

# Cookbook

## Recipe

<b>1</b>	<b>Random Variable Control</b>	<b>5</b>
1.1	Maximum of i.i.d gaussian . . . . .	5
1.2	Union bound for partial sums . . . . .	5
1.3	Random Singed Summation Bound . . . . .	5
1.4	Summation Bound . . . . .	5
1.5	Concentraion ineqaulity . . . . .	6
1.5.1	Bernstein's ineqaulity . . . . .	6
1.5.2	Matrix version Bernstein's inequality . . . . .	6
1.5.3	Hoeffding's inequality . . . . .	6
1.5.4	Paley-Zygmund inequality . . . . .	7
1.5.5	Max of independent Gaussians . . . . .	7
1.5.6	DKW inequality . . . . .	7
1.5.7	Quadratic form . . . . .	7
1.5.8	Gaussian CCDF bound . . . . .	8
1.5.9	McDiarmid's inequality . . . . .	8
1.6	Decoupling lemma . . . . .	8
1.6.1	Quadratic form . . . . .	8
1.7	Global variance control . . . . .	8
1.7.1	Efron-Stein inequality . . . . .	8
1.8	Information . . . . .	9
1.8.1	Fano's inequality . . . . .	9
1.9	Do you like martingale? . . . . .	9
1.9.1	Tail Distribution . . . . .	9
1.10	Stochastic Dominance . . . . .	9
1.10.1	SD is preserved under sums/convolutions . . . . .	9
<b>2</b>	<b>Privacy</b>	<b>9</b>
2.1	Privacy Loss . . . . .	9
2.2	alternative . . . . .	9
2.2.1	alternative . . . . .	9
2.2.2	Moment difference bound . . . . .	10
2.3	zCDP . . . . .	10
2.3.1	Key properties . . . . .	10
2.4	Approximate Rényi Differential Privacy . . . . .	10
2.4.1	Properties . . . . .	12
2.5	Composition . . . . .	12
2.5.1	alternative . . . . .	12
2.6	Joint Differential Privacy . . . . .	12
2.7	Lower bound tools . . . . .	12
2.7.1	Query moment w.r.t binary data . . . . .	12
2.8	Decomposition . . . . .	13
2.8.1	Basic decomposition . . . . .	13
2.8.2	Bayesian version [2, Lemma 5.6] . . . . .	13
<b>3</b>	<b>Probability and Statistics</b>	<b>13</b>
3.1	Moment . . . . .	13
3.1.1	Moment and CDF . . . . .	13
3.1.2	alternative . . . . .	14
3.1.3	Clipped random variable . . . . .	14
3.2	Conditional Distribution . . . . .	14

3.3	Weak convergence . . . . .	15
3.3.1	Delta method . . . . .	15
3.3.2	Slutsky's theorem . . . . .	15
3.4	Central limit theorem . . . . .	15
3.4.1	Salem-Zygmund . . . . .	15
3.5	Conditional independence . . . . .	15
3.5.1	Bayes' Theorem . . . . .	16
3.5.2	Conditional expectation under independence . . . . .	16
3.6	Regular conditional distribution . . . . .	16
3.6.1	Disintegration formula . . . . .	16
3.7	Markov kernel . . . . .	17
3.8	Mutual Information . . . . .	17
3.8.1	Concavity of Mutual information . . . . .	17
3.9	Fisher Information . . . . .	17
3.9.1	Cramér–Rao bound . . . . .	17
3.10	MLE estimation . . . . .	18
3.10.1	Asymptotic normality . . . . .	18
3.11	Famous Family . . . . .	18
3.11.1	Poission . . . . .	18
3.11.2	Binomial . . . . .	18
3.11.3	Negative Binomial . . . . .	18
3.12	M-estimator . . . . .	18
3.12.1	Consistency . . . . .	18
<b>4</b>	<b>Convexity</b> . . . . .	<b>19</b>
4.1	Characterization . . . . .	19
4.1.1	Characterization of Strict Convexity . . . . .	19
4.1.2	When can we assume equal variables? . . . . .	19
4.2	Dual Problem . . . . .	19
<b>5</b>	<b>Integral Technique</b> . . . . .	<b>20</b>
5.1	Gaussian Integral . . . . .	20
5.1.1	Gauss-Hermite quadrature . . . . .	20
5.1.2	Stein's lemma . . . . .	20
5.1.3	Change of measure . . . . .	20
<b>6</b>	<b>Stochastic process</b> . . . . .	<b>20</b>
6.1	Predictable process . . . . .	20
6.1.1	How to understand the predictable process? . . . . .	20
6.2	Local martingale . . . . .	20
6.2.1	Quadratic variation . . . . .	20
6.2.2	Stochastic integral . . . . .	20
6.3	Doob's h transform . . . . .	21
<b>7</b>	<b>Fundamental Algebra</b> . . . . .	<b>21</b>
7.1	Series . . . . .	21
7.1.1	Geometric . . . . .	21
7.1.2	Quotient Stack . . . . .	21
7.2	Inequalities . . . . .	21
7.2.1	Ratio of Summation . . . . .	21
7.2.2	Weighted AM-GM Inequality . . . . .	21
7.2.3	Titu's Lemma . . . . .	21
7.2.4	Cauchy-Schwarz Inequality . . . . .	22
7.2.5	Chebyshev's Sum Inequality . . . . .	22
7.2.6	Symmetric Parametric Inequality . . . . .	22

7.2.7	We love Jensen . . . . .	22
7.3	Bounds and Approximations . . . . .	22
7.3.1	Exponential Bound on Hyperbolic Ratio . . . . .	22
7.3.2	Exponential Inequalities . . . . .	22
7.3.3	Softplus Quadratic Bound . . . . .	22
7.3.4	Miscellaneous upper Bounds . . . . .	22
7.3.5	Order of Rademacher Sums . . . . .	23
7.3.6	Sum Approximation . . . . .	23
7.4	Combinatorics . . . . .	23
7.4.1	Expansion . . . . .	23
7.4.2	Vandermonde's Identity . . . . .	23
<b>8</b>	<b>Useful real analysis results</b>	<b>23</b>
8.1	Leibniz's Integral Rule . . . . .	23
<b>9</b>	<b>Linear Algebra</b>	<b>24</b>
9.1	Determinant . . . . .	24
9.1.1	Expansion formula . . . . .	24
9.1.2	Rearrangement . . . . .	24
9.1.3	Weinstein-Aronszajn identity . . . . .	24
9.1.4	DPP related . . . . .	24
9.1.5	Ratio . . . . .	25
9.1.6	Inverse of trace . . . . .	25
9.2	Vectorization . . . . .	25
9.3	Trace . . . . .	25
9.3.1	Von Neumann's trace inequality . . . . .	25
9.4	Inversion . . . . .	25
9.4.1	Woodbury identity . . . . .	25
9.4.2	Schur Complement . . . . .	25
9.5	Hadamard Product . . . . .	26
9.5.1	Quadratic Relation . . . . .	26
9.5.2	Rank Relation . . . . .	26
9.5.3	Spectrum Relation . . . . .	26
9.5.4	Determinant . . . . .	26
9.6	Matrix Calculus . . . . .	26
9.6.1	Matrix Chain rule . . . . .	26
9.6.2	Differentials . . . . .	26
9.6.3	Useful first derivatives . . . . .	27
9.6.4	Quadratic form . . . . .	27
9.6.5	Hessian product rule . . . . .	27
9.6.6	Integration by parts . . . . .	27
9.7	Eigenvalues and Eigenvectors . . . . .	27
9.7.1	General Properties . . . . .	27
9.7.2	Symmetric . . . . .	28
9.7.3	Singular Value Decomposition . . . . .	28
9.7.4	LU decomposition . . . . .	28
9.7.5	Cholesky decomposition . . . . .	28
9.7.6	Eigenvalues of its reverse . . . . .	28
9.7.7	Row stochastic matrix . . . . .	28
9.8	Inverses . . . . .	29
9.8.1	Rank-1 update of the inverse of inner product . . . . .	29
9.8.2	Approximations . . . . .	29
9.8.3	Block matrix . . . . .	29
9.9	PSD matrix . . . . .	29

9.9.1	Decomposition . . . . .	29
9.9.2	Sylvester's characterization . . . . .	29
9.9.3	Equation with zeros . . . . .	29
9.9.4	Rank of product . . . . .	30
9.9.5	Outer product . . . . .	30
9.9.6	Small perturbations . . . . .	30
9.9.7	Hadamard inequality . . . . .	30
9.9.8	Loewner order . . . . .	30
9.9.9	Inverse of PSD . . . . .	30
9.10	Symmetric and skew-symmetric matrix . . . . .	30
9.10.1	Properties of symmetric matrix . . . . .	30
9.10.2	Youla decomposition . . . . .	30
9.11	Some techniques . . . . .	31
9.11.1	Binary analysis . . . . .	31

## Open problem

### Tail of hypergeometric distribution

If we estimate the maximum from the sampled subset of the population, what is the tail probability of the cardinality of a set whose value is greater than the estimated maximum?

2025-06-17

## 1. Random Variable Control

### 1.1 Maximum of i.i.d gaussian

Let  $\xi_1, \dots, \xi_k$  be  $k$  independent samples from  $\mathcal{N}(0, 1)$ . Then

$$\mathbb{E} [\max \{\xi_1^2, \dots, \xi_k^2\}] \leq 2 \log(2k) \quad (1.1)$$

### 1.2 Union bound for partial sums

**Etemadi's inequality.** Let  $X_1, \dots, X_n$  be independent random variables. For  $i \in [n]$ , let  $Y_i = \sum_{j=1}^i X_j$  denote the partial sum up to  $i$ . Then for all  $\alpha \geq 0$ ,

$$\Pr[\max_{i=1}^n |Y_i| > 3 \cdot \alpha] \leq 3 \cdot \max_{i=1}^n \Pr[|Y_i| > \alpha]. \quad (1.2)$$

*Proof Sketch.*  $\mathbb{P}[|Y_i| > \alpha]$  term을 얻기 위해서  $|Y_i - Y_n|$ 과  $|Y_i|$  사이의 independence를 사용함. 그리고 partial sum의 maximum과 각 partial sum을 연결하기 위해서  $i$ 번 째 partial sum이 처음으로  $3\alpha$  보다 큰 event로 분해함. (Detail)

### 1.3 Random Singed Summation Bound

**Khinchine inequalilty.** Let  $\{\varepsilon_n\}_{n=1}^N$  be i.i.d. Rademacher random variables. Let  $0 < p < \infty$  and let  $x_1, \dots, x_N \in \mathbb{C}$ . Then

$$A_p \left( \sum_{n=1}^N |x_n|^2 \right)^{1/2} \leq \left( \mathbb{E} \left| \sum_{n=1}^N \varepsilon_n x_n \right|^p \right)^{1/p} \leq B_p \left( \sum_{n=1}^N |x_n|^2 \right)^{1/2} \quad (1.3)$$

### 1.4 Summation Bound

**Marcinkiewicz-Zygmund inequalilty.** If  $X_i, i = 1, \dots, n$  are independent random variables with  $\mathbb{E}[X_i] = 0$  and  $\mathbb{E}[|X_i|^p], 1 < p < +\infty$ , then

$$A_p \mathbb{E} \left[ \left( \sum_{i=1}^n |X_i|^2 \right)^{p/2} \right] \leq \mathbb{E} \left[ \left| \sum_{i=1}^n X_i \right|^p \right] \leq B_p \mathbb{E} \left[ \left( \sum_{i=1}^n |X_i|^2 \right)^{p/2} \right] \quad (1.4)$$

where  $A_p$  and  $B_p$  are positive constants, which depend only on  $p$ . for some constants  $A_p, B_p$  depending only on  $p$ .

**Latala's inequality.** If  $p \geq 2$  and  $X, X_1, \dots, X_n$  are i.i.d. mean 0 random variables, then we have

$$\left\| \sum_{i=1}^n X_i \right\|_{L^p} \sim \sup \left\{ \frac{p}{s} \left( \frac{n}{p} \right)^{1/s} \|X\|_{L^s} \mid \max \left\{ 2, \frac{p}{n} \right\} \leq s \leq p \right\} \quad (1.5)$$

## 1.5 Concentraion ineqauality

### 1.5.1 Bernstein's ineqauality

Let  $X_1, \dots, X_n$  be independent random variables. Assume  $\mathbb{E}[X_i] = 0, \mathbb{E}[X_i^2] = \sigma_i^2$ , and  $\Pr[|X_i| \leq 1] = 1$  for each  $i \in [n]$ . Let  $\sigma^2 := \sum_{i=1}^n \sigma_i^2$ . Then for all  $t \geq 0$ ,

$$\Pr \left[ \sum_{i=1}^n X_i \geq t \right] \leq \exp \left( \frac{-3t^2}{6\sigma^2 + 2t} \right) \quad (1.6)$$

*Proof Sketch.* First bound the MGF of each  $X_i$  using taylor expansion. Then use Markov inequality for  $\Pr [\sum_{i=1}^n X_i \geq t]$  with  $\exp(\lambda \cdot)$  and minimize the upper bound with  $\lambda$ . Then use the following lemma to finish the proof.

