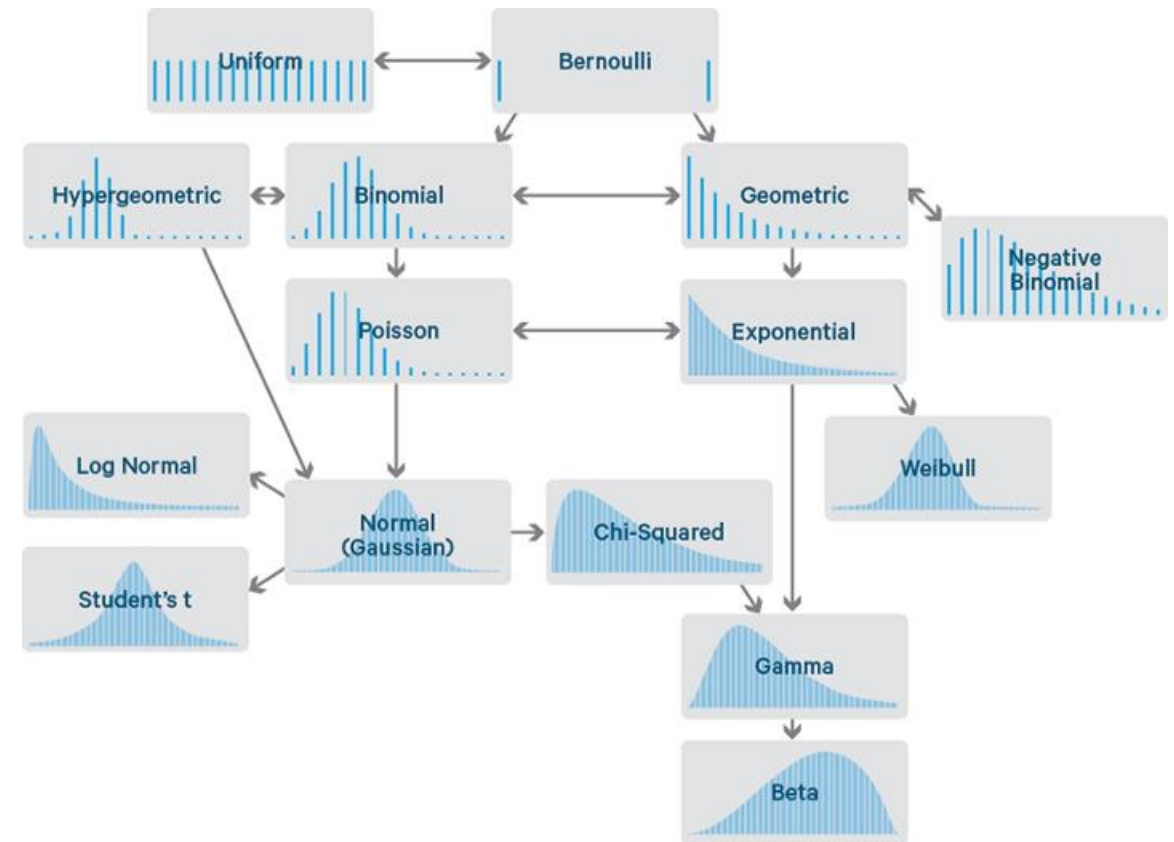


# GAUSSIAN DISTRIBUTION





# FAMILY OF PROBABILITY DENSITY





# ESTIMATION FRAMEWORK

## Parameter $\theta$

A statistical parameter is a quantity characterize probability distribution or statistical model.

## Estimator $\hat{\theta}$

An estimator is a quantity as an estimate of a given quantity based on observed data.

Unbiased estimator  $\theta = E[\hat{\theta}]$

Consistent estimator  $\hat{\theta} \xrightarrow{p} \theta, \text{ as } n \rightarrow \infty$

# EXAMPLE

Sample Mean

$$\hat{\mu} = \frac{\sum x_i}{n}$$

One can show

$$E[\hat{\mu}] = \mu$$
$$var[\hat{\mu}] = \frac{\sigma_x^2}{\sqrt{n}}$$

$$\hat{\mu} \xrightarrow{p} \mu, \text{ as } n \rightarrow \infty$$

One can see that estimator of sample average  $\hat{\mu}$  is an unbiased and consistent estimator.

# SUMMARY

1. Framework of optimization
2. Probability distribution
3. Bayesian probability
4. Least square method

# 下一課...

迴歸分析及例子:

1. 線性迴歸基礎
2. 線性迴歸的統計特質
3. 非線性迴歸 (Nonlinear regression)
4. 正規化迴歸 (Regularized regression)