



# Machine Learning and Advanced Topics (Lecture Notes)

University of Copenhagen

**Yijun Bian**

SID: ??

Courses: ML, ATML

Department of Computer Science  
University of Copenhagen (DIKU)

Postdoc Training  
Prof. Yevgeny Seldin

23. November 2021

# Contents

<b>1</b>	<b>Preliminaries</b>	<b>1</b>
1.1	Self-Assessment Assignment for the Machine Learning course . . . . .	1
1.1.1	Vectors and matrices . . . . .	1
1.1.2	Derivatives . . . . .	3
1.1.3	Probability Theory: Sample Space . . . . .	3
1.1.4	Probability Theory: Properties of Expectation . . . . .	4
1.1.5	Probability Theory: Complements of Events . . . . .	8
1.1.6	Induction . . . . .	8
1.1.7	Programming . . . . .	9
1.2	ATML Self-Preparation Assignment . . . . .	9
1.2.1	Illustration of Hoeffding's Inequality . . . . .	9
1.2.2	The effect of scale (range) and normalization of random variables in Hoeffding's inequality . . . . .	9
1.2.3	Probability in Practice . . . . .	9
1.2.4	Occam's Razor . . . . .	9
1.2.5	Train-Validation Split Trade-off . . . . .	9
1.2.6	Optional Assignment Covering the VC-Analysis . . . . .	9
1.2.7	The growth function . . . . .	9
1.2.8	VC-dimension . . . . .	9
1.2.9	Airline Revisited . . . . .	9
<b>2</b>	<b>Chapter 2. Concentration of Measure Inequalities</b>	<b>10</b>
2.1	Markov's Inequality . . . . .	10
2.2	Chebyshev's Inequality . . . . .	10
2.3	Hoeffding's Inequality . . . . .	11
2.3.1	Understanding Hoeffding's Inequality . . . . .	14
2.3.2	Example Appendix . . . . .	16
2.4	Basics of Information Theory: Entropy, Relative Entropy, and the Method of Types	17
2.5	kl Inequality . . . . .	20
2.5.1	Relaxations of the kl-inequality: Pinsker's and refined Pinsker's equations .	21
2.6	Sampling Without Replacement . . . . .	23
<b>3</b>	<b>Appendix B: Probability Theory Basics</b>	<b>25</b>

# 1 Preliminaries

## 1.1 Self-Assessment Assignment for the Machine Learning course

Needed:

- basic linear algebra
- calculus
- probability theory

### 1.1.1 Vectors and matrices

**Problem 1** (Vectors and matrices). Consider the two vectors  $\mathbf{a} = (1, 2, 2)^\top$  and  $\mathbf{b} = (3, 2, 1)^\top$  and the matrix

$$M = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} \quad (1.1)$$

**A 1–3:** No, the inner product  $\mathbf{a}^\top \mathbf{b}$  (i.e., scalar product, dot product) and the outer product  $\mathbf{a} \mathbf{b}^\top$ , they are not equal.

$$\langle \mathbf{a}, \mathbf{b} \rangle \stackrel{\text{def}}{=} \mathbf{a}^\top \mathbf{b} \stackrel{\text{def}}{=} \mathbf{a} \cdot \mathbf{b} = 1 \times 3 + 2 \times 2 + 2 \times 1 = 3 + 2 + 2 = 7, \quad (1.2a)$$

$$\|\mathbf{a}\| = \sqrt{1^2 + 2^2 + 2^2} = \sqrt{1 + 4 + 4} = \sqrt{9} = 3, \quad (1.2b)$$

$$\|\mathbf{b}\| = \sqrt{3^2 + 2^2 + 1^2} = \sqrt{9 + 4 + 1} = \sqrt{14}, \quad (1.2c)$$

$$\mathbf{a} \mathbf{b}^\top = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} \begin{pmatrix} 3 & 2 & 1 \end{pmatrix} = \begin{pmatrix} 3 & 2 & 1 \\ 6 & 4 & 2 \\ 6 & 4 & 2 \end{pmatrix}. \quad (1.2d)$$

$$(1.2e)$$

**A4:** The inner product is symmetric ( $\mathbf{a}^\top \mathbf{b} = \mathbf{b}^\top \mathbf{a}$ ), but the outer product is not ( $\mathbf{a} \mathbf{b}^\top \neq \mathbf{b} \mathbf{a}^\top$ ).

$$\mathbf{a} \mathbf{b}^\top = \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} \begin{pmatrix} 3 & 2 & 1 \end{pmatrix} = \begin{pmatrix} 3 & 2 & 1 \\ 6 & 4 & 2 \\ 6 & 4 & 2 \end{pmatrix}, \quad (1.3a)$$

$$\mathbf{a}^\top \mathbf{b} = \begin{pmatrix} 1 & 2 & 2 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix} = 1 \times 3 + 2 \times 2 + 2 \times 1 = 3 + 4 + 2 = 9, \quad (1.3b)$$

$$\mathbf{b} \mathbf{a}^\top = \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 2 \end{pmatrix} = \begin{pmatrix} 3 & 6 & 6 \\ 2 & 4 & 4 \\ 1 & 2 & 2 \end{pmatrix}, \quad (1.3c)$$

$$\mathbf{b}^\top \mathbf{a} = \begin{pmatrix} 3 & 2 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} = 3 \times 1 + 2 \times 2 + 1 \times 2 = 3 + 4 + 2 = 9. \quad (1.3d)$$

**Q5:** How to find inverse matrix?:

(1)  $(A \ I)$ , (2) elementary row transformation  $(I \ B)$ , (3)  $B = A^{-1}$ . e.g.,

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, \quad A^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}. \quad (1.4)$$

**A5:**

$$M^{-1}M = MM^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} M^{-1} = I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}. \quad (1.5)$$

$$(MI) = \begin{pmatrix} 1 & 0 & 0 & \vdots & 1 & 0 & 0 \\ 0 & 4 & 0 & \vdots & 0 & 1 & 0 \\ 0 & 0 & 2 & \vdots & 0 & 0 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 & \vdots & 1 & 0 & 0 \\ 0 & 1 & 0 & \vdots & 0 & \frac{1}{4} & 0 \\ 0 & 0 & 1 & \vdots & 0 & 0 & \frac{1}{2} \end{pmatrix} = (IM^{-1}), \quad (1.6a)$$

$$M^{-1}M = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{4} & 0 \\ 0 & 0 & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = I, \quad (1.6b)$$

$$MM^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{4} & 0 \\ 0 & 0 & \frac{1}{2} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = I. \quad (1.6c)$$

**A6:**

$$M\mathbf{a} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 \\ 8 \\ 4 \end{pmatrix} \quad (1.7)$$

**A7:** No, the matrix is not symmetric ( $A \neq A^T$ ).

$$A = \mathbf{a}\mathbf{b}^T = \begin{pmatrix} 3 & 2 & 1 \\ 6 & 4 & 2 \\ 6 & 4 & 2 \end{pmatrix}, \quad (1.8a)$$

$$A^T = \begin{pmatrix} 3 & 6 & 6 \\ 2 & 4 & 4 \\ 1 & 2 & 2 \end{pmatrix} = \mathbf{b}\mathbf{a}^T = (\mathbf{a}\mathbf{b}^T)^T. \quad (1.8b)$$

**A8:** The rank of  $A$  is 1.

$$(AI) = \begin{pmatrix} 3 & 2 & 1 & \vdots & 1 & 0 & 0 \\ 6 & 4 & 2 & \vdots & 0 & 1 & 0 \\ 6 & 4 & 2 & \vdots & 0 & 0 & 1 \end{pmatrix} \sim \begin{pmatrix} 3 & 2 & 1 & | & 1 & 0 & 0 \\ 0 & 0 & 0 & | & -2 & 1 & 0 \\ 0 & 0 & 0 & | & 0 & -1 & 1 \end{pmatrix}, \quad (1.9a)$$

$$\begin{pmatrix} A \\ I \end{pmatrix} = \begin{pmatrix} 3 & 2 & 1 \\ 6 & 4 & 2 \\ 6 & 4 & 2 \\ \cdots & \cdots & \cdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 \\ 2 & 0 & 0 \\ 2 & 0 & 0 \\ \cdots & \cdots & \cdots \\ \frac{1}{3} & -\frac{2}{3} & -\frac{1}{3} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}. \quad (1.9b)$$

**A9:** A square matrix is invertible if and only if its rank is equal to the number of its columns/rows.  $A = \mathbf{a}\mathbf{b}^T$  is not invertible.

**A10:** The projection of vector **a** onto vector **b** is

$$\|\vec{a}\| = \sqrt{1^2 + 2^2 + 2^2} = \sqrt{1 + 4 + 4} = \sqrt{9} = 3, \quad (1.10a)$$

$$\|\vec{b}\| = \sqrt{3^2 + 2^2 + 1^2} = \sqrt{9 + 4 + 1} = \sqrt{14}, \quad (1.10b)$$

$$\cos(\theta) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \cdot \|\vec{b}\|} = \frac{(1, 2, 2) \cdot (3, 2, 1)}{\sqrt{1 + 4 + 4} \sqrt{9 + 4 + 1}} = \frac{3 + 4 + 2}{\sqrt{9} \sqrt{14}} = \frac{9}{3\sqrt{14}} = \frac{3}{\sqrt{14}}, \quad (1.10c)$$

$$\Pi_{\vec{b}}(\vec{a}) = \|\vec{a}\| \cos(\theta) \vec{e}_b = 3 \frac{3}{\sqrt{14}} \frac{(3, 2, 1)}{\sqrt{14}} = \frac{9}{14} (3, 2, 1) \quad (1.10d)$$

$$= \frac{\vec{a} \cdot \vec{b}}{\vec{b} \cdot \vec{b}} \vec{b} = \frac{3 + 4 + 2}{9 + 4 + 1} (3, 2, 1) = \frac{9}{14} (3, 2, 1) \quad (1.10e)$$

$$= \left( \vec{a} \cdot \frac{\vec{b}}{\|\vec{b}\|} \right) \frac{\vec{b}}{\|\vec{b}\|} = \|\vec{a}\| \cos(\theta) \frac{\vec{b}}{\|\vec{b}\|} = \frac{\vec{a} \cdot \vec{b}}{\vec{b} \cdot \vec{b}} \vec{b} \quad (1.10f)$$

### 1.1.2 Derivatives

**Problem 2** (Derivatives). We denote the derivative of a univariate function  $f(x)$  with respect to the variable  $x$  by  $\frac{df(x)}{dx}$ . We denote the partial derivative of a multivariate function  $f(x_1, \dots, x_n)$  with respect to the variable  $x_i$ , where  $1 \leq i \leq n$ , by  $\frac{\partial f(x_1, \dots, x_n)}{\partial x_i}$ . The partial derivative  $\frac{\partial f(x_1, \dots, x_n)}{\partial x_i}$  is the derivative of  $f(x_1, \dots, x_n)$  with respect to  $x_i$  when we treat all other variables  $x_j$  for  $j \neq i$  in  $f$  as constants.

Please recall the basic rules for derivatives, namely the sum rule, the chain rule, and the product rule.

