

Chen Tang's Knowledge Database

`ChenTang@link.cuhk.edu.cn`

Last updated on September 17, 2023

Preface

The following is a compendium of my academic notes spanning various domains. I present these notes publicly to share my methodological framework for managing and structuring an individual's knowledge networks.

The inevitability of encountering occasional errors is acknowledged.

This notebook will undergo continuous updates.

Contents

Preface	1
1 Mathematics and Optimization	5
1.1 Calculus and Linear Algebra	5
1.1.1 Keys of Linear Algebra	5
1.2 Analysis and Algebra	8
1.3 Probability Theory	9
1.3.1 Basics of Probability	9
1.4 Stochastic Process	13
1.5 Linear Programming	14
1.5.1 Simplex Method	19
1.6 Convex Optimization	20
1.7 Non-Convex Optimization	21
2 Statistics and Econometrics	22
2.1 Statistical Inference	22
2.1.1 Exponential Families	22
2.1.2 Scale and Location	28
2.1.3 Data Reduction	29
2.2 Causal Model	31
2.2.1 Rubin's Causal Model	31
2.3 Reduced-Form Identification	32
2.3.1 Randomized Trials	32
2.3.2 Regression and Matching	33
2.3.3 Asymptotic Analysis	37
2.4 Advanced Econometrics	38
2.5 Machine Learning Interface	39
3 Economics	40
3.1 Microeconomics	40

3.2	Macroeconomics	41
3.3	Game Theory	42
3.4	Development Economics	43
3.4.1	Models of Development Economics	43
3.4.2	Clan Culture	43
3.5	Data Science Interface	46
4	Computer Science and Data Science	47
4.1	Cloud Computing	47
4.1.1	Principles of Cloud Computing System	47
4.1.2	Virtual Machines	52
4.2	Machine Learning	55
4.3	Deep Learning	56
4.4	Reinforcement Learning	57
4.4.1	Introduction and MDP	57
5	Information Systems and Operations Management	62
5.1	Empirical Operations Management	62
5.1.1	12 Papers of [Roth, 2007]	63
5.2	Revenue Management	66
5.2.1	Traditional RM	66
5.2.2	Overbooking	68
5.2.3	Traditional Consumer Choice Model	68
5.2.4	Current Consumer Choice Model	68
5.3	Platform Operations Management	69
5.3.1	Platform Owner's Entry	69
5.3.2	Consumer Polarization	70
5.3.3	Network Effect	71
5.3.4	Online Gaming	72
5.4	Behavioral Operations Management	75
5.5	Data-Driven Operations Management	76
6	Miscellaneous	77
6.1	Notes on Tools	77
6.1.1	LaTeX Shortcuts	77
6.2	Important Proofs	81
6.3	Beautiful Phrases	82
6.4	Eureka Ideas	83

6.4.1	RNN and Causal Model	83
6.4.2	Earthquakes on Immigration and Investment	83

Chapter 1

Mathematics and Optimization

1.1 Calculus and Linear Algebra

1.1.1 Keys of Linear Algebra

The properties of Matrix Multiplication:

- Associativity: $(AB)C = A(BC)$;
- Distributivity: $A(C + D) = AC + AD$;
- Identity Multiplication: $I_m A = A, A I_n = A$

Only square matrix has an inverse matrix, and the inverse is unique. If a matrix doesn't have an inverse, then it's called **regular/invertible/nonsingular**. The properties of inverses and transposes:

- $(AB)^{-1} = B^{-1}A^{-1}$;
- $(A + B)^{-1} \neq A^{-1} + B^{-1}$;
- $(AB)^T = B^T A^T$;
- $(A + B)^T = A^T + B^T$.

To find the inverse matrix, gaussian elimination can be applied: $[A \mid I] = [I \mid A^{-1}]$.

Solutions of Linear Systems

Reduced Row-Echelon Form Matrix:

$$A = \begin{bmatrix} 1 & 3 & 0 & 0 & 3 \\ 0 & 0 & 1 & 0 & 9 \\ 0 & 0 & 0 & 1 & -4 \end{bmatrix} \quad (1.1.1.1)$$

There are *pivot (basic variables)* and *free variable*, and the column of free variable is dependent on pivots. The steps to find solutions of linear systems:

1. Find a particular solution for $Ax = \mathbf{b}$ by setting all free variable zero;
2. Find all solutions for $Ax = 0$;
3. Add them up.

There is iterative method to solve large-scale linear equations: define error $= \|x^{(k+1)} - x_*\|$, then optimize the function $x^{(k+1)} = Cx^{(k)} + d$ and iterate it.

Vector Spaces, Basis and Rank

Definition 1.1.1.1: Vector Space and Subspace

A **Vector Space** $V = (\mathcal{V}, +, \cdot)$ is a set \mathcal{V} with two operations:

1. $+: \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$;
2. $\cdot: \mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V}$

For $\mathcal{U} \subseteq \mathcal{V}$ and $\mathcal{U} \neq \emptyset$, then $U = (\mathcal{U}, +, \cdot)$ is a vector subspace.

Remark 1.1.1

- $V = (\mathcal{V}, +)$ is an Abelian group;
- The subspace needs to satisfy *closure*, $\lambda x \in U, x + y \in U$.

If $v = \sum_{i=1}^k \lambda_i x_i$, then v is a linear combination of (x_1, x_2, \dots, x_k) . If **not** all values of a solution are 0, then it's called a non-trivial solution. For $\sum_{i=1}^k \lambda_i x_i = 0$, if the non-trivial solution exists, then the vectors are called **linearly dependent**.

Definition 1.1.1.2: Basis and Rank

- **Generating Set**: if all vectors in V can be expressed as a linear combination of $\mathcal{A} = \{x_1, \dots, x_k\} \subseteq \mathcal{V}$, then \mathcal{A} is a generating set of V ;
- **Span**: The set of linear combinations of \mathcal{A} is its span;
- **Basis**: The minimal generating set (linearly independent) of a vector space V is called its basis;
- **Rank**: the number of linearly independent column vectors in a matrix $\mathcal{A} \in \mathbb{R}^m \times n$.

Remark 1.1.2

- $\text{rk}(A) = \text{rk}(A^\top)$;
- A matrix $A \in \mathbb{R}^{n \times n}$ is invertible if and only if $\text{rk}(A) = n$;
- The span of a matrix is also called its **image**, $\dim(U) = \text{rk}(A)$;
- $Ax = b$ only has solution if and only if $\text{rk}(A) = \text{rk}(A \mid b)$;
- The solution (kernel, null space) to $Ax = 0$ has a dimension of $n - \text{rk}(A)$;

Definition 1.1.1.3: Linear Mappings

For vector spaces V, W , a mapping $\Phi : V \rightarrow W$ is called linear mapping if:

$$\forall x, y \in V \forall \lambda, \psi \in \mathbb{R} : \Phi(\lambda x + \psi y) = \lambda \Phi(x) + \psi \Phi(y) \quad (1.1.1.2)$$

1.2 Analysis and Algebra

1.3 Probability Theory

1.3.1 Basics of Probability

Common Discrete Distribution

Bernoulli(p)

pmf: $P(X = x | p) = p^x(1 - p)^{1-x}; \quad x = 0, 1; \quad 0 \leq p \leq 1$

mean: $EX = p$; *variance*: $VarX = p(1 - p)$

- A Bernoulli trial (named after James Bernoulli) is an experiment with only two possible outcomes;
- Bernoulli random variable $X = 1$ if “success” occurs and $X = 0$ if “failure” occurs where the probability of a “success” is p .

Binomial(n, p)

pmf: $P(X = x | n, p) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, \dots, n; \quad 0 \leq p \leq 1$

mean: $EX = np$; *variance*: $np(1 - p)$

- A Binomial experiment consists of n independent identical Bernoulli trials;
- $X = \sum_{i=1}^n Y_i$, where Y_1, \dots, Y_n are n identical, independent Bernoulli random variables.

Poisson(λ)

pmf: $P(X = x | \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}; \quad x = 0, 1, \dots; \quad 0 \leq \lambda < \infty$

mean: $EX = \lambda$; *variance*: $VarX = \lambda$

- A Poisson distribution is typically used to model the probability distribution of the number of occurrences (with λ being the intensity rate) per unit time or per unit area;
- Binomial pmf approximates Poisson pmf. Poisson pmf is also a limiting distribution of a negative binomial distribution;
- A useful result: By Taylor series expansion: $e^\lambda = \sum_{x=0}^{\infty} \frac{\lambda^x}{x!}$.

Assumption 1.3.1.1: Poisson Process

1. Let $X(\Delta)$ be the number of events that occur during an interval Δ ;
2. The events are independent: if $\Delta_1, \dots, \Delta_n$ are disjoint intervals, then $X(\Delta_1), \dots, X(\Delta_n)$ are independent;

3. $X(\Delta)$ only depends on the length of Δ ;
4. The probability that exactly one event occurs in a small interval of length Δt equals $\lambda \Delta t + o(\Delta t)$;
5. Poisson distribution is **not** memoryless, but its interval (exponential distribution) is memoryless.

Geometric(p)

pmf: $P(X = x | p) = p(1 - p)^{x-1}; \quad x = 1, 2, \dots; \quad 0 \leq p \leq 1$

mean: $\frac{1}{p}$; variance: $\frac{1-p}{p^2}$

- The experiment consists of a sequence of independent trials;
- 🌲 The property of memoryless: $P(X > s | X.t) = P(X > s - t)$.

Negative Binomial(r, p)

pmf: $P(X = x | r, p) = \binom{r+x-1}{x} p^r (1-p)^x, \quad x = 0, 1, \dots; \quad 0 \leq p \leq 1$

mean: $EX = \frac{r(1-p)}{p}$; variance: $\frac{r(1-p)}{p^2}$

- assume there are many independent and identical experiments, to observe the r th success, X is the number of games to see the failure;

Hypergeometric(N, M, K)

pmf: $P(X = x | N, M, K) = \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}}; \quad x = 0, 1, 2, \dots, K;$

$M - (N - K) \leq x \leq M; \quad N, M, K \geq 0$

mean: $EX = \frac{KM}{N}$; variance: $\frac{KM}{N} \frac{(N-M)(N-K)}{N(N-1)}$

Common Continuous Distribution

Uniform(a, b)

pdf: $f(x | a, b) = \frac{1}{b-a};$ mean: $EX = \frac{b+a}{2};$ variance: $Var X = \frac{(b-a)^2}{12}.$

Exponential(β)

pdf: $f(x | \beta) = \frac{1}{\beta} e^{-x/\beta}, 0 \leq x < \infty, \beta > 0;$ mean: $EX = \beta;$ variance: $Var X = \beta^2.$

Gamma(α, β)

pdf: $f(x | \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}, \quad 0 \leq x < \infty, \quad \alpha, \beta > 0;$ mean: $\alpha\beta;$ variance: $\alpha\beta^2.$

- The *gamma function* is defined as $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt;$

- $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha), \alpha > 0;$
- $\Gamma(n) = (n - 1)!, \quad \text{for any integer } n > 0.$

Normal (μ, σ^2)

pdf: $f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/(2\sigma^2)}, \quad -\infty < x < \infty; \text{ mean: } \mu; \text{ variance: } \sigma^2.$

Example 1.3.1.1

for $f(x) = \frac{\beta\alpha^\beta}{x^{\beta+1}}, \quad \alpha < x < \infty, \quad \alpha > 0, \quad \beta > 0:$

- (a) Verify $f(x)$ is a pdf;
- (b) Derive the mean and variance of this distribution;
- (c) Prove that the variance does not exist if $\beta \leq 2$.

Solution:

(a)

$$\begin{aligned} \int_{\alpha}^{\infty} f(x) dx &= \int_{\alpha}^{\infty} \frac{\beta \cdot \alpha^\beta}{x^{\beta+1}} dx = -x^{\beta \cdot \alpha^\beta |_{\alpha}} \\ &= 0 + \alpha^{-\beta} \cdot \alpha^\beta = 1 \end{aligned} \quad (1.3.1.1)$$

(b)

$$\begin{aligned} EX &= \int_{\alpha}^{\infty} x f(x) dx = \int_{\alpha}^{\infty} \frac{\beta \cdot \alpha^\beta}{X^\beta} dx \\ &= \frac{\beta}{-\beta + 1} \cdot \alpha^\beta \cdot x^{-\beta+1} \Big|_{\alpha}^{\infty} = \frac{\beta \cdot \alpha}{\beta - 1} \end{aligned} \quad (1.3.1.2)$$

$$\begin{aligned} EX^2 &= \int_{\alpha}^{\infty} \frac{\beta \cdot \alpha^\beta}{x^{\beta-1}} dx = \frac{\beta}{-\beta + 2} \cdot \alpha^\beta \cdot x^{-\beta+2} \Big|_{\alpha}^{\infty} \\ &= \frac{\alpha^2 \beta}{\beta - 2} \end{aligned} \quad (1.3.1.3)$$

$$\text{Var } X = EX^2 - (EX)^2 = \frac{\beta \alpha^2}{(\beta - 1)^2 (\beta - 2)} \quad (1.3.1.4)$$

(c)

If $\beta < 2$, then the variance is negative.

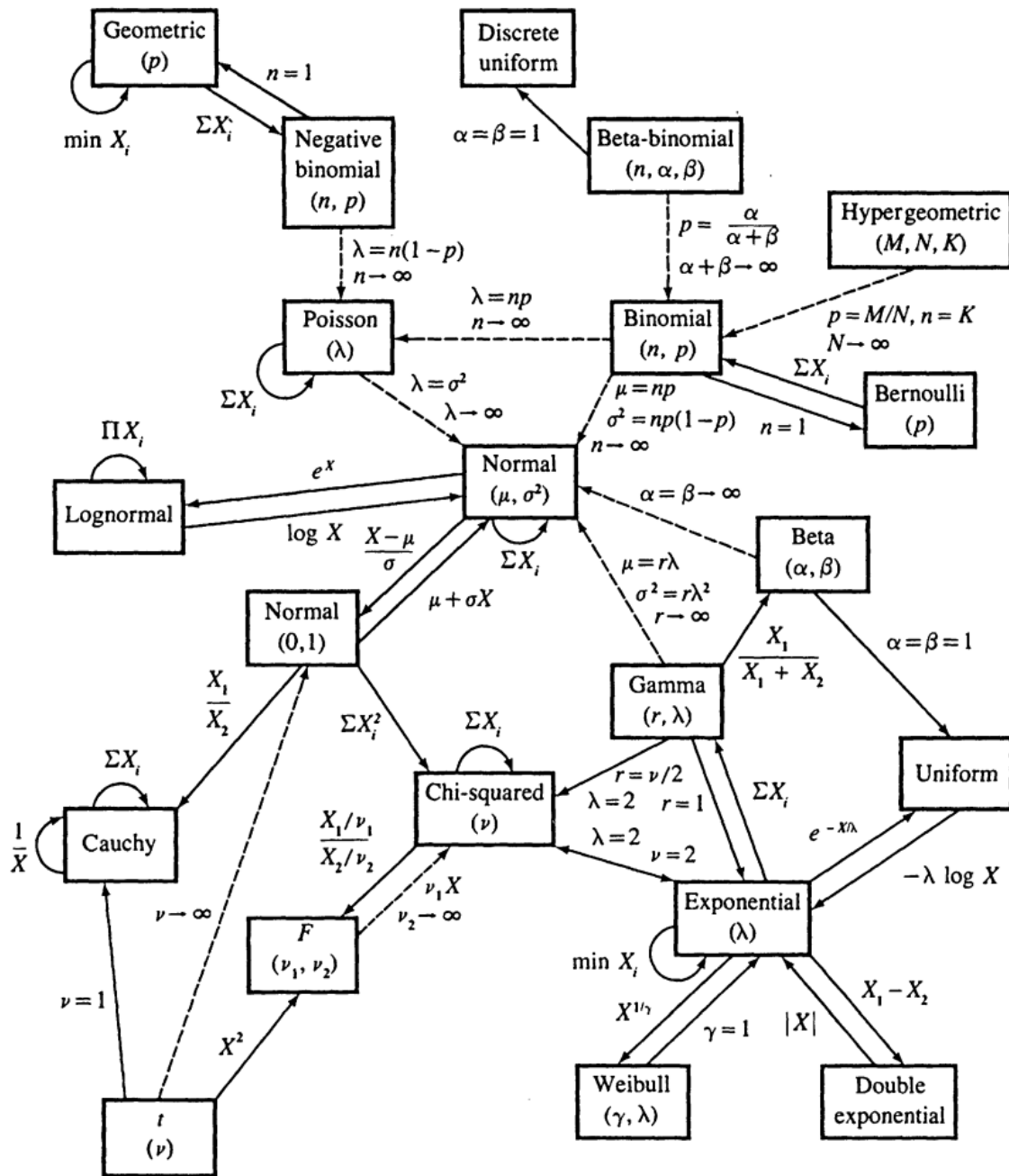


Figure 1.1: Type of distribution

1.4 Stochastic Process

1.5 Linear Programming

The main contents of this section are notes from [Luenberger et al., 1984], [Bertsimas and Tsitsiklis, 1997].