**Lemma.** Let  $v > -1$ . Then  $(1+v) \log(1+v) \geq v + \frac{3v^2}{2v+6}$

### 1.5.2 Matrix version Bernstein's inequality

Let  $\mathbf{B}$  a fixed  $q \times d$  matrix. Construct  $q \times d$  matrix  $\mathbf{R}$  such that

$$\mathbb{E}[\mathbf{R}] = \mathbf{B}, \quad \|\mathbf{R}\|_{\text{op}} \leq L \quad (1.7)$$

Form the matrix sampling estimator

$$\bar{\mathbf{R}}_m = \frac{1}{m} \sum_{i=1}^m \mathbf{R}_i, \quad (1.8)$$

where each  $\mathbf{R}_i$  is an independent copy of  $\mathbf{R}$ . Then for every  $t > 0$ , the estimator satisfies

$$\mathbb{P} \left[ \|\bar{\mathbf{R}}_m - \mathbf{B}\|_{\text{op}} \geq t \right] \leq (q+d) \cdot \exp \left( \frac{-mt^2}{m_2(\mathbf{R}) + 2Lt/3} \right), \quad (1.9)$$

where  $m_2(R)$  is the second moment  $m_2(\mathbf{R}) = \max \left\{ \|\mathbb{E}[\mathbf{R}^* \mathbf{R}]\|_{\text{op}}, \|\mathbb{E}[\mathbf{R} \mathbf{R}^*]\|_{\text{op}} \right\}$ .

### 1.5.3 Hoeffding's inequality

Let  $X_1, \dots, X_n$  be independent random variables such that  $a_i \leq X_i \leq b_i$  a.s. Then for all  $t > 0$ ,

$$\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) \leq \exp \left( -\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right), \quad (1.10)$$

where  $S_n = X_1 + \dots + X_n$ . Also, consider a set of  $r$  i.i.d. random variables  $X_1, \dots, X_r$  such that  $-\Delta \leq X_i \leq \Delta$  and  $\mathbb{E}[X_i] = 0$  for each  $i \in [r]$ . Let  $\sum_{i=1}^r X_i$ . Then for any  $\alpha \in (0, 1/2)$

$$\mathbb{P}[|M| > \alpha] \leq 2 \exp \left( -\frac{\alpha^2}{2r\Delta^2} \right) \quad (1.11)$$

The proof uses the following:

**Lemma.** Let  $X$  be any real-valued random variable such that  $a \leq X \leq b$  a.s. Then. for all  $\lambda \in \mathbb{R}$ ,

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \exp \left( \frac{\lambda^2(b-a)^2}{8} \right) \quad (1.12)$$

**COOL FACT:** conditional expectation also works for the lemma.

**Variance-only form.** Consider a set of  $r$  independent random variables  $X_1, \dots, X_r$ . Let  $M = \sum_{i=1}^r X_i$ . Then for  $\alpha \in (0, 2\text{Var}[M]/(\max_i |X_i - \mathbb{E}[X_i]|))$

$$\mathbb{P}[|M - \mathbb{E}[M]| > \alpha] \leq 2 \exp \left( \frac{-\alpha^2}{4 \sum_{i=1}^r \text{Var}[X_i]} \right). \quad (1.13)$$

#### 1.5.4 Paley-Zygmund inequality

If  $Z \geq 0$  is a random variable with finite variance, and if  $0 \leq \theta \leq 1$ , then

$$\mathbb{P}(Z > \theta \mathbb{E}[Z]) \geq (1 - \theta)^2 \frac{\mathbb{E}[Z]^2}{\mathbb{E}[Z^2]} \quad (1.14)$$

#### 1.5.5 Max of independent Gaussians

Let  $X_1, X_2, \dots, X_n$  i.i.d.  $\mathcal{N}(0, 1)$ , then

$$\mathbb{E}[\max(X_1, \dots, X_n)] = \sqrt{2 \log(n)} + o(\sqrt{\log(n)}) \quad (1.15)$$

#### 1.5.6 DKW inequality

DKW inequality provides a bound on the worst-case distance of empirical CDF and the true CDF:

$$\mathbb{P}\left(\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| > \varepsilon\right) \leq C e^{-2n\varepsilon^2} \quad \text{for every } \varepsilon > 0. \quad (1.16)$$

For multivariate case, let  $X_1, X_2, \dots, x_n$  be an i.i.d. sequence of  $k$ -dimensional vectors,

$$\mathbb{P}\left(\sup_{t \in \mathbb{R}^k} |F_n(t) - F(t)| > \varepsilon\right) \leq (n+1)k e^{-2n\varepsilon^2} \quad (1.17)$$

for every  $\varepsilon, n, k > 0$ .

Also see [local DKW inequality](#)

**Steinke version** Let  $X_1, \dots, X_n$  be independent random variables with CDF  $f(v) := \mathbb{P}[X_i \leq v]$  for all  $i \in [n]$  and  $v \in \mathbb{R}$ . Let the empirical CDF be  $F_x(v) := \frac{1}{n} \sum_{i=1}^n 1[X_i \leq v]$  for all  $v \in \mathbb{R}$ . Then, for all  $\beta > 0$ ,

$$\mathbb{P}_X \left[ \sup_{v \in \mathbb{R}} F_x(v) - f(v) \leq \sqrt{\frac{2 \log(1/\beta)}{n}} + \frac{\log(1/\beta)}{2n} \right] \geq 1 - \beta. \quad (1.18)$$

**Lemma.** For all  $t, \lambda > 0$ ,

$$\mathbb{P} \left[ \sup_{v \in \mathbb{R}} F_x(v) \log \left( 1 + \frac{t}{f(v)} \right) > \frac{\lambda}{n} \right] \leq (1+t)^n e^{-\lambda} \leq e^{tn-\lambda}. \quad (1.19)$$

Note: maximum bound되는 event 확률 구할 때는 martingale construction해서 optional stopping theorem 적용하는 것도 좋음  $\Rightarrow$  Lemma에서는 binomial exponent에 놓아서 martingale 만듬.

#### 1.5.7 Quadratic form

**Definition (Subgaussian random variable).** A centered random variable  $X$  is said to be  $v$ -subgaussian if its cumulant generating function is subquadratic:

$$\xi_X(t) \leq \frac{1}{2} v t^2 \quad \forall t \in \mathbb{R} \quad (1.20)$$

**Hanson-Wright tail bound.** Let  $\mathbf{x}$  be a random vector with independent centered  $v$ -subgaussian entries and let  $\mathbf{A}$  be a square matrix. Then

$$\mathbb{P}(|\mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x}]| \geq t) \leq 2 \exp \left( - \frac{c \cdot t^2}{v^2 \|\mathbf{A}\|_F^2 + v \|\mathbf{A}\| t} \right), \quad (1.21)$$

where  $c > 0$  is a constant independent of  $v, \mathbf{x}, t$  or  $\mathbf{A}$ .

### 1.5.8 Gaussian CCDF bound

$$1 - \Phi(w) \leq \min \left\{ \frac{1}{2}, \frac{1}{w\sqrt{2\pi}} \right\} e^{-w^2/2}, \quad w > 0 \quad (1.22)$$

### 1.5.9 McDiarmid's inequality

A function  $f : \mathcal{X} \times \mathcal{X} \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$  satisfies the bounded differences property if substituting the value of the  $i$ th coordinate  $x_i$  changes the value of  $f$  by at most  $c_i$ . More formally, if there are constants  $c_1, c_2, \dots, c_n$  such that for all  $i \in [n]$ , and all  $x_i \in \mathcal{X}_1, x_2 \in \mathcal{X}_2, \dots, x_n \in \mathcal{X}_3$ ,

$$\sup_{x'_i \in \mathcal{X}_i} |f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i \quad (1.23)$$

Let  $f : \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$  satisfy the bounded differences property with bounds  $c_1, c_2, \dots, c_n$ .

Consider independent random variables  $X_1, X_2, \dots, X_n$  where  $X_i \in \mathcal{X}_i$  for all  $i$ . Then, for any  $\varepsilon > 0$ ,

$$\mathbb{P}(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] \geq \varepsilon) \leq \exp \left( -\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2} \right), \quad (1.24)$$

$$\mathbb{P}(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] \leq -\varepsilon) \leq \exp \left( -\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2} \right) \quad (1.25)$$

and as an immediate consequence,

$$\mathbb{P}(|f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)]| \geq \varepsilon) \leq 2 \exp \left( -\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2} \right) \quad (1.26)$$

## 1.6 Decoupling lemma

### 1.6.1 Quadratic form

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be convex. Let  $X_1, \dots, X_n \in \mathbb{R}$  be independent mean-zero random variables. For  $i, j \in [n]$ , let  $a_{i,j} \in \mathbb{R}$  be a constant. Then

$$\mathbb{E} \left[ f \left( \sum_{i \neq j} a_{ij} X_i X_j \right) \right] \leq \mathbb{E} \left[ f \left( 4 \sum_{i \neq j} a_{ij} X_i X'_j \right) \right], \quad (1.27)$$

where  $X'_1, \dots, X'_n$  are independent copies of  $X_1, \dots, X_n$ .

Note: We can analyze  $\mathbf{x}^* \mathbf{A} \mathbf{x}$  by  $\mathbf{x}^* \mathbf{A} \mathbf{x}'$  with independent  $\mathbf{x}'$ .

## 1.7 Global variance control

### 1.7.1 Efron-Stein inequality

For  $i \in [n]$  and tuple  $Z = (Z_1, \dots, Z_n)$ , let  $Z^{(i)}$  denote the tuple  $(Z_1, \dots, Z_{i-1}, \tilde{Z}_i, Z_{i+1}, \dots, Z_n)$ , where  $\tilde{Z}_i$  is an independent copy of  $Z_i$ . For a scalar function  $f(Z)$ , the Efron-Stein inequality states that

$$\text{Var}[f(Z)] = \mathbb{E}[(f(Z) - \mathbb{E}[f(Z)])^2] \leq \frac{1}{2} \cdot \sum_{i \in [n]} \mathbb{E} \left[ \left( f(Z) - f(Z^{(i)}) \right)^2 \right] \quad (1.28)$$

$$\stackrel{\dagger}{=} \underbrace{\sum_{i \in [n]} \mathbb{E} \left[ \left( f(Z) - \mathbb{E}_i[f(Z^{(i)})] \right)^2 \right]}_{\text{sum of conditional variance}} \quad (1.29)$$



† : Note that  $\mathbb{E} = \mathbb{E}_{-i} \mathbb{E}_i$ ,  $\mathbb{E}[f(Z^{(i)}) \mid Z] = \mathbb{E}_i[f(Z^{(i)})]$  and

$$\mathbb{E}_i[(Z - \mathbb{E}_i[Z])^2] = \frac{1}{2} \mathbb{E}_i[(Z - Z^{(i)})^2] \quad (1.30)$$

## 1.8 Information

### 1.8.1 Fano's inequality

Let  $X \in \{0, 1\}^d$  be uniformly random and let  $Y \in \mathbb{R}^d$  be a random variable that depends on  $X$ .

If  $\mathbb{E}[\|X - Y\|_1] \leq \alpha \cdot d$  for  $\alpha \leq \frac{1}{2}$ , then

$$I(X; Y) \geq d \cdot D_{\text{KL}}\left(\text{Ber}(\alpha) \parallel \text{Ber}\left(\frac{1}{2}\right)\right). \quad (1.31)$$

## 1.9 Do you like martingale?

### 1.9.1 Tail Distribution

Let  $X$  be a nonnegative cadlag submartingale. Then, for each  $K, t > 0$ ,

$$K \mathbb{P}(X_t^* \geq K) \leq \mathbb{E}[1_{\{X_t^* \geq K\}} X_t] \quad (1.32)$$

## 1.10 Stochastic Dominance

**Definition.** Let  $X, Y \in \mathbb{R}$  be random variables. We say  $X$  is *stochastically dominated* by  $Y$  if  $\mathbb{P}[X > t] \leq \mathbb{P}[Y > t]$  for all  $t \in \mathbb{R}$ . Equivalently,  $X$  is stochastically dominated by  $Y$  if there exists a coupling such that  $\mathbb{P}[X \leq Y] = 1$ .

### 1.10.1 SD is preserved under sums/convolutions

**Lemma.** Suppose  $X_1$  is stochastically dominated by  $Y_1$ . Suppose that, for all  $x \in \mathbb{R}$ , the conditional distribution  $X_2 \mid X_1 = x$  is stochastically dominated by  $Y_2$ . Assume that  $Y_1$  and  $Y_2$  are independent. Then  $X_1 + X_2$  is stochastically dominated by  $Y_1 + Y_2$ .

## 2. Privacy

### 2.1 Privacy Loss

**Definition.** Let  $Y$  and  $Z$  be two random variables. The privacy loss random variables  $\mathcal{L}_{Y \parallel Z}$  is distributed by drawing  $t \sim \text{Law}(Y)$ , and outputting  $\log \left( \frac{\mathbb{P}[Y=t]}{\mathbb{P}[Z=t]} \right)$ . If the support of  $Y$  and  $Z$  are not equal, then the privacy loss random variable is undefined.

### 2.2 $\varepsilon - \delta$ DP

#### 2.2.1 4 ways to see $\delta$

**Proposition.** Let  $P$  and  $Q$  be two probability distributions on  $\mathcal{Y}$  such that the privacy loss distribution  $\text{PrivLoss}(P \parallel Q)$  is well-defined. Fix  $\varepsilon \geq 0$  and define

$$\delta := \sup_{S \subset \mathcal{Y}} P(S) - e^\varepsilon Q(S). \quad (2.1)$$

Then

$$\begin{aligned}
 \delta &= \mathbb{P}_{Z \sim \text{PrivLoss}(P||Q)}[Z > \varepsilon] - e^\varepsilon \cdot \mathbb{P}_{Z' \sim \text{PrivLoss}(Q||P)}[-Z' > \varepsilon] \\
 &= \mathbb{E}_{Z \sim \text{PrivLoss}(P||Q)}[\max\{0, 1 - \exp(\varepsilon - Z)\}] \\
 &= \int_\varepsilon^\infty e^{\varepsilon - z} \mathbb{P}_{Z \sim \text{PrivLoss}(P||Q)}[Z > z] dz \\
 &\leq \mathbb{P}_{Z \sim \text{PrivLoss}(P||Q)}[Z > \varepsilon].
 \end{aligned}$$

### 2.2.2 Moment difference bound

Let  $X$  and  $Y$  be a random variable supported on  $[-\Delta, \Delta]$  satisfying  $\mathbb{P}[X \in S] \leq e^\varepsilon \mathbb{P}[Y \in S] + \delta$  for all measurable  $S$  and vice versa. Then

$$\mathbb{E}[X] - \mathbb{E}[Y] \leq (e^\varepsilon - 1) \mathbb{E}[|Y|] + 2\delta\Delta \quad (2.2)$$

## 2.3 zCDP

**Definition.** A randomised mechanism  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  is  $(\xi, \rho)$ -zero-concentrated differentially private if, for all  $x, x' \in \mathcal{X}^n$  differing on a single entry and all  $\alpha \in (1, \infty)$ ,

$$\mathbb{E}[e^{(\alpha-1)Z}] \leq e^{(\alpha-1)(\xi+\rho\alpha)}, \quad (2.3)$$

where  $Z = \text{PrivLoss}(M(x)||M(x'))$  is the privacy loss random variable.