**Q1:** What is the derivative of  $f(x) = \frac{1}{1 + \exp(-x)}$  with respect to  $x$ ?

**Q2:** What is the partial derivative of  $f(w, x) = 2(wx + 5)^2$  with respect to  $w$ ?

**A 1–2:**

$$\frac{df(x)}{dx} = (1 + e^{-x})^{-1} = -(1 + e^{-x})^{-2} e^{-x} (-1) = \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^x}. \quad (1.11a)$$

$$\frac{\partial f(w, x)}{\partial x} = 2 \cdot 2(wx + 5)^1 w = 4w(wx + 5), \quad (1.11b)$$

$$\frac{\partial f(w, x)}{\partial w} = 2 \cdot 2(wx + 5)^1 x = 4x(wx + 5). \quad (1.11c)$$

### 1.1.3 Probability Theory: Sample Space

**Problem 3** (Probability theory: sample space). An urn contains five red, three orange, and one blue ball. Two balls are randomly selected (without replacement).

**A1:** The sample space is balls in the urn, that is, 5 red, 3 orange, and 1 blue ball.

**A2:**

$$\mathbf{P}_{\text{red}}(x) = \frac{5}{5 + 3 + 1} = \frac{5}{9}, \quad (1.12a)$$

$$\mathbf{P}_{\text{orange}}(x) = \frac{3}{5 + 3 + 1} = \frac{3}{9} = \frac{1}{3}, \quad (1.12b)$$

$$\mathbf{P}_{\text{blue}}(x) = \frac{1}{5 + 3 + 1} = \frac{1}{9}. \quad (1.12c)$$

**A 3–4:** The possible values of  $X$  are 0, 1, 2, and 3. ( $X = 3$  will never be true.)

$$\mathbf{P}[X = 3] = 0, \quad (1.13a)$$

$$\mathbf{P}[X = 2] = \frac{C_3^2}{C_9^2} = \frac{3}{1} \cdot \frac{2 \times 1}{9 \times 8} = \frac{1}{12}, \quad (1.13b)$$

$$\begin{aligned} \mathbf{P}[X = 1] &= \frac{C_3^1 C_5^1 + C_3^1 C_1^1}{C_9^2} = \frac{3 \times 5 + 3 \times 1}{1} \cdot \frac{2 \times 1}{9 \times 8} = \frac{6}{3 \times 4} = \frac{1}{2} \\ &= \frac{C_3^1 C_{5+1}^1}{C_9^2} = \frac{3 \times 6}{1} \cdot \frac{2 \times 1}{9 \times 8} = \frac{1}{2}, \end{aligned} \quad (1.13c)$$

$$\begin{aligned} \mathbf{P}[X = 0] &= 1 - 0 - \frac{1}{12} - \frac{1}{2} = 1 - \frac{7}{12} = \frac{5}{12} \\ &= \frac{C_1^0 C_5^2 + C_1^1 C_5^1}{C_9^2} = \frac{1 \times \frac{5 \times 4}{2 \times 1} + 1 \times 5}{1} \cdot \frac{2 \times 1}{9 \times 8} = \frac{10 + 5}{1} \cdot \frac{1}{9 \times 4} = \frac{15}{36} = \frac{5}{12}. \end{aligned} \quad (1.13d)$$

Table 1: Probability of orange balls. Let  $Y$  represent the number of selected red balls. There are  $C_9^2 = \frac{9 \times 8}{2 \times 1} = 36$  probabilities.

(a) No blue ball								(b) One blue ball							
$X \setminus Y$	0	1	2	3	4	5		$X \setminus Y$	0	1	2	3	4	5	
0	0	0	10	0	0	0	10	0	0	5	0	0	0	0	5
1	0	15	0	0	0	0	15	1	3	0	0	0	0	0	3
2	3	0	0	0	0	0	3	2	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	3	0	0	0	0	0	0	0
	3	15	10	0	0	0	28		3	5	0	0	0	0	8

**A 5–6:**

$$\mathbf{E}[X] = \sum_{x=0}^3 x \mathbf{P}[X = x] = 3 \times 0 + 2 \times \frac{1}{12} + 1 \times \frac{1}{2} + 0 \times \frac{5}{12} = 0 + \frac{1}{6} + \frac{1}{2} + 0 = \frac{4}{6} = \frac{2}{3}, \quad (1.14a)$$

$$\begin{aligned} \mathbf{V}[X] &= \mathbf{E}[(X - \mu)^2] = \sum_{x=0}^3 (x - \mu)^2 \mathbf{P}(X = x) \\ &= \left(-\frac{2}{3}\right)^2 \frac{5}{12} + \left(\frac{1}{3}\right)^2 \frac{1}{2} + \left(\frac{4}{3}\right)^2 \frac{1}{12} + \left(\frac{7}{3}\right)^2 \cdot 0 \\ &= \frac{5}{9 \times 3} + \frac{1}{9 \times 2} + \frac{4}{27} + 0 = \frac{9}{27} + \frac{1}{18} = \frac{1}{3} + \frac{1}{18} = \frac{7}{18}. \end{aligned} \quad (1.14b)$$

### 1.1.4 Probability Theory: Properties of Expectation

#### Basic Theorems [ref](#)

**Remark 1 (Linear Function).** If  $Y = aX + b$ , where  $a$  and  $b$  are finite constants, then

$$\mathbf{E}(Y) = a\mathbf{E}(X) + b. \quad (1.15)$$

*Proof.*

$$\begin{aligned}\mathbf{E}(Y) &= \mathbf{E}(aX + b) = \int_{-\infty}^{\infty} (ax + b)f(x)dx \\ &= a \int_{-\infty}^{\infty} xf(x)dx + b \int_{-\infty}^{\infty} f(x)dx = a\mathbf{E}(X) + b.\end{aligned}\quad (1.16)$$

□

**Corollary 1.1.** *If  $X = c$  with probability 1, then  $\mathbf{E}(X) = c$ .*

*Proof.*

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} cf(x)dx = c \int_{-\infty}^{\infty} f(x)dx = c. \quad (1.17)$$

□

**Theorem 1.2.** *If there exists a constant such that  $\mathbf{Pr}(X \geq a) = 1$ , then  $\mathbf{E}(X) \geq a$ . If there exists a constant  $b$  such that  $\mathbf{Pr}(X \leq b) = 1$ , then  $\mathbf{E}(X) \leq b$ .*

*Proof.*

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} xf(x)dx = 0 + \int_a^{\infty} xf(x)dx \geq \int_a^{\infty} af(x)dx = a \mathbf{Pr}(X \geq a) = a, \quad (1.18)$$

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} xf(x)dx = \int_{-\infty}^b xf(x)dx + 0 \leq \int_{-\infty}^b bf(x)dx = b \mathbf{Pr}(X \leq b) = b. \quad (1.19)$$

□

**Theorem 1.3.** *Suppose that  $\mathbf{E}(x) = a$  and that either  $\mathbf{Pr}(X \geq a) = 1$  or  $\mathbf{Pr}(X \leq a) = 1$ . Then  $\mathbf{Pr}(X = a) = 1$ .*

*Proof.* **I don't know why?** Discrete situations here. (1) When  $\mathbf{Pr}(X \geq a) = 1$ , suppose  $x_1, x_2, \dots$  are all situations of  $x > a$ , then  $\mathbf{Pr}(X = x) > 0$ . Let  $p_0 = \mathbf{Pr}(X = a)$ , then

$$\mathbf{E}(X) = p_0a + \sum_{j=1}^{\infty} x_j \mathbf{Pr}(X = x_j) \geq p_0a + \sum_{j=1}^{\infty} a \mathbf{Pr}(X = x_j) = a \quad (1.20)$$

(2) When  $\mathbf{Pr}(X \leq a) = 1$ , suppose  $\dots, x_{-2}, x_{-1}$  are all situations of  $x < a$ , then  $\mathbf{Pr}(X = x) > 0$ . Let  $p_0 = \mathbf{Pr}(X = a)$ , then

$$\mathbf{E}(X) = p_0a + \sum_{j=-\infty}^{-1} x_j \mathbf{Pr}(X = x_j) \leq p_0a + \sum_{j=-\infty}^{-1} a \mathbf{Pr}(X = x_j) = a \quad (1.21)$$

□

**Theorem 1.4.** *If  $X_1, \dots, X_n$  are  $n$  random variables such that each expectation  $\mathbf{E}(X_i)$  is finite ( $i = 0, \dots, n$ ), then*

$$\mathbf{E}(X_1 + \dots + X_n) = \mathbf{E}(X_1) + \dots + \mathbf{E}(X_n). \quad (1.22)$$

*Proof.* Two continuous variables

$$\begin{aligned}
\mathbf{E}(X_1 + X_2) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_1 + x_2) f(x_1, x_2) dx_1 dx_2 \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1, x_2) dx_1 dx_2 + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_2 f(x_1, x_2) dx_1 dx_2 \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1, x_2) dx_2 dx_1 + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_2 f(x_1, x_2) dx_1 dx_2 \tag{1.23a}
\end{aligned}$$

$$= \int_{-\infty}^{\infty} x_1 \left[ \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 \right] dx_1 + \int_{-\infty}^{\infty} x_2 \left[ \int_{-\infty}^{\infty} f(x_1, x_2) dx_1 \right] dx_2 \tag{1.23b}$$

$$\begin{aligned}
&= \int_{-\infty}^{\infty} x_1 f_1(x_1) dx_1 + \int_{-\infty}^{\infty} x_2 f_2(x_2) dx_2 \\
&= \mathbf{E}(X_1) + \mathbf{E}(X_2). \tag{1.23c}
\end{aligned}$$

□

**Corollary 1.5.** Assume that  $\mathbf{E}(x_i)$  is finite for  $i = 1, \dots, n$ . For all constants  $a_1, \dots, a_n$  and  $b$ ,

$$\mathbf{E}(a_1 X_1 + \dots + a_n X_n + b) = a_1 \mathbf{E}(X_1) + \dots + a_n \mathbf{E}(X_n) + b. \tag{1.24}$$

**Definition 1.6** (Convex Functions). A function  $g$  of a vector argument is convex if, for every  $\alpha \in (0, 1)$ , and every  $x$  and  $y$ ,

$$g[\alpha x + (1 - \alpha)y] \geq \alpha g(x) + (1 - \alpha)g(y). \tag{1.25}$$

**Theorem 1.7** (Jensen's Inequality). Let  $g$  be a convex function, and let  $X$  be a random vector with finite mean. Then  $\mathbf{E}[g(X)] \geq g(\mathbf{E}[X])$ . The equality is satisfied/met if and only if  $g$  is a linear function.

**Remark 2.** If  $X_1, \dots, X_n$  are  $n$  independent random variables such that each expectation  $\mathbf{E}(X_i)$  is finite ( $i = 1, \dots, n$ ), then

$$\mathbf{E}(\prod_{i=1}^n X_i) = \prod_{i=1}^n \mathbf{E}(X_i). \tag{1.26}$$

*Proof.* The p.d.f. are  $f_i$  for the variable  $X_i$  in this group, then

$$f(x_1, \dots, x_i) = \prod_{i=1}^n f_i(x_i), \tag{1.27}$$

therefore,

$$\begin{aligned}
\mathbf{E}(\prod_{i=1}^n X_i) &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (\prod_{i=1}^n x_i) f(x_1, \dots, x_n) dx_1, \dots, x_n \\
&= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \prod_{i=1}^n [x_i f_i(x_i)] dx_1, \dots, x_n \tag{1.28a}
\end{aligned}$$

$$\begin{aligned}
&= \mathbf{E}(X_1) \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \prod_{i=2}^n [x_i f_i(x_i)] dx_2, \dots, x_n \\
&= \mathbf{E}(X_1) \mathbf{E}(X_2) \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \prod_{i=3}^n [x_i f_i(x_i)] dx_3, \dots, x_n \\
&= \prod_{i=1}^n \int_{-\infty}^{\infty} x_i f_i(x_i) dx_i = \prod_{i=1}^n \mathbf{E}(X_i). \tag{1.28b}
\end{aligned}$$

□

**Problem 4** (Probability theory: properties of expectation). Let  $X$  and  $Y$  be two discrete random variables taking values in  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Starting from the definitions, prove the following identities:



1.  $\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y]$
2. If  $X$  and  $Y$  are independent then  $\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]$ . (Mark the step where you are using the independence assumption. Note that this assumption was not required in point 1.)
3. Provide an example of two random variables  $X$  and  $Y$  for which  $\mathbf{E}[XY] \neq \mathbf{E}[X]\mathbf{E}[Y]$ . (Describe how you define the random variables, provide a joint probability distribution table [see comment below], and calculate  $\mathbf{E}[XY]$  and  $\mathbf{E}[X]\mathbf{E}[Y]$ .)
4.  $\mathbf{E}[\mathbf{E}[X]] = \mathbf{E}[X]$ .
5. Variance of a random variable is defined as  $\mathbf{V}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2]$ . Show that  $\mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2$ .