The *standard form* of the linear programming:

$$\begin{aligned} & \text{minimize} && c_1x_1 + c_2x_2 + \dots + c_nx_n \\ & \text{subject to} && a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ & && a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ & && \vdots \quad \vdots \\ & && a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \\ & \text{and} && x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0, \end{aligned} \tag{1.5.0.1}$$

or the above equations can be concisely written in:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^T \mathbf{x} \\ & \text{subject to} && \mathbf{Ax} = \mathbf{b} \quad \text{and} \quad \mathbf{x} \geq \mathbf{0}. \end{aligned} \tag{1.5.0.2}$$

Conversion to standard form LP

Slack Variable

For this kind of LP formation:

$$\begin{aligned} & \text{minimize} && c_1x_1 + c_2x_2 + \dots + c_nx_n \\ & \text{subject to} && a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \leq b_1 \\ & && a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \leq b_2 \\ & && \vdots \quad \vdots \\ & && a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \leq b_m \\ & \text{and} && x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0, \end{aligned} \tag{1.5.0.3}$$

The above formulation can be transformed into the following standard form:

$$\begin{aligned} & \text{minimize} && c_1x_1 + c_2x_2 + \dots + c_nx_n \\ & \text{subject to} && a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n + y_1 = b_1 \\ & && a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n + y_2 = b_2 \\ & && \vdots \quad \vdots \\ & && a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n + y_m = b_m \\ & \text{and} && x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0, \\ & \text{and} && y_1 \geq 0, y_2 \geq 0, \dots, y_m \geq 0. \end{aligned} \tag{1.5.0.4}$$

Now the constraint would be modified from \mathbf{A} to $[\mathbf{A}, \mathbf{I}]$, and the number of unknowns is changed from n to $n + m$.

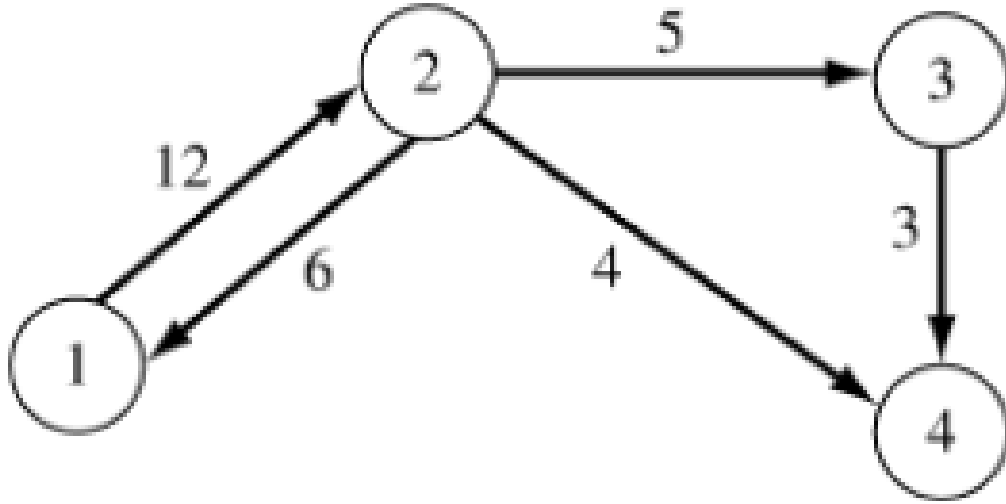


Figure 1.2: A network with capacities

Surplus Variable

Similar to the slack variable, formulation like:

$$a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n \geq b_i \quad (1.5.0.5)$$

can be transformed into:

$$a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n - y_i = b_i \quad (1.5.0.6)$$

Free Variable

For some variables without the constraint of the sign, there are two methods to transform it into the standard form. One is to $x_i = u_i - v_i$, where both u_i and v_i are larger or equal to zero. Another method can be used when the following condition holds:

$$a_1x_1 + \cdots + a_ix_i + \cdots + a_nx_n = b_i \quad (1.5.0.7)$$

where x_i is the free variable, thus we can replace x_i by $b_i - \sum_{j \neq i}^n a_j x_j$, which can eliminate one variable and one constraint at the same time.

One important example of the linear programming model is the **maximal flow problem**: This

problem can be formulated into:

$$\begin{aligned}
& \text{minimize} \\
& \text{subject to } \sum_{j=1}^n x_{1j} - \sum_{j=1}^n x_{j1} - f = 0 \\
& \sum_{j=1}^n x_{ij} - \sum_{j=1}^n x_{ji} = 0, \quad i \neq 1, m \\
& \sum_{j=1}^n x_{mj} - \sum_{j=1}^n x_{jm} + f = 0 \\
& 0 \leq x_{ij} \leq k_{ij}, \quad \text{for all } i, j
\end{aligned} \tag{1.5.0.8}$$

Basic Solutions

The system of linear equalities:

$$\mathbf{Ax} = \mathbf{b} \tag{1.5.0.9}$$

where \mathbf{A} is a $m \times n$ constraint matrix and \mathbf{x} is the $n \times 1$ decision variables.

Assumption 1.5.0.1: Full Rank Assumption

The $m \times n$ \mathbf{A} has $m < n$, and the m rows of \mathbf{A} are linearly independent

. This assumption makes the linear equalities always have at least one basic solution. At least m linearly independent columns $\rightarrow \mathbf{B}$, then would get:

$$\mathbf{Bx}_B = \mathbf{b} \tag{1.5.0.10}$$

Definition 1.5.0.1: Basic Feasible Solution

$\mathbf{x} = (\mathbf{x}_B, \mathbf{0})$ is the **basic solution**, the components of \mathbf{x} related to the columns of \mathbf{B} are **basic variables**

If the solution also satisfies $\mathbf{x} \geq 0$, then it's called a **basic feasible solution**.

If some of the basic variables have value zero, then it's called a **degenerate basic solution**.

Similar to the definition above, we have **degenerate basic feasible solution**.

Theorem 1.5.0.1: Fundamental Theorem of Linear Programming

Given a linear program in standard form, where \mathbf{A} is an $m \times n$ matrix of rank m :

1. if there is a feasible solution, there is a basic feasible solution;
2. if there is an optimal feasible solution, there is an optimal basic feasible solution.

Proof 1.5.1

(1)

If there is a feasible solution \mathbf{x} , the solution satisfy:

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b} \quad (1.5.0.11)$$

Assume there are p of variable $x_i > 0$, then the following equation holds:

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_p\mathbf{a}_p = \mathbf{b} \quad (1.5.0.12)$$

Then there are two cases:

1. if $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$ are linearly independent, then $p \leq m$, which means \mathbf{x} is already a basic solution;
2. otherwise, there would be a non-trivial solution for $y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \cdots + y_p\mathbf{a}_p = \mathbf{0}$, which means $(x_1 - \epsilon y_1)\mathbf{a}_1 + (x_2 - \epsilon y_2)\mathbf{a}_2 + \cdots + (x_p - \epsilon y_p)\mathbf{a}_p = \mathbf{b}$.

Then for any value of ϵ , $\mathbf{x} - \epsilon\mathbf{y}$ is a solution but may violate the signal constraint. Mention that there is at least one y_i is negative or positive, thus there is at least one x_i decreasing when we increase the ϵ ($\epsilon > 0$), thus we set $\epsilon = \min\{x_i/y_i : y_i > 0\}$, which would bring us to a new feasible solution but the number of zero is larger.

Through such iteration, we can get at least one basic feasible solution.

(2)

The idea is the same as the first one. In case one it's obvious. in case two, we need to prove for any $\mathbf{x} - \epsilon\mathbf{y}$ is still optimal.

Note the new value is $\mathbf{c}^T\mathbf{x} - \epsilon\mathbf{c}^T\mathbf{y}$.

Q.E.D.

Remark 1.5.1

This theorem reduce the original problem to the size of $\binom{n}{m} = \frac{n!}{m!(n-m)!}$.

Definition 1.5.0.2: Extreme Point

A point \mathbf{x} in a convex set \mathbf{C} is an *extreme point* if there are **no** two distinct $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{C}$ such that $\mathbf{x} = \alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2$ for some $\alpha, 0 < \alpha < 1$

Theorem 1.5.0.2: Equivalence of extreme points and basic solutions

Let \mathbf{A} be an $m \times n$ matrix with rank m , Let \mathbf{K} denote the *convex polytope* consisting all vector \mathbf{x} satisfying $\mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}$.

A vector \mathbf{x} is an extreme point of \mathbf{K} if and only if \mathbf{x} is a basic feasible solution.

Proof 1.5.2

(1) BFS \rightarrow extreme point:

Suppose $\mathbf{x} = (x_1, x_2, \dots, x_m, 0, 0, \dots, 0)$ is a BFS, it satisfies $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \dots + x_m\mathbf{a}_m = \mathbf{b}$.

b. If $\mathbf{x} = \alpha\mathbf{y} + (1 - \alpha)\mathbf{z}$, since the value in \mathbf{y}, \mathbf{z} is larger than 0, then we have:

$$y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \dots + y_m\mathbf{a}_m = \mathbf{b} \quad z_1\mathbf{a}_1 + z_2\mathbf{a}_2 + \dots + z_m\mathbf{a}_m = \mathbf{b} \quad (1.5.0.13)$$

Because the vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m$ are linearly independent, then we can get $\mathbf{x} = \mathbf{y} = \mathbf{z}$, which means that \mathbf{z} is an extreme point. (2) Extreme point \rightarrow BFS:

Assume that \mathbf{x} has k components larger than zero, then we have: $y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \dots + y_k\mathbf{a}_k = \mathbf{0}$. Assume that $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$ are linearly dependent, which would leads to:

$$y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \dots + y_k\mathbf{a}_k = \mathbf{0} \quad (1.5.0.14)$$

Define $\mathbf{y} = (y_1, y_2, \dots, y_k, 0, 0, \dots, 0)$, it's obvious to see the following can exist:

$$\mathbf{x} + \epsilon\mathbf{y} \geq \mathbf{0}, \quad \mathbf{x} - \epsilon\mathbf{y} \geq \mathbf{0} \quad (1.5.0.15)$$

Contradicts that \mathbf{x} is an extreme point $\rightarrow \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$ are linearly independent $\rightarrow \mathbf{x}$ is a BFS.

Q.E.D.

Corollary 1.5.0.1

1. If the convex set \mathbf{K} is nonempty, there is at least one extreme point;
2. If there is a finite optimal solution to a linear programming problem, there is a finite optimal solution which is an extreme point of the constraint set.;
3. The constraint set \mathbf{K} possesses at most a finite number of extreme points;

4. If K is bounded, then it's a *convex polyhedron*.

1.5.1 Simplex Method

1.6 Convex Optimization

1.7 Non-Convex Optimization

Chapter 2

Statistics and Econometrics

2.1 Statistical Inference

This section is mainly the notes from [Casella and Berger, 2021] Statistics is using observed data to **inference** the statistical model.

2.1.1 Exponential Families

Definition 2.1.1.1: Exponential Families

A family of pdfs or pmfs is called an *exponential family* if:

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp \left(\sum_{i=1}^k w_i(\boldsymbol{\theta})t_i(x) \right) \quad (2.1.1.1)$$

Where $h(x) \geq 0$, $c(\boldsymbol{\theta}) \geq 0$, and $h(x), t_i(x)$ don't depend on $\boldsymbol{\theta}$. $c(\boldsymbol{\theta}), w_i(\boldsymbol{\theta})$ don't depend on x .

- Continuous: normal, gamma, beta, exponential;
- Discrete: binomial, poisson, negative binomial;
- $\boldsymbol{\theta} = \theta_1, \theta_2, \dots, \theta_d$, k must $\geq d$;
- If $k = d$, then it's a *full exponential family*, if $k > d$, then it's a *curved exponential family* (For example, most normal distributions are *full exponential family*, but normal distribution satisfy $\mu = \sigma^2$ is a *curved exponential family*).

To verify a pdf is an exponential family, identify the function $h(x), c(\boldsymbol{\theta}), t_i(x), w_i(\boldsymbol{\theta})$, then verify these functions satisfy the condition above.

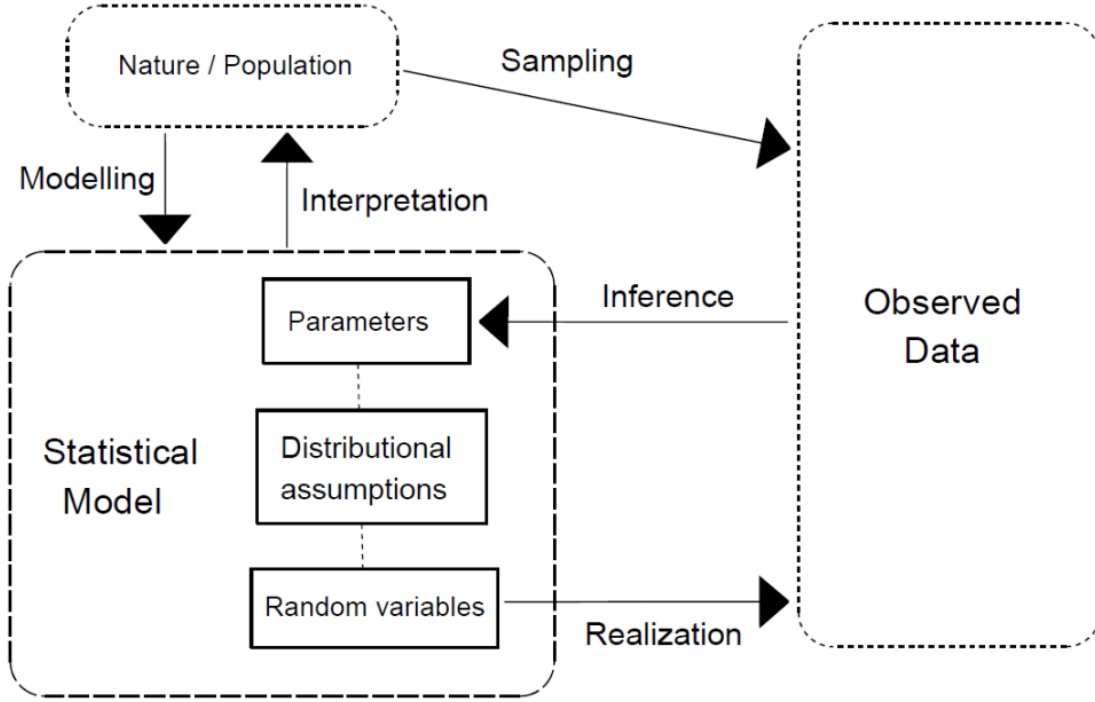


Figure 2.1: The scope of statistical inference

Example 2.1.1.1

Show that Binomial, Poisson, Exponential and Normal distribution belongs to the exponential families.