### 2.3.1 Key properties

1. Pure  $\varepsilon$ -DP implies  $\frac{1}{2}\varepsilon^2$ -zCDP
2. The composition of  $k$  independent  $\frac{1}{2}\varepsilon^2$ -zCDP algorithms satisfies  $\frac{1}{2}\varepsilon^2 k$ -zCDP.
3.  $\frac{1}{2}\varepsilon^2 k$ -zCDP implies approximate  $(\varepsilon', \delta)$ -DP with  $\delta \in (0, 1)$  arbitrary and  $\varepsilon' = \varepsilon \cdot \sqrt{2k \log(1/\delta)} + \frac{1}{2}\varepsilon^2 k$ .

## 2.4 Approximate Rényi Differential Privacy

“Rényi differential privacy was introduced by Minorov and was motivated by analyzing privacy amplification by subsampling interleaved with composition, which arises in differentially private deep learning”

— Thomas Steinke **Definition (RDP).** An algorithm  $M$  is said to be  $(\lambda, \varepsilon)$ -RDP with  $\lambda \geq 1$  and  $\varepsilon \geq 0$ , if for any adjacent inputs  $x, x'$

$$D_\lambda(M(x)||M(x')) := \frac{1}{\lambda-1} \log_{Y \leftarrow M(x)} \left[ \left( \frac{\mathbb{P}[M(x) = Y]}{\mathbb{P}[M(x') = Y]} \right)^{\lambda-1} \right] \leq \varepsilon \quad (2.4)$$

**Tip:** The  $\varepsilon$  should be thought of as a function  $\varepsilon(\lambda)$ , rather than a single number.

### Properties

Let  $P, Q$  be probability distributions over  $\mathcal{Y}$  with a common sigma-algebra such that  $P$  is absolutely continuous with respect to  $Q$ .

1. **Postprocessing (a.k.a. data processing inequality) & non-negativity:**

Let  $f : \mathcal{Y} \rightarrow \mathcal{Z}$  be a measurable function. Let  $f(P)$  denote the distribution on  $\mathcal{Z}$

obtained by applying  $f$  to a sample from  $P$ ; define  $f(Q)$  similarly. Then

$$0 \leq D_\alpha(f(P)\|f(Q)) \leq D_\alpha(P\|Q) \quad \text{for all } \alpha \in [1, \infty].$$

2. **Composition:** If  $P = P' \times P''$  and  $Q = Q' \times Q''$  are product distributions, then

$$D_\alpha(P\|Q) = D_\alpha(P'\|Q') + D_\alpha(P''\|Q'') \quad \text{for all } \alpha \in [1, \infty].$$

More generally, suppose  $P$  and  $Q$  are distributions on  $\mathcal{Y} = \mathcal{Y}' \times \mathcal{Y}''$ . Let  $P'$  and  $Q'$  be the marginal distributions on  $\mathcal{Y}'$  induced by  $P$  and  $Q$  respectively. For  $y' \in \mathcal{Y}'$ , let  $P_{y'}''$  and  $Q_{y'}''$  be the conditional distributions on  $\mathcal{Y}''$  induced by  $P$  and  $Q$  respectively. That is, we can generate a sample  $Y = (Y', Y'') \leftarrow P$  by first sampling  $Y' \leftarrow P'$  and then sampling  $Y'' \leftarrow P_{Y'}''$ , and similarly for  $Q$ . Then

$$D_\alpha(P\|Q) \leq D_\alpha(P'\|Q') + \sup_{y' \in \mathcal{Y}'} D_\alpha(P_{y'}''\|Q_{y'}'') \quad \text{for all } \alpha \in [1, \infty].$$

3. **Monotonicity:** For all  $1 \leq \alpha \leq \alpha' \leq \infty$ ,

$$D_\alpha(P\|Q) \leq D_{\alpha'}(P\|Q).$$

4. **Gaussian divergence:** For all  $\mu, \mu' \in \mathbb{R}$  with  $\sigma > 0$  and all  $\alpha \in [1, \infty)$ ,

$$D_\alpha(\mathcal{N}(\mu, \sigma^2)\|\mathcal{N}(\mu', \sigma^2)) = \alpha \cdot \frac{(\mu - \mu')^2}{2\sigma^2}.$$

5. **Pure DP to Concentrated DP:** For all  $\alpha \in [1, \infty)$ ,

$$D_\alpha(P\|Q) \leq \frac{\alpha}{8} \cdot (D_\infty(P\|Q) + D_\infty(Q\|P))^2.$$

6. **Quasi-convexity:** Let  $P'$  and  $Q'$  be probability distributions over  $\mathcal{Y}$  such that  $P'$  is absolutely continuous with respect to  $Q'$ . For  $s \in [0, 1]$ , let  $(1-s) \cdot P + s \cdot P'$  denote the convex combination of the distributions  $P$  and  $P'$  with weighting  $s$ . For all  $\alpha \in (1, \infty)$  and all  $s \in [0, 1]$ ,

$$\begin{aligned} & D_\alpha((1-s) \cdot P + s \cdot P' \parallel (1-s) \cdot Q + s \cdot Q') \\ & \leq \frac{1}{\alpha-1} \log((1-s) \cdot \exp((\alpha-1)D_\alpha(P\|Q)) + s \cdot \exp((\alpha-1)D_\alpha(P'\|Q'))) \\ & \leq \max\{D_\alpha(P\|Q), D_\alpha(P'\|Q')\}, \end{aligned}$$

and

$$D_1((1-s) \cdot P + s \cdot P' \parallel (1-s) \cdot Q + s \cdot Q') \leq (1-s) \cdot D_1(P\|Q) + s \cdot D_1(P'\|Q').$$

7. **Triangle-like inequality (a.k.a. group privacy):** Let  $R$  be a distribution on  $\mathcal{Y}$  and assume that  $Q$  is absolutely continuous with respect to  $R$ . For all  $1 < \alpha < \alpha' < \infty$ ,

$$D_\alpha(P\|R) \leq \frac{\alpha'}{\alpha' - 1} \cdot D_{\alpha', \frac{\alpha' - 1}{\alpha}}(P\|Q) + D_{\alpha'}(Q\|R).$$

In particular, if  $D_\alpha(P\|Q) \leq \rho_1 \cdot \alpha$  and  $D_\alpha(Q\|R) \leq \rho_2 \cdot \alpha$  for all  $\alpha \in (1, \infty)$ , then

$$D_\alpha(P\|R) \leq (\sqrt{\rho_1} + \sqrt{\rho_2})^2 \cdot \alpha \quad \text{for all } \alpha \in (1, \infty).$$

8. **Conversion to approximate DP:** For all measurable  $S \subset \mathcal{Y}$ , all  $\alpha \in (1, \infty)$ , and all  $\tilde{\epsilon} \geq D_\alpha(P\|Q)$ ,

$$P(S) \leq e^{\tilde{\epsilon}} \cdot Q(S) + e^{-(\alpha-1)(\tilde{\epsilon} - D_\alpha(P\|Q))} \cdot \frac{1}{\alpha} \left(1 - \frac{1}{\alpha}\right)^{\alpha-1} \leq e^{\tilde{\epsilon}} \cdot Q(S) + e^{-(\alpha-1)(\tilde{\epsilon} - D_\alpha(P\|Q))}.$$

**Definition (Approximate RDP).** A randomized algorithm  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  is  $\delta$ -approximately  $(\lambda, \varepsilon)$ -Rényi differentially private if, for all neighboring pairs of inputs  $x, x' \in \mathcal{X}^n$ ,  $D_\lambda^\delta(M(x) || M(x')) \leq \varepsilon$ .

#### 2.4.1 Properties

1.  $(\varepsilon, \delta)$ -DP is equivalent to  $\delta$ -approximate  $(\infty, \delta)$ -RDP.
2.  $(\varepsilon, \delta)$ -DP implies  $\delta$ -approximate  $(\lambda, \frac{1}{2}\varepsilon^2\delta)$ -RDP for all  $\lambda \in (1, \infty)$ .
3.  $\delta$ -approximate  $(\lambda, \varepsilon)$ -RDP implies  $(\hat{\varepsilon}, \hat{\delta})$ -DP for

$$\hat{\delta} = \delta + \frac{\exp((\lambda - 1)(\hat{\varepsilon} - \varepsilon))}{\lambda} \cdot \left(1 - \frac{1}{\lambda}\right)^{\lambda-1}. \quad (2.5)$$

4.  $\delta$ -approximate  $(\lambda, \varepsilon)$ -RDP is closed under postprocessing.
5. If  $M_1$  is  $\delta_1$ -approximately  $(\lambda, \varepsilon_1)$ -RDP and  $M_2$  is  $\delta_2$ -approximately  $(\lambda, \varepsilon_2)$ -RDP, then their composition is  $(\delta_1 + \delta_2)$ -approximately  $(\lambda, \varepsilon_1 + \varepsilon_2)$ -RDP.

## 2.5 Composition

#### 2.5.1 Advanced composition $(\varepsilon, \delta)$

**Theorem.** If each mechanism  $m_i$  is in a  $k$ -fold adaptive composition  $m_1, \dots, m_k$  satisfies  $\varepsilon$ -differential privacy, then for any  $\delta' \geq 0$ , the entire  $k$ -fold adaptive composition satisfies  $(\varepsilon', \delta')$ -differential privacy, where

$$\varepsilon' = \varepsilon \sqrt{2k \log(1/\delta')} + k\varepsilon(e^\varepsilon - 1) \quad (2.6)$$

**Theorem.** For  $j \in [k]$ , let  $M_j \in \mathcal{X}^n \times \mathcal{Y}_{i-1} \rightarrow \mathcal{Y}_i$  be randomized algorithms. Suppose  $M_j$  is  $(\varepsilon_j, \delta_j)$ -DP for each  $j \in [k]$ . For  $j \in [k]$ , inductively define  $M_{1\dots j} : \mathcal{X}^n \rightarrow \mathcal{Y}_j$  by  $M_{1\dots j}(x) = M_j(x, M_{1\dots(j-1)}(x))$ , where each algorithm is run independently and  $M_{1\dots 0} = y$  for some fixed  $y_0 \in \mathcal{Y}_0$ . Then  $M_{1\dots k}$  is  $(\varepsilon, \delta)$ -DP for any  $\delta > \sum_{j=1}^k \delta_j$  with

$$\varepsilon = \min \left\{ \sum_{j=1}^k \varepsilon_j, \frac{1}{2} \sum_{j=1}^k \varepsilon_j^2 + \sqrt{2 \log(1/\delta') \sum_{k=1}^k \varepsilon_j^2} \right\} \quad (2.7)$$

## 2.6 Joint Differential Privacy

**Definition.** For  $\varepsilon, \delta \geq 0$ , a randomized algorithm  $\mathcal{M} : \mathbb{N}^{\mathcal{X}} \rightarrow \mathcal{Y}^N$  is  $(\varepsilon, \delta)$ -joint differentially private if for every possible pair of  $z, z' \in \mathcal{X}$ , for every  $i \in [N]$ , and for every subset of possible outputs  $E \subseteq \mathcal{Y}^{N-1}$ , we have

$$\mathbb{P}_{\mathcal{M}}[\mathcal{M}(z \cup D_{-z})_{-i} \in E] \leq e^\varepsilon \mathbb{P}_{\mathcal{M}}[\mathcal{M}(z' \cup D_{-z})_{-i} \in E] + \delta \quad (2.8)$$

where  $\mathcal{M}_{-i}$  denotes the output of  $\mathcal{M}$  that excludes the  $i$ th dimension.

## 2.7 Lower bound tools

#### 2.7.1 Query moment w.r.t binary data

**Correlation-Variance Dichotomy.** Let  $f : \{0, 1\}^d \rightarrow [0, 1]$  be an arbitrary function. Let  $P \in [0, 1]$  be uniformly random and, conditioned on  $P$ , let  $X_1, X_2, \dots, X_n$  be independent with  $\mathbb{E}[X_i] = P$  for each  $i \in [n]$ . Then

$$\underbrace{\mathbb{E}_{X, P} \left[ (f(X) - P) \cdot \sum_{i=1}^n (X_i - P) \right]}_{\text{일종의 total correlation}} + \mathbb{E}_P \left[ \mathbb{E}_X [f(X) - \overline{X}]^2 \right] \geq \frac{1}{12} \quad (2.9)$$

## 2.8 Decomposition

### 2.8.1 Basic decomposition

Let  $P$  and  $Q$  be probability distributions over  $\mathcal{Y}$ . Fix  $\varepsilon, \delta \geq 0$ . Suppose that, for all measurable  $S \subset \mathcal{Y}$ , we have  $P(S) \leq e^\varepsilon Q(S) + \delta$  and vice versa. Then there exist  $\delta' \in [0, \delta]$  and distributions  $P', Q', P''$  and  $Q''$  over  $\mathcal{Y}$  such that the following three properties are all satisfied.