**A0:** We have already known that

$$\mathbf{E}[X] = \sum_{k=1}^{\infty} x_k p_k, \quad \mathbf{E}[X] = \int_{-\infty}^{\infty} x f(x) dx, \quad (1.29a)$$

$$\mathbf{E}[Y] = \sum_{k=1}^{\infty} y_k p_k, \quad \mathbf{E}[Y] = \int_{-\infty}^{\infty} y f(y) dy, \quad (1.29b)$$

$$\mathbf{E}[X + Y] = . \quad (1.29c)$$

**A1:** therefore,

$$\begin{aligned} \mathbf{E}[X + Y] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + y) f(x + y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x + y) dy dx + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(x + y) dx dy \end{aligned} \quad (1.30a)$$

$$= \int_{-\infty}^{\infty} x f_x(x) dx + \int_{-\infty}^{\infty} y f_y(y) dy \quad (1.30b)$$

$$= \mathbf{E}[X] + \mathbf{E}[Y]. \quad (1.30c)$$

**A2:** therefore,

$$\begin{aligned} \mathbf{E}[XY] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (xy) f(xy) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_x(x) f_y(y) dx dy \end{aligned} \quad (1.31a)$$

$$= \int_{-\infty}^{\infty} x f_x(x) dx \int_{-\infty}^{\infty} y f_y(y) dy \quad (1.31b)$$

$$= \mathbf{E}[X]\mathbf{E}[Y]. \quad (1.31c)$$

**A4:**

$$\mathbf{E}[X] = \text{a constant} \quad (1.32a)$$

$$\mathbf{E}[\mathbf{E}[X]] = \int_{-\infty}^{\infty} \mathbf{E}[X] \mathbf{Pr}(x) dx = \mathbf{E}[X] \int_{-\infty}^{\infty} \mathbf{Pr}(x) dx = \mathbf{E}[X] \quad (1.32b)$$

**A5:**

$$\begin{aligned} \mathbf{V}[X] &= \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2 - 2X\mathbf{E}[X] + (\mathbf{E}[X])^2] \\ &= \mathbf{E}[X^2] - 2\mathbf{E}[X]\mathbf{E}[X] + (\mathbf{E}[X])^2 = \mathbf{E}[X^2] - (\mathbf{E}[X])^2. \end{aligned} \quad (1.33)$$

### 1.1.5 Probability Theory: Complements of Events

**Problem 5** (1). The complement of event  $A$  is denoted by  $\bar{A}$  and defined by  $\bar{A} = \Omega \setminus A$ . Starting from probabilities axioms prove that  $\Pr\{A\} = 1 - \Pr\{\bar{A}\}$ .

**Problem 6** (2). In many cases it is easier to calculate the probability of a complement of an event than to calculate the probability of the event itself. Use this to solve the following question. We flip a fair coin 10 times.

- What is the probability to observe at least one tail?
- What is the probability to observe at least two tails?

**A2:**

$$\Pr[x \geq 1] = 1 - \Pr[x = 0] = 1 - \frac{1^{10}}{2^{10}} = \frac{1023}{1024} \quad (1.34)$$

$$\Pr[x \geq 2] = 1 - \Pr[x = 0] - \Pr[x = 1] = 1 - \frac{1^{10}}{2^{10}} - A_{10}^1 \frac{1^9}{2^{10}} = 1 - \frac{1}{1024} - \frac{10}{1024} = \frac{1024 - 11}{1024} = \frac{1013}{1024} \quad (1.35)$$

### 1.1.6 Induction

**Problem 7.** Prove by induction that for all integer  $n$  and  $d$ , such that  $n \geq d \geq 1$ :

$$\sum_{i=0}^d \binom{n}{i} \leq n^d + 1. \quad (1.36)$$

*Proof.* **I don't know why?** (a1) When  $d = 1$ ,

$$Eq.^{\text{left}} = \sum_{i=0}^1 \binom{n}{i} = \binom{n}{0} + \binom{n}{1} = \frac{n!}{n!0!} + \frac{n!}{(n-1)!1!} = 1 + n = n^d + 1 = Eq.^{\text{right}} \quad (1.37)$$

(a2) When  $d = j$ , the inequality is met, then for  $d = j + 1$ ,

$$Eq.^{\text{left}} = \sum_{i=0}^{j+1} \binom{n}{i} = \sum_{i=0}^j \binom{n}{i} + \binom{n}{j+1} \leq n^j + 1 + \frac{n!}{(n-j-1)!(j+1)!}, \quad (1.38a)$$

$$Eq.^{\text{right}} = n^{j+1} + 1 = \quad (1.38b)$$

(a3) In conclusion.

(b1) When  $n = 1$ , then  $d = 1$  as well,

$$Eq.^{\text{left}} = \sum_{i=0}^1 \binom{1}{i} = \binom{1}{0} + \binom{1}{1} = \frac{1!}{1!0!} + \frac{1!}{0!1!} = 1 + 1 = n^d + 1 = Eq.^{\text{right}} \quad (1.39)$$

(b2) When  $n = t$ , the inequality is met, then for  $n = t + 1$ ,

$$\begin{aligned} Eq.^{\text{left}} &= \sum_{i=0}^d \binom{t+1}{i} = \sum_{i=0}^d \frac{(t+1)!}{(t+1-i)!i!} = \sum_{i=0}^d \frac{t+1}{t+1-i} \frac{t!}{(t-i)!i!} = \sum_{i=0}^d \binom{t}{i} \frac{t+1}{t+1-i} \\ &\leq (t^d + 1) \frac{t+1}{t+1-i} \end{aligned} \quad (1.40a)$$

$$Eq.^{\text{right}} = (t+1)^d + 1 = \quad (1.40b)$$

(b3) In conclusion. □

### 1.1.7 Programming

## 1.2 ATML Self-Preparation Assignment

### 1.2.1 Illustration of Hoeffding's Inequality

#### Q1:

1. Make 1,000,000 repetitions of the experiment of drawing 20 i.i.d. Bernoulli random variables  $X_1, \dots, X_{20}$  (20 coins) with bias  $\frac{1}{2}$ .
2. Plot the empirical frequency of observing  $\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha$  for  $\alpha \in \{0.5, 0.55, 0.6, \dots, 0.95, 1\}$ .
3. Explain why the above granularity of  $\alpha$  is sufficient. I.e., why, e.g., taking  $\alpha = 0.51$  will not provide any extra information about the experiment.
4. In the same figure plot the Hoeffding's bound<sup>1</sup> on  $\Pr(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$  for the same values of  $\alpha$ .
5. In the same figure plot the Markov's bound<sup>2</sup> on  $\Pr(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$ .
6. Compare the three plots.
7. For  $\alpha = 1$  and  $\alpha = 0.95$  calculate the exact probability  $\Pr(\frac{1}{20} \sum_{i=1}^{20} X_i \geq \alpha)$  and compare it with the Hoeffding's bound. (No need to add this one to the plot.)

Do not forget to put axis labels and a legend in your plot!

### 1.2.2 The effect of scale (range) and normalization of random variables in Hoeffding's inequality

### 1.2.3 Probability in Practice

### 1.2.4 Occam's Razor

### 1.2.5 Train-Validation Split Trade-off

### 1.2.6 Optional Assignment Covering the VC-Analysis

### 1.2.7 The growth function

### 1.2.8 VC-dimension

### 1.2.9 Airline Revisited

---

<sup>1</sup>Hoeffding's bound is the right hand side of Hoeffding's inequality.

<sup>2</sup>Markov's bound is the right hand side of Markov's inequality.

## 2 Chapter 2. Concentration of Measure Inequalities

### 2.1 Markov's Inequality

**Theorem 2.1** (Markov's Inequality). *For any non-negative random variable  $X$  and  $\varepsilon > 0$ :*

$$\Pr(X \geq \varepsilon) \leq \frac{\mathbf{E}[X]}{\varepsilon}. \quad (2.1)$$

*Proof.* Define a random variable  $Y = \mathbb{I}(X \geq \varepsilon)$  to be the indicator function of whether  $X$  exceeds  $\varepsilon$ . Then  $Y \leq \frac{X}{\varepsilon}$ . Since  $Y$  is a Bernoulli random variable,  $\mathbf{E}[Y] = \Pr(Y = 1)$ . We have:

$$\Pr(X \geq \varepsilon) = \Pr(Y = 1) = \mathbf{E}[Y] \leq \mathbf{E}\left[\frac{X}{\varepsilon}\right] = \frac{\mathbf{E}[X]}{\varepsilon}. \quad (2.2)$$

□

**A Question:** Because  $Y = \mathbb{I}[\frac{X}{\varepsilon} \geq 1] \in [0, 1]$ , and  $\frac{X}{\varepsilon} \in [0, +\infty)$  due to the non-negativity of  $X$ .

1. If  $\frac{X}{\varepsilon} \in [0, 1)$ , then  $Y = 0$ , and  $Y \leq \frac{X}{\varepsilon}$ ;
2. If  $\frac{X}{\varepsilon} = 1$ , then  $Y = 1$ , and  $Y = \frac{X}{\varepsilon}$ ;
3. If  $\frac{X}{\varepsilon} \in (1, +\infty)$ , then  $Y = 1$ , and  $Y < \frac{X}{\varepsilon}$ ;

Overall,  $Y \leq \frac{X}{\varepsilon}$ . Then  $\mathbf{E}[Y] = \Pr(Y = 1)$ .

$$\Pr(X \geq \varepsilon) = \Pr(Y = 1) = \mathbf{E}[Y] \leq \mathbf{E}\left[\frac{X}{\varepsilon}\right] = \frac{\mathbf{E}[X]}{\varepsilon} \quad (2.3a)$$

$$\Pr(X \geq \frac{1}{\delta} \mathbf{E}[X]) \leq \frac{\mathbf{E}[X]}{\frac{1}{\delta} \mathbf{E}[X]} = \delta \quad (2.3b)$$

By denoting the right hand side of Markov's inequality by  $\delta$  we obtain the following equivalent statement. For any non-negative random variable  $X$ :

$$\Pr(X \geq \frac{1}{\delta} \mathbf{E}[X]) \leq \delta. \quad (2.4)$$

### 2.2 Chebyshev's Inequality

This one exploits variance to obtain tighter concentration.