Solution:

Binomial:

$$\begin{aligned}
 f(x|p) &= \binom{n}{x} p^x (1-p)^{n-x} = \binom{n}{x} (1-p)^n \left(\frac{p}{1-p} \right)^x \\
 &= \binom{n}{x} (1-p)^n \exp \left(x \log \left(\frac{p}{1-p} \right) \right),
 \end{aligned} \tag{2.1.1.2}$$

Among which $\binom{n}{x}$ is $h(x)$, $(1-p)^n$ is $c(\theta)$, x is $t_1(x)$, and $\log \left(\frac{p}{1-p} \right)$ is $w_i(\theta)$. Note that $f(x | p)$ is only in exponential families when $0 < p < 1$.

Poisson:

$$f(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{1}{x!} e^{-\lambda} \exp(x \log(\lambda))$$

(2.1.1.3)

then

$$h(x) = \frac{1}{x!}, c(\lambda) = e^{-\lambda}, t(x) = x \text{ and } w(\lambda) = \log(\lambda).$$

Exponential:

$$f(x|\beta) = \frac{1}{\beta} \exp\left(-\frac{x}{\beta}\right)$$

(2.1.1.4)

then

$$h(x) = 1, c(\beta) = \frac{1}{\beta}, t(x) = x \text{ and } w(\beta) = -\frac{1}{\beta}.$$

Normal:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right)$$

(2.1.1.5)

then

$$h(x) = 1, c(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right),$$

$$t_1(x) = -\frac{x^2}{2}, w_1(\mu, \sigma) = \frac{1}{\sigma^2}, t_2(x) = x \text{ and } w_2(\mu, \sigma) = \frac{\mu}{\sigma^2}.$$

Example 2.1.1.2

Show the following are exponential families:

- Gamma family with either α, β is unknown or both unknown;
- Beta family with either α, β is unknown or both unknown;
- Negative Binomial family when r is unknown.

Solution:

Gamma α unknown:

$$f(x|\alpha; \beta) = e^{-x/\beta} \frac{1}{\Gamma(\alpha)\beta^\alpha} \exp((\alpha - 1) \log x)$$

(2.1.1.6)

thus $h(x) = e^{-x/\beta}, x > 0; c(\alpha) = \frac{1}{\Gamma(\alpha)\beta^\alpha}; w_1(\alpha) = \alpha - 1; t_1(x) = \log x.$

Gamma beta unknown:

$$f(x|\beta; \alpha) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-\frac{x}{\beta}} \quad (2.1.1.7)$$

thus $h(x) = \frac{x^{\alpha-1}}{\Gamma(\alpha)}, x > 0; c(\beta) = \frac{1}{\beta^\alpha}; w_1(\beta) = \frac{1}{\beta}; t_1(x) = -x.$

Gamma both unknown:

$$f(x|\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \exp((\alpha - 1) \log x - \frac{x}{\beta}) \quad (2.1.1.8)$$

thus $h(x) = I_{\{x>0\}}(x); c(\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha}; w_1(\alpha) = \alpha - 1; t_1(x) = \log x; w_2(\alpha, \beta) = -1/\beta; t_2(x) = x.$

Beta α unknown:

$h(x) = (1 - x)^{\beta-1} I_{[0,1]}(x), \quad c(\alpha) = \frac{1}{B(\alpha, \beta)}, \quad w_1(\alpha) = \alpha - 1, \quad t_1(x) = \log x$

Beta β unknown:

$h(x) = x^{\alpha-1} I_{[0,1]}(x), \quad c(\beta) = \frac{1}{B(\alpha, \beta)}, \quad w_1(\beta) = \beta - 1, \quad t_1(x) = \log(1 - x)$

Beta both unknown:

$h(x) = I_{[0,1]}(x), \quad c(\alpha, \beta) = \frac{1}{B(\alpha, \beta)}, \quad w_1(\alpha) = \alpha - 1, \quad t_1(x) = \log x, \quad w_2(\beta) = \beta - 1, \quad t_2(x) = \log(1 - x).$

Negative Binomial:

$$h(x) = \binom{r+x-1}{x} I_{\mathbb{N}}(x), \quad c(p) = \left(\frac{p}{1-p} \right)^r, \quad w_1(p) = \log(1-p), \quad t_1(x) = x. \quad (2.1.1.9)$$

Theorem 2.1.1.1: Expectation and Variance of Exponential Families

If X is a random variable that satisfies any distribution from the exponential families, then:

$$\begin{aligned} 1. \quad & E \left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) = - \frac{\partial}{\partial \theta_j} \log(c(\boldsymbol{\theta})); \\ 2. \quad & \text{Var} \left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) = - \frac{\partial^2}{\partial \theta_j^2} \log(c(\boldsymbol{\theta})) - E \left(\sum_{i=1}^k \frac{\partial^2 w_i(\boldsymbol{\theta})}{\partial \theta_j^2} t_i(X) \right). \end{aligned} \quad (2.1.1.10)$$

Example 2.1.1.3

Derive the mean and variance for binomial and normal distribution using the above theorem.

Solution:

Binomial:

$$h(x) = \binom{n}{x}, c(p) = (1-p)^n, t(x) = x \text{ and } w(p) = \log\left(\frac{p}{1-p}\right).$$

Then,

$$\begin{aligned}\frac{d}{dp}w(p) &= \frac{d}{dp} \log\left(\frac{p}{1-p}\right) = \frac{1}{p(1-p)}, \\ \frac{d^2}{dp^2}w(p) &= -\frac{1}{p^2} + \frac{1}{(1-p)^2} = \frac{2p-1}{p^2(1-p)^2}, \\ \frac{d}{dp} \log(c(p)) &= \frac{d}{dp} n \log(1-p) = -\frac{n}{1-p}, \\ \frac{d^2}{dp^2} \log(c(p)) &= -\frac{n}{(1-p)^2}.\end{aligned}\tag{2.1.1.11}$$

Therefore, from Theorem 3.4.2, we have

$$\begin{aligned}E\left(\frac{1}{p(1-p)}X\right) &= \frac{n}{1-p} \Rightarrow E(X) = np, \\ \text{Var}\left(\frac{1}{p(1-p)}X\right) &= \frac{n}{(1-p)^2} - E\left(\frac{2p-1}{p^2(1-p)^2}X\right) \Rightarrow \text{Var}(X) = np(1-p)\end{aligned}$$

Normal:

For Normal Distribution, we have

$$h(x) = 1, c(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right),$$

$$t_1(x) = -\frac{x^2}{2}, w_1(\mu, \sigma) = \frac{1}{\sigma^2}, t_2(x) = x \text{ and } w_2(\mu, \sigma) = \frac{\mu}{\sigma^2}.$$

Then,

$$\begin{aligned}\frac{\partial w_1(\mu, \sigma)}{\partial \mu} &= \frac{\partial(1/\sigma^2)}{\partial \mu} = 0, \\ \frac{\partial w_2(\mu, \sigma)}{\partial \mu} &= \frac{\partial(\mu/\sigma^2)}{\partial \mu} = \frac{1}{\sigma^2}, \\ \frac{\partial w_1(\mu, \sigma)}{\partial \sigma} &= \frac{\partial(1/\sigma^2)}{\partial \sigma} = -\frac{2}{\sigma^3}, \\ \frac{\partial w_2(\mu, \sigma)}{\partial \sigma} &= \frac{\partial(\mu/\sigma^2)}{\partial \sigma} = -\frac{2\mu}{\sigma^3}, \\ \frac{\partial}{\partial \mu} \log(c(\mu, \sigma)) &= \frac{\partial}{\partial \mu} \left(-\frac{\log(2\pi)}{2} - \log(\sigma) - \frac{\mu^2}{2\sigma^2} \right) = -\frac{\mu}{\sigma^2}, \\ \frac{\partial}{\partial \sigma} \log(c(\mu, \sigma)) &= \frac{\partial}{\partial \sigma} \left(-\frac{\log(2\pi)}{2} - \log(\sigma) - \frac{\mu^2}{2\sigma^2} \right) = -\frac{1}{\sigma} + \frac{\mu^2}{\sigma^3}, \\ E\left(\frac{1}{\sigma^2}X\right) &= \frac{\mu}{\sigma^2} \text{ and } E\left(-\frac{2}{\sigma^3}\left(-\frac{X^2}{2}\right) - \frac{2\mu}{\sigma^3}X\right) = \frac{1}{\sigma} - \frac{\mu^2}{\sigma^3},\end{aligned}\tag{2.1.1.12}$$

which implies

$$E(X) = \mu, E(X^2) = \mu^2 + \sigma^2 \text{ and } \text{Var}(X) = E(X^2) - (EX)^2 = \sigma^2$$

Definition 2.1.1.2: The indicator function

$$I_A(x) = \begin{cases} 1, & \text{if } x \in A \\ 0, & \text{if } x \notin A \end{cases}$$

Example 2.1.1.4

Show that $f(x | \theta) = \frac{1}{\theta} \exp(1 - \frac{x}{\theta})$ is **NOT** an exponential family.

Solution:

$$f(x|\theta) = \frac{1}{\theta} \exp\left(1 - \frac{x}{\theta}\right) I_{[\theta, \infty)}(x)\tag{2.1.1.13}$$

At here $h(x) = I_{[\theta, \infty)}(x)$, which is not independent with θ .

Definition 2.1.1.3: Reparameterization of Exponential Families

$$f(x|\boldsymbol{\eta}) = h(x)c^*(\boldsymbol{\eta}) \exp \left(\sum_{i=1}^k \eta_i t_i(x) \right) \quad (2.1.1.14)$$

Where $h(x)$ and $t_i(x)$ is identical with the original of parameterization. $\boldsymbol{\eta} = (\eta_1, \eta_2, \dots, \eta_n)$ and $\eta_i = w_i(\theta)$. And to make it a pdf (integrates to 1):

$$c^*(\boldsymbol{\eta}) = \left[\int_{-\infty}^{\infty} h(x) \exp \left(\sum_{i=1}^k \eta_i t_i(x) \right) dx \right]^{-1} \quad (2.1.1.15)$$

2.1.2 Scale and Location

step 1: define the standard pdf $f(Z)$, for example: $f(Z) = e^{-Z}, Z \geq 0$;

step 2: find the relationship between X and Z : $\sigma Z + \mu = X$, where σ is the scale parameter and μ is the location parameter;

step 3: replace Z using X : $f(x) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right)$ is a distribution transformed from the standard pdf.

Theorem 2.1.2.1: Sacle-Location Family

If $f(x)$ is any pdf, for $\mu, \sigma > 0$, then the function $g(x | \mu, \sigma) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right)$ is a valid pdf.

Proof 2.1.1

$$g(x|\mu, \sigma) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right) \geq 0 \quad (2.1.2.1)$$

$$\int_{-\infty}^{\infty} g(x|\mu, \sigma) dx = \int_{-\infty}^{\infty} \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right) dx \xrightarrow{(y=\frac{x-\mu}{\sigma})} \int_{-\infty}^{\infty} f(y) dy = 1.$$

Q.E.D.

Remark 2.1.1

- For discrete R.V.(pmf), the above theorem doesn't hold. (ignore the $\frac{1}{\sigma}$);
- The σ is the scale parameter, the μ is the location parameter.

Example 2.1.2.1

$$\int_{\mathbf{X}} (x_1 | \theta) = \frac{1}{\theta} \exp\{1 - \frac{x}{\theta}\}, \quad x \geq \theta \quad (2.1.2.2)$$

$$= \frac{1}{\theta} \exp\{-\frac{x - \theta}{\theta}\} \cdot I_{[0, \infty)}(x) \quad (2.1.2.3)$$

Is θ a scale parameter or a location parameter?

Solution:

It depends on the standard pdf:

If the standard pdf is $f_Z(z) = \exp\{-z\} \cdot I_{[0, +\infty)}(z)$:

then θ serves as both parameters because $f_X(x) = \frac{1}{\theta} \exp\{-\frac{x-\theta}{\theta}\} \cdot I_{[0, +\infty)}(\frac{x-\theta}{\theta})$.

Else if the standard pdf is $f_Z(z) = \exp\{1 - z\} \cdot I_{[0, +\infty)}(z)$:

then θ is only a scale parameter because $f_X(x) = \frac{1}{\theta} \exp\{1 - \frac{x}{\theta}\} \cdot I_{[0, +\infty)}(\frac{x-\theta}{\theta})$.

Theorem 2.1.2.2: For any pdf $f(\cdot)$, and $\sigma, \mu > 0$. Then X is a random variable with pdf $\frac{1}{\sigma} f(\frac{x-\mu}{\sigma})$ **if and only if** there is a random variable Z with pdf $f(z)$ and $X = \sigma Z + \mu$.

Proof 2.1.2: Rigorous Proof

The necessity of the proof can be found in 2.1.2, here the sufficiency need
Q.E.D.

Proof 2.1.3: Intuitive Proof

Q.E.D.

2.1.3 Data Reduction

The idea of data reduction is to summarize or reduce the data X_1, X_2, \dots, X_n to get the information of the unknown parameter θ .

There are many sample point $\mathbf{x} = (x_1, x_2, \dots, x_n)$, which are realizations (observations) of the random variable $\mathbf{X} = (X_1, X_2, \dots, X_n)$. A **Statistic** $T(\mathbf{X})$ is a form of data reduction, or a summary of the data, $T(\mathbf{x})$ is an observation of $T(\mathbf{X})$. \mathcal{X} is the **sample space**. $\mathcal{T} = \{t : t = T(\mathbf{x}), \text{ for } \mathbf{x} \in \mathcal{X}\}$ is the image of cX under $T(\mathbf{X})$. $T(\mathbf{X})$ partition the sample space \mathcal{X} into sets $A_t = \{\mathbf{x} : T(\mathbf{x}) = t, \mathbf{x} \in \mathcal{X}\}$.

Example 2.1.3.1

Give a two-dimensional example for random variable, sample point, Statistic, observation of Statistic, sample space, image and A_t .

Solution:

Assume $\mathbf{X} = X_1, X_2$, among which X_1, X_2 are Bernoulli R.V. with $p = 0.5$. $\mathbf{x} = (0, 1)$ is a sample point.

The sample space \mathcal{X} is $(0, 0), (0, 1), (1, 0), (1, 1)$, set $T(\mathbf{X})$ as the statistic, then the image \mathcal{T} is $(0, 1, 2)$.

$A_1(\mathbf{x}) = ((1, 0), (0, 1))$, $A_2(\mathbf{x}) = (1, 1)$, $A_0(\mathbf{x}) = (0, 0)$.

2.2 Causal Model

2.2.1 Rubin's Causal Model

2.3 Reduced-Form Identification

The main contents of this chapter are the notes of [Angrist and Pischke, 2014], [Angrist and Pischke, 2009].

2.3.1 Randomized Trials

Keywords 2.1

ATT, ATE, Counter-factual World, Potential Outcome

The outcome is Y_i , the potential outcome is Y_{1i}, Y_{0i} :

$$Y_i = \begin{cases} Y_{1i} & \text{if } D_i = 1 \\ Y_{0i} & \text{if } D_i = 0 \end{cases} \quad Y_i = Y_{0i} + (Y_{1i} - Y_{0i})D_i \quad (2.3.1.1)$$

Naive Comparison:

$$\begin{aligned} \{\mathbf{E}[Y_i|D_i = 1] - \mathbf{E}[Y_i|D_i = 0]\} &= \{\mathbf{E}[Y_{1i}|D_i = 1] - \mathbf{E}[Y_{0i}|D_i = 1]\} \\ &\quad + \{\mathbf{E}[Y_{0i}|D_i = 1] - \mathbf{E}[Y_{0i}|D_i = 0]\} \end{aligned} \quad (2.3.1.2)$$

- observed difference: $\mathbf{E}[Y_i|D_i = 1] - \mathbf{E}[Y_i|D_i = 0]$;
- ATT: $\mathbf{E}[Y_{1i}|D_i = 1] - \mathbf{E}[Y_{0i}|D_i = 1]$;
- selection bias: $\mathbf{E}[Y_{0i}|D_i = 1] - \mathbf{E}[Y_{0i}|D_i = 0]$.