1. We can express  $P$  and  $Q$  as convex combinations:

$$\begin{aligned} P &= (1 - \delta')P' + \delta'P'' \\ Q &= (1 - \delta')Q' + \delta'Q'' \end{aligned}$$

2. Second, for all measurable  $S \subset \mathcal{Y}$ , we have  $e^{-\varepsilon}P'(S) \leq Q'(S) \leq e^\varepsilon P'(S)$
3. There exists measurable  $S, T \subset \mathcal{Y}$  such that  $P''(S) = 1, Q''(T) = 1, \forall S' \subset S P(S') \geq Q(S')$ , and  $\forall T' \subset T Q(T') \geq P(T')$

**Corollary.** Let  $P$  and  $Q$  be probability distribution over  $\mathcal{Y}$ . Fix  $\varepsilon, \delta$ . Suppose that for all measurable  $S \subset \mathcal{Y}$ , we have  $P(S) \leq e^\varepsilon Q(S) + \delta$  and  $Q(S) \leq e^\varepsilon P(S) + \delta$ . Then there exist distributions  $A, B, P''$ , and  $Q''$  over  $\mathcal{Y}$  such that

$$\begin{aligned} P &= (1 - \delta) \frac{e^\varepsilon}{e^\varepsilon + 1} A + (1 - \delta) \frac{1}{e^\varepsilon + 1} B + \delta P'', \\ Q &= (1 - \delta) \frac{e^\varepsilon}{e^\varepsilon + 1} B + (1 - \delta) \frac{1}{e^\varepsilon + 1} A + \delta Q'' \end{aligned}$$

**Interpretation:** All  $(\varepsilon, \delta)$  DP distributions can be represented as a postprocessing of the  $(\varepsilon, \delta)$  randomized response with the postprocessing  $F$  such that  $F(0, \perp) = A, F(1, \perp) = B, F(0, \top) = P''$  and  $F(1, \top) = Q''$

### 2.8.2 Bayesian version [2, Lemma 5.6]

#### Question

Suppose we observe a sample from either  $P$  or  $Q$  and we have a prior on these two possibilities, what is the posterior distribution of possibilities? We need to account for the event with  $\delta$  where things "fail" arbitrarily.

Let  $P$  and  $Q$  be probability distributions over  $\mathcal{Y}$ . Fix  $\varepsilon, \delta \geq 0$ . Suppose that, for all measurable  $S \subset \mathcal{Y}$ , we have  $P(S) \leq e^\varepsilon Q(S) + \delta$  and vice versa.

Then there exists a randomized function  $E_{P,Q} : \mathcal{Y} \rightarrow \{0, 1\}$  with the following properties:

1. Fix  $p \in [0, 1]$  and suppose  $X \sim \text{Bernoulli}(p)$ . If  $X = 1$ , sample  $Y \sim P$  else  $Y \sim Q$ . Then for all  $Y \in \mathcal{Y}$ , we have

$$\mathbb{P}_{\substack{X \sim \text{Bernoulli}(p) \\ Y \sim XP + (1-X)Q}} [X = 1 \wedge E_{P,Q}(Y) = 1 | Y = y] \leq \frac{p}{p + (1-p)e^{-\varepsilon}}$$

2. Under each hypothesis  $Y \sim P$  and  $Y \sim Q$ , the expected value  $E_{P,Q}(Y)$  is equal or greater than  $1 - \delta$ .

## 3. Probability and Statistics

### 3.1 Moment

#### 3.1.1 Moment and CDF

For random variable  $X$ ,

$$\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X \geq t) dt - \int_{-\infty}^0 \mathbb{P}(X \leq t) dt \quad (3.1)$$

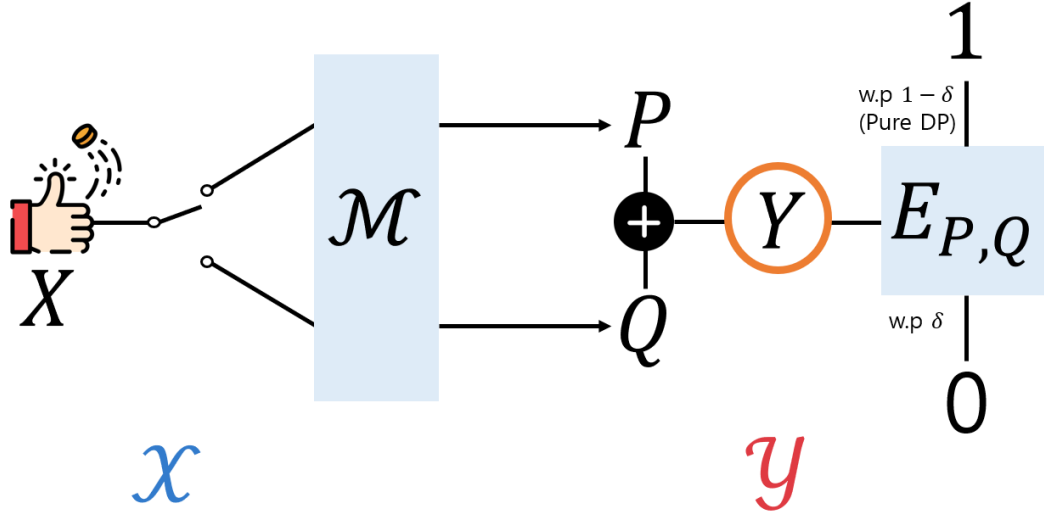


Figure 1: Visualization of the Bayesian version decomposition

In general,  $X \geq 0$  and a smooth function  $g$  with  $g(0) = 0$

$$\mathbb{E}[g(X)] = \int_0^\infty g'(t)P(X > t) dt \quad (3.2)$$

### 3.1.2 $k$ -th moment in the lens of CDF

$$\frac{1}{k} \mathbb{E}[X^k] = - \int_{-\infty}^0 x^{k-1} F(x) dx + \int_0^\infty x^{k-1} (1 - F(x)) dx \quad (3.3)$$

### 3.1.3 Clipped random variable

For complementary CDF  $\bar{F}$  of random variable  $X \geq 0$ ,

$$\mathbb{E}[\min(X, k)] = \int_0^k x f(x) dx + k \bar{F}(k), \quad (3.4)$$

## 3.2 Conditional Distribution

**Tip:** ML/Statistics 분야에서 흔히 쓰는 notation  $p(x|y), p(x, y)$  같은 것은  $X, Y$  들의 pdf/pmf라고 생각하자. 거의 measure는 안나옴.

**Conditioning on event.** Let  $f_{X,Y}$  be the joint density of  $X$  and  $Y$ , and  $f_X(x)$  is the marginal density of  $X$ .

1. Single point conditioning  $X = 1$

$$f_{Y|X=1}(y) = \frac{f_{X,Y}(1, y)}{f_X(1)}$$

2. Set conditioning  $X \in S$

$$f_{Y|X \in S}(y) = \frac{\int_S f_{X,Y}(x, y) dx}{\int_S f_X(x) dx} = \frac{\int_S f_{Y|X=x}(y) f_X(x) dx}{\int_S f_X(x) dx}$$

3. Event conditioning

$$f_{Y|A}(y) = \frac{f_Y(y) 1[y \in A]}{\int_A f_Y(y) dy}$$

### 3.3 Weak convergence

#### 3.3.1 Delta method

If there is a sequence of random variables  $X_n$  satisfying

$$\sqrt{n}[X_n - \theta] \xrightarrow{w} \mathcal{N}(0, \sigma^2) \quad (3.5)$$

then

$$\sqrt{n}[g(X_n) - g(\theta)] \xrightarrow{w} \mathcal{N}(0, \sigma^2 \cdot [g'(\theta)]^2) \quad (3.6)$$

given that  $g'(\theta)$  exists and is non-zero value.

#### 3.3.2 Slutsky's theorem

If  $X_n$  converges in distribution to a random element  $X$  and  $Y_n$  converges in probability to a constant  $c$ , then

1.  $X_n + Y_n \xrightarrow{w} X + c$
2.  $X_n Y_n \xrightarrow{w} Xc$
3.  $X_n / Y_n \xrightarrow{w} X/c$

where  $\xrightarrow{w}$  denotes convergence in distribution.

### 3.4 Central limit theorem

#### 3.4.1 Salem-Zygmund

**Theorem.** Let  $U$  be a uniform random variable with support  $(0, 2\pi)$ , and let  $X_k = r_k \cos(n_k U + a_k)$  ( $0 \leq a_k < 2\pi$ ), where

1.  $n_k$  satisfy the *lacunarity condition*: there exists  $q > 1$  such that  $n_{k+1} \geq qn_k$  for all  $k$
2.  $\sum_{i=1}^{\infty} r_i^2 = \infty$  and  $\frac{r_k^2}{r_1^2 + \dots + r_k^2} \rightarrow 0$

Then,

$$\frac{X_1 + \dots + X_k}{\sqrt{r_1^2 + \dots + r_k^2}} \quad (3.7)$$

converges in distribution to  $\mathcal{N}(0, 1/2)$ .

### 3.5 Conditional independence

**Theorem.** Let  $p_{XYZ}$  be the joint PDF/PMF of  $X, Y$  and  $Z$ . Then the following are equivalent with up to almost-everywhere equivalence:

1.  $X \perp Y \mid Z$
2.  $p_{XYZ}(x, y, |z) = p_{X|Z}(x|z)p_{Y|Z}(y|z)$
3.  $p_{X|YZ}(x|y, z) = p_{X|Z}(x|z)$
4.  $p_{XYZ}(x, y, z) = \frac{p_{XZ}(x, z)p_{YZ}(y, z)}{p_Z(z)}$
5.  $p_{XYZ}(x, y, z) = g(x, z)h(y, z)$  for some measurable functions  $g$  and  $h$
6.  $p_{X|YZ}(x|y, z) = w(x, z)$  for some measurable function  $w$

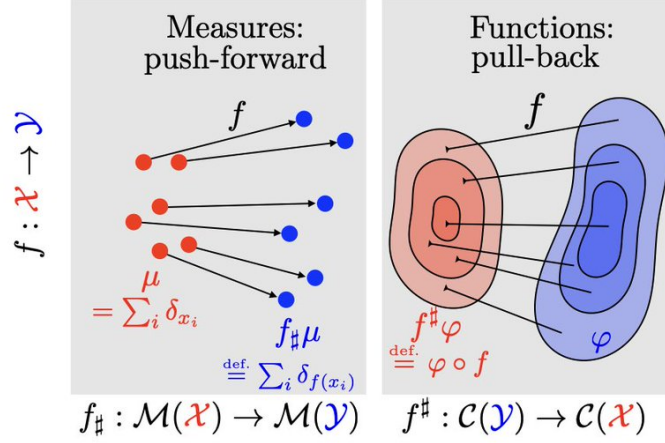
**Properties.** Let  $X, Y, Z, W$  be RVs

1. (symmetry)  $X \perp Y \mid Z \iff Y \perp X \mid Z$
2. (decomposition)  $X \perp Y \mid Z \Rightarrow h(X) \perp Y \mid Z$  for any measurable function  $h$
3. (weak union)  $X \perp Y \mid Z \Rightarrow X \perp Y \mid Z, h(X)$  for any measurable function  $h$
4. (contraction)

$$X \perp Y \mid Z \text{ and } X \perp W \mid (Y, Z) \iff X \perp (W, Y) \mid Z \quad (3.8)$$

5. If the joint PDF  $P_{XYZW}(x, y, z, w)$  satisfies  $f_{YZW}(y, z, w) > 0$  almost everywhere. Then

$$X \perp Y \mid (W, Z) \text{ and } X \perp W \mid (Y, Z) \iff X \perp (W, Y) \mid Z \quad (3.9)$$



Remark:  $f^{\#}$  and  $f_{\#}$  are adjoints

$$\int_{\mathcal{Y}} \varphi d(f_{\#}\mu) = \int_{\mathcal{X}} (f^{\#}\varphi) d\mu$$

Figure 2: Pullback of functions and pushforward of measures are dual one with each other!

### 3.5.1 Bayes' Theorem

Assume that  $X$  is a random variable on  $(\Omega, \mathcal{F}, P)$ , and let  $Q$  be another probability measure on  $(\Omega, \mathcal{F})$  with Radon-Nikodym derivative

$$L = \frac{dQ}{dP} \text{ on } \mathcal{F} \quad (3.10)$$

Assume that  $X \in L^1(\Omega, \mathcal{F}, Q)$  and that  $\mathcal{G}$  is a sigma-algebra with  $\mathcal{G} \subseteq \mathcal{F}$ . Then

$$\mathbb{E}_Q[X | \mathcal{G}] = \frac{\mathbb{E}_P[L \cdot X | \mathcal{G}]}{\mathbb{E}_P[L | \mathcal{G}]}, \quad Q - a.s. \quad (3.11)$$

### 3.5.2 Conditional expectation under independence

**Proposition.** Let  $(\Omega, \mathcal{B}, \mathbb{P})$  be a probability space,  $(\mathbb{X}, \mathcal{M}), (\mathbb{Y}, \mathcal{N})$  be measurable spaces,  $X: \Omega \rightarrow \mathbb{X}$  and  $Y: \Omega \rightarrow \mathbb{Y}$  be measurable functions. If  $X$  and  $Y$  are independent and  $f \in (\mathcal{M} \otimes \mathcal{N})_b$  then

$$\mathbb{E}[f(X, Y) | X] = \mathbb{E}[f(x, Y)]|_{x=X} \text{ a.s.} \quad (3.12)$$

## 3.6 Regular conditional distribution

**Theorem.** If  $X$  is a real random variable defined on probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  then for every  $\sigma$ -algebra  $\mathcal{G} \subset \mathcal{F}$  there is a regular conditional distribution for  $X$  given  $\mathcal{G}$ .