**Theorem 2.2** (Chebyshev's inequality). *For any  $\varepsilon > 0$ ,*

$$\Pr(|X - \mathbf{E}[X]| \geq \varepsilon) \leq \frac{\mathbf{V}[X]}{\varepsilon^2}. \quad (2.5)$$

*Proof.* We use a transformation of a random variable. We have that  $\Pr(|X - \mathbf{E}[X]| \geq \varepsilon) = \Pr((X - \mathbf{E}[X])^2 \geq \varepsilon^2)$ , because the first statement holds if and only if the second holds. In addition, using Markov's inequality and the fact that  $(X - \mathbf{E}[X])^2$  is a non-negative random variable we have

$$\Pr(|X - \mathbf{E}[X]| \geq \varepsilon) = \Pr((X - \mathbf{E}[X])^2 \geq \varepsilon^2) \leq \frac{\mathbf{E}[(X - \mathbf{E}[X])^2]}{\varepsilon^2} = \frac{\mathbf{V}[X]}{\varepsilon^2}. \quad (2.6)$$

□

In order to illustrate the relative advantage of Chebyshev's inequality compared to Markov's consider the following example. Let  $X_1, \dots, X_n$  be  $n$  independent identically distributed Bernoulli random variables and let  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$  be their average. We would like to bound the probability that  $\hat{\mu}_n$  deviates from  $\mathbf{E}[\hat{\mu}_n]$  by more than  $\varepsilon$  (this is the central question in machine learning). We have  $\mathbf{E}[\hat{\mu}_n] = \mathbf{E}[X_1] = \mu$  and by independence of  $X_i$ -s and Theorem B.26<sup>3</sup> we have  $\mathbf{V}[\hat{\mu}_n] = \frac{1}{n^2} \mathbf{V}[n\hat{\mu}_n] = \frac{1}{n^2} \sum_{i=1}^n \mathbf{V}[X_i] = \frac{1}{n} \mathbf{V}[X_1]$ . By Markov's inequality,

$$\Pr(\hat{\mu}_n - \mathbf{E}[\hat{\mu}_n] \geq \varepsilon) = \Pr(\hat{\mu}_n \geq \mathbf{E}[\hat{\mu}_n] + \varepsilon) \leq \frac{\mathbf{E}[\hat{\mu}_n]}{\mathbf{E}[\hat{\mu}_n] + \varepsilon} = \frac{\mathbf{E}[X_1]}{\mathbf{E}[X_1] + \varepsilon}. \quad (2.9)$$

Note that as  $n$  grows the inequality stays the same. By Chebyshev's inequality we have

$$\Pr(\hat{\mu}_n - \mathbf{E}[\hat{\mu}_n] \geq \varepsilon) \leq \Pr(|\hat{\mu}_n - \mathbf{E}[\hat{\mu}_n]| \geq \varepsilon) \leq \frac{\mathbf{V}[\hat{\mu}_n]}{\varepsilon^2} = \frac{\mathbf{V}[X_1]}{n\varepsilon^2}. \quad (2.10)$$

Note that as  $n$  grows the right hand side of the inequality decreases at the rate of  $\frac{1}{n}$ . Thus, in this case, Chebyshev's inequality is much tighter than Markov's and it illustrates that as the number of random variables grows the probability that their average significantly deviates from the expectation decreases. In the next section we show that this probability actually decreases at an exponential rate.

**A Question:** I don't know why  $\mathbf{E}[\hat{\mu}_n] = \mathbf{E}[X_1] = \mu$ ?

## 2.3 Hoeffding's Inequality

Hoeffding's inequality is a much more powerful concentration result.

**Theorem 2.5** (Hoeffding's Inequality). *Let  $X_1, \dots, X_n$  be independent real-valued random variables, such that for each  $i \in \{1, \dots, n\}$  there exist  $a_i \leq b_i$ , such that  $X_i \in [a_i, b_i]$ . Then for every  $\varepsilon > 0$ :*

$$\Pr\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right] \geq \varepsilon\right) \leq e^{-2\varepsilon^2 / \sum_{i=1}^n (b_i - a_i)^2}, \quad (2.11)$$

and

$$\Pr\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right] \leq -\varepsilon\right) \leq e^{-2\varepsilon^2 / \sum_{i=1}^n (b_i - a_i)^2}. \quad (2.12)$$

By taking a union bound of the events, we obtain the following corollary.

---

3

**Theorem 2.3** (B.26). *If  $X_1, \dots, X_n$  are independent random variables then*

$$\mathbf{V}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbf{V}[X_i]. \quad (2.7)$$

**Theorem 2.4** (B.23). *If  $X$  and  $Y$  are independent random variables, then*

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]. \quad (2.8)$$

*We emphasize that in contrast with Theorem B.22, this property does not hold in the general case (if  $X$  and  $Y$  are not independent).*

**Corollary 2.6.** *Under the assumptions of Theorem 2.5:*

$$\Pr \left( \left| \sum_{i=1}^n X_i - \mathbf{E} \left[ \sum_{i=1}^n X_i \right] \right| \geq \varepsilon \right) \leq 2e^{-2\varepsilon^2 / \sum_{i=1}^n (b_i - a_i)^2}. \quad (2.13)$$

Equations (2.11) and (2.12) are known as “one-sided Hoeffding’s inequalities” and 2.13 is known as “two-sided Hoeffding’s inequality”. If we assume that  $X_i$ -s are identically distributed and belong to the  $[0, 1]$  interval we obtain the following corollary.

**A Proof of Corollary 2.6:** Let the event  $A$  be  $\sum_{i=1}^n X_i - \mathbf{E} [\sum_{i=1}^n X_i]$ , then

$$\Pr(A_+) = \Pr(A \geq \varepsilon) \quad (2.14a)$$

$$\Pr(A_-) = \Pr(A \leq -\varepsilon) \quad (2.14b)$$

$$\begin{aligned} \Pr(|A| \geq \varepsilon) &= \Pr((A \geq \varepsilon) \vee (A \leq -\varepsilon)) \\ &= \Pr(A \geq \varepsilon) + \Pr(A \leq -\varepsilon) \\ &\leq 2 \exp \left( \frac{-2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \end{aligned} \quad (2.14c)$$

**Corollary 2.7.** *Let  $X_1, \dots, X_n$  be independent random variables, such that  $X_i \in [0, 1]$  and  $\mathbf{E}[X_i] = \mu$  for all  $i$ , then for every  $\varepsilon > 0$ :*

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \varepsilon \right) \leq e^{-2n\varepsilon^2} \quad (2.15)$$

and

$$\Pr \left( \mu - \frac{1}{n} \sum_{i=1}^n X_i \geq \varepsilon \right) \leq e^{-2n\varepsilon^2}. \quad (2.16)$$

**A Proof of Corollary 2.7:** According to Theorem 2.5, we get two inequalities. Let  $Y_i = \frac{1}{n}X_i \in [0, 1]$ , we can get

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \mathbf{E} \left[ \sum_{i=1}^n X_i \right] \geq \varepsilon \right) \leq \exp \left( \frac{-2\varepsilon^2}{\sum_{i=1}^n \frac{1}{n^2} (b_i - a_i)^2} \right) = \exp \left( \frac{-2\varepsilon^2}{\frac{1}{n} 1^2} \right) = e^{-2n\varepsilon^2} \quad (2.17a)$$

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \mathbf{E} \left[ \frac{1}{n} \sum_{i=1}^n X_i \right] \leq -\varepsilon \right) \leq \exp \left( \frac{-2\varepsilon^2}{\sum_{i=1}^n \frac{1}{n^2} (b_i - a_i)^2} \right) = \exp \left( \frac{-2\varepsilon^2}{\frac{1}{n} 1^2} \right) = e^{-2n\varepsilon^2} \quad (2.17b)$$

$$\mathbf{E} \left[ \sum_{i=1}^n \frac{1}{n} X_i \right] = \frac{1}{n} \sum_{i=1}^n \mathbf{E}[X_i] = \mu \quad (2.17c)$$

Recall that by Chebyshev’s inequality  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$  converges to  $\mu$  at the rate of  $n^{-1}$ . Hoeffding’s inequality demonstrates that the convergence is actually much faster, at least at the rate of  $e^{-n}$ . The proof of Hoeffding’s inequality is based on Hoeffding’s lemma.

**Lemma 2.8** (Hoeffding’s Lemma). *Let  $X$  be a random variable, such that  $X \in [a, b]$ . Then for any  $\lambda \in \mathbb{R}$ :*

$$\mathbf{E} \left[ e^{\lambda X} \right] \leq e^{\lambda \mathbf{E}[X] + \frac{\lambda^2 (b-a)^2}{8}}. \quad (2.18)$$

The function  $f(\lambda) = \mathbf{E}[e^{\lambda X}]$  is known as the *moment generating function* of  $X$ , since  $f'(0) = \mathbf{E}[X]$ ,  $f''(0) = \mathbf{E}[X^2]$ , and, more generally,  $f^{(k)}(0) = \mathbf{E}[X^k]$ . We provide the proof of the lemma immediately after the proof of Theorem 2.5.

*Proof of Theorem 2.5.* We prove the first inequality in Theorem 2.5. The second inequality follows by applying the first inequality to  $-X_1, \dots, -X_n$ . The proof is based on Chernoff's bounding technique. For any  $\lambda > 0$  the following holds:

$$\Pr\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right] \geq \varepsilon\right) = \Pr\left(e^{\lambda\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right]\right)} \geq e^{\lambda\varepsilon}\right) \leq \frac{\mathbf{E}\left[e^{\lambda\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right]\right)}\right]}{e^{\lambda\varepsilon}}, \quad (2.19)$$

where the first step holds since  $e^{\lambda x}$  is a monotonously increasing function for  $\lambda > 0$  and the second step holds by Markov's inequality. We now take a closer look at the nominator:

$$\begin{aligned} \mathbf{E}\left[e^{\lambda\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right]\right)}\right] &= \mathbf{E}\left[e^{\sum_{i=1}^n \lambda(X_i - \mathbf{E}[X_i])}\right] \\ &= \mathbf{E}\left[\prod_{i=1}^n e^{\lambda(X_i - \mathbf{E}[X_i])}\right] \\ &= \prod_{i=1}^n \mathbf{E}\left[e^{\lambda(X_i - \mathbf{E}[X_i])}\right] \end{aligned} \quad (2.20a)$$

$$\leq \prod_{i=1}^n e^{\lambda^2(b_i - a_i)^2/8} \quad (2.20b)$$

$$= e^{(\lambda^2/8)\sum_{i=1}^n (b_i - a_i)^2}, \quad (2.20c)$$

where (2.20a) holds since  $X_1, \dots, X_n$  are independent and (2.20b) holds by Hoeffding's lemma applied to a random variable  $Z_i = X_i - \mathbf{E}[X_i]$  (note that  $\mathbf{E}[Z_i] = 0$  and that  $Z_i \in [a_i - \mu_i, b_i - \mu_i]$  for  $\mu_i = \mathbf{E}[X_i]$ ). *Put attention to the crucial role that independence of  $X_1, \dots, X_n$  plays in the proof! Without independence we would not have been able to exchange the expectation with the product and the proof would break down!* To complete the proof we substitute the bound on the expectation into the previous calculation and obtain:

$$\Pr\left(\sum_{i=1}^n X_i - \mathbf{E}\left[\sum_{i=1}^n X_i\right] \geq \varepsilon\right) \leq e^{(\lambda^2/8)\left(\sum_{i=1}^n (b_i - a_i)^2\right) - \lambda\varepsilon}. \quad (2.21)$$

This expression is minimized<sup>4</sup> by

$$\lambda^* = \underset{\lambda}{\operatorname{argmin}} e^{(\lambda^2/8)\left(\sum_{i=1}^n (b_i - a_i)^2\right) - \lambda\varepsilon} = \underset{\lambda}{\operatorname{argmin}} \left( (\lambda^2/8) \left( \sum_{i=1}^n (b_i - a_i)^2 \right) - \lambda\varepsilon \right) = \frac{4\varepsilon}{\sum_{i=1}^n (b_i - a_i)^2}. \quad (2.23)$$