The existence of the selection bias is due to the dependence between D_i and the **potential** outcome Y_{1i}, Y_{0i} . Randomization can make $D_i \perp Y_{1i}, Y_{0i}$:

$$E[Y_{1i}|D_i = 1] = E[Y_{1i}|D_i = 0] = E[Y_{1i}]E[Y_{0i}|D_i = 1] = E[Y_{0i}|D_i = 0] = E[Y_{0i}] \quad (2.3.1.3)$$

So selection = $\mathbf{E}[Y_{0i}|D_i = 1] - \mathbf{E}[Y_{0i}|D_i = 0] = 0$. 😊 Besides, by randomization, $ATT = ATE$. Randomization example:

- Rand HIE experiment: whether the insurance program makes people healthier;
- STAR: the effects of class size on education;
- OHP: this group experiment is not perfect because the group is not a determinant of whether to receive the treatment, but the treatment group does have a higher probability to get the treatment (**Instrumental Variable** can handle this situation);

Remark 2.3.1

- By randomization, the individual differences still exists;
- Checking for balance is an important step in randomization;
- The most critical idea of randomization is **Other Things Equal**(*ceteris paribus*);
- Randomization was invented by *Ronald Aylmer Fisher* in 1925.

2.3.2 Regression and Matching

Keywords 2.2

- CEF, CEF decomposition, ANOVA;
- regression justification,

Conditional Expectation Function (CEF) is a **population** concept:

$$\begin{aligned} E[Y_i|X_i = x] &= \int t f_y(t|X_i = x) dt \\ E[Y_i|X_i = x] &= \sum_t t P(Y_i = t|X_i = x) \end{aligned} \tag{2.3.2.1}$$

Lemma 2.3.2.1: The law of iterated expectations

$$E[y_i|X_i = x] = \int t f_y(t|X_i = x) dt. \tag{2.3.2.2}$$

Proof 2.3.1

$$\begin{aligned}
 E\{E[y_i|X_i]\} &= \int E[y_i|X_i = u] g_x(u) du \\
 &= \int \left[\int t f_y(t|X_i = u) dt \right] g_x(u) du \\
 &= \int \int t f_y(t|X_i = u) g_x(u) du dt \\
 &= \int t \left[\int f_y(t|X_i = u) g_x(u) du \right] dt = \int t \left[\int f_{xy}(u, t) du \right] dt \\
 &= \int t g_y(t) dt.
 \end{aligned}
 \tag{2.3.2.3}$$

Q.E.D.

♥ 3 important property of CEF:

Theorem 2.3.2.1: CEF Decomposition Property

$$Y_i = E[Y_i|X_i] + \epsilon_i \tag{2.3.2.4}$$

where ϵ_i is mean independent of X_i , and X_i is uncorrelated with any function of X_i .

Proof 2.3.2

Take the expectation of X_i at both sides:

$$\begin{aligned}
 E[Y_i|X_i] &= E[E[Y_i|X_i]|X_i] + E[\epsilon_i|X_i] \\
 E[\epsilon_i|X_i] &= E[Y_i|X_i] - E[Y_i|X_i] = 0
 \end{aligned}
 \tag{2.3.2.5}$$

$$E[\epsilon_i] = \int_{X_i} f_x(t) E[\epsilon_i|X_i] dt = \int_{X_i} 0 dt = 0 = E[\epsilon_i|X_i] \tag{2.3.2.6}$$

Q.E.D.

This means that Y_i can be decomposed into 2 parts: explained by X_i and terms uncorrelated with X_i .

Theorem 2.3.2.2: CEF Prediction Property

CEF is the best estimator of Y_i in MMSE sense, which means:

$$E[Y_i|X_i] = \arg \min_{m(X_i)} E[(Y_i - m(X_i))^2] \quad (2.3.2.7)$$

Proof 2.3.3

$$\begin{aligned} (Y_i - m(X_i))^2 &= ((Y_i - E[Y_i|X_i]) + (E[Y_i|X_i] - m(X_i)))^2 \\ &= (Y_i - E[Y_i|X_i])^2 + 2(E[Y_i|X_i] - m(X_i))(Y_i - E[Y_i|X_i]) \\ &\quad + (E[Y_i|X_i] - m(X_i))^2 \end{aligned} \quad (2.3.2.8)$$

By setting $m(X_i) = \text{CEF}$, the formula has the lowest constant value.

Q.E.D.

Theorem 2.3.2.3: ANOVA Theorem

$$V(Y_i) = E[V(Y_i|X_i)] + V(E[Y_i|X_i]) \quad (2.3.2.9)$$

This indicates that the variance of Y_i can be decomposed into two parts:

1. the variance of the CEF;
2. the variance of the residual;

Remark 2.3.2

- The CEF property **dosen't rely on any assumption!** It has nothing to do with regression right now;
- If X_i is not mean independent of Y_i , then by ANOVA theorem, the variance of the outcome variable controlled by X_i could be smaller;

CEF and (Population) Regression

$$\begin{aligned} \beta &= \arg \min_b E[(Y_i - X_i'b)^2] \\ 1storder : E[X_i(Y_i - X_i'b)] &= 0 \\ solution : \beta &= E[X_i X_i']^{-1} E[X_i Y_i] \end{aligned} \quad (2.3.2.10)$$

Theorem 2.3.2.4: Regression Anatomy

$$\beta_k = \frac{Cov(Y_i, \tilde{X}_{ki})}{V(\tilde{X}_{ki})} \quad (2.3.2.11)$$

Corollary 2.3.2.1: Bivariate Case

$$\beta = \frac{Cov(Y_i, \tilde{X}_i)}{V(\tilde{X}_i)} \quad (2.3.2.12)$$

Proof 2.3.4

Substitute

$$Y_i = \alpha + \beta_1 x_{1i} + \cdots + \beta_k x_{ki} + e_i \quad (2.3.2.13)$$

\tilde{x}_{ki} is uncorrelated with e_i and other covariates by construction, thus

$$Cov(\tilde{x}_{ki}, x_{ki}) = Var(\tilde{x}_{ki}), \text{ thus } Cov(Y_i, \tilde{x}_{ki}) = \beta_k Cov(\tilde{x}_{ki}, x_{ki}).$$

Q.E.D.

Remark 2.3.3

The regression anatomy shows that each β_k in multi-regression is the bivariate slope after "partialing out" all the other regressors.

☺ Why the population regression coefficient is what we are interested in (Link with CEF):

Theorem 2.3.2.5: Regression Justification

1. Suppose the CEF is linear, then the population regression function is it;
2. In any condition, $X_i'\beta$ is the best predictor of Y_i in a MMSE sense;
3. The function $X_i'\beta$ provides the MMSE linear approximation to $E[Y_i|X_i]$.

$$\beta = \arg \min_b E\{(E[Y_i|X_i] - X_i'b)^2\} \quad (2.3.2.14)$$

Proof 2.3.5

Suppose $E[Y_i|X_i] = X_i'\beta^*$. By regression decomposition theorem:

$$\begin{aligned} E[X_i(Y_i - X_i'\beta^*)] &= 0 \\ \beta^* &= E[X_i X_i']^{-1} E[X_i Y_i] = \tilde{\beta} \end{aligned} \tag{2.3.2.15}$$

$$\begin{aligned} (Y_i - X_i'b)^2 &= \{(y_i - E[y_i|X_i]) + (E[y_i|X_i] - X_i'b)\}^2 \\ &= (y_i - E[y_i|X_i])^2 + (E[y_i|X_i] - X_i'b)^2 \\ &\quad + 2(y_i - E[y_i|X_i])(E[y_i|X_i] - X_i'b). \end{aligned} \tag{2.3.2.16}$$

Q.E.D.

Corollary 2.3.2.2

- For saturated model, the population linear regression is the CEF;
- For single dummy variable, the coefficient is the mean probability of receiving treatment;

2.3.3 Asymptotic Analysis

2.4 Advanced Econometrics

2.5 Machine Learning Interface

Chapter 3

Economics

3.1 Microeconomics

3.2 Macroeconomics

3.3 Game Theory

3.4 Development Economics

3.4.1 Models of Development Economics

3.4.2 Clan Culture

Definition 3.4.2.1: Clan Culture

A clan is a consolidated kin group made up of component families that trace their patrilineal descent from a common ancestor.

History of Clans

”Modern” clan originated in the Song Dynasty (860-1279 CE). At that time Neo-Confucian ideology was formed, which provided the theoretical basis as well as clan organization structure design. The characteristics of ”modern” clan culture:

1. The families of a clan lived in the same or several nearby communities;
2. Common properties and organized routine group activities, resource pooling;
3. Compilation of genealogies;
4. Own internal governance structures.

Currently although China has been transitioning for a long time from a traditional society to a modern society, clan culture is still prevalent and has a broad impact on the lives of Chinese people, especially in rural areas.

[Bertrand and Schoar, 2006] shows the positive correlation between the fraction of family control among listed firms and family ties using cross-country level data. [Cheng et al., 2021] uses IV (the minimum distance to two prominent neoConfucian academies, the Kaoting Academy (Kaoting Shuyuan) and the Xiangshan Academy (Xiangshan Shuyuan)) to identify that clan culture causes higher firm ownership concentration. Potential reasons that culture affects the concentration of family ownership:

1. Clan culture fosters high trust within the family and low trust in outsiders(**short-radius trust attitude**). According to agency-cost based theories, family ownership can be concentrated in such a situation.
2. *Resources Pooling*: common property ownership.
3. *Amenity Potential*: other things constant, owners subject to stronger influences of clan culture could have a higher utility.

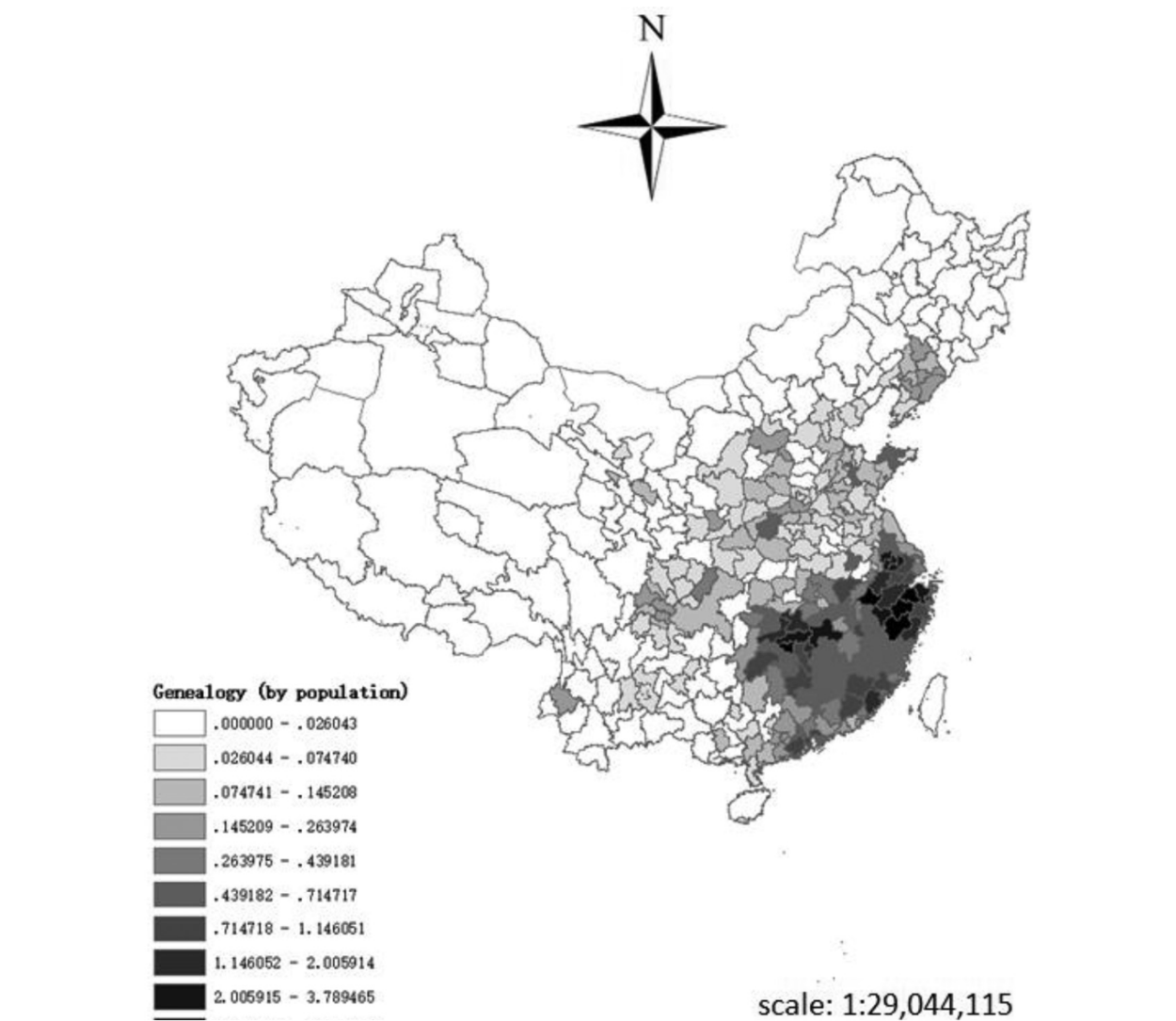


Figure 3.1: Clan Culture Intensity

[Zhang, 2020] estimates the effect of clan on entrepreneurship. He finds that clan leads to a higher occurrence of entrepreneurship by helping overcome financing constraints and escape from local governments' "grabbing hand". [Zhang, 2019] investigates the relationship between the low take-up rate of social pensions and the clan culture intensity. In his article, dummy variable $temple_c$ (whether community c has ancestral temple) is constructed as the proxy variable for the strength of clan culture. Some interesting insights are obtained:

1. Clan culture is positively related to adults raising children for support in their old age;
2. Clan culture is associated with a larger number of children being born and a higher probability of having sons;
3. Clan culture is associated with a higher coresidence rate between old parents and adult sons;
4. Clan culture is associated with a higher likelihood of receiving financial transfers from non-coresident children;
5. Clan culture is associated with a lower likelihood of participating in rural pension programs.

[Cao et al., 2022] There are much research about clan culture outside Mainland China. [Yang, 2019] found the concave relationship the between the **heterogeneity** of clan family and the provision of public goods. This finding implies that group homogeneity yields not only benefits, but also some possible costs. reCommon Control Variables in Clan Research Identification(Individual Level):

1. *Hukou* Status;

Regional Level:

1. Distance to the sea;

Data resources in Clan Research

Remark 3.4.1

- *Comprehensive Catalogue of the Chinese Genealogy* can be used to construct the strength of clan culture;
-

3.5 Data Science Interface

Chapter 4

Computer Science and Data Science

4.1 Cloud Computing

This section is mainly the notes of [Hwang, 2017].

4.1.1 Principles of Cloud Computing System

Keywords 4.1

- HTC, HPC, physical machine, virtual machine, VM clusters;
- IaaS, PaaS, SaaS,

Definition 4.1.1.1: HPC and HTC

- **HTC** refers to *highly-throughput computing* system built with parallel and distributed computing technologies;
- **HPC** refers to *highly-performance computing* system in terms of raw speed in batch processing.

Definition 4.1.1.2: Cloud

A cloud is a pool of virtualized computer resources. A cloud can host a variety of different workloads, including batch-style backend jobs and interactive, user-facing applications. Some view the clouds as computing clusters with modest changes in *virtualization*.