Regular conditional distributions are useful in part because they allow one to reduce many problems concerning conditional expectations to problems concerning only ordinary expectations. For such applications the following disintegration formula for conditional expectations is essential.

### 3.6.1 Disintegration formula

**Theorem.** Let  $\mu_w(dx)$  be a regular conditional distribution for  $X$  given  $\mathcal{G}$ , let  $Y$  be  $\mathcal{G}$ -measurable, and let  $f(x, y)$  be a jointly measurable real-valued function such that  $\mathbb{E}[|f(X, Y)|] < \infty$ . Then,

$$\mathbb{E}[f(X, Y) | \mathcal{G}] = \int f(x, Y(w)) \mu_w(dx) \quad \text{a.s.} \quad (3.13)$$



**Theorem2.** Let  $Y$  and  $X$  be two Radon spaces. Let  $\mu \in P(Y)$ , let  $\pi : Y \rightarrow X$  be a Borel-measurable function, and let  $\nu \in P(X)$  be the pushforward measure from  $Y$  to  $X$  by  $\pi$ . Then there exists a  $\nu$ -almost everywhere uniquely determined family of probability measures  $\{\mu_x\}_{x \in X} \subseteq P(Y)$  such that

1. the function  $x \mapsto \mu_x$  is Borel measurable
2.  $\mu_x$  lives on the fiber  $\pi^{-1}(x)$
3. for every Borel-measurable function  $f : Y \rightarrow [0, +\infty]$ ,

$$\int_Y f(y) d\mu(y) = \int_X \int_{\pi^{-1}(x)} f(y) d\mu_x(y) d\nu(x) \quad (3.14)$$

### 3.7 Markov kernel

A Markov kernel (also called transition kernel, stochastic kernel, or probability kernel) is a mathematical formalization of a “function with random outcomes”.

### 3.8 Mutual Information

#### 3.8.1 Concavity of Mutual information

Let  $\alpha$  be the law of  $X$  and  $\pi$  be the conditions law of  $Y|X$ . Let  $I_1$  be  $I(X; Y)$  where  $(X, Y) \sim (\alpha_1, \pi)$ , let  $I_2$  be  $I(X; Y)$  where  $(X, Y) \sim (\alpha_2, \pi)$ , let  $I$  be  $I(X; Y)$  where  $(X, Y) \sim (\lambda\alpha_1 + (1 - \lambda)\alpha_2, \pi)$ , for some  $0 \leq \lambda \leq 1$ , then

$$I \geq \lambda I_1 + (1 - \lambda) I_2.$$

### 3.9 Fisher Information

Given the score function  $\log p(\theta; X)$ , the Fisher Information is defined as

$$I(\theta) := \mathbb{E} \left[ -\frac{\partial^2}{\partial \theta^2} \log p(\theta; X) \right] \quad (3.15)$$

It gives you uncertainty about the estimation since

$$\underbrace{\text{Var} \left[ \frac{\partial \ell(\theta; X)}{\partial \theta} \right]}_{\text{variance of score}} = -\mathbb{E} \left[ \frac{\partial^2 \ell(\theta; X)}{\partial \theta^2} \right] \quad (3.16)$$

holds.

#### 3.9.1 Cramér–Rao bound

Unbiased means there

$$\mathbb{E}[\hat{\theta}(X) - \theta \mid \theta] = \int (\hat{\theta}(x) - \theta) f(x; \theta) dx = 0 \text{ regardless of the value of } \theta \quad (3.17)$$

Then, the following holds

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)} \quad (3.18)$$

The precision to which we can estimate  $\theta$  is fundamentally limited by the Fisher information of the likelihood function.

### 3.10 MLE estimation

#### 3.10.1 Asymptotic normality

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} \mathcal{N}(0, I^{-1}(\theta)) \quad (3.19)$$

The mean square error (MSE) of  $\hat{\theta}_n$  is

$$\text{MSE}(\hat{\theta}_n, \theta_0) = \text{bias}^2(\hat{\theta}) + \text{Var}(\hat{\theta}_n) \approx \frac{1}{nI(\theta_0)} \quad (3.20)$$

Moreover, if we know about  $I(\theta_0)$ , we can construct a  $1 - \alpha$  confidence interval using

$$\left[ \hat{\theta}_n - \frac{z_{1-\alpha/2}}{\sqrt{n\hat{I}(\theta_0)}}, \hat{\theta}_n + \frac{z_{1-\alpha/2}}{\sqrt{n\hat{I}(\theta_0)}} \right] \quad (3.21)$$

1

### 3.11 Famous Family

#### 3.11.1 Poission

**Binomial of Poisson trials is Poisson.** Let  $\lambda \geq 0, p \in [0, 1]$ . Suppose  $(X_i)_{i=1}^\infty$  are i.i.d Bernoulli random variables with parameter  $p$ , and  $N$  is a Poisson( $\lambda$ ) random variable independent of the  $X_i$ 's. Then  $\sum_{i=1}^N X_i \sim \text{Poi}(\lambda p)$ .

**Tail Distribution.** Let  $X \sim \text{Poi}(\lambda)$ , for some parameter  $\lambda > 0$ . Then for any  $x > 0$ , we have

$$\mathbb{P}[X \geq \lambda + x] \leq e^{-\frac{x^2}{2\lambda} h(\frac{x}{\lambda})} \quad (3.22)$$

, and, for any  $0 < x < \lambda$ ,

$$\mathbb{P}[X \leq \lambda - x] \leq e^{-\frac{x^2}{2\lambda} h(\frac{x}{\lambda})}. \quad (3.23)$$

where  $h(u) := 2 \frac{(1+u) \log(1+u) - u}{u^2}$ . In particular, this implies that for every  $x > 0$ ,

$$\mathbb{P}[|X - \lambda| \geq x] \leq 2e^{-\frac{x^2}{2(\lambda+x)}}. \quad (3.24)$$

#### 3.11.2 Binomial

**Fact.** If  $X$  is Binomial( $n, p$ ), then  $\mathbb{E}[1/(X+1)] \leq 1/((n+1) \cdot p)$

#### 3.11.3 Negative Binomial

**Negative binomial as a Poisson + logarithmic.**

1. Draw  $T$  from a Poisson distribution and draw  $K_1, K_2, \dots, K_T$  independently from a logarithmic distribution.
2. Then  $K = \sum_{t=1}^T K_t$  follows a negative binomial distribution.

### 3.12 M-estimator

#### 3.12.1 Consistency

**Theorem.** Suppose that

1.  $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \rightarrow 0$  in probability (ULLN)
2. For all  $\varepsilon > 0$ ,  $\sup \{M(\theta) : d(\theta, \theta_0) \geq \varepsilon\} < M(\theta_0)$  (identifiability)
3.  $M_n(\hat{\theta}_n) \geq M_n(\theta_0) - o_p(1)$

Then  $\hat{\theta}_n \rightarrow \theta_0$  in probability.

---

<sup>1</sup>CI: estimator  $\pm$  z-value \* (SD of estimator)

## 4. Convexity

### 4.1 Characterization

Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice differentiable over an open domain. Then, the following are equivalent

1.  $f$  is convex
2.  $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ , for all  $x, y \in \text{dom}(f)$
3.  $\nabla^2 f(x) \succeq 0$ , for all  $x \in \text{dom}(f)$

여기서 Condition 3은 모든 점에서 non-negative curvature를 가지고 있다는 의미.

#### 4.1.1 Characterization of Strict Convexity

1.  $\nabla^2 f(x) \succ 0, \forall x \in \Omega$  ( The converse is not true)
2. A function  $f$  is strictly convex on  $\Omega \subseteq \mathbb{R}^n$  if and only if

$$f(y) > f(x) + \nabla^\top f(x)(y - x), \forall x, y \in \Omega, x \neq y$$

3.  $f$  is strongly convex if and only if there exists  $m > 0$  such that

$$\begin{aligned} f(y) &\geq f(x) + \nabla^\top f(x)(y - x) + m \|y - x\|^2, \forall x, y \in \text{dom}(f) \\ \iff \nabla^2 f(x) &\succeq m \mathbf{I}, \forall x \in \text{dom}(f) \end{aligned}$$

#### 4.1.2 When can we assume equal variables?

If the **constraints and the function to be optimized are both symmetric with respect to a group of permutations of the variables**, then the solution set will also be symmetric with respect to this group.

### 4.2 Dual Problem

$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu)$  is concave and lower bound of the optimal value.

**Example.**

$$\begin{aligned} \text{maximize} \quad & \sum_{i=0}^m \mathbb{P}_F[F(t) = i] \cdot \mathbb{P}_{\tilde{W}}[\tilde{W}(t) \geq v - i] \\ \text{subject to} \quad & \sum_{i=0}^m \mathbb{P}[F(t) = i] \cdot i \leq 2m \cdot \delta \\ & \sum_{i=0}^m \mathbb{P}_F[F(t) = i] = 1, \text{ and } \mathbb{P}_F[F(t) = i] \geq 0 \forall i \in \{0, 1, \dots, m\} \end{aligned}$$

**Step 1. Find Lagrangian**

$$L(F, \alpha, \beta, \lambda) = \sum_{i=0}^m \mathbb{P}_F[F(t) = i] (\mathbb{P}_{\tilde{W}}[\tilde{W}(t) \geq v - i] - \alpha i - \beta - \lambda_i) + 2m\delta\alpha + \beta \quad (4.1)$$

with  $\alpha \geq 0, \lambda \geq 0$ . Then,  $g(\alpha, \beta, \lambda) := \sup_F L(F, \alpha, \beta, \lambda) \geq \sup_{F \in \mathcal{C}} L(F, \alpha, \beta, \lambda) \geq f^*$ . To drop the first term, add the constraint:  $\alpha \cdot i + \beta \geq \mathbb{P}_{\tilde{W}}[\tilde{W}(t) \geq v - i] \forall i \in \{0, 1, \dots, m\}$ .

**Step 2. Optimize the Dual function**

$$\begin{aligned} \text{minimize} \quad & 2m\delta\alpha + \beta \\ \text{subject to} \quad & \alpha \cdot i + \beta \geq \mathbb{P}_{\tilde{W}}[\tilde{W}(t) \geq v - i] \quad \forall i \in \{0, 1, \dots, m\} \\ & \alpha \geq 0 \end{aligned}$$

Make  $\alpha, \beta$  as small as possible!

$$\beta = \mathbb{P}_{\tilde{W}^*}[\tilde{W} \geq v],$$

$$\alpha = \max \left( \{0\} \cup \left\{ \frac{1}{i} \left( \mathbb{P}_{\tilde{W}^*}[\tilde{W} \geq v - i] - \beta \right) : i \in \{1, 2, \dots, m\} \right\} \right)$$

where  $\tilde{W}^*$  is a distribution on  $\mathbb{R}$  such that  $\mathbb{P}[\tilde{W}^* \geq v - i] \geq \mathbb{P}[\tilde{W}(t) \geq v - i]$  for all  $i \in \{0, 1, \dots, m\}$  and all  $t$  in the support of  $T$ .

**Remark.** Dual problem을 구하는 과정에서 constraint 추가해도 괜찮음.

## 5. Integral Technique

### 5.1 Gaussian Integral

#### 5.1.1 Gauss-Hermite quadrature

$$\int_{-\infty}^{\infty} e^{-x^2/2} f(x) dx \approx \sum_{i=1}^n w_i f(x_i) \quad (5.1)$$

where  $n$  is the number of sampled points used. The  $x_i$  are the roots of the physicists' version of the Hermite polynomial  $H_n(x)$  ( $i = 1, 2, \dots, n$ ) and the associated weights  $w_i$  are given by

$$w_i = \frac{2^{n-1} n! \sqrt{\pi}}{n^2 [H_{n-1}(x_i)]^2} \quad (5.2)$$

#### 5.1.2 Stein's lemma

For a differentiable function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  with  $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[|\phi'(x)|] < \infty$

$$\mathbb{E}_{x \sim \mathcal{N}(0,1)}[\phi(x)x] = \mathbb{E}_{x \sim \mathcal{N}(0,1)}[\phi'(x)] \quad (5.3)$$

#### 5.1.3 Change of measure

$$\mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) d\mathbf{w} = \mathcal{N}(\mathbf{x}^\top \mathbf{w} | \mathbf{x}^\top \boldsymbol{\mu}_n, \mathbf{x}^\top \boldsymbol{\Sigma}_n \mathbf{x}) d(\mathbf{x}^\top \mathbf{w}) \quad (5.4)$$

## 6. Stochastic process

### 6.1 Predictable process

#### 6.1.1 How to understand the predictable process?

**Fact 1.** The predictable sigma-algebra is generated by continuous and adapted processes.

**Fact 2.** The predictable sigma-algebra is generated by the sets of the form

$$\{(s, t] \times A : t > x \geq 0, A \in \mathcal{F}_s\} \cup \{\{0\} \times A : A \in \mathcal{F}_0\} \quad (6.1)$$

### 6.2 Local martingale

#### 6.2.1 Quadratic variation

**Fact.** If  $X$  is a continuous local martingale, then  $[X]_t < \infty$  a.s. for every  $t \geq 0$ , where  $[X]$  denote the quadratic variation of the process  $X$ .

#### 6.2.2 Stochastic integral

**Proposition.** For any continuous  $L^2$ -martingale  $M$  where  $M_0 = 0$ , and any predictable step process  $V$  where  $|V| \leq 1$ , the process  $(V \cdot M)$  is an  $L^2$ -martingale with  $\mathbb{E}(V \cdot M)_t^2 \leq \mathbb{E} M_t^2$ .