---

4

$$\begin{aligned} \lambda^* &= \underset{\lambda}{\operatorname{argmin}} \frac{1}{8} \left( \left( \sum_{i=1}^n (b_i - a_i)^2 \right) \lambda^2 - 8\varepsilon\lambda \right) = \underset{\lambda}{\operatorname{argmin}} \frac{1}{8} \sum_{i=1}^n (b_i - a_i)^2 \left( \lambda^2 - \frac{8\varepsilon}{\sum_{i=1}^n (b_i - a_i)^2} \lambda \right) \\ &= \underset{\lambda}{\operatorname{argmin}} \frac{1}{8} \sum_{i=1}^n (b_i - a_i)^2 \left( \left( \lambda - \frac{4\varepsilon}{\sum_{i=1}^n (b_i - a_i)^2} \right)^2 - \frac{16\varepsilon^2}{\left( \sum_{i=1}^n (b_i - a_i)^2 \right)^2} \right) \\ &= \frac{4\varepsilon}{\sum_{i=1}^n (b_i - a_i)^2} \quad \text{s.t.} \quad \min_{\lambda} = -\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \end{aligned} \quad (2.22a)$$

It is important to note that the best choice of  $\lambda$  does not depend on the sample. In particular, it allows to fix  $\lambda$  before observing the sample. By substituting  $\lambda^*$  into the calculation we obtain the result of the theorem.  $\square$

*Proof of Lemma 2.8.* Note that

$$\mathbf{E} \left[ e^{\lambda X} \right] = \mathbf{E} \left[ e^{\lambda(X - \mathbf{E}[X]) + \lambda \mathbf{E}[X]} \right] = e^{\lambda \mathbf{E}[X]} \times \mathbf{E} \left[ e^{\lambda(X - \mathbf{E}[X])} \right]. \quad (2.24)$$

Hence, it is sufficient to show that for any random variable  $Z$  with  $\mathbf{E}[Z] = 0$  and  $Z \in [a, b]$  we have

$$\mathbf{E} \left[ e^{\lambda Z} \right] \leq e^{\lambda^2(b-a)^2/8}. \quad (2.25)$$

By convexity of the exponential function, for  $z \in [a, b]$  we have:

$$e^{\lambda z} \leq \frac{z-a}{b-a} e^{\lambda b} + \frac{b-z}{b-a} e^{\lambda a}. \quad (2.26)$$

Let  $p = -a/(b-a)$ . Then:

$$\begin{aligned} \mathbf{E} \left[ e^{\lambda Z} \right] &\leq \mathbf{E} \left[ \frac{Z-a}{b-a} e^{\lambda b} + \frac{b-Z}{b-a} e^{\lambda a} \right] \\ &= \frac{\mathbf{E}[Z] - a}{b-a} e^{\lambda b} + \frac{b - \mathbf{E}[Z]}{b-a} e^{\lambda a} \end{aligned} \quad (2.27a)$$

$$\text{nonumber} \quad (2.27b)$$

$$\begin{aligned} &= \frac{-a}{b-a} e^{\lambda b} + \frac{b}{b-a} e^{\lambda a} \\ &= \left( 1 - p + p e^{\lambda(b-a)} \right) e^{-p\lambda(b-a)} \\ &= e^{\phi(u)}, \end{aligned} \quad (2.27c)$$

where  $u = \lambda(b-a)$  and  $\phi(u) = -pu + \ln(1 - p + p e^u)$  and we used the fact that  $\mathbf{E}[Z] = 0$ . It is easy to verify that the derivative of  $\phi$  is

$$\phi'(u) = -p + \frac{p}{p + (1-p)e^{-u}}, \quad (2.28)$$

and, therefore,  $\phi(0) = \phi'(0) = 0$ . Furthermore,

$$\phi''(u) = \frac{p(1-p)e^{-u}}{(p + (1-p)e^{-u})^2} \leq \frac{1}{4}. \quad (2.29)$$

By Taylor's theorem,  $\phi(u) = \phi(0) + u\phi'(0) + \frac{u^2}{2}\phi''(\theta)$  for some  $\theta \in [0, u]$ . Thus, we have:

$$\phi(u) = \phi(0) + u\phi'(0) + \frac{u^2}{2}\phi''(\theta) = \frac{u^2}{2}\phi''(\theta) \leq \frac{u^2}{8} = \frac{\lambda^2(b-a)^2}{8}. \quad (2.30)$$

$\square$

### 2.3.1 Understanding Hoeffding's Inequality

Hoeffding's inequality involves three interconnected terms:  $n$ ,  $\varepsilon$ , and  $\delta = 2e^{-2n\varepsilon^2}$ , which is the bound on the probability that the event under  $\mathbf{Pr}(\cdot)$  holds (for the purpose of the discussion we consider two-sided Hoeffding's inequality for random variables bounded in  $[0, 1]$ ). We can fix



any two of the three terms  $n$ ,  $\varepsilon$ , and  $\delta$  and then the relation  $\delta = e^{-2n\varepsilon^2}$  provides the value of the third. Thus, we have

$$\delta = 2e^{-2n\varepsilon^2}, \quad (2.31a)$$

$$\varepsilon = \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2n}}, \quad (2.31b)$$

$$n = \frac{\ln\left(\frac{2}{\delta}\right)}{2\varepsilon^2}. \quad (2.31c)$$

Overall, Hoeffding's inequality tells by how much the empirical average  $\frac{1}{n} \sum_{i=1}^n X_i$  can deviate from its expectation  $\mu$ , but the interplay between the three parameters provides several ways of seeing and using Hoeffding's inequality. For example, if the number of samples  $n$  is fixed (we have made a fixed number of experiments and now analyze what we can get from them), there is an interplay between the precision  $\varepsilon$  and confidence  $\delta$ . We can request higher precision  $\varepsilon$ , but then we have to compromise on the confidence  $\delta$  that the desired bound  $\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \leq \varepsilon$  holds. And the other way around: we can request higher confidence  $\delta$ , but then we have to compromise on precision  $\varepsilon$ , i.e., we have to increase the allowed range  $\pm\varepsilon$  around  $\mu$ , where we expect to find the empirical average  $\frac{1}{n} \sum_{i=1}^n X_i$ .

As another example, we may have target precision  $\varepsilon$  and confidence  $\delta$  and then the inequality provides us the number of experiments  $n$  that we have to perform in order to achieve the target. It is often convenient to write the inequalities (2.15) and (2.16) with a fixed confidence in mind, thus we have

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2n}} \right) \leq \delta, \quad (2.32a)$$

$$\Pr \left( \mu - \frac{1}{n} \sum_{i=1}^n X_i \geq \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2n}} \right) \leq \delta, \quad (2.32b)$$

$$\Pr \left( \left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2n}} \right) \leq \delta. \quad (2.32c)$$

(Put attention that the  $\ln 2$  factor in the last inequality comes from the union bound over the first two inequalities: if we want to keep the same confidence we have to compromise on precision.)

**A1:** I don't know why, that is, Corollary 2.6.

In many situations we are interested in the complimentary events. Thus, for example, we have

$$\Pr \left( \mu - \frac{1}{n} \sum_{i=1}^n X_i \leq \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2n}} \right) \geq 1 - \delta. \quad (2.33)$$

Careful reader may point out that the inequalities above should be strict (" $<$ " and " $>$ "). This is true, but if it holds for strict inequalities it also holds for non-strict inequalities (" $\leq$ " and " $\geq$ "). Since strict inequalities provide no practical advantages we will use the non-strict inequalities to avoid the headache of remembering which inequalities should be strict and which should not.

The last inequality essentially says that with probability at least  $(1 - \delta)$  we have

$$\mu \leq \frac{1}{n} \sum_{i=1}^n X_i + \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2n}}, \quad (2.34)$$

and this is how we will occasionally use it. Note that the random variable is  $\frac{1}{n} \sum_{i=1}^n X_i$  and the right way of interpreting the above inequality is actually that with probability at least  $(1 - \delta)$

$$\frac{1}{n} \sum_{i=1}^n X_i \geq \mu - \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{2n}}, \quad (2.35)$$

that is, the probability is over  $\frac{1}{n} \sum_{i=1}^n X_i$  and not over  $\mu$ . However, many generalisation bounds that we study in Chapter 3 are written in the first form in the literature and we follow the tradition.

### 2.3.2 Example Appendix

**Q Example:** We would like to bound the probability that we flip a fair coin 10 times and obtain 8 or more heads. Let  $X_1, \dots, X_{10}$  be i.i.d. Bernoulli random variables with bias  $\frac{1}{2}$ . The question is equivalent to asking what is the probability that  $\sum_{i=1}^{10} X_i \geq 8$ . We have  $\mathbf{E} \left[ \sum_{i=1}^{10} X_i \right] = 5$  and by Markov's inequality,

$$\Pr \left( \sum_{i=1}^{10} X_i \geq 8 \right) \leq \frac{\mathbf{E} \left[ \sum_{i=1}^{10} X_i \right]}{8} = \frac{5}{8}. \quad (2.36)$$

**A Example:** Let the event  $A$  denote the value/number of  $\sum_{i=1}^{10} X_i$ .

$$\Pr(A \geq 8) \leq \frac{\mathbf{E}[A]}{8} = \frac{5}{8}. \quad (2.37a)$$

$$\Pr(A \geq 8) = \frac{A_{10}^8 + A_{10}^9 + A_{10}^{10}}{2^{10}} = \frac{\frac{10!}{2!} + \frac{10!}{1!} + \frac{10!}{0!}}{2^{10}} = \frac{10!}{2^{10}}, \quad (2.37b)$$

$$\Pr(A \geq 8) = \frac{C_{10}^8 + C_{10}^9 + C_{10}^{10}}{2^{10}} = \frac{\frac{10!}{8!2!} + \frac{10!}{9!1!} + \frac{10!}{10!0!}}{2^{10}} = \frac{45+10+1}{1024} = \frac{56}{1024} = \frac{7}{128}. \quad (2.37c)$$

According to the Markov's inequality, we have

$$C_{10}^i = C_{10}^i \text{ for } i \in \{0, 1, \dots, 10\} \quad (2.38a)$$

$$= \frac{10!}{10!0!} \frac{10!}{9!1!}, \frac{10!}{8!2!}, \frac{10!}{7!3!}, \frac{10!}{6!4!}, \frac{10!}{5!5!}, \frac{10!}{4!6!}, \frac{10!}{3!7!}, \frac{10!}{2!8!}, \frac{10!}{1!9!}, \frac{10!}{0!10!}, \quad (2.38b)$$

$$= 1, 10, 45, 120, 210, 252, 210, 120, 45, 10, 1, \quad (2.38c)$$

$$\text{in summary, } 1024 \quad (2.38d)$$

$\Pr(A \geq 0) \leq \frac{5}{0} \approx +\infty$	(2.39a)	$\Pr(A \geq 0) = 1$	(2.40a)
$\Pr(A \geq 1) \leq \frac{5}{1} = 5$	(2.39b)	$\Pr(A \geq 1) =$	(2.40b)
$\Pr(A \geq 2) \leq \frac{5}{2} = 2.5$	(2.39c)	$\Pr(A \geq 2) =$	(2.40c)
$\Pr(A \geq 3) \leq \frac{5}{3} \approx 1.6667$	(2.39d)	$\Pr(A \geq 3) =$	(2.40d)
$\Pr(A \geq 4) \leq \frac{5}{4} = 1.25$	(2.39e)	$\Pr(A \geq 4) =$	(2.40e)
$\Pr(A \geq 5) \leq \frac{5}{5} = 1$	(2.39f)	$\Pr(A \geq 5) =$	(2.40f)
$\Pr(A \geq 6) \leq \frac{5}{6} \approx 0.8333$	(2.39g)	$\Pr(A \geq 6) =$	(2.40g)
$\Pr(A \geq 7) \leq \frac{5}{7} \approx 0.7143$	(2.39h)	$\Pr(A \geq 7) =$	(2.40h)
$\Pr(A \geq 8) \leq \frac{5}{8} = 0.625$	(2.39i)	$\Pr(A \geq 8) =$	(2.40i)
$\Pr(A \geq 9) \leq \frac{5}{9} \approx 0.5556$	(2.39j)	$\Pr(A \geq 9) =$	(2.40j)
$\Pr(A \geq 10) \leq \frac{5}{10} = 0.5$	(2.39k)	$\Pr(A \geq 10) = \frac{C_{10}^{10}}{2^{10}} = \frac{1}{1024}$	(2.40k)
			(2.40l)

## 2.4 Basics of Information Theory: Entropy, Relative Entropy, and the Method of Types

In this section we briefly introduce a number of basic concepts from information theory that are very useful for deriving concentration inequalities. Specifically, we introduce the notions of entropy and relative entropy and some basic tools from the method of types. We start with some definitions.