Basic cloud service models are *infrastructure as a service* (IaaS) or infrastructure cloud, *platform as a service* (PaaS) or platform cloud, and *software as a service* (SaaS) or application cloud. Their differences with on-premise computing are listed below:

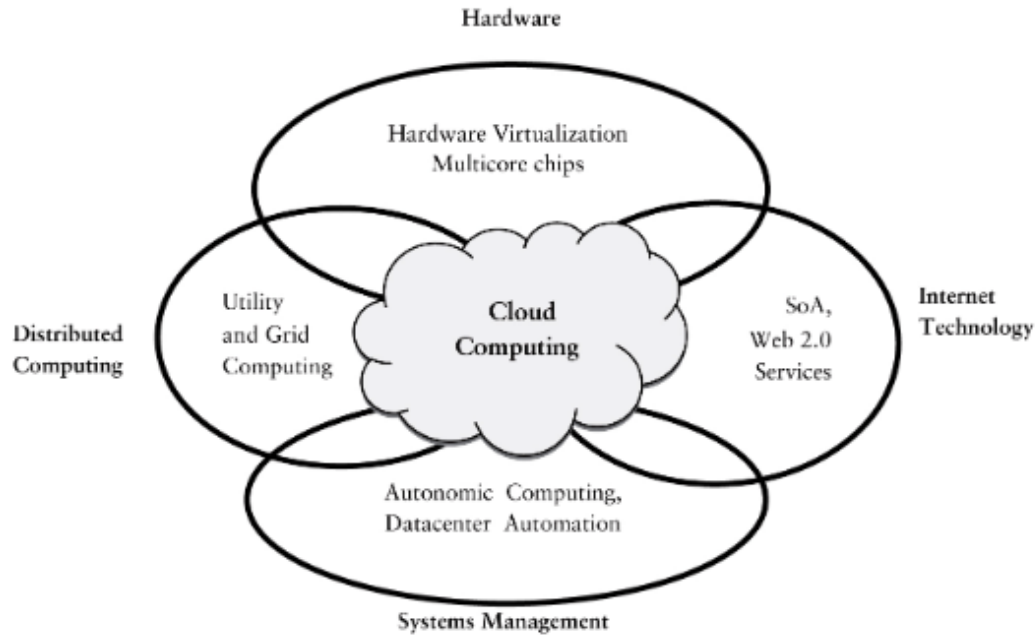


Figure 4.1: Cloud Computing Technology Convergence

- **IaaS:** AWS, GoGrid;
- **PaaS:** Google App Engine, Microsoft Azure, Salesforce;
- **SaaS:** CRM, ERP, HR, Hadoop, Google Docs.

A physical cluster is a collection of servers (PMs) interconnected by a physical network. **virtual clusters** are built with multiple VMs installed at PM belong to one or more physical clusters. The virtual clusters have the following interesting properties:

1. Multiple VMs running with different OSs can be deployed on the same physical node;
2. A VM runs with a guest OS, which is often different than the host OS that manages the resources in the PM, where the VM is implemented;
3. Using VMs can greatly enhance server utilization and application flexibility;
4. VMs can be colonized (replicated) in multiple servers for fault tolerance and disaster recovery;
5. The size (number of nodes) of a virtual cluster can grow or shrink dynamically;
6. the failure of any physical nodes may disable some VMs installed on the failing nodes but the failure of VMs will not pull down the host system.

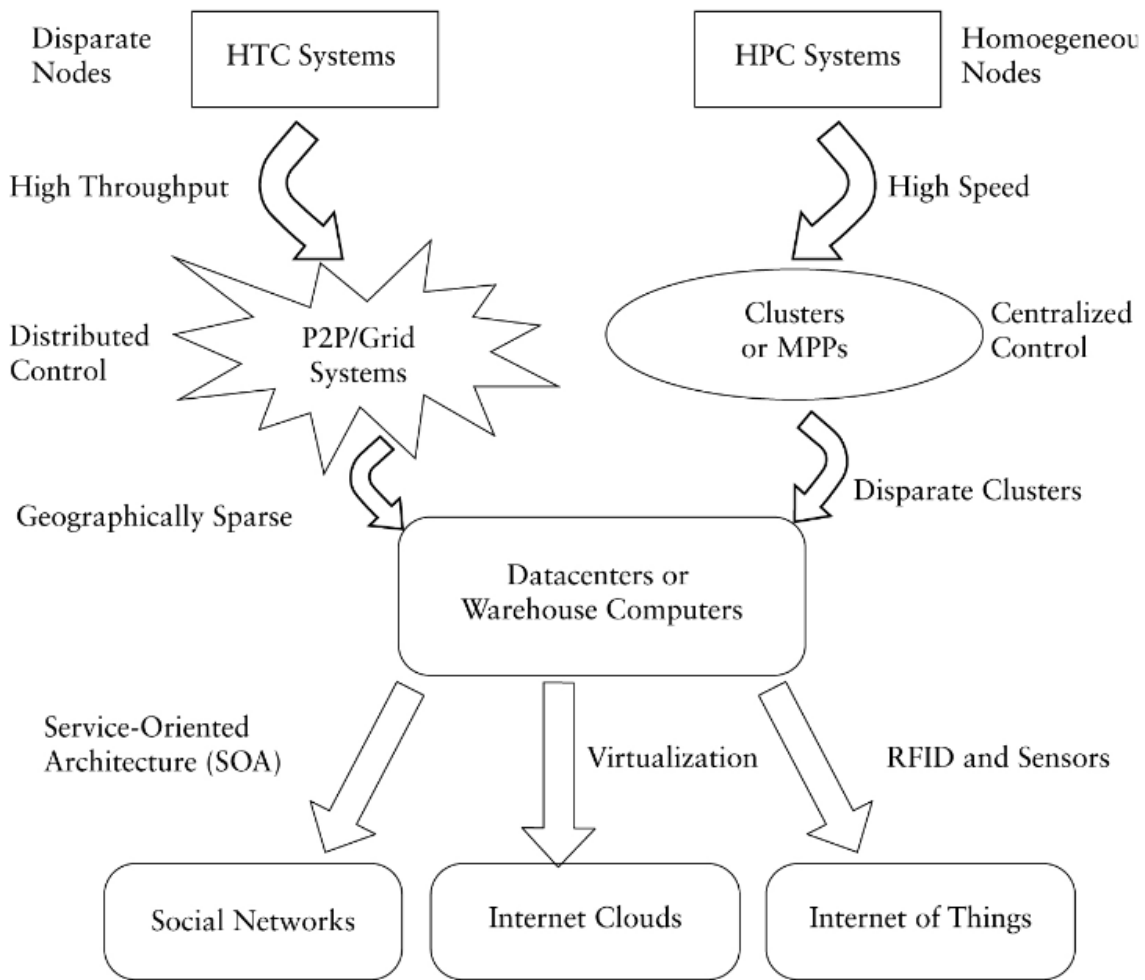


Figure 4.2: From HPC systems and clusters to grids, p2p networks, clouds and IoT

Resource Types	On-Premise Computing	IaaS Model	PaaS Model	SaaS Model
App Software	User	User	Shared	Vendor
Virtual Machines	User	Shared	Shared	Vendor
Servers	User	Vendor	Vendor	Vendor
Storage	User	Vendor	Vendor	Vendor
Networking	Shared	Vendor	Vendor	Vendor

Figure 4.3: Comparing three cloud service models with on-premise computing

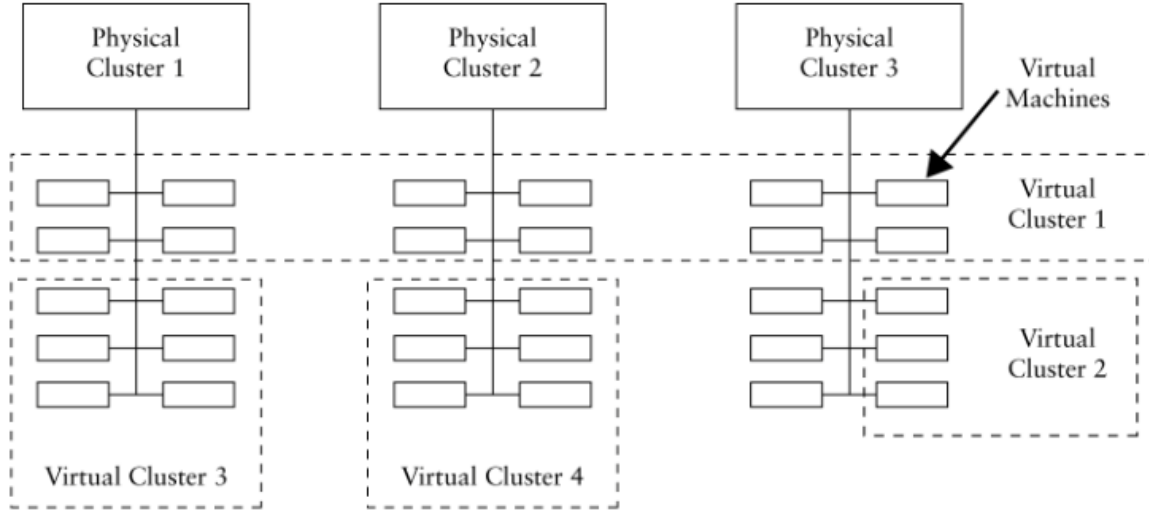


Figure 4.4: physical cluster and virtual cluster

Cloud Scalability

☺ There are two fundamental issues on cloud performance: **Scalability** and **Availability**. The total execution time of the program is calculated by $\alpha T + (1 - \alpha)T/n$.

Theorem 4.1.1.1: Amdahl's Law

The total execution time of the program is calculated by $\alpha T + (1 - \alpha)T/n$.

$$\text{Speedup factor} = S = T / [\alpha T + (1 - \alpha)T/n] = 1 / [\alpha + (1 - \alpha)/n] \quad (4.1.1.1)$$

- The communication time and I/O time are excluded in this formula;
- When $n \rightarrow \infty$, the S can have the upper bound $\frac{1}{\alpha}$, thus α is the *sequential bottleneck* here;
- This speedup is called *fixed-workload speedup*;
- The *cluster efficiency* is defined by $E = S/n = \frac{1}{\alpha n + 1 - \alpha}$; (efficiency means that one more cluster can reduce how much time to process)
- Large sequential bottleneck would lead to many idle servers in the cluster;

Theorem 4.1.1.2: Gustafson's Law

By fixing the parallel execution time at level W :

$$S' = W'/W = [\alpha W + (1 - \alpha)nW]/W = \alpha + (1 - \alpha)n \quad (4.1.1.2)$$

- The efficiency is obtained by: $E' = S'/n = \alpha/n + (1 - \alpha)$, this means one more cluster can increase how much workload;
- For a fixed workload, use Amdahl's Law; for a scaled problem, apply Gustafson's Law.

Cloud Availability

Theorem 4.1.1.3: Cluster Availability

$$\text{Cluster Availability} = MTTF / (MTTF + MTTR) \quad (4.1.1.3)$$

- $MTTF$: mean time to failure;
- $MTTR$: mean time to repair.

Example 4.1.1.1

There is a double-redundancy cluster, the MTTF is 200 units while MTTR is 5 units, calculate the availability of this system.

Solution:

The availability of each server is $200/(200 + 5) = 97.5\%$, the failure rate of the whole system is $1 - (1 - 97.5\%)^2 = 0.625\%$

Consider the use of a cluster of n homogeneous servers in a system, the system is available when more than k machine is running, then the system availability can be expressed by:

$$\begin{aligned} A &= \sum_{i=k}^n \binom{n}{i} p^i (1-p)^{n-i} \\ &= \binom{n}{k} p^k (1-p)^{n-k} + \binom{n}{k+1} p^{k+1} (1-p)^{n-k-1} \\ &\quad + \cdots + \binom{n}{n-1} p^{n-1} (1-p)^1 + \binom{n}{n} p^n (1-p)^0, \end{aligned} \quad (4.1.1.4)$$

CPU and GPU

- CPU has high flexibility for different applications. *Von Neumann bottleneck*: memory access is slow compared to calculation.

CPU vs GPU Architectures

• CPU

- Few Cores
- Lots of Cache
- Handful of Threads
- Independent Processes

• GPU

- Hundreds of Cores
- Thousands of Threads
- Single Process Execution

Systolic SIMD Execution

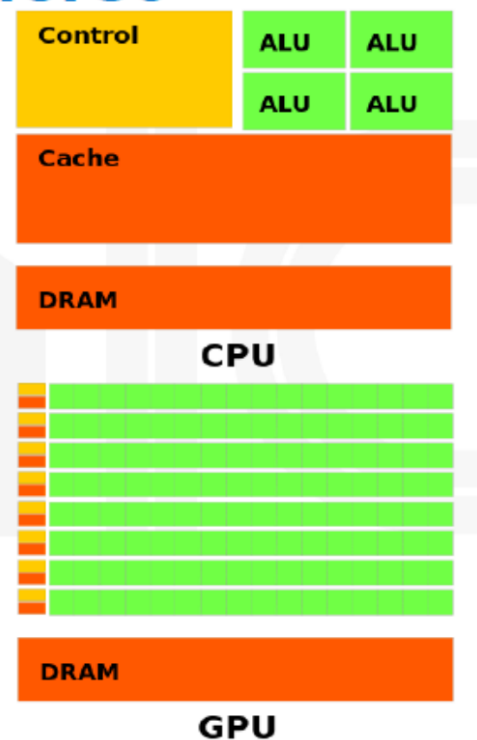


Figure 4.5: CPU and GPU Architecture

- GPU has high throughput, it works well on applications with *massive parallelism*.

4.1.2 Virtual Machines

A VM is essentially built as a software package that can be loaded into a host computer to execute certain user applications. Once the jobs are done, the VM package can be removed from the host computer. The host acts like a “hotel” to accommodate different “guests” at different timeframes.

VMM (virtual machine monitor) can be used to build virtual machine for the guest OS, it can have the following operations: There are five levels of virtualization, among which only 2 are valuable:

- **Hypervisor:** virtualization on top of bare-metal hardware;
- **Container:** virtualization on operating system level, isolated containers of user app with isolated resources.

Hypervisor

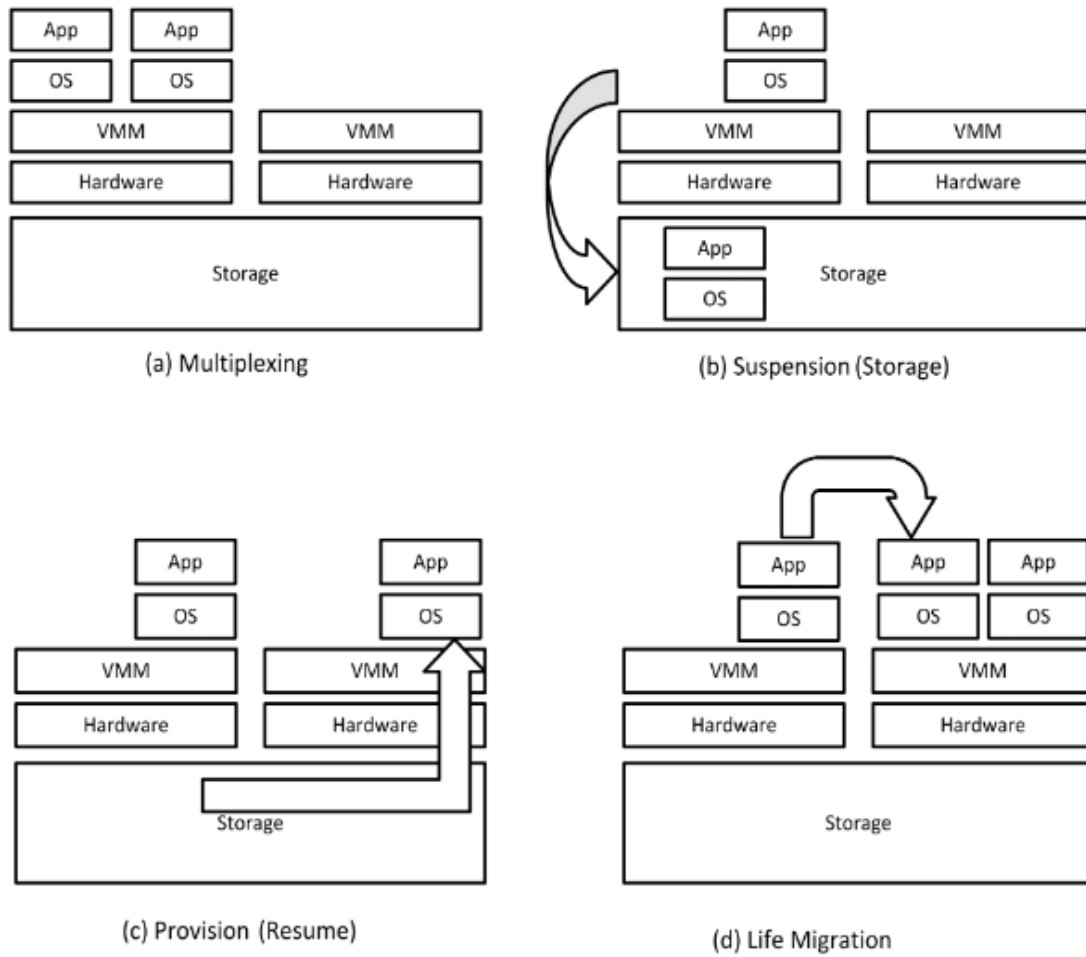


Figure 4.6: VMM Operations

4.2 Machine Learning

4.3 Deep Learning

4.4 Reinforcement Learning

This section is mainly the notes of [Thrun and Littman, 2000].