## 6.3 Doob's h transform

Set

$$h(x) = \mathbb{P}_x(\tau_A < \tau_B) \quad (6.2)$$

Then,  $h(x)$  is the probability, starting from  $x$  to hit  $A$  before hitting  $B$ . Then  $h$  is positive on  $\mathcal{X} \setminus (A \cup B)$ . Furthermore, for  $x \notin A \cup B$

$$\hat{P}(x, y) = \mathbb{P}_x[X_1 = y | \tau_A < \tau_B] \quad (6.3)$$

Finally,  $h(x) = \mathbb{P}_x(\tau_A < \tau_B)$  satisfies both

1.  $h(x) = 1$  for  $x \in A$  and  $h(z) = 0$  for  $z \in B$
2.  $h$  is harmonic at  $x$  for every  $x \notin A \cup B$

and  $h(\cdot)$  is the unique solution of linear system given by 1 and 2 above. (Good source)

## 7. Fundamental Algebra

### 7.1 Series

#### 7.1.1 Geometric

$$\sum_{n=a}^b r^n = r^a \frac{1 - r^{b-a+1}}{1 - r} = \frac{r^a - r^{b+1}}{1 - r} \quad (7.1)$$

$$\sum_{j=1}^{\infty} \frac{j^2}{\rho^j} = \frac{\rho(\rho+1)}{(\rho-1)^3} \quad (7.2)$$

#### 7.1.2 Quotient Stack

(source)

$$\sum_{k=1}^n \lfloor \frac{n}{k} \rfloor = 2 \sum_{k=1}^{\lfloor \sqrt{n} \rfloor} \lfloor \frac{n}{k} \rfloor - \lfloor \sqrt{n} \rfloor^2 \quad (7.3)$$

### 7.2 Inequalities

#### 7.2.1 Ratio of Summation

$$\frac{\sum_i a_i}{\sum_i b_i} \leq \max_i \frac{a_i}{b_i} \quad (7.4)$$

#### 7.2.2 Weighted AM-GM Inequality

Let the nonnegative numbers  $x_1, x_2, \dots, x_n$  and the nonnegative weights  $w_1, w_2, \dots, w_n$  be given. Set  $w = w_1 + w_2 + \dots + w_n$ . If  $w > 0$ , then the inequality

$$\frac{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}{w} \geq \sqrt[w]{x_1^{w_1} x_2^{w_2} \dots x_n^{w_n}} \quad (7.5)$$

holds.

#### 7.2.3 Titu's Lemma

**Summation Form.** For any real numbers  $a_1, a_2, \dots, a_n$  and positive reals  $b_1, b_2, b_3, \dots, b_n$ , we have

$$\frac{a_1^2}{b_1} + \frac{a_2^2}{b_2} + \dots + \frac{a_n^2}{b_n} \geq \frac{(a_1 + a_2 + \dots + a_n)^2}{b_1 + b_2 + \dots + b_n} \quad (7.6)$$

**Probabilistic Form.** Let  $X$  be a real random variable and  $Y$  be a positive random variable such that  $\mathbb{E}[|X|]$  and  $\mathbb{E}[Y]$  are well defined. Then

$$\mathbb{E}[X^2/Y] \geq \mathbb{E}[|X|^2]/\mathbb{E}[Y] \geq \mathbb{E}[X]^2/\mathbb{E}[Y] \quad (7.7)$$

### 7.3. Bounds and Approximations

---

#### 7.2.4 Cauchy-Schwarz Inequality

For any non-zero vector  $\mathbf{x}$ ,

$$\|\mathbf{x}\|_2^2 \leq \|\mathbf{x}\|_1^2 \leq \|\mathbf{x}\|_0 \|\mathbf{x}\|_2^2 \quad (7.8)$$

(Note: Useful in binary matrix multiplication.)

#### 7.2.5 Chebyshev's Sum Inequality

If  $a_1 \geq a_2 \geq \dots \geq a_n$  and  $b_1 \geq b_2 \geq \dots \geq b_n$ , then

$$\frac{1}{n} \sum_{k=1}^n a_k b_k \geq \left( \frac{1}{n} \sum_{k=1}^n a_k \right) \left( \frac{1}{n} \sum_{k=1}^n b_k \right) \quad (7.9)$$

#### 7.2.6 Symmetric Parametric Inequality

$$\left(1 - p + \frac{p}{x}\right)^{\alpha-2} \cdot (1 - p + px)^{\alpha-2} \geq 1 \quad (7.10)$$

#### 7.2.7 We love Jensen

By convexity of  $(u, v) \mapsto u^\lambda v^{1-\lambda}$

$$\mathbb{E}[U]^\lambda \mathbb{E}[V]^{1-\lambda} \leq \mathbb{E}[U^\lambda V^{1-\lambda}] \quad (7.11)$$

### 7.3 Bounds and Approximations

#### 7.3.1 Exponential Bound on Hyperbolic Ratio

For  $0 \leq y < x \leq 2$ ,

$$\frac{\sinh(x) - \sinh(y)}{\sinh(x - y)} \leq e^{\frac{1}{2}xy} \quad (7.12)$$

#### 7.3.2 Exponential Inequalities

For all  $t \in \mathbb{R}$  and  $0 \leq p \leq 1$ ,

$$1 - p + p \cdot x \leq e^{p(x-1)} \quad (7.13)$$

Additionally:

$$\forall x, y \geq 0 \quad \frac{1 + e^{x+y}}{e^x + e^y} \leq e^{xy/2} \quad (7.14)$$

$$\frac{1}{2}(e^x + e^{-x}) \leq e^{x^2/2} \quad (7.15)$$

#### 7.3.3 Softplus Quadratic Bound

For all  $a, x \in \mathbb{R}$  with  $a \neq 0$ , we have

$$\log(1 + e^x) \leq \log(1 + e^a) + \frac{x - a}{1 + e^{-a}} + \frac{(e^a - 1) \cdot (x - a)^2}{4 \cdot a \cdot (e^a + 1)} \quad (7.16)$$

#### 7.3.4 Miscellaneous upper Bounds

**Linear-Rational function.** For all  $x > 0$ ,  $t \in [0, 1]$ :

$$\frac{1}{1 - t + tx} \leq 1 - t(1 - t)(1 + 3t)(x - 1) + t^2 \left( (1 - t)x^2 + \frac{t}{x} - 1 \right) \quad (7.17)$$

**Inverse logarithmic.** For all  $u > 0$ :

$$\frac{1}{\log(1 + 1/u)} \leq u + \frac{1}{2} \quad (7.18)$$

$x \log_+ y$  **decoupling**. For non negative reals  $x$  and  $y$ ,

$$x \log_+ y + 1 \wedge +1 \leq x \log x + e^{-1}y + 1 \quad (7.19)$$

**Weighted reciprocal**. For all  $p \in [0, 1]$  and  $x \in (0, \infty)$ ,

$$\frac{1}{1 - p + p/x} \leq 1 - p + p \cdot x \quad (7.20)$$

### 7.3.5 Order of Rademacher Sums

Let  $\sigma \in \{-1, 1\}^n$  be a random Rademacher sequence and let  $a \in \mathbb{R}^n$  be an arbitrary real vector with sorted entries  $|a_1| \geq |a_2| \geq \dots \geq |a_n|$ . Then

$$\|\langle a, \sigma \rangle\|_{L^p} \sim \sum_{i \leq p} a_i + \sqrt{p} \left( \sum_{i > p} a_i^2 \right)^{1/2} \quad (7.21)$$

### 7.3.6 Sum Approximation

$$\max(a, b) \leq a + b \leq 2 \max(a, b), \quad a, b \geq 0 \quad (7.22)$$

## 7.4 Combinatorics

### 7.4.1 Expansion

$$\binom{2n}{m} = \sum_{j=0}^{\lfloor \frac{m}{2} \rfloor} \binom{n}{j} \binom{n-j}{m-2j} 2^{m-2j} \quad (7.23)$$

### 7.4.2 Vandermonde's Identity

$$\sum_{i=0}^r \binom{m}{i} \binom{n}{r-i} = \binom{n+m}{r} \quad (7.24)$$

$$\sum_{m=0}^n \binom{m}{j} \binom{n-m}{k-j} = \binom{n+1}{k+1} \quad (7.25)$$

Special form (**hockey-stick identity**):

$$\sum_{m=k}^n \binom{m}{k} = \binom{n+1}{k+1} \quad (7.26)$$

$$(\text{c.f. } \prod_{\ell=0}^{k-1} \binom{\ell+\eta}{\ell+1}) = \binom{k+\eta-1}{k})$$

## 8. Useful real analysis results

### 8.1 Leibniz's Integral Rule

Let  $\mu$  be a probability distribution with support  $\Omega$ , let  $I \subset \mathbb{R}$  be a nontrivial open interval, also let  $f : \Omega \times I \rightarrow \mathbb{R}$  be a map with the following properties:

1. For any  $x \in I$ ,  $\mathbb{E}_{w \sim \mu}[|f(w, x)|] < \infty$  (**Uniformly finite moment**)
2. For almost all  $w \in \Omega$ , the map  $x \mapsto f(w, x)$  is differentiable with derivative  $\frac{\partial}{\partial x} f(w, x)$  (**differentiability**)
3. There is a map  $h : \Omega \rightarrow \mathbb{R}$  with the property that  $\mathbb{E}_{w \sim \mu}[|h(w)|] < \infty$ , such that  $|\frac{\partial}{\partial x} f(\cdot, x)| \leq h$  (**derivative bound**).

Then, for any  $x \in I$ ,  $\mathbb{E}_{w \sim \mu} \left[ \left| \frac{\partial}{\partial x} f(w, x) \right| \right] < \infty$  and the function  $F : x \rightarrow \mathbb{E}_{w \sim \mu}[f(w, x)]$  is differentiable with derivative

$$F'(x) = \mathbb{E}_{w \sim \mu} \left[ \frac{\partial}{\partial x} f(w, x) \right] \quad (8.1)$$

## 9. Linaer Algebra

### 9.1 Determinant

#### 9.1.1 Expansion formula

For any  $A \subseteq \mathcal{Y}$ ,

$$\sum_{A \subseteq Y \subseteq \mathcal{Y}} \det(\mathbf{L}_Y) = \det(\mathbf{L} + \mathbf{I}_{\bar{A}}), \quad (9.1)$$

#### 9.1.2 Rearrangement

$$\sum_{(I', J') \in \mathcal{S}(I, J)} \det(\mathbf{Z}_{Y, I'}) \det(\mathbf{Z}_{Y, J'}) \leq \sum_{(I', *) \in \mathcal{S}(I, J)} \det(\mathbf{Z}_{Y, I'})^2 \quad (9.2)$$

where  $\mathbf{I}_{\bar{A}}$  is the diagonal matrix with ones in the diagonal positions with indices in  $\bar{A}$  and zeros elsewhere.

#### 9.1.3 Weinstein-Aronszajn identity

If  $\mathbf{A}$  and  $\mathbf{B}$  are matrices of size  $m \times n$  and  $n \times m$  respectively, given that  $\mathbf{AB}$  is of trass class, then

$$\det(\mathbf{I} + \mathbf{AB}) = \det(\mathbf{I} + \mathbf{BA}) \quad (9.3)$$

#### 9.1.4 DPP related

**Proposal matrix for NDPP.** Given  $\mathbf{V}, \mathbf{B}, \mathbf{D}$  such that  $\mathbf{L} = \mathbf{V}\mathbf{V}^\top + \mathbf{B}(\mathbf{D} - \mathbf{D}^\top)\mathbf{B}^\top$ , let  $\{\rho_i, \mathbf{v}_i\}_{i=1}^K$  be the eigendecomposition of  $\mathbf{V}\mathbf{V}^\top$  and  $\{(\sigma_j, \mathbf{y}_{2j-1}, \mathbf{y}_{2j})\}$  be the Youla decomposition of  $\mathbf{B}(\mathbf{D} - \mathbf{D}^\top)\mathbf{B}^\top$ . Denote  $\mathbf{Z} := [\mathbf{v}_1, \dots, \mathbf{v}_K, \mathbf{y}_1, \dots, \mathbf{y}_K] \in \mathbb{R}^{M \times 2K}$  and

$$\begin{aligned} \mathbf{X} &:= \text{diag} \left( \rho, \dots, \rho_K, \begin{bmatrix} 0 & \sigma_1 \\ -\sigma_1 & 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 & \sigma_{K/2} \\ -\sigma_{K/2} & 0 \end{bmatrix} \right), \\ \hat{\mathbf{X}} &:= \text{diag} \left( \rho, \dots, \rho_K, \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_1 \end{bmatrix}, \dots, \begin{bmatrix} \sigma_{K/2} & 0 \\ 0 & \sigma_{K/2} \end{bmatrix} \right), \end{aligned}$$

so that  $\mathbf{L} = \mathbf{Z}\mathbf{X}\mathbf{Z}^\top$  and  $\hat{\mathbf{L}} = \mathbf{Z}\hat{\mathbf{X}}\mathbf{Z}^\top$ . Then, for every subset  $\mathbf{Y} \subseteq [M]$ , it holds that

$$\det(\mathbf{L}_Y) \leq \det(\hat{\mathbf{L}}_Y) \quad (9.4)$$

and the equality holds when the size of  $\mathbf{Y}$  is equal to the rank of  $\mathbf{L}$ .

**Proposal matrix for NDPP II.** Given  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\mathbf{W}^A \in \mathbb{R}^{d \times d}$ . Then,

$$\det([\mathbf{X}\mathbf{W}^A\mathbf{X}^\top]_S) \leq \det([\mathbf{X}\hat{\mathbf{W}}^A\mathbf{X}^\top]_S) \quad (9.5)$$

for every  $S \subseteq [n]$ . In addition, equality holds when  $|S| \geq d$ .