**Definition 2.9** (Entropy). *Let  $p(x)$  be a distribution of a discrete random variable  $X$  taking values in a finite set  $\mathcal{X}$ . We define the entropy of  $p$  as:*

$$\mathbf{H}(p) = - \sum_{x \in \mathcal{X}} p(x) \ln p(x). \quad (2.41)$$

We use the convention that  $0 \ln 0 = 0$  (which is justified by continuity of  $z \ln z$ , since  $z \ln z \rightarrow 0$  as  $z \rightarrow 0$ ).

We have special interest in Bernoulli random variables.

**Definition 2.10** (Bernoulli random variable).  *$X$  is a Bernoulli random variable with bias  $p$  if  $X$  accepts values in  $\{0, 1\}$  with  $\Pr(X = 0) = 1 - p$  and  $\Pr(X = 1) = p$ .*

Note that expectation of a Bernoulli random variable is equal to its bias:

$$\mathbf{E}[X] = 0 \times \Pr(X = 0) + 1 \times \Pr(X = 1) = \Pr(X = 1) = p. \quad (2.42)$$

With a slight abuse of notation we specialise the definition of entropy to Bernoulli random variables.

**Definition 2.11** (Binary entropy). *Let  $p$  be a bias of Bernoulli random variable  $X$ . We define the entropy of  $p$  as*

$$\mathbf{H}(p) = -p \ln p - (1 - p) \ln(1 - p). \quad (2.43)$$

Note that when we talk about Bernoulli random variables  $p$  denotes the bias of the random variable and when we talk about more general random variables  $p$  denotes the complete distribution.

Entropy is one of the central quantities in information theory and it has numerous applications. We start by using binary entropy to bound binomial coefficients.

**Lemma 2.12.**

$$\frac{1}{n+1} e^{n\mathbf{H}(\frac{k}{n})} \leq \binom{n}{k} \leq e^{n\mathbf{H}(\frac{k}{n})}. \quad (2.44)$$

(Note that  $\frac{k}{n} \in [0, 1]$  and  $\mathbf{H}(\frac{k}{n})$  in the lemma is the binary entropy.)

*Proof.* By the binomial formula we know that for any  $p \in [0, 1]$ :

$$\sum_{i=0}^n \binom{n}{i} p^i (1-p)^{n-i} = 1. \quad (2.45)$$

We start with the upper bound. Take  $p = \frac{k}{n}$ . Since the sum is larger than any individual term, for the  $k$ -th term of the sum we get:

$$\begin{aligned} 1 &\geq \binom{n}{k} p^k (1-p)^{n-k} = \binom{n}{k} \left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k} = \binom{n}{k} \left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} \\ &= \binom{n}{k} e^{k \ln \frac{k}{n} + (n-k) \ln \frac{n-k}{n}} = \binom{n}{k} e^{n \left( \frac{k}{n} \ln \frac{k}{n} + \frac{n-k}{n} \ln \frac{n-k}{n} \right)} = \binom{n}{k} e^{-n\mathbf{H}(\frac{k}{n})}. \end{aligned} \quad (2.46a)$$

By changing sides of the inequality we obtain the upper bound.

For the lower bound it is possible to show that if we fix  $p = \frac{k}{n}$  then  $\binom{n}{k} p^k (1-p)^{n-k} \geq \binom{n}{i} p^i (1-p)^{n-i}$  for any  $i \in \{0, \dots, n\}$ . We also note that there are  $(n+1)$  elements in the sum in Eq. (2.45). Again, take  $p = \frac{k}{n}$ , then

$$1 \leq (n+1) \max_i \binom{n}{i} \left(\frac{k}{n}\right)^i \left(\frac{n-k}{n}\right)^{n-i} = (n+1) \binom{n}{k} \left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} = (n+1) \binom{n}{k} e^{-n\mathbf{H}(\frac{k}{n})}, \quad (2.47)$$

where the last step follows the same steps as in the derivation of the upper bound.  $\square$

**A Lemma 2.12:**

$$e^{k \ln \frac{k}{n} + (n-k) \ln \frac{n-k}{n}} = e^{\ln \left(\frac{k}{n}\right)^k + \ln \left(\frac{n-k}{n}\right)^{n-k}} = e^{\ln \left(\frac{k}{n}\right)^k} e^{\ln \left(\frac{n-k}{n}\right)^{n-k}} = \left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} \quad (2.48a)$$

$$\left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k} = \left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} = e^{n \left( \frac{k}{n} \ln \frac{k}{n} + \frac{n-k}{n} \ln \frac{n-k}{n} \right)} = e^{-\mathbf{H}(\frac{k}{n}) \cdot n} \quad (2.48b)$$

Lemma 2.12 shows that the number of configurations of choosing  $k$  out of  $n$  objects is directly related to the entropy of the imbalance  $\frac{k}{n}$  between the number of objects that are selected ( $k$ ) and the number of objects that are left out ( $n-k$ ).

We now introduce one additional quantity, the *Kullback-Leibler (KL) divergence*, also known as *Kullback-Leibler distance* and as *relative entropy*.

**Definition 2.13** (Relative entropy or Kullback-Leibler divergence). Let  $p(x)$  and  $q(x)$  be two probability distributions of a random variable  $X$  (or two probability density functions, if  $X$  is a continuous random variable), the Kullback-Leibler divergence or relative entropy is defined as:

$$\mathbf{KL}(p\|q) = \mathbb{E}_p \left[ \ln \frac{p(X)}{q(X)} \right] = \begin{cases} \sum_{x \in \mathcal{X}} p(x) \ln \frac{p(x)}{q(x)}, & \text{if } \mathcal{X} \text{ is discrete;} \\ \int_{x \in \mathcal{X}} p(x) \ln \frac{p(x)}{q(x)} dx, & \text{if } \mathcal{X} \text{ is continuous.} \end{cases} \quad (2.49)$$

We use the convention that  $0 \ln \frac{0}{0} = 0$  and  $0 \ln \frac{0}{q} = 0$  and  $p \ln \frac{p}{0} = \infty$ .

We specialize the definition to Bernoulli distributions.

**Definition 2.14** (Binary kl-divergence). Let  $p$  and  $q$  be biases of two Bernoulli random variables. The binary kl divergence is defined as:

$$\mathbf{kl}(p\|q) = \mathbf{KL}([1-p, p]\| [1-q, q]) = p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q}. \quad (2.50)$$

KL divergence is the central quantity in information theory. Although it is not a distance measure, because it does not satisfy the triangle inequality, it is the right way of measuring distances between probability distributions. This is illustrated by the following example.

**Example 2.15.** Let  $X_1, \dots, X_n$  be an i.i.d. sample of  $n$  Bernoulli random variables with bias  $p$  and let  $\frac{1}{n} \sum_{i=1}^n X_i$  be the empirical bias of the sample. (Note that  $\frac{1}{n} \sum_{i=1}^n X_i \in \{0, \frac{1}{n}, \frac{2}{n}, \dots, \frac{n}{n}\}$ .) Then by Lemma 2.12:

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i = \frac{k}{n} \right) = \binom{n}{k} p^k (1-p)^{n-k} \leq e^{n\mathbf{H}(\frac{k}{n})} e^{n(\frac{k}{n} \ln p + \frac{n-k}{n} \ln(1-p))} = e^{-n\mathbf{kl}(\frac{k}{n}\|p)}, \quad (2.51)$$

and

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i = \frac{k}{n} \right) \geq \frac{1}{n+1} e^{-n\mathbf{kl}(\frac{k}{n}\|p)}. \quad (2.52)$$

Thus,  $\mathbf{kl}(\frac{k}{n}\|p)$  governs the probability of observing empirical bias  $\frac{k}{n}$  when the true bias is  $p$ . It is easy to verify that  $\mathbf{kl}(p\|p) = 0$  and it is also possible to show that  $\mathbf{kl}(\hat{p}\|p)$  is convex in  $\hat{p}$  and that  $\mathbf{kl}(\hat{p}\|p) \geq 0$ . Thus, the probability of empirical bias is maximised when it coincides with the true bias.

**A Example 2.15:** According to Lemma 2.12, we get  $\frac{1}{n+1} e^{n\mathbf{H}(\frac{k}{n})} \leq \binom{n}{k} \leq e^{n\mathbf{H}(\frac{k}{n})}$ , then we already know  $\left(\frac{k}{n}\right)^k \left(\frac{n-k}{n}\right)^{n-k} = e^{k \ln \frac{k}{n} + (n-k) \ln \frac{n-k}{n}} = e^{-n\mathbf{H}(\frac{k}{n})}$ . Let  $p = \frac{k}{n}$ , then

$$\binom{n}{k} \leq e^{n\mathbf{H}(\frac{k}{n})} = e^{-n\left(\frac{k}{n} \ln \frac{k}{n} + \frac{n-k}{n} \ln \frac{n-k}{n}\right)} \quad (2.53a)$$

$$e^{n\left(\frac{k}{n} \ln p + \frac{n-k}{n} \ln(1-p)\right)} = e^{\ln p^k + \ln(1-p)^{n-k}} = p^k (1-p)^{n-k} \quad (2.53b)$$

$$\mathbf{kl}\left(\frac{k}{n}\|p\right) = \frac{k}{n} \left(\ln \frac{k}{n} - \ln p\right) + \frac{n-k}{n} \left(\ln \frac{n-k}{n} - \ln(1-p)\right) \quad (2.53c)$$

$$\binom{n}{k} p^k (1-p)^{n-k} \leq e^{-n\left(\frac{k}{n} \ln \frac{k}{n} + \frac{n-k}{n} \ln \frac{n-k}{n}\right) + n\left(\frac{k}{n} \ln p + \frac{n-k}{n} \ln(1-p)\right)} = e^{-n\mathbf{kl}\left(\frac{k}{n}\|p\right)} \quad (2.53d)$$

$$\binom{n}{k} p^k (1-p)^{n-k} \geq \frac{1}{n+1} e^{-n\left(\frac{k}{n} \ln \frac{k}{n} + \frac{n-k}{n} \ln \frac{n-k}{n}\right) + n\left(\frac{k}{n} \ln p + \frac{n-k}{n} \ln(1-p)\right)} = \frac{1}{n+1} e^{-n\mathbf{kl}\left(\frac{k}{n}\|p\right)} \quad (2.53e)$$

## 2.5 kl Inequality

Example 2.15 shows that **kl** can be used to bound the empirical bias when the true bias is known. But in machine learning we are usually interested in the inverse problem — how to infer the true bias  $p$  when the empirical bias  $\hat{p}$  is known. Next we demonstrate that this is also possible and that it leads to an inequality, which in most cases is tighter than Hoeffding's inequality. We start with the following lemma.