4.4.1 Introduction and MDP

Keywords 4.2

- Exploration, Exploitation, Reward, Policy, Value Function;
-

The four elements of reinforcement learning:

1. *Policy*: agent's way to interact with the environment;
2. *Reward Signal*: on each time step, the environment sends to the agent a single number;
3. *Value Function*: specify what is good in the long run, which is the discount value of the rewards;
4. *Model*: mimics behaviors of the environment;

Remark 4.4.1

"Deadly Trials"

1. The balance between **Exploration** and **Exploitation**;
2. Reinforcement learning is difficult to generalize;
3. Delayed consequences may cause RL algorithm to perform poorly.

MDP is a mathematical framework to model *discrete-time* sequential decision process, denoted by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \rho_0, \gamma)$:

- \mathcal{S} : the state space, which is the states for the **entire** environment(MOBA games may can not directly be modeled by MDP);
- \mathcal{A} : the action space. \mathcal{A} can depend on $s \in \mathcal{S}$;
- $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$:the environment transition probability function;
- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathbb{R})$: the reward function;
- $\rho_0 \in \Delta(\mathcal{S})$: the initial state distribution;
- $\gamma \in [0, 1]$ is the discount factor.

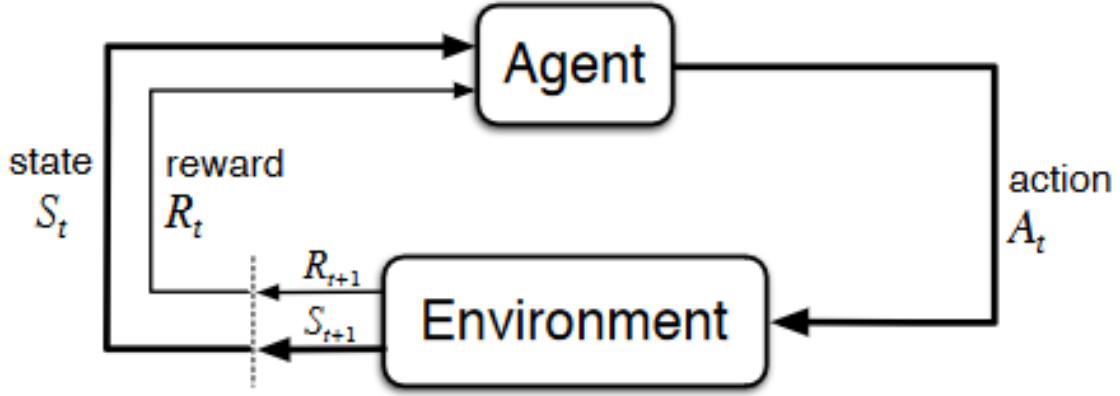


Figure 4.7: The agent–environment interaction in a Markov decision process

Remark 4.4.2

- The $\Delta(\cdot)$ may not be deterministic, but some random distribution;
- Among the above tuple, $\mathcal{S}, \mathcal{T}, \mathcal{R}, \rho_0$ can not be modified by the agent, to train a good policy, \mathcal{A}, γ is the key;
- RL is more like infants, rather than adults;
- The reward function is the way of communicating with the agent *what* to do, not *how* to do;
- The *trajectory* of the MDP sequence: $S_0, A_0, R_1, S_1, A_1, \dots$.

Theorem 4.4.1.1: Dynamics of MDP

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\} \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r | s, a) = 1, \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}(s) \quad (4.4.1.1)$$

Corollary 4.4.1.1: Some formula derived from the dynamics theorem

$$p(s' | s, a) \doteq \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a) \quad (4.4.1.2)$$

$$r(s, a) \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) \quad (4.4.1.3)$$

$$r(s, a, s') \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathbb{R}} r \frac{p(s', r \mid s, a)}{p(s' \mid s, a)} \quad (4.4.1.4)$$

Definition 4.4.1.1: Some useful function:

The act value function given policy π :

$$Q^\pi(s, a) = \mathbb{E}_{s_t, a_t, r_t, t \geq 0} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right] \quad (4.4.1.5)$$

The expected return at state s given policy π :

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} [Q^\pi(s, a)] \quad (4.4.1.6)$$

The **advantage function**:

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) \quad (4.4.1.7)$$

The temporal-difference error:

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t) \quad (4.4.1.8)$$

If we denote $G_t = R_{t+1} + \gamma V(s_{t+1})$, then we have:

Theorem 4.4.1.2: Bellman Equation

$$\begin{aligned} v_\pi(s) &\doteq \mathbb{E}_\pi[G_t \mid S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma V(s_{t+1}) \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) \left[r + \gamma \mathbb{E}_\pi[V(s_{t+1}) \mid S_{t+1} = s'] \right] \\ &= \sum_a \pi(a \mid s) \sum_{s'} p(s', r \mid s, a) \left[r + \gamma v_\pi(s') \right], \quad \text{for all } s \in \mathcal{S}, \end{aligned} \quad (4.4.1.9)$$

A policy π is better π' if and only if $v_\pi(s) \leq v_{\pi'}(s)$ for all $s \in \mathcal{S}$. There is always at list one *optimal policy* denoted by π_* (doesn't hold for partially observed MDP). They share the same state-value function, called the *optimal state-value function*: $v_*(s) \doteq \max_\pi v_\pi(s)$. Optimal policy also shares the same *optimal action-value function*: $q_*(s, a) \doteq \max_\pi q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$. Then we can get the **Bellman optimality function** in an action-value function sense:

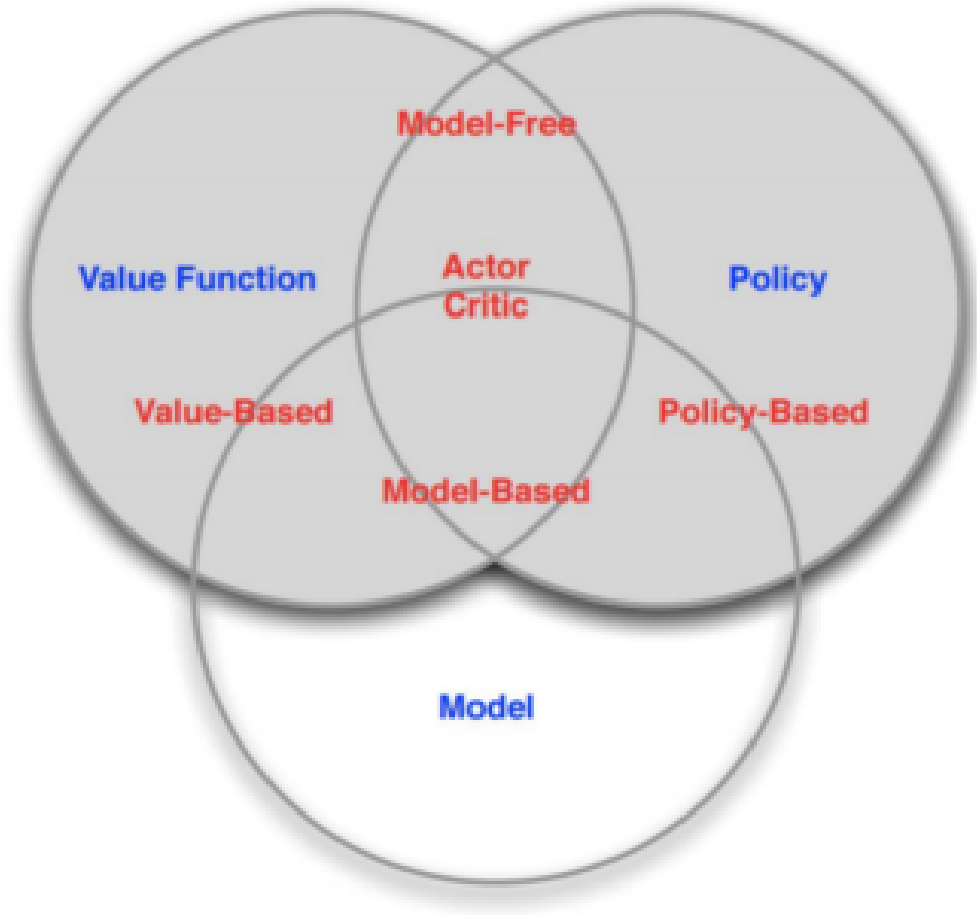


Figure 4.8: Classification of different reinforcement learning agents

Theorem 4.4.1.3: Bellman optimality function

$$\begin{aligned}
 q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \middle| S_t = s, A_t = a \right] \\
 &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right].
 \end{aligned} \tag{4.4.1.10}$$

Definition 4.4.1.2: The optimal policy π^*

A policy π^* is an optimal policy if for every policy π and every $s \in \mathcal{S}$, we have:

$$V^{\pi^*}(s) \geq V^{\pi}(s) \tag{4.4.1.11}$$

Remark 4.4.3

- If we have the full information of the game (MDP framework), then the optimal policy is always deterministic;
- The optimal policy is always stochastic when there are **minimax** structure (e.g. perfect information);
- Whether the optimal policy is stochastic or deterministic has nothing to do with the stochasticity of the game.

Chapter 5

Information Systems and Operations Management

5.1 Empirical Operations Management

[Roth, 2007] describes the evolution of empirical OM from 1980 to 2007, the author selects 12 profounding papers in this domain. [Brusco et al., 2017] reviewed the clustering methods applied in 6 OM journals.

[Choi et al., 2016], multi-methodological OM is advocated, which includes the empirical methodology.

Definition 5.1.0.1: Multi-methodological OM

an approach for OM research in which at least two distinct OM research methods are employed nontrivially to meet the research goals.

[Roth and Singhal, 2022] classified 75 papers as empirical among the top 200 cited papers in *POM*, these papers are mainly from 3 topical areas:

1. Responsibility Operations: covers environmental management, sustainability, humanitarian efforts;
2. Supply Chain Management: bullwhip effect, risk management, supply chain finance;
3. Manufacturing Strategy and Quality Management.

Among these studies, primary data (surveys, experiments, interviews) are used, followed by secondary data (public database, firm's data). Roth's suggestions:

1. For some topic which is very intuitive and not surprising, focus on the **size** of effect rather than sign;

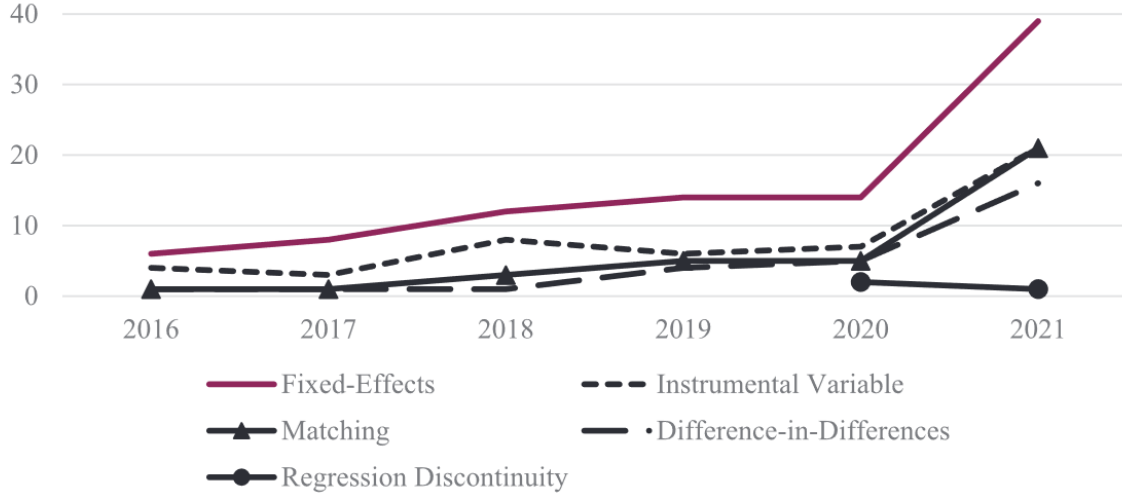


Figure 5.1: identification strategies from 2016 to 2021

2. Avoid the confirmation bias and focus on consistency;
3. Focus on endogeneity and causality, using common sense at the same time.

[Kumar and Tang, 2022] reviewed different domains of OM publications on *POM*, within which a section about empirical OM is covered.

[Mithas et al., 2022] reviewed 411 empirical papers from 2016-2021 on *POM*, with a causal inference and counterfactual perspective.

Remark 5.1.1

Two challenges in assessing causality:

$$T = ATE + [E(Y(0)|Z = 1) - E(Y(0)|Z = 0)] + (1 - \pi) * [(ATT - ATU)]. \quad (5.1.0.1)$$

- $[E(Y(0)|Z = 1) - E(Y(0)|Z = 0)]$ is the **baseline bias**, coming from *OV*B or *simultaneity*;
- $[(ATT - ATU)]$ is the **differential treatment bias**.

[Fisher and Raman, 2022] especially investigate the empirical research in retail operations from traditional ones like forecasting and inventory planning, to new technologies, like radio frequency identification (RFID) and e-commerce.