**DPP probability expansion [1, Lemma 2.6].**

$$\mathbb{P}_{\hat{\mathbf{L}}}(Y) = \frac{\det(\hat{\mathbf{L}}_Y)}{\det(\hat{\mathbf{L}} + \mathbf{I})} = \sum_{E \subseteq [2K], |E|=|Y|} \det(\underbrace{\mathbf{Z}_{Y,E}\mathbf{Z}_{Y,E}^\top}_{\text{elementary DPP}}) \prod_{i \in E} \frac{\lambda_i}{\lambda_i + 1} \prod_{i \notin E} \frac{1}{\lambda_i + 1} \quad (9.6)$$

1. Choose an elementary DPP according to its mixture weight
2. Sample a subset from the selected elementary DPP

**DPP probability expansion II.** The probability of sampling  $S \in \binom{[n]}{k}$  from the  $k$ -DPP with  $\hat{\mathbf{L}}$  can be decomposed into the following

$$\frac{\det(\hat{\mathbf{L}}_S)}{e_k(\{\lambda_i\}_{i=1}^d)} = \sum_{E \in \binom{[d]}{k}} \frac{\prod_{i \in E} \lambda_i}{e_k(\{\lambda_i\}_{i=1}^d)} \cdot \det(\mathbf{K}_S^E) \quad (9.7)$$

where  $\mathbf{K}^E$  is a rank- $k$  projection matrix consisting of eigenvalues of  $\hat{\mathbf{L}}$ .



## 9.1.5 Ratio

Given that  $\det(\mathbf{Q}\mathbf{S}\mathbf{Q}^\top) \neq 0$

$$\frac{\det(\mathbf{Q}(\mathbf{S} + \mathbf{R})\mathbf{Q}^\top)}{\det(\mathbf{Q}\mathbf{S}\mathbf{Q}^\top)} \leq \det(\mathbf{I}_2 + (\mathbf{Q}\mathbf{S}\mathbf{Q}^\top)^{-1/2}\mathbf{Q}\mathbf{R}\mathbf{Q}^\top(\mathbf{Q}\mathbf{S}\mathbf{Q}^\top)^{-1/2}) \quad (9.8)$$

## 9.1.6 Inverse of trace

For an invertible matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\text{tr}(\mathbf{A}^{-1}) = \sum_{i=1}^n \det(\mathbf{A}_{-i}) / \det(\mathbf{A}), \quad (9.9)$$

where  $\mathbf{A}_{-i} \in \mathbb{R}^{(n-1) \times (n-1)}$  is the submatrix of  $\mathbf{A}$  where the  $i$ th row and column of  $\mathbf{A}$  are removed.

## 9.2 Vectorization

$$\text{vec}(\mathbf{A}\mathbf{X}\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{A})\text{vec}(\mathbf{X}) \quad (9.10)$$

**Lyapunov Equation.**

$$\mathbf{A}\mathbf{X} + \mathbf{X}\mathbf{B} = \mathbf{C} \quad (9.11)$$

$$\mathbf{A}\mathbf{X}\mathbf{I} + \mathbf{I}\mathbf{X}\mathbf{B} = \mathbf{C} \quad (9.12)$$

$$(\mathbf{I} \otimes \mathbf{A})\text{vec}(\mathbf{X}) + (\mathbf{B}^\top \otimes \mathbf{I})\text{vec}(\mathbf{X}) = \text{vec}(\mathbf{C}) \quad (9.13)$$

$$\text{vec}(\mathbf{X}) = (\mathbf{I} \otimes \mathbf{A} + \mathbf{B}^\top \otimes \mathbf{I})^{-1}\text{vec}(\mathbf{C}) \quad (9.14)$$

## 9.3 Trace

$$\text{tr}(\mathbf{A}\mathbf{B}\mathbf{C}) = \text{vec}(\mathbf{A}^\top)^\top (\mathbf{I} \otimes \mathbf{B})\text{vec}(\mathbf{C}) \quad (9.15)$$

$$\text{tr}(\mathbf{A}^\top \mathbf{B}\mathbf{C}\mathbf{D}^\top) = \text{vec}(\mathbf{A})^\top (\mathbf{D} \otimes \mathbf{B})\text{vec}(\mathbf{C}) \quad (9.16)$$

## 9.3.1 Von Neumann's trace inequality

**Theorem.** If  $\mathbf{A}, \mathbf{B}$  are complex  $n \times n$  matrices with singular values

$$\alpha_1 \geq \dots \geq \alpha_n, \quad \beta_1 \geq \dots \geq \beta_n, \quad (9.17)$$

respectively, then

$$|\text{tr}(\mathbf{A}\mathbf{B})| \leq \sum_{i=1}^n \alpha_i \beta_i \quad (9.18)$$

## 9.4 Inversion

## 9.4.1 Woodbury identity

$$[\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D}]^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}[\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B}]^{-1}\mathbf{D}\mathbf{A}^{-1} \quad (9.19)$$

given that  $\mathbf{A}^{-1}$  and  $\mathbf{C}^{-1}$  exist. If  $\mathbf{B} = \mathbf{x}, \mathbf{C} = \mathbf{I}, \mathbf{D} = \mathbf{y}^\top$

$$(\mathbf{A} + \mathbf{x}\mathbf{y}^\top)^{-1} = \mathbf{A}^{-1} - \frac{(\mathbf{A}^{-1}\mathbf{x})(\mathbf{y}\mathbf{A}^{-1})}{1 + \mathbf{y}^\top \mathbf{A}^{-1} \mathbf{x}} \quad (9.20)$$

## 9.4.2 Schur Complement

Schur Complement essentially is a block Cholesky factorization of a matrix

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix} \quad (9.21)$$

$\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$  is called the *Schur complement* of  $\mathbf{D}$ .

## 9.5 Hadamard Product

### 9.5.1 Quadratic Relation

$$\mathbf{x}^\top (\mathbf{A} \odot \mathbf{B}) \mathbf{y} = \text{tr}(\text{Diag}(\mathbf{x}) \mathbf{A} \text{Diag}(\mathbf{y}) \mathbf{B}^\top) \quad (9.22)$$

By setting  $\mathbf{x} = \mathbf{y}$ , it shows that the Hadamard product of two PSD matrices is PSD.

### 9.5.2 Rank Relation

$$\text{rank}(\mathbf{A} \odot \mathbf{B}) \leq \text{rank}(\mathbf{A}) \text{rank}(\mathbf{B}) \quad (9.23)$$

### 9.5.3 Spectrum Relation

$$\prod_{i=k}^n \lambda_i(\mathbf{A} \odot \mathbf{B}) \geq \prod_{i=k}^n \lambda_i(\mathbf{A} \mathbf{B}), \quad \forall k = 1, \dots, n \quad (9.24)$$

with  $\lambda_i(\cdot)$  denotes PD matrix.

### 9.5.4 Determinant

$$|\mathbf{A} \odot \mathbf{B}| \geq |\mathbf{A}| |\mathbf{B}| \quad (9.25)$$

## 9.6 Matrix Calculus

### 9.6.1 Matrix Chain rule

$$[\nabla_{\mathbf{X}} f(g(\mathbf{X}))]_{ij} = \sum_{k=1}^p \sum_{\ell=1}^q \frac{\partial f(G)}{\partial g_{k\ell}} \frac{\partial g_{k\ell}}{\partial x_{ij}} \quad (9.26)$$

### 9.6.2 Differentials

$$d(\text{tr } \mathbf{X}) = \text{tr } d\mathbf{X} \quad (9.27)$$

$$d(\mathbf{X} \otimes \mathbf{Y}) = (d\mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes (d\mathbf{Y}) \quad (9.28)$$

$$d\mathbf{X}^{-1} = -\mathbf{X}^{-1} \cdot d\mathbf{X} \cdot \mathbf{X}^{-1} \quad (9.29)$$

$$d(\det(\mathbf{X})) = \text{tr}(\text{adj}(\mathbf{X}) d\mathbf{X}) \quad (9.30)$$

$$d \det(\mathbf{X}) = \det(\mathbf{X}) \text{tr}(\mathbf{X}^{-1} d\mathbf{X}) \quad (9.31)$$

$$d \log(\det(\mathbf{X})) = \text{tr}(\mathbf{X}^{-1} d\mathbf{X}) \quad (9.32)$$

$$d\sigma(a) = (\text{Diag}(\sigma) - \text{Diag}(\sigma)^2) da \quad (9.33)$$

$$d(\text{softmax}(\theta)) = (\text{Diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top) d\theta \quad (9.34)$$

Note: Elementwise function은 일단 Diagonal 형태로 바꿔서 생각해 보삼 ㅋ

## 9.6.3 Useful first derivatives

$$\frac{\partial \operatorname{tr} \mathbf{X}}{\partial \mathbf{X}} = \mathbf{I} \quad (9.35)$$

$$\frac{\partial \operatorname{tr} \mathbf{X}^{-1}}{\partial \mathbf{X}} = -\mathbf{X}^{-2} \quad (9.36)$$

$$\frac{\partial \operatorname{tr} (\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A}^\top \quad (9.37)$$

$$\frac{\partial \operatorname{tr} (\mathbf{X}^k)}{\partial \mathbf{X}} = k \cdot (\mathbf{X}^\top)^{k-1} \quad (9.38)$$

$$\frac{\partial \operatorname{tr} (\mathbf{X}\mathbf{A}\mathbf{X}\mathbf{B})}{\partial \mathbf{X}} = \mathbf{B}^\top \mathbf{X}^\top \mathbf{A}^\top + \mathbf{A}^\top \mathbf{X}^\top \mathbf{B}^\top \quad (9.39)$$

$$\frac{\partial \mathbf{A}\mathbf{X}^{-1}\mathbf{B}}{\partial \mathbf{X}} = -\mathbf{X}^\top \mathbf{A}^\top \mathbf{B}^\top \mathbf{X}^{-\top} \quad (9.40)$$

$$\frac{\partial \log \det(\mathbf{X})}{\partial \mathbf{X}} = \mathbf{X}^{-\top} \quad (9.41)$$

$$\frac{\partial \det(\mathbf{X}^{-1})}{\partial \mathbf{X}} = \frac{\mathbf{X}^{-\top}}{\det \mathbf{X}} \quad (9.42)$$

$$\frac{\partial \det(\mathbf{X}^k)}{\partial \mathbf{X}} = k \det(\mathbf{X}^k) \mathbf{X}^{-\top} \quad (9.43)$$

$$\frac{\partial \log \det(\mathbf{X}\mathbf{X}^\top)}{\partial \mathbf{X}} = 2\mathbf{X}[\mathbf{X}^\top \mathbf{X}]^{-1} \cdot \det(\mathbf{X}\mathbf{X}^\top) \quad (9.44)$$

$$\frac{\partial \det(\mathbf{A}\mathbf{X}\mathbf{B})}{\partial \mathbf{X}} = \det(\mathbf{A}\mathbf{X}\mathbf{B}) \mathbf{A}^\top (\mathbf{A}\mathbf{X}\mathbf{B})^{-\top} \mathbf{B}^\top \quad (9.45)$$

## 9.6.4 Quadratic form

$$\frac{\partial (\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s})}{\partial \mathbf{s}} = -2\mathbf{A}^\top \mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s}) \quad (9.46)$$

$$\frac{\partial (\mathbf{x} - \mathbf{s})^\top \mathbf{W}(\mathbf{x} - \mathbf{s})}{\partial \mathbf{x}} = 2\mathbf{W}(\mathbf{x} - \mathbf{s}) \quad (9.47)$$

$$\frac{\partial (\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W}(\mathbf{s} - \mathbf{A}\mathbf{s})}{\partial \mathbf{x}} = 2\mathbf{W}(\mathbf{s} - \mathbf{A}\mathbf{s}) \quad (9.48)$$

$$\frac{\partial (\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s})}{\partial \mathbf{A}} = -2\mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s}) \mathbf{s}^\top \quad (9.49)$$

## 9.6.5 Hessian product rule

Given two functions  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ ,

$$H_c(fg) = (H_c f)g(c) + \nabla_c f^\top \nabla_c g + \nabla_c g^\top \nabla_c f + f(c)H_c g \quad (9.50)$$

## 9.6.6 Integration by parts

Given vector valued function  $\varphi$  and scalar function  $f$  with vanishing condition,

$$\int_{\mathbb{R}^d} \varphi(\mathbf{x}) \cdot \nabla f(\mathbf{x}) \, d\mathbf{x} = - \int_{\mathbb{R}^d} (\nabla \cdot \varphi(\mathbf{x})) f(\mathbf{x}) \, d\mathbf{x} \quad (9.51)$$

## 9.7 Eigenvalues and Eigenvectors

## 9.7.1 General Properties

Assume that  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and  $\mathbf{B} \in \mathbb{R}^{m \times n}$ ,

$$\operatorname{eig}(\mathbf{A}\mathbf{B}) = \operatorname{eig}(\mathbf{B}\mathbf{A}) \quad (9.52)$$

$$\operatorname{rank}(\mathbf{A}) = r \Rightarrow \text{At most } r \text{ non-zero } \lambda_i \quad (9.53)$$

## 9.7.2 Symmetric

Assume  $\mathbf{A}$  is symmetric, then

$$\mathbf{V}\mathbf{V}^\top = \mathbf{I} \quad (9.54)$$

$$\lambda_i \in \mathbb{R} \quad (9.55)$$

$$\text{tr}(\mathbf{A}^p) = \sum_i \lambda_i^p \quad (9.56)$$

$$\text{eig}(\mathbf{I} + c\mathbf{A}) = 1 + c\lambda_i \quad (9.57)$$

$$\text{eig}(\mathbf{A} - c\mathbf{I}) = \lambda_i - c \quad (9.58)$$

$$\text{eig}(\mathbf{A}^{-1}) = \lambda_i^{-1} \quad (9.59)$$

For a symmetric, positive matrix  $\mathbf{A}$

$$\text{eig}(\mathbf{A}^\top \mathbf{A}) = \text{eig}(\mathbf{A}\mathbf{A}^\top) = \text{eig}(\mathbf{A}) \circ \text{eig}(\mathbf{A}) \quad (9.60)$$

## 9.7.3 Singular Value Decomposition

Any  $n \times m$  matrix  $\mathbf{A}$  can be written as

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top \quad (9.61)$$

where

$$\mathbf{U} = \text{eigenvectors of } \mathbf{A}\mathbf{A}^\top \quad n \times n$$

$$\mathbf{D} = \sqrt{\text{diag}(\text{eig}(\mathbf{A}\mathbf{A}^\top))} \quad n \times m$$

$$\mathbf{V} = \text{eigenvectors of } \mathbf{A}^\top \mathbf{A} \quad m \times m$$

**Square decomposed into rectangular.** Assume  $\mathbf{V}_* \mathbf{D}_* \mathbf{U}_*^\top = 0$  then we can expand the SVD of  $\mathbf{A}$  into

$$\mathbf{A} = \left[ \begin{array}{c|c} \mathbf{V} & \mathbf{V}_* \end{array} \right] \left[ \begin{array}{c|c} \mathbf{D} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{D}_* \end{array} \right] \left[ \begin{array}{c} \mathbf{U}^\top \\ \hline \mathbf{U}_*^\top \end{array} \right] \quad (9.62)$$

where the SVD of  $\mathbf{A}$  is  $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$

## 9.7.4 LU decomposition

Assume  $\mathbf{A}$  is a square matrix with non-zero leading principal minors, then

$$\mathbf{A} = \mathbf{L}\mathbf{U} \quad (9.63)$$

where  $\mathbf{L}$  is a unique unit lower triangular matrix and  $\mathbf{U}$  is a unique upper triangular matrix.