**Lemma 2.16.** *Let  $X_1, \dots, X_n$  be i.i.d. Bernoulli with bias  $p$  and  $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$  be the empirical bias. Then*

$$\mathbf{E} \left[ e^{n \mathbf{kl}(\hat{p} \| p)} \right] \leq n + 1. \quad (2.54)$$

*Proof.*

$$\mathbf{E} \left[ e^{n \mathbf{kl}(\hat{p} \| p)} \right] = \sum_{k=0}^n \mathbf{Pr} \left( \hat{p} = \frac{k}{n} \right) e^{n \mathbf{kl}(\frac{k}{n} \| p)} \leq \sum_{k=0}^n e^{-n \mathbf{kl}(\frac{k}{n} \| p)} e^{n \mathbf{kl}(\frac{k}{n} \| p)} = \sum_{k=0}^n 1 = n + 1, \quad (2.55)$$

where the inequality was derived in Eq. (2.51).  $\square$

**A Lemma 2.16:** We have already known that  $\mathbf{E}[X] = \sum_{k=1}^{\infty} x_k p_k = \int_{-\infty}^{\infty} x f(x) dx$ .

$$\mathbf{E} \left[ e^{n \mathbf{kl}(\hat{p} \| p)} \right] = \sum_{k=0}^n \mathbf{Pr} \left( \hat{p} = \frac{k}{n} \right) e^{n \mathbf{kl}(\frac{k}{n} \| p)} \stackrel{\text{def}}{=} A \quad (2.56a)$$

$$\frac{1}{n+1} e^{-n \mathbf{kl}(\frac{k}{n} \| p)} \leq \mathbf{Pr} \left( \frac{1}{n} \sum_{i=1}^n X_i = \frac{k}{n} \right) \leq e^{-n \mathbf{kl}(\frac{k}{n} \| p)} \quad (2.56b)$$

$$A \leq \sum_{k=0}^n e^{-n \mathbf{kl}(\frac{k}{n} \| p)} e^{n \mathbf{kl}(\frac{k}{n} \| p)} = \sum_{k=0}^n e^0 = \sum_{k=0}^n 1 = n + 1 \quad (2.56c)$$

$$A \geq \sum_{k=0}^n \frac{1}{n+1} e^{-n \mathbf{kl}(\frac{k}{n} \| p)} e^{n \mathbf{kl}(\frac{k}{n} \| p)} = \sum_{k=0}^n \frac{1}{n+1} e^0 = \frac{1}{n+1} \sum_{k=0}^n 1 = \frac{1}{n+1} (n + 1) = 1 \quad (2.56d)$$

We combine this lemma with Markov's inequality to obtain the following result.

**Theorem 2.17 (kl inequality).** *Let  $X_1, \dots, X_n$  be i.i.d. Bernoulli with bias  $p$  and let  $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$  be the empirical bias. Then*

$$\mathbf{Pr}(\mathbf{kl}(\hat{p} \| p) \geq \varepsilon) \leq (n + 1) e^{-n\varepsilon}. \quad (2.57)$$

*Proof.* By Markov's inequality and Lemma 2.16:

$$\mathbf{Pr}(\mathbf{kl}(\hat{p} \| p) \geq \varepsilon) = \mathbf{Pr} \left( e^{n \mathbf{kl}(\hat{p} \| p)} \geq e^{n\varepsilon} \right) \leq \frac{\mathbf{E} \left[ e^{n \mathbf{kl}(\hat{p} \| p)} \right]}{e^{n\varepsilon}} \leq \frac{n + 1}{e^{n\varepsilon}}. \quad (2.58)$$

$\square$

**A Theorem 2.17:** Markov's inequality:  $\mathbf{Pr}[X \geq \varepsilon] \leq \frac{\mathbf{E}[X]}{\varepsilon}$ . And then it's obviously demonstrated.

### 2.5.1 Relaxations of the kl-inequality: Pinsker's and refined Pinsker's equations

By denoting the right hand side of kl inequality (2.57) by  $\delta$ , we obtain that with probability greater than  $(1 - \delta)$ :

$$\mathbf{kl}(\hat{p} \| p) \leq \frac{\ln \frac{n+1}{\delta}}{n}. \quad (2.59)$$

This leads to an implicit bound on  $p$ , which is not very intuitive and not always convenient to work with. In order to understand the behavior of the kl inequality better, we use a couple of its relaxations. The first relaxation is known as Pinsker's inequality.

**A:** Because  $\Pr(\mathbf{kl}(\hat{p} \| p) \geq \varepsilon) \leq (n+1)e^{-n\varepsilon} \stackrel{\text{def}}{=} \delta$ , we get that with the probability greater than  $(1 - \delta)$ , we have  $e^{-n\varepsilon} = \frac{\delta}{n+1}$  and  $-n\varepsilon = \ln \frac{\delta}{n+1}$ , and then

$$\mathbf{kl}(\hat{p} \| p) \leq \varepsilon = -\frac{1}{n} \ln \frac{\delta}{n+1} = \frac{1}{n} \ln \frac{n+1}{\delta}. \quad (2.60)$$

**Lemma 2.18** (Pinsker's inequality).

$$\mathbf{KL}(p \| q) \geq \frac{1}{2} \|p - q\|_1^2, \quad (2.61)$$

where  $\|p - q\|_1 = \sum_{x \in \mathcal{X}} |p(x) - q(x)|$  is the  $L_1$ -norm.

**Corollary 2.19** (Pinsker's inequality for the binary kl divergence).

$$\mathbf{kl}(p \| q) \geq \frac{1}{2} (|p - q| + |(1 - p) - (1 - q)|)^2 = 2(p - q)^2. \quad (2.62)$$

**A:** We have already known that

$$\mathbf{KL}(p \| q) = \mathbb{E}_p \left[ \ln \frac{p(X)}{q(X)} \right] = \begin{cases} \sum_{x \in \mathcal{X}} p(x) \ln \frac{p(x)}{q(x)}, & \text{discrete} \\ \int_{x \in \mathcal{X}} p(x) \ln \frac{p(x)}{q(x)} dx, & \text{continuous} \end{cases} \quad (2.63)$$

$$\mathbf{kl}(p \| q) = \mathbf{KL}([1 - p, p] \| [1 - q, q]) = p \ln \frac{p}{q} + (1 - p) \ln \frac{1-p}{1-q} \quad (2.64)$$

$$\mathbf{kl}(p \| q) \geq \frac{1}{2} (|p - q| + |(1 - p) - (1 - q)|)^2 = \frac{1}{2} (|p - q| + |q - p|)^2 = \frac{1}{2} (2|p - q|)^2 = 2(p - q)^2 \quad (2.65a)$$

$$\mathbf{kl}(\hat{p} \| p) = \hat{p} \ln \frac{\hat{p}}{p} + (1 - \hat{p}) \ln \frac{1-\hat{p}}{1-p} \leq \frac{1}{n} \ln \frac{n+1}{\delta} \quad (2.65b)$$

$$2(\hat{p} - p)^2 \leq \mathbf{kl}(\hat{p} \| p) \leq \frac{1}{n} \ln \frac{n+1}{\delta} \quad (2.65c)$$

$$|\hat{p} - p|^2 \leq \frac{\mathbf{kl}(\hat{p} \| p)}{2} \leq \frac{1}{2n} \ln \frac{n+1}{\delta} \quad (2.65d)$$

$$|\hat{p} - p| \leq \sqrt{\frac{\mathbf{kl}(\hat{p} \| p)}{2}} \leq \sqrt{\frac{1}{2n} \ln \frac{n+1}{\delta}} \quad (2.65e)$$

By applying Corollary 2.19 to inequality (2.59), we obtain that with probability greater than  $(1 - \delta)$ ,

$$|p - \hat{p}| \leq \sqrt{\frac{\mathbf{kl}(\hat{p} \| p)}{2}} \leq \sqrt{\frac{\ln \frac{n+1}{\delta}}{2n}}. \quad (2.66)$$

Recall that Hoeffding's inequality assures that with probability greater than  $(1 - \delta)$ ,

$$p \leq \hat{p} + \sqrt{\frac{\ln \frac{1}{\delta}}{2n}}. \quad (2.67)$$

Thus, in the worst case the kl inequality is only weaker by the  $\ln(n+1)$  factor and in fact the  $\ln(n+1)$  factor can be reduced by a more careful analysis. Next we show that the kl inequality can actually be significantly tighter than Hoeffding's inequality. For this we use refined Pinsker's inequality.

**A:** We know that  $\delta = (n+1)e^{-n\epsilon}$ , and  $\Pr(\mathbf{kl}(\hat{p}||p) \geq \epsilon) \leq \delta$ , then with the probability greater than  $(1 - \delta)$ , we have  $\mathbf{kl}(\hat{p}||p) \leq \frac{\ln \frac{n+1}{\delta}}{n}$ . We also know that Hoeffding's inequality

$$\Pr\left(\sum_{i=1}^n X_i - \mathbb{E}\left[\sum_{i=1}^n X_i\right] \geq \epsilon\right) \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad \Pr\left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \epsilon\right) \leq e^{-2n\epsilon^2} \quad (2.68a)$$

$$\Pr\left(\sum_{i=1}^n X_i - \mathbb{E}\left[\sum_{i=1}^n X_i\right] \leq -\epsilon\right) \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad \Pr\left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \leq -\epsilon\right) \leq e^{-2n\epsilon^2} \quad (2.68b)$$

$$\Pr\left(\left|\sum_{i=1}^n X_i - \mathbb{E}\left[\sum_{i=1}^n X_i\right]\right| \geq \epsilon\right) \leq 2 \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \quad \Pr\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq \epsilon\right) \leq 2e^{-2n\epsilon^2} \quad (2.68c)$$

We know that

$$\hat{p} - \sqrt{\frac{\ln \frac{n+1}{\delta}}{2n}} \leq p \leq \hat{p} + \sqrt{\frac{\ln \frac{n+1}{\delta}}{2n}} \quad (2.69)$$

**Is it right in the following?** If let  $\delta = e^{-2n\epsilon^2}$ , then with the probability greater than  $(1 - \delta)$ , we have  $\ln \delta = -2n\epsilon^2$  and  $\epsilon^2 = -\frac{\ln \delta}{2n} = \frac{\ln \frac{1}{\delta}}{2n}$ , then we get

$$\hat{p} - p \leq \epsilon = \sqrt{\frac{\ln \frac{1}{\delta}}{2n}} \quad (2.70a)$$

$$\hat{p} - p \geq -\epsilon = -\sqrt{\frac{\ln \frac{1}{\delta}}{2n}} \quad (2.70b)$$

$$\hat{p} - \sqrt{\frac{\ln \frac{1}{\delta}}{2n}} \leq p \leq \hat{p} + \sqrt{\frac{\ln \frac{1}{\delta}}{2n}} \quad (2.70c)$$

Note that  $\ln \frac{1}{\delta} = \ln \delta^{-1} = \ln e^{2n\epsilon^2} = 2n\epsilon^2$ , and  $\sqrt{\frac{\ln \frac{1}{\delta}}{2n}} = \epsilon$ .