5.1.1 12 Papers of [Roth, 2007]

[Fisher, 2007] is a conceptual paper, in which a matrix is proposed to navigate conducting research: Fisher suggests that a good empirical OM research should include the following:

Approach	How baseline bias is addressed	How differential treatment effect bias is addressed	What is estimated	Limitations
1. Regression-based	Assumes strong ignorability (i.e., selection on observables only)	Ruled out by constant-coefficient models	ATE	Assumption of correct specification of the relationship between Y and the controls. No distinction between the treatment and covariates
2. Matching and weighting	Assumes strong ignorability (i.e., selection on observables only)	Possible to assess differential treatment effects across strata	ATT, ATU, or ATE (depends on the assumptions)	It is possible to conduct sensitivity analyses (e.g., Rosenbaum's gamma) for assessing selection bias and for the violation of the strong ignorability assumption
3. Instrumental variable (IV)	Relaxes the assumption of selection on observables. Exploits randomization induced by the IV	Ignores treatment heterogeneity by estimation for only a subgroup	LATE (Only for compliers or defiers, but not both)	Difficult to find strong and relevant IVs. Exclusion restriction (IV affects the outcome only via treatment) is not testable. Yields large standard errors if the sample sizes are small
4. Regression discontinuity (RD)	Sharp RD uses the assumption of selection on observables, but fuzzy RD relaxes that. Exploits the randomization induced by the cutoff score on the running variable	Ignores treatment heterogeneity by estimation for only a subgroup	LATE	Running variable perfectly (or fuzzily) determines the treatment assignment. Assumes no discontinuity in other factors that could also affect the outcome other than the treatment
5. Differences-in-differences (DID)	Exploits within-unit variation over time, assuming that all unobservables are time-invariant	Ignores treatment heterogeneity between ATT and ATU	ATT	Assumes parallel trends: Had there been no treatment, the trend in the outcomes would be parallel between the treated and the control. Mostly useful for sharp binary interventions. Requires longitudinal data on both the treated and control units
6. Fixed effects (FEs)	Exploits within-unit variation over time, assuming that all unobservables are time-invariant. Uses changes over time in the control group as a counterfactual for the changes in the treated group	Ignores treatment heterogeneity (assumes that the heterogeneity of the unit-specific causal effect across the population is random)	ATT	Needs strong assumptions or long time series for modeling the counterfactual trajectory. Assumes that past outcomes do not influence the treatment and no lagged treatment effects

Figure 5.2: How identification techniques works

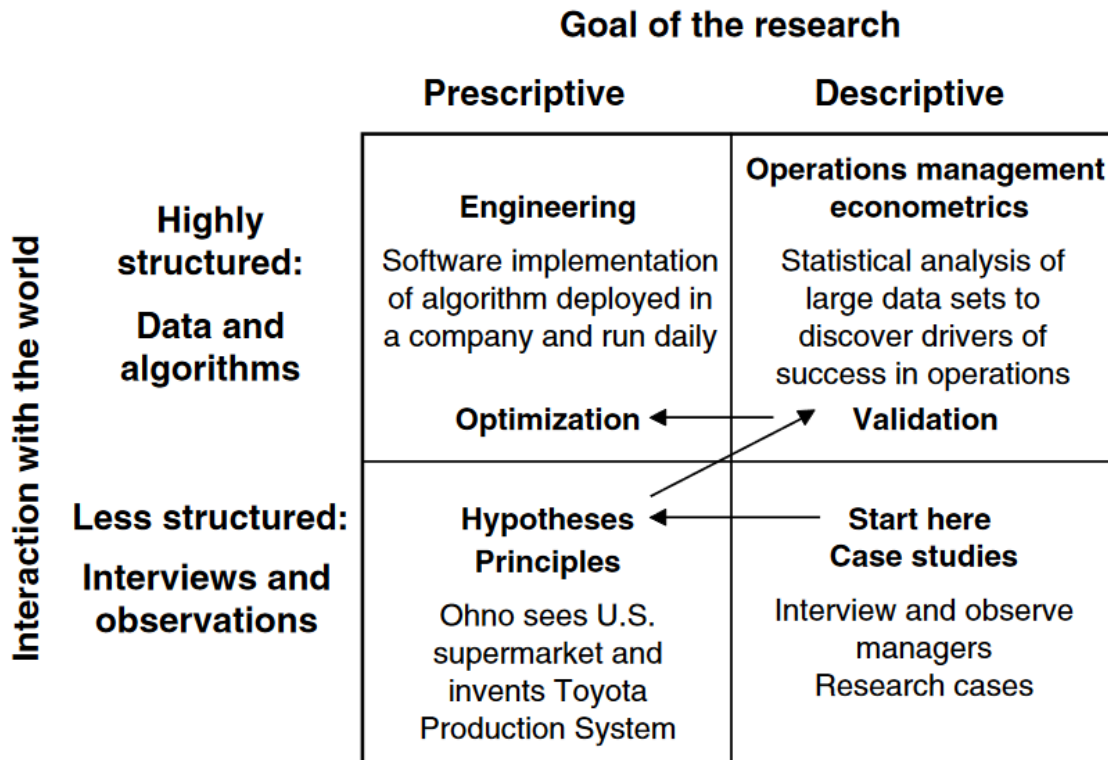


Figure 5.3: Navigating Matrix Cells

1. Identifying and verifying important phenomena;
2. Identifying and characterizing important questions on which we can do useful research;
3. Validating models and assumptions that have been made;
4. Establishing the relevance of our research by demonstrating how the research outputs apply to practice.

[Gans et al., 2007] investigates the posit of some bandit consumer choice models:

5.2 Revenue Management

Revenue management is a data-driven system to price perishable assets tactically at the micro-market level to maximize expected revenue or profit. Some important reviews before 2009 can be found in the book [Gallego and Topaloglu, 2019]. There are some extra summary papers like [Strauss et al., 2018], [Klein et al., 2020].

5.2.1 Traditional RM

Keywords 5.1

- Protection Level, Booking Limit, Littlewood's Rule

Assumption 5.2.1.1: What does "traditional" means in RM?

1. The traditional RM system doesn't consider the choice model, in particular, it assumes the demands are independent random variables;
2. Further assumption: consumer will leave without purchasing if preferred fare class is unavailable (holds when gaps in fares are large enough);
3. The capacity is fixed, the capacity's marginal profit is zero(can be relaxed);
4. All booked consumers would arrive (another circumstance see 5.2.2).

Assumption 5.2.1.2: Single Resource RM

1. The units of capacity is c , pricing at multiple different level $p_n < \dots < p_1$;
2. Low-before-high fare class arrival order: D_2 before D_1 for example (this is the worst case for revenue);
3. **Protection level** for customer j : leave $y \in \{0, 1, \dots, c\}$ for D_{j-1}, \dots, D_1 ; $c - y$ is the **booking limit** which serves D_j ;

☺ So the problem is to solve the optimal protection level given the current consumer level j .

Let $V_j(x)$ be the optimal revenue given D_j coming in, x units remained. $V_0(x) = 0$ by design. Let y be the protection level for D_{j-1}, \dots, D_1 : sales at $p_j = \min\{x - y, D_j\}$. The remaining capacity for D_{j-1}, \dots, D_1 is $x - \min\{x - y, D_j\} = \max\{y, x - D_j\}$. Now let $W_j(y, x)$ be the optimal solution. We have:

$$W_j(y, x) = p_j \mathbb{E}\{\min\{x - y, D_j\}\} + \mathbb{E}\{V_{j-1}(\max\{y, x - D_j\})\} \quad (5.2.1.1)$$

$$V_j(x) = \max_{y \in \{0, \dots, x\}} W_j(y, x) = \max_{y \in \{0, \dots, x\}} \{p_j \mathbb{E}\{\min\{x - y, D_j\}\} + \mathbb{E}\{V_{j-1}(\max\{y, x - D_j\})\}\} \quad (5.2.1.2)$$

Proposition 5.2.1.1: Structure of the Optimal Policy

$$y_{j-1}^* = \max\{y \in \mathbb{N}_+ : \Delta V_{j-1}(y) > p_j\}. \quad (5.2.1.3)$$

The maximizer of $W_j(y, x)$ is given by y_j^*, y_1^*

Remark 5.2.1

The optimal solution for y_j is independent of the distribution of D_j ;

Corollary 5.2.1.1: When $j = 2$:

Theorem 5.2.1.1: Littlewood's rule

$$y_1^* = \max\{y \in \mathbb{N}_+ : \mathbb{P}\{D_1 \geq y\} > r\} \quad (5.2.1.4)$$

;

Remark 5.2.2

The Littlewood's Rule:

1. The solution depends on the **fare ratio**: $r := p_2/p_1$;
2. When the distribution of D_2 is continuous: $F_1(y) = \mathbb{P}\{D_1 \leq y\}$. The optimal protection level is $y_1^* = F_1^{-1}(1 - r) = \mu_1 + \sigma_1 \Phi^{-1}(1 - r)$:
 - (a) if $r > \frac{1}{2}$, $y_1^* < \mu_1$ and y_1^* decreases with σ_1 ;
 - (b) if $r < \frac{1}{2}$, $y_1^* < \mu_1$ and y_1^* increases with σ_1 ;
 - (c) if $r = \frac{1}{2}$, $y_1^* = \mu_1$;
3. Using Littlewood's rule would result some D_1 served by competitors (high spill rates).
Solution: add penalty to save more seats for the high fare consumers:

$$y_1^* = \max\left\{y \in \mathbb{N}_+ : \mathbb{P}\{D_1 \geq y\} > \frac{p_2}{p_1 + \rho}\right\} \quad (5.2.1.5)$$

5.2.2 Overbooking

5.2.3 Traditional Consumer Choice Model

5.2.4 Current Consumer Choice Model

5.3 Platform Operations Management

[Wang et al., 2023] uses game theory framework to analyze the cross-licensing policy initiated by Qualcomm, which provide some insights for the up-stream manufacturing company:

1. The supplier shouldn't adopt the cross-licensing policy if the inferior manufacturer's cost of innovation is high;
2. Cross-licensing may achieve a higher level of total innovation if the superior manufacturer's cost of innovation is low;
3. The superior manufacturer can benefit from cross-licensing, if innovation is costly but the manufacturers' costs of innovation are similar.

5.3.1 Platform Owner's Entry

Keywords 5.2

Complementary Markets,

[Zhu and Liu, 2018] studies the entry of Amazon platform. Logit regression is adopted to verify the following **Hypothesis**:

- Platform owners are more likely to compete with a complementor when its products are successful;
- Platform owners are less likely to compete with a complementor when its products require significant platform-specific investments to grow.

Identification Techniques:

- The sales ranking is used as proxy variable for the sales of the products;
- To overcome the impact of the referral rates by category-level fixed effects;
- To measure the seller's platform-specific investment, they calculate the seller's average answers;

[Shi et al., 2023] investigates the timing of the platform owner's entry on the value creation.

Definition 5.3.1.1: Timing of Owner's Entry

- **Platform Complementors:** actors that offer an application that brings additional value to platform users when used in combination with the platform;
- **Early-entry:** the entry occurs when the ratio between the current and the eventual

complementary market's size is low;

- **Late-entry**: entry to a relatively mature complementary market;
- **Value creation**: the activities geared toward increasing the perceived attractiveness of the platform ecosystem among customers and measure it as changes to complement popularity among customers. (proxy variable)

Identification Techniques:

- Use the *the number of reviews* as the proxy for the popularity of complements (dependent variables);
- *functional specificity* measures the heterogeneity of a complement based on the complexity of services offered by the complement;
- Follows [Zhu and Liu, 2018] to account for platform-specific investments by *interface coupling*: whether the complement connects with the platform core;
- To verify the exogeneity of Amazon's entry decision, a logit regression is conducted to test the number of reviews (popularity) does not influence Amazon's decision;
- PSM method is adopted.

$$Reviews_{it} = \alpha + \beta Treated_i \times After_t + \delta Controls_{it} + C_i + T_i + \epsilon_{it} \quad (5.3.1.1)$$

[Hagiu et al., 2020] examine the

5.3.2 Consumer Polarization

Consumer polarization is a topic in consumer research. **Group-Polarization Hypothesis** suggests that group discussion generally produces attitudes that are more extreme in the direction of the average of prediscussion attitudes in a variety of situations. Works like [Rao and Steckel, 1991] provides a mathematical presentation for this phenomenon (in the domain of preference):

$$U_s = \sum_{i=1}^m \lambda_i u_i + \phi(\bar{u} - K) \quad 0 \leq \lambda_i \leq 1, \sum_{i=1}^m \lambda_i = 1, \phi \geq 0 \quad (5.3.2.1)$$

In this model, the \bar{u} is the algebraic mean of all consumers' utility, and K is the **Pivot Point**. Rewrite this formula:

$$\begin{aligned}
 U_g &= \sum_{i=1}^m \left(\lambda_i + \frac{\phi}{m} \right) u_i - \phi K & \begin{aligned} w_0 &\leq 0; \\ \sum_{i=1}^m w_i &\geq 1; \end{aligned} \\
 &= w_0 + w_1 u_1 + w_2 u_2 + \cdots + w_m u_m & 0 \leq \frac{w_0}{1 - \sum w_i} \leq 1.
 \end{aligned} \tag{5.3.2.2}$$

Remark 5.3.1

- [Zhao et al., 2023] use experiment results to suggest that eWOM (electronic word of mouth) polarization (the degree of eWOM to which positive and negative sentiments are simultaneously strong) would decrease the consumers' intention to purchase, mediating by the enhancement of attitude ambivalence. (Ambivalence is a psychological state where a person endorses both positive and negative attitudinal positions)
- [Iyer and Yoganasimhan, 2021] use game theory framework, to get the conclusion that sequential decision making could reduce the polarization.

5.3.3 Network Effect

Network Effect and **Network Externality**:

[Narayan et al., 2011] verifies that peer influence affects attribute preferences via a bayesian updating mechanism. In their model, the utility is given as follows:

$$U_{ijp}^R = X_{jp} \beta_i^R + \lambda_i \varepsilon_{ijp}^R \tag{5.3.3.1}$$

Where U_{ijp} is the utility of consumer i for product j given choice set p , X_{jp} is the attribute of product j in the choice set p , β_i^R is the customers' weights. The bayesian updating process is given below:

$$\begin{aligned}
 \beta_{ik}^R &= \rho_{ik} \beta_{ik}^l + (1 - \rho_{ik}) \frac{\sum_{i=1, i \neq i}^N w_{ii} \beta_{ik}^l}{\max \left[\left(\sum_{i=1, i \neq i}^N w_{ii} \right), 1 \right]}, \\
 &\text{where } 0 \leq \rho_{ik} \leq 1.
 \end{aligned} \tag{5.3.3.2}$$

Other research on peer-influence:

Remark 5.3.2

The consumers' interaction and social connections have a proposition proposed in [Zhang et al., 2017] for their goal attainment and spending: a positive linear term plus a negative squared term;

5.3.4 Online Gaming

Many industrial news about online gaming can be found in [Chen et al., 2017]. In [Lei, 2022], the dissertation fully discussed loot box pricing, matchmaking, and price discrimination with fairness constraints.

Play-Duration and Spending

[Zhang et al., 2017]'s work shows that there is a nonlinear effect of social connections and interactions on consumers' goal attainment and spending: A positive linear terms and a negative squared term. Mechanism: functional in providing useful information or tips that can facilitate goal attainment, but would raise information overload problem.

Player engagement can be embodied by many specific metrics, such as time or money spent in the game, the number of matches played within a time window, or churn risk. [Chen et al., 2017] define churn risk as the proportion of total players stopping playing the game over a period of time.

Matchmaking

Matchmaking connects multiple players to participate in online PvP games. (PvP(Player-versus-Player) games, which cover many popular genres, such as multiplayer online battle arena (MOBA), first-person shooting (FPS), and e-sports, have increased worldwide popularity in recent years.)

The past matchmaking strategy matches similar skilled players in the same round (SBMM), the current MM system focuses on improving the players' engagement and decreasing the churn rate. For example, in [Chen et al., 2017] EOMM (Engagement Optimization MatchMaking) is proposed to minimize the churn rate.

[Chen et al., 2021] propose an algorithm to maximize the cumulative active players.

Assumption 5.3.4.1: Chen 2021 MatchMaking

1. players can have heterogeneous skill levels: level 1 to level K ;
2. the outcome of each match is a Bernoulli random variable: $p_{kj} = 1 - p_{jk}$, $p_{kk} = 0.5$, $p_{kj} > 0.5$ if $k > j$;
3. player's skill level is fixed: *relative* level;
4. and their state depends on the win-loss outcomes of the past m matches: $g \in \mathcal{G}$ ($2^m + 1$ possible cardinality);
5. A geometric losing churn model: players churn with a fixed probability, starting from the second loss in a row;
6. $P_{win}^k, P_{lose}^k \in [0, 1]^{|\mathcal{G}| \times |\mathcal{G}|}$ is the transition matrix of level k player's engagement state;
7. $M_{kj} = p_{kj}P_{win}^k + (1 - p_{kj})P_{lose}^k$ is the aggregate transition matrix. (\bar{G} is the reduced aggregate transition matrix);
8. using the fluid matching model and assume players are infinitely divisible;

The **Dynamic Programming** formulation: $f_{kg,jg'}$ is the amount of kg players matched with jg' players, s_{kg}^t is the number of kg players at time t .

FB flow balance constraints:

$$\begin{aligned}
 & \sum_{j=1}^K \sum_{g' \in \bar{\mathcal{G}}} f_{kg,jg'}^t = s_{kg}^t, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, \\
 & \sum_{j=1}^K \sum_{g' \in \bar{\mathcal{G}}} f_{jg',kg}^t = s_{kg}^t, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, \\
 & f_{kg,jg'}^t = f_{jg',kg}^t, j = 1, \dots, K, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, g' \in \bar{\mathcal{G}} \\
 & f_{kg,jg'}^t \geq 0, j = 1, \dots, K, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, g' \in \bar{\mathcal{G}}
 \end{aligned} \tag{5.3.4.1}$$

ED evolution of demographics:

$$\mathbf{s}_k^{t+1} = \sum_{j=1, \dots, K} \left(\mathbf{f}_{kj}^t \mathbf{1} \right)^\top (\bar{M}_{kj} + N_k) \quad k = 1, \dots, K \tag{5.3.4.2}$$

The value-to-go function is:

$$V^\pi(\mathbf{s}^t) = \sum_{k=1}^K \sum_{g \in \bar{\mathcal{G}}} s_{kg}^{t+1} + \gamma V^\pi(\mathbf{s}^{t+1}) \quad (5.3.4.3)$$

subject to (FB), (ED).