## 9.7.5 Cholesky decomposition

Assume  $\mathbf{A}$  is a symmetric positive definite square matrix, then

$$\mathbf{A} = \mathbf{U}^\top \mathbf{U} = \mathbf{L}\mathbf{L}^\top \quad (9.64)$$

where  $\mathbf{U}$  is a unique upper triangular matrix and  $\mathbf{L}$  is a lower triangular matrix.

## 9.7.6 Eigenvalues of its reverse

**Proposition.** Given  $M \times K$  matrix  $\mathbf{A}, \mathbf{B}$ , the nonzero eigenvalues of  $\mathbf{A}\mathbf{B}^\top \in \mathbb{C}^{M \times M}$  and  $\mathbf{B}^\top \mathbf{A} \in \mathbb{C}^{K \times K}$  are identical. In addition, if  $(\lambda, \mathbf{v})$  is an eigenpair of  $\mathbf{B}^\top \mathbf{A}$  with  $\lambda \neq 0$ , then  $(\lambda, \mathbf{A}\mathbf{v} / \|\mathbf{A}\mathbf{v}\|_2)$  is an eigenpair of  $\mathbf{A}\mathbf{B}^\top$ .

## 9.7.7 Row stochastic matrix

**Fact.** The operator norm of a row-stochastic matrix is 1.

## 9.8 Inverses

### 9.8.1 Rank-1 update of the inverse of inner product

Denote  $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1}$  and that  $\mathbf{X}$  is extended to include a new column vector in the end  $\tilde{\mathbf{X}} = [\mathbf{X}, \mathbf{v}]$ , let  $N = \mathbf{v}^\top (\mathbf{I} - \mathbf{X} \mathbf{A} \mathbf{X}^\top) \mathbf{v}$  then

$$(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} = N^{-1} \begin{bmatrix} N\mathbf{A} + \mathbf{A} \mathbf{X}^\top \mathbf{v} (\mathbf{A} \mathbf{X}^\top \mathbf{v})^\top & -\mathbf{A} \mathbf{X}^\top \mathbf{v} \\ -\mathbf{v}^\top \mathbf{X} \mathbf{A}^\top & 1 \end{bmatrix} \quad (9.65)$$

### 9.8.2 Approximations

The following identity is known as the *Neuman series* of a matrix, which holds when  $|\lambda_i| < 1$  for all eigenvalues  $\lambda_i$

$$(\mathbf{I} - \mathbf{A})^{-1} = \sum_{n=0}^{\infty} \mathbf{A}^n \quad (9.66)$$

$$(\mathbf{I} + \mathbf{A})^{-1} = \sum_{n=0}^{\infty} (-1)^n \mathbf{A}^n \quad (9.67)$$

$$\mathbf{A} - \mathbf{A}(\mathbf{I} + \mathbf{A})^{-1}\mathbf{A} = \mathbf{A} - \mathbf{A}(\mathbf{I} + \mathbf{A}^{-1})^{-1} \quad (9.68)$$

$$= \mathbf{A}(\mathbf{I} - (\mathbf{I} + \mathbf{A}^{-1})^{-1}) \quad (9.69)$$

$$\approx \mathbf{A}(\mathbf{I} - \mathbf{I} + \mathbf{A}^{-1} - \mathbf{A}^{-2}) \quad (9.70)$$

$$= \mathbf{I} - \mathbf{A}^{-1} \quad (9.71)$$

### 9.8.3 Block matrix

Using Schur complements

$$\mathbf{C}_1 = \mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} \quad (9.72)$$

$$\mathbf{C}_2 = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} \quad (9.73)$$

as

$$\left[ \begin{array}{c|c} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \hline \mathbf{A}_{21} & \mathbf{A}_{22} \end{array} \right]^{-1} = \left[ \begin{array}{c|c} \mathbf{C}_1^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{C}_2^{-1} \\ \hline -\mathbf{C}_2^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} & \mathbf{C}_2^{-1} \end{array} \right] \quad (9.74)$$

## 9.9 PSD matrix

### 9.9.1 Decomposition

1. The matrix is PSD with rank  $r \iff$  there exists a matrix  $\mathbf{B}$  of rank  $r$  such that  $\mathbf{A} = \mathbf{B} \mathbf{B}^\top$
2. The matrix is PD  $\iff$  there exists an invertible matrix  $\mathbf{B}$  such that  $\mathbf{A} = \mathbf{B} \mathbf{B}^\top$
3. Given  $\mathbf{A}$  is an  $n \times n$  PSD matrix, there exists an  $n \times r$  matrix  $\mathbf{B}$  of rank  $r$  such that  $\mathbf{B}^\top \mathbf{A} \mathbf{B} = \mathbf{I}$ .

### 9.9.2 Sylvester's characterization

$$\mathbf{A} \succeq 0 \iff \text{All } 2^n - 1 \text{ principal minors are nonnegative.} \quad (9.75)$$

$$\mathbf{A} \succ 0 \iff \text{All } n \text{ leading principal minors are positive.} \quad (9.76)$$

### 9.9.3 Equation with zeros

Assume  $\mathbf{A}$  is PSD, then  $\mathbf{X}^\top \mathbf{A} \mathbf{X} = \mathbf{0} \Rightarrow \mathbf{A} \mathbf{X} = \mathbf{0}$

## 9.9.4 Rank of product

Assume  $\mathbf{A}$  is positive definite, then  $\text{rank}(\mathbf{B}\mathbf{A}\mathbf{B}^\top) = \text{rank}(\mathbf{B})$

## 9.9.5 Outer product

If  $\mathbf{X} \in n \times r$ , where  $n \leq r$  and  $\text{rank}(\mathbf{X}) = n$ , then  $\mathbf{X}\mathbf{X}^\top$  is positive definite.

## 9.9.6 Small perturbations

If  $\mathbf{A}$  is positive definite, and  $\mathbf{B}$  is symmetric, then  $\mathbf{A} - t\mathbf{B}$  is positive definite for sufficiently small  $t$ .

## 9.9.7 Hadamard inequality

If  $\mathbf{A}$  is a positive definite or semi-definite matrix, then

$$\det(\mathbf{A}) \leq \prod_i A_{ii} \quad (9.77)$$

## 9.9.8 Loewner order

**Fact.** Let  $\mathbf{A}$  and  $\mathbf{B}$  be hermitian positive definite. Then

$$\mathbf{A} \succeq \mathbf{B} \iff \mathbf{I} \succeq \mathbf{A}^{-1/2} \mathbf{B} \mathbf{A}^{-1/2} \quad (9.78)$$

## 9.9.9 Inverse of PSD

**Fact.** Suppose that  $\mathbf{A}, \mathbf{B}$  and  $\mathbf{A} - \mathbf{B}$  are all positive definite, then  $\mathbf{B}^{-1} - \mathbf{A}^{-1}$  is also positive definite.

## 9.10 Symmetric and skew-symmetric matrix

## 9.10.1 Properties of symmetric matrix

1. Every real symmetric matrix can be orthogonally diagonalizable.<sup>2</sup>
2. The rank of a symmetric matrix  $\mathbf{A}$  is equal to the number of non-zero eigenvalues of  $\mathbf{A}$ .
3. If  $\mathbf{A}$  and  $\mathbf{B}$  are  $n \times n$  real symmetric matrices that commute, then they can be simultaneously diagonalized by an orthogonal matrix.

## 9.10.2 Youla decomposition

Given  $\mathbf{B} \in \mathbb{R}^{M \times K}$  and  $\mathbf{D} \in \mathbb{R}^{K \times K}$ , consider a rank- $K$  skew-symmetric matrix  $\mathbf{B}^\top(\mathbf{D} - \mathbf{D}^\top)\mathbf{B}^\top$ . Then, we can write

$$\mathbf{B}(\mathbf{D} - \mathbf{D}^\top)\mathbf{B} = \sum_{j=1}^{K/2} i\sigma_j(\mathbf{a}_j + i\mathbf{b}_j)(\mathbf{a}_j + i\mathbf{b}_j)^H - i\sigma_j(\mathbf{a}_j - i\mathbf{b}_j)(\mathbf{a}_j - i\mathbf{b}_j)^H \quad (9.79)$$

$$= \sum_{j=1}^{K/2} 2\sigma_j(\mathbf{a}_j\mathbf{b}_j^\top - \mathbf{b}_j\mathbf{a}_j^\top) \quad (9.80)$$

$$= \sum_{j=1}^{K/2} [\mathbf{a}_j - \mathbf{b}_j \quad \mathbf{a}_j + \mathbf{b}_j] \begin{bmatrix} 0 & \sigma_j \\ -\sigma_j & 0 \end{bmatrix} \begin{bmatrix} \mathbf{a}_j^\top - \mathbf{b}_j^\top \\ \mathbf{a}_j^\top + \mathbf{b}_j^\top \end{bmatrix} \quad (9.81)$$

Note that  $\mathbf{a}_1 \pm \mathbf{b}_1, \dots, \mathbf{a}_{K/2} \pm \mathbf{b}_{K/2}$  are real-valued orthonormal vectors. The pair  $\{(\sigma_j, \mathbf{a}_j - \mathbf{b}_j, \mathbf{a}_j + \mathbf{b}_j)\}_{j=1}^{K/2}$  is often called the **Youla decomposition** of  $\mathbf{B}(\mathbf{D} - \mathbf{D}^\top)\mathbf{B}^\top$ .

---

<sup>2</sup>Think of it this way: every symmetric matrix can be triangulated and normality is preserved under a similar transform. When is the triangular matrix normal? Of course, it is the diagonal matrix.

## 9.11 Some techniques

### 9.11.1 Binary analysis

어떤 matrix의 operator norm을 분석하기 위해 matrix를 binary matrix로 decomposition하는 것은 유용할 수 있다.

#### Example

Let  $\mathbf{v}$  be the unit-normed vector that realizes the operator norm of  $\mathbf{D}^{-1}\mathbf{A}$ . We define the sequence of binary matrices  $\mathbf{B}^0, \mathbf{B}^1, \mathbf{B}^2$  as follows:

$$B_{i,j}^t := 1_{\{2^{-t-1}\sqrt{\alpha/n} < [\mathbf{D}^{-1}\mathbf{A}]_{i,j} \leq 2^{-t}\sqrt{\alpha/n}\}} \text{ for every integers } t \geq 0, \quad (9.82)$$

where  $\sqrt{\alpha/n}$  is the upper bound for entries of  $\mathbf{D}^{-1}\mathbf{A}$ . Then we have the following inequalities for the entries and the  $l_2$ -norm:

$$[\mathbf{D}^{-1}\mathbf{A}]_{i,j} \leq \sqrt{\alpha/n} \sum_{t=0}^{\infty} 2^{-t} \cdot [\mathbf{B}^t]_{i,j} \quad (9.83)$$

$$\|\mathbf{D}^{-1}\mathbf{A} \cdot \mathbf{v}\|_2 \leq \sqrt{\alpha/n} \sum_{t=0}^{\infty} 2^{-t} \cdot \|\mathbf{B}^t \mathbf{v}\|_2 \quad (9.84)$$

만약  $\mathbf{B}^t$  matrix의 row, column들의 non-zero elements를 estimate 하면  $\|\mathbf{B}^t \mathbf{v}\|_2^2$ 도 estimate 할 수 있고  $\mathbf{D}^{-1}\mathbf{A}$ 의 operator norm bound도 estimate 할 수 있다.

## Deferred proof

Section 9.7.7 Suppose  $\mathbf{Ax} = \lambda \mathbf{x}$  for some  $\lambda > 1$ . Since the rows of  $\mathbf{A}$  are nonnegative and sum to 1, each element of vector  $\mathbf{Ax}$  is a convex combination of the components of  $\mathbf{x}$ , which can be no greater than  $\mathbf{x}_{max}$ . On the other hand, at least one element of  $\lambda \mathbf{x}$  is greater than  $\mathbf{x}_{max}$ , which proves that  $\mathbf{x}_{max}$ , which shows that  $\lambda > 1$  is impossible.

## References

- [1] A. Kulesza and B. Taskar. Determinantal point processes for machine learning. *Found. Trends Mach. Learn.*, 5:123–286, 2012.
- [2] T. Steinke, M. Nasr, and M. Jagielski. Privacy auditing with one (1) training run. *ArXiv*, abs/2305.08846, 2023.