The above model can be formulated in a linear programming style:

Theorem 5.3.4.1: Chen 2021 MM LP Formulation

$$\begin{aligned} V^*(\mathbf{s}^0) &= \max \sum_{t=1}^{\infty} \gamma^{t-1} \sum_k \sum_{g \in \bar{\mathcal{G}}} s_{kg}^t \\ \text{s.t. } &\sum_{j=1}^K \sum_{g' \in \bar{\mathcal{G}}} f_{kg,jg'}^t = s_{kg}^t, \forall k, \forall g \in \bar{\mathcal{G}}, t = 0, 1, \dots \\ &\sum_{j=1}^K \sum_{g' \in \bar{\mathcal{G}}} f_{jg',kg}^t = s_{kg}^t, \forall k, \forall g \in \bar{\mathcal{G}}, t = 0, 1, \dots \\ &f_{kg,jg'}^t = f_{jg',kg}^t, j = 1, \dots, K, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, g' \in \bar{\mathcal{G}}, t = 0, 1, \dots \\ &f_{kg,jg'}^t \geq 0, j = 1, \dots, K, k = 1, \dots, K, \forall g \in \bar{\mathcal{G}}, g' \in \bar{\mathcal{G}}, t = 0, 1, \dots \\ &\mathbf{s}_k^{t+1} = \sum_{j=1, \dots, K} \left(\mathbf{f}_{kj}^t \mathbf{1} \right)^\top (\bar{M}_{kj} + N_k), \forall k, t = 0, 1, \dots \end{aligned} \quad (5.3.4.4)$$

Remark 5.3.3

- Using an optimal matchmaking policy instead of SBMM may reduce the required bot ratio significantly while maintaining the same level of engagement.

5.4 Behavioral Operations Management

5.5 Data-Driven Operations Management

Chapter 6

Miscellaneous

6.1 Notes on Tools

6.1.1 LaTeX Shortcuts

There are 6x6 colors in the preset preamble:

aa	ab	ac	ad	ae	af
ba	bb	bc	bd	be	bf
ca	cb	cc	cd	ce	cf
da	db	dc	dd	de	df
ea	eb	ec	ed	ee	ef
fa	fb	fc	fd	fe	ff

Using `\href{URL}{text}` to refer a [website](#).

Using `\eq` to write equation, `\tab` to get an unordered list, `\lis` to get an ordered list.

$$E = mc^2 \tag{6.1.1.1}$$

- item 1;
- item 2;
- item 3.

1. item 1;
2. item 2;

3. item 3.

Format of Words

```
\hl{highlighted},
\u{underlined},
\st{striketrough}\
\rt{red}, \yt{yellow},
\bt{blue}, \gt{green}
```

highlighted, underlined, ~~striketrough~~
red, yellow, blue, green

Shortcuts

```
\RR, \NN, \ZZ, \QQ\
\bA, \bB, \bC, \bD
```

\mathbb{R} , \mathbb{N} , \mathbb{Z} , \mathbb{Q}
 \mathbb{A} , \mathbb{B} , \mathbb{C} , \mathbb{D}

Shortcuts

```
\tbf{Text}, \tit{Text}\
\cA, \cB, \cC, \cD
```

Text, *Text*
A, B, C, D

Emoji

```
\emogood, \emobad, \emocool,
\emoheart, \emotree
```

😊, 😞, 😴, ❤️, 🌳

Use `\ass`, `\ax`, `\thm`, `\co`, `\pro`, `\defi`, `\re`, `\key`, `\ex`, `\proo` to use preset `tcolorboxes` template.

Assumption 6.1.1.1: Example

1. item 1;
2. item 2;
3. item 3.

Axiom 6.1.1.1: Exmaple

Test

Theorem 6.1.1.1: Exmaple

Test

Corollary 6.1.1.1: Exmaple

Test

Proposition 6.1.1.1: Exmaple

Test

Definition 6.1.1.1: Exmaple

Test

Remark 6.1.1

Expample

Keywords 6.1

Example

Example 6.1.1.1

Problem example

Solution:

The solution is omitted.

Proof 6.1.1: proposition 6.1.1

The Formal Proof

Q.E.D.

Using `\label`, `\ref` to refer to chapters 6, sections 6.1, equations 6.1.1, and boxes 6.1.1.

Using `\cite` to cite the literature in apa style. For example: [Klein et al., 2020]

Using `\sep` to insert a horizontal line with words in the middle:

Compilation

Clearning all the auxiliary files: LaTeXmk \rightarrow BibTeX \rightarrow LaTeXmk \rightarrow LaTeXmk. Or, zip the main files and upload to Overleaf.

Put photos in the *pic* file and use `\fig` to show it.

$$h(x) = \binom{n}{x}, \quad c(p) = (1-p)^n, \quad t(x) = x \quad \text{and} \quad w(p) = \log\left(\frac{p}{1-p}\right).$$

Then,

$$\begin{aligned} \frac{d}{dp} w(p) &= \frac{d}{dp} \log\left(\frac{p}{1-p}\right) = \frac{1}{p(1-p)}, \\ \frac{d^2}{dp^2} w(p) &= -\frac{1}{p^2} + \frac{1}{(1-p)^2} = \frac{2p-1}{p^2(1-p)^2}, \\ \frac{d}{dp} \log(c(p)) &= \frac{d}{dp} n \log(1-p) = -\frac{n}{1-p}, \\ \frac{d^2}{dp^2} \log(c(p)) &= -\frac{n}{(1-p)^2}. \end{aligned}$$

Therefore, from Theorem 3.4.2, we have

$$\begin{aligned} \mathbb{E}\left(\frac{1}{p(1-p)}X\right) &= \frac{n}{1-p} \Rightarrow \mathbb{E}(X) = np, \\ \text{Var}\left(\frac{1}{p(1-p)}X\right) &= \frac{n}{(1-p)^2} - \mathbb{E}\left(\frac{2p-1}{p^2(1-p)^2}X\right) \Rightarrow \text{Var}(X) = np(1-p). \end{aligned}$$

Figure 6.1: Example.png

6.2 Important Proofs

6.3 Beautiful Phrases

Introduction

🌱 To the limit of our knowledge, this study is among the first to unveil ...

PUBG Paper

😄 The bot can be designed so that it is competitive but still loses to the human player, which may result in the human breaking their losing streak and remaining in the system longer. On the other hand, due to the limitations of technology, AI-powered bots can be identified by experienced players. If a human player is frequently matched with bots, they may find out that their opponents are not human and perhaps be discouraged from playing the game.

6.4 Eureka Ideas

6.4.1 RNN and Causal Model

On 1st Aug 2023.

The RNN structure and rubin's causal model under a changing environment looks similar;

To-do:

1. Inspect the papers containing keywords "RNN" and "Causal";
2. Explorations of Causal Models in Bayesian Causal Framework, with a Focus Beyond Rubin's Work;
3. explore the further connection between RNN and causal model.

6.4.2 Earthquakes on Immigration and Investment

On 19th Aug 2023.

Since 2019, numerous enterprises have engaged in shale gas extraction across various cities and counties in southern Sichuan, leading to frequent yet non-hazardous seismic activities. The geographical locations subjected to shale gas extraction remain independent of the local economic conditions, and the introduction of these enterprises has shown no discernible positive impact on the fiscal health of the local governments or the regional employment scenario. Essentially, this situation represents a discontinuity, wherein the extraction of shale gas corresponds to a quasi-random selection of regions experiencing seismic events. This phenomenon can be effectively studied using the Differences-in-Differences (DiD) methodology to investigate the causal effects of this unstable geological activity on the migration rates of local residents and the influx of external capital.

To Do:

1. Reviewing Literature to Investigate the Independence of Shale Gas from Local Economic Levels;
2. Assessing the Accessibility of Relevant Data for Investigation.

List of Figures

1.1	Type of distribution	12
1.2	A network with capacities	15
2.1	The scope of statistical inference	23
3.1	Clan Culture Intensity	44
4.1	Cloud Computing Technology Convergence	48
4.2	From HPC systems and clusters to grids, p2p networks, clouds and IoT	49
4.3	Comparing three cloud service models with on-premise computing	49
4.4	physical cluster and virtual cluster	50
4.5	CPU and GPU Architecture	52
4.6	VMM Operations	53
4.7	The agent–environment interaction in a Markov decision process	58
4.8	Classification of different reinforcement learning agents	60
5.1	identification strategies from 2016 to 2021	63
5.2	How identification techniques works	64
5.3	Navigating Matrix Cells	65
6.1	Example.png	80

Bibliography

- [Angrist and Pischke, 2009] Angrist, J. D. and Pischke, J.-S. (2009). *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton university press.
- [Angrist and Pischke, 2014] Angrist, J. D. and Pischke, J.-S. (2014). *Mastering'metrics: The Path from Cause to Effect*. Princeton university press.
- [Bertrand and Schoar, 2006] Bertrand, M. and Schoar, A. (2006). The role of family in family firms. *Journal of economic perspectives*, 20(2):73–96.
- [Bertsimas and Tsitsiklis, 1997] Bertsimas, D. and Tsitsiklis, J. N. (1997). *Introduction to Linear Optimization*, volume 6. Athena scientific Belmont, MA.
- [Brusco et al., 2017] Brusco, M. J., Singh, R., Cradit, J. D., and Steinley, D. (2017). Cluster analysis in empirical OM research: Survey and recommendations. *International Journal of Operations & Production Management*, 37(3):300–320.
- [Cao et al., 2022] Cao, J., Xu, Y., and Zhang, C. (2022). Clans and calamity: How social capital saved lives during China's Great Famine. *Journal of Development Economics*, 157:102865.
- [Casella and Berger, 2021] Casella, G. and Berger, R. L. (2021). *Statistical Inference*. Cengage Learning.
- [Chen et al., 2021] Chen, M., Elmachoub, A. N., and Lei, X. (2021). Matchmaking strategies for maximizing player engagement in video games. *Available at SSRN 3928966*.
- [Chen et al., 2017] Chen, Z., Xue, S., Kolen, J., Aghdaie, N., Zaman, K. A., Sun, Y., and Seif El-Nasr, M. (2017). EOMM: An Engagement Optimized Matchmaking Framework. In *Proceedings of the 26th International Conference on World Wide Web*, pages 1143–1150, Perth Australia. International World Wide Web Conferences Steering Committee.
- [Cheng et al., 2021] Cheng, J., Dai, Y., Lin, S., and Ye, H. (2021). Clan culture and family ownership concentration: Evidence from China. *China Economic Review*, 70:101692.
- [Choi et al., 2016] Choi, T.-M., Cheng, TCE., and Zhao, X. (2016). Multi-methodological research in operations management. *Production and Operations Management*, 25(3):379–389.

- [Fisher, 2007] Fisher, M. (2007). Strengthening the empirical base of operations management. *Manufacturing & Service Operations Management*, 9(4):368–382.
- [Fisher and Raman, 2022] Fisher, M. and Raman, A. (2022). Innovations in retail operations: Thirty years of lessons from Production and Operations Management. *Production and Operations Management*, 31(12):4452–4461.
- [Gallego and Topaloglu, 2019] Gallego, G. and Topaloglu, H. (2019). *Revenue Management and Pricing Analytics*, volume 279 of *International Series in Operations Research & Management Science*. Springer New York, New York, NY.
- [Gans et al., 2007] Gans, N., Knox, G., and Croson, R. (2007). Simple models of discrete choice and their performance in bandit experiments. *Manufacturing & Service Operations Management*, 9(4):383–408.
- [Hagiu et al., 2020] Hagiu, A., Jullien, B., and Wright, J. (2020). Creating platforms by hosting rivals. *Management Science*, 66(7):3234–3248.
- [Hwang, 2017] Hwang, K. (2017). *Cloud Computing for Machine Learning and Cognitive Applications*. Mit Press.
- [Iyer and Yoganarasimhan, 2021] Iyer, G. and Yoganarasimhan, H. (2021). Strategic Polarization in Group Interactions. *Journal of Marketing Research*, 58(4):782–800.
- [Klein et al., 2020] Klein, R., Koch, S., Steinhardt, C., and Strauss, A. K. (2020). A review of revenue management: Recent generalizations and advances in industry applications. *European Journal of Operational Research*, 284(2):397–412.
- [Kumar and Tang, 2022] Kumar, S. and Tang, C. S. (2022). Expanding the boundaries of the discipline: The 30th-anniversary issue of Production and Operations Management. *Production and Operations Management*, 31(12):4257–4261.
- [Lei, 2022] Lei, X. (2022). *Revenue Management in Video Games and with Fairness*. PhD thesis, Columbia University.
- [Luenberger et al., 1984] Luenberger, D. G., Ye, Y., et al. (1984). *Linear and Nonlinear Programming*, volume 2. Springer.
- [Mithas et al., 2022] Mithas, S., Chen, Y., Lin, Y., and De Oliveira Silveira, A. (2022). On the causality and plausibility of treatment effects in operations management research. *Production and Operations Management*, 31(12):4558–4571.

- [Narayan et al., 2011] Narayan, V., Rao, V. R., and Saunders, C. (2011). How Peer Influence Affects Attribute Preferences: A Bayesian Updating Mechanism. *Marketing Science*, 30(2):368–384.
- [Rao and Steckel, 1991] Rao, V. R. and Steckel, J. H. (1991). A polarization model for describing group preferences. *Journal of Consumer Research*, 18(1):108–118.
- [Roth and Singhal, 2022] Roth, A. M. and Singhal, V. R. (2022). Pioneering role of the Production and Operations Management in promoting empirical research in operations management. *Production and Operations Management*, 31(12):4529–4543.
- [Roth, 2007] Roth, A. V. (2007). Applications of empirical science in manufacturing and service operations. *Manufacturing & Service Operations Management*, 9(4):353–367.
- [Shi et al., 2023] Shi, R., Aaltonen, A., Henfridsson, O., and Gopal, R. D. (2023). Comparing platform owners’ early and late entry into complementary markets. *MIS Quarterly*.
- [Strauss et al., 2018] Strauss, A. K., Klein, R., and Steinhardt, C. (2018). A review of choice-based revenue management: Theory and methods. *European Journal of Operational Research*, 271(2):375–387.
- [Thrun and Littman, 2000] Thrun, S. and Littman, M. L. (2000). Reinforcement learning: An introduction. *AI Magazine*, 21(1):103–103.
- [Wang et al., 2023] Wang, J., Huang, T., and Lee, J. (2023). Cross-licensing in a Supply Chain with Asymmetric Manufacturers.
- [Yang, 2019] Yang, H. (2019). Family clans and public goods: Evidence from the New Village beautification project in South Korea. *Journal of Development Economics*, 136:34–50.
- [Zhang, 2019] Zhang, C. (2019). Family support or social support? The role of clan culture. *Journal of Population Economics*, 32:529–549.
- [Zhang, 2020] Zhang, C. (2020). Clans, entrepreneurship, and development of the private sector in China. *Journal of Comparative Economics*, 48(1):100–123.
- [Zhang et al., 2017] Zhang, C., Phang, C. W., Wu, Q., and Luo, X. (2017). Nonlinear Effects of Social Connections and Interactions on Individual Goal Attainment and Spending: Evidences from Online Gaming Markets. *Journal of Marketing*, 81(6):132–155.
- [Zhao et al., 2023] Zhao, P., Ma, Z., Gill, T., and Ranaweera, C. (2023). Social media sentiment polarization and its impact on product adoption. *Marketing Letters*.

[Zhu and Liu, 2018] Zhu, F. and Liu, Q. (2018). Competing with complementors: An empirical look at Amazon. com. *Strategic management journal*, 39(10):2618–2642.