**Lemma 2.20** (Refined Pinsker's inequality).

$$\mathbf{kl}(p||q) \geq \frac{(p-q)^2}{2 \max\{p, q\}} + \frac{(p-q)^2}{2 \max\{(1-p), (1-q)\}}. \quad (2.71)$$

**Corollary 2.21** (Refined Pinsker's inequality). *If  $q > p$  then*

$$\mathbf{kl}(p||q) \geq \frac{(p-q)^2}{2q}. \quad (2.72)$$

**Corollary 2.22** (Refined Pinsker's inequality). *If  $\mathbf{kl}(p||q) \leq \epsilon$ , then*

$$q \leq p + \sqrt{2p\epsilon} + 2\epsilon. \quad (2.73)$$

By applying Corollary 2.22 to inequality (2.59) we obtain that with probability greater than  $(1 - \delta)$ :

$$p \leq \hat{p} + \sqrt{\frac{2\hat{p} \ln \frac{n+1}{\delta}}{n}} + \frac{2 \ln \frac{n+1}{\delta}}{n}. \quad (2.74)$$



Note that when  $\hat{p}$  is close to zero, the latter inequality is much tighter than Hoeffding's inequality. Finally, we note that although there is no analytic inversion of  $\mathbf{kl}(\hat{p}||p)$  it is possible to invert it numerically to obtain even tighter bounds than the relaxations above. Additionally, the bound in Theorem 2.17 can be improved slightly.

**A:** According to Eq. (2.59) we know that  $\delta = (n+1)e^{-n\varepsilon}$  and with probability greater than  $(1-\delta)$ , we have  $\mathbf{kl}(\hat{p}||p) \leq \frac{\ln \frac{n+1}{\delta}}{n} = \frac{\ln e^{n\varepsilon}}{n} = \varepsilon$ . Therefore, we have  $e^{-n\varepsilon} = \frac{\delta}{n+1}$  and  $-n\varepsilon = \ln \frac{\delta}{n+1}$  and  $\varepsilon = -\frac{1}{n} \ln \frac{\delta}{n+1} = \frac{1}{n} \ln \frac{n+1}{\delta}$ . Thus, if  $\mathbf{kl}(\hat{p}||q) \leq \varepsilon$ , then

$$p \leq \hat{p} + \sqrt{2\hat{p}\frac{1}{n} \ln \frac{n+1}{\delta}} + 2\frac{1}{n} \ln \frac{n+1}{\delta} = \hat{p} + \sqrt{\frac{2\hat{p} \ln \frac{n+1}{\delta}}{n}} + \frac{2 \ln \frac{n+1}{\delta}}{n}. \quad (2.75)$$

We have known that

$$\mathbf{KL}(p||q) \geq \frac{1}{2} \|p - q\|_1^2 \quad (2.76a)$$

$$\mathbf{kl}(p||q) \geq 2(p-q)^2 = \frac{(p-q)^2}{\frac{1}{2}} \quad (2.76b)$$

$$\because x \in [0, 1] \quad \therefore \frac{1}{x} \in [1, +\infty) \quad (2.76c)$$

$$\frac{1}{\max\{p, q\}}, \frac{1}{\max\{(1-p), (1-q)\}} \in [1, +\infty) \quad (2.76d)$$

$$\frac{(p-q)^2}{2\max\{p, q\}} + \frac{(p-q)^2}{2\max\{(1-p), (1-q)\}} \geq 2\frac{(p-q)^2}{2} = (p-q)^2 \quad (2.76e)$$

due to  $p, q, (1-p), (1-q) \in [0, 1]$ .

## 2.6 Sampling Without Replacement

Let  $X_1, \dots, X_n$  be a sequence of random variables *sampled without replacement* from a finite set of values  $\mathcal{X} = \{x_1, \dots, x_N\}$  of size  $N$ . The random variables  $X_1, \dots, X_n$  are dependent. For example, if  $\mathcal{X} = \{-1, +1\}$  and we sample two values then  $X_1 = -X_2$ . Since  $X_1, \dots, X_n$  are dependent, the concentration results from previous sections do not apply directly. However, the following result by Hoeffding, which we cite without a proof, allows to extend results for sampling with replacement to sampling without replacement.

**Lemma 2.23.** *Let  $X_1, \dots, X_n$  denote a random sample without replacement from a finite set  $\mathcal{X} = \{x_1, \dots, x_N\}$  of  $N$  real values. Let  $Y_1, \dots, Y_n$  denote a random sample with replacement from  $\mathcal{X}$ . Then for any continuous and convex function  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,*

$$\mathbf{E} \left[ f \left( \sum_{i=1}^n X_i \right) \right] \leq \mathbf{E} \left[ f \left( \sum_{i=1}^n Y_i \right) \right]. \quad (2.77)$$

In particular, the lemma can be used to prove Hoeffding's inequality for sampling without replacement.

**Theorem 2.24** (Hoeffding's inequality for sampling without replacement). *Let  $X_1, \dots, X_n$  denote a random sample without replacement from a finite set  $\mathcal{X} = \{x_1, \dots, x_N\}$  of  $N$  values, where each element  $x_i$  is in the  $[0, 1]$  interval. Let  $\mu = \frac{1}{N} \sum_{i=1}^N x_i$  be the average of the values in  $\mathcal{X}$ . Then for all  $\varepsilon > 0$ ,*

$$\Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \varepsilon \right) \leq e^{-2n\varepsilon^2}, \quad (2.78a)$$

$$\Pr \left( \mu - \frac{1}{n} \sum_{i=1}^n X_i \geq \varepsilon \right) = \Pr \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \leq -\varepsilon \right) \leq e^{-2n\varepsilon^2}. \quad (2.78b)$$

The proof is a minor adaptation of the proof of Hoeffding's inequality for sampling with replacement using Lemma 2.23 and is left as an exercise. (Note that it requires a small modification inside the proof, because Lemma 2.23 cannot be applied directly to the statement of Hoeffding's inequality.)

While formal proof requires a bit of work, intuitively the result is quite expected. Imagine the process of sampling without replacement. If the average of points sampled so far starts deviating from the mean of the values in  $\mathcal{X}$ , the average of points that are left in  $\mathcal{X}$  deviates in the opposite direction and "applies extra force" to new samples to bring the average back to  $\mu$ . In the limit when  $n = N$  we are guaranteed to have the average of  $X_i$ -s being equal to  $\mu$ .

**A:** I don't quite understand that. What does it mean? opposite direction

### 3 Appendix B: Probability Theory Basics

**Theorem 3.1** (Theorem B.26). *If  $X_1, \dots, X_n$  are independent random variables then*

$$\mathbf{V} \left[ \sum_{i=1}^n X_i \right] = \sum_{i=1}^n \mathbf{V}[X_i]. \quad (3.1)$$

*Proof.* We have already known that

$$\mathbf{V}[X] = \mathbf{E} \left[ (X - \mathbf{E}[X])^2 \right] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2, \quad (3.2a)$$

$$\mathbf{E}[cX] = c\mathbf{E}[X], \quad (3.2b)$$

$$\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y], \quad (3.2c)$$

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]. \quad (3.2d)$$

Then we could obtain that

$$\begin{aligned} \mathbf{V} \left[ \sum_{i=1}^n X_i \right] &= \mathbf{E} \left[ \left( \sum_{i=1}^n X_i - \mathbf{E} \left[ \sum_{i=1}^n X_i \right] \right)^2 \right] \\ &= \mathbf{E} \left[ \left( \sum_{i=1}^n X_i - \sum_{i=1}^n \mathbf{E}[X_i] \right)^2 \right] = \mathbf{E} \left[ \left( \sum_{i=1}^n (X_i - \mathbf{E}[X_i]) \right)^2 \right] \end{aligned} \quad (3.3a)$$

$$= \mathbf{E} \left[ (g)^2 \right] \quad (3.3b)$$

$$\begin{aligned} \mathbf{V} \left[ \sum_{i=1}^n X_i \right] &= \mathbf{E} \left[ \left( \sum_{i=1}^n X_i \right)^2 \right] - \left( \mathbf{E} \left[ \sum_{i=1}^n X_i \right] \right)^2 = \mathbf{E} \left[ \left( \sum_{i=1}^n X_i \right)^2 \right] - \left( \sum_{i=1}^n \mathbf{E}[X_i] \right)^2 \\ \sum_{i=1}^n \mathbf{V}[X_i] &= \sum_{i=1}^n \mathbf{E} \left[ (X_i - \mathbf{E}[X_i])^2 \right] = \sum_{i=1}^n \mathbf{E} \left[ X_i^2 - 2X_i\mathbf{E}[X_i] + (\mathbf{E}[X_i])^2 \right] \\ &= \sum_{i=1}^n (\mathbf{E}[X_i^2] - (\mathbf{E}[X_i])^2) \end{aligned} \quad (3.3c)$$

Eq. (3.3a) not necessarily independent.

$$\mathbf{V}[X_i + X_i] = \mathbf{E} \left[ (X_i + X_i - \mathbf{E}[X_i] - \mathbf{E}[X_i])^2 \right] = \mathbf{E} \left[ 4(X_i - \mathbf{E}[X_i])^2 \right] = 4\mathbf{V}[X_i] \quad (3.4a)$$

$$\begin{aligned} \mathbf{V}[X_i + X_j] &= \mathbf{E} \left[ (X_i + X_j - \mathbf{E}[X_i] - \mathbf{E}[X_j])^2 \right] = \mathbf{E} \left[ ((X_i - \mathbf{E}[X_i]) + (X_j - \mathbf{E}[X_j]))^2 \right] \\ &= \mathbf{E} \left[ (X_i - \mathbf{E}[X_i])^2 + (X_j - \mathbf{E}[X_j])^2 + 2(X_i - \mathbf{E}[X_i])(X_j - \mathbf{E}[X_j]) \right] \end{aligned} \quad (3.4b)$$

$$\begin{aligned} &= \mathbf{E} \left[ (X_i - \mathbf{E}[X_i])^2 \right] + \mathbf{E} \left[ (X_j - \mathbf{E}[X_j])^2 \right] + 2\mathbf{E} [X_i X_j - X_i \mathbf{E}[X_j] - X_j \mathbf{E}[X_i] + \mathbf{E}[X_i] \mathbf{E}[X_j]] \\ &= \mathbf{V}[X_i] + \mathbf{V}[X_j] + 2(\mathbf{E}[X_i X_j] - 2\mathbf{E}[X_i] \mathbf{E}[X_j] + \mathbf{E}[X_i] \mathbf{E}[X_j]) \end{aligned} \quad (3.4c)$$

$$= \mathbf{V}[X_i] + \mathbf{V}[X_j] + 2(\mathbf{E}[X_i X_j] - \mathbf{E}[X_i] \mathbf{E}[X_j]) = \mathbf{V}[X_i] + \mathbf{V}[X_j] \quad (3.4d)$$

**Induction** (1) When  $n = 1$  and  $n = 2$ , the equality is satisfied. (2) If the equality satisfied when  $n = t$ , that is,

$$\mathbf{V} \left[ \sum_{i=1}^t X_i \right] = \sum_{i=1}^t \mathbf{V}[X_i], \quad (3.5)$$

then when  $n = t + 1$ ,

$$Eq.^{\text{left}} = \mathbf{V} \left[ \sum_{i=1}^t X_i + X_{t+1} \right] = \mathbf{E} \left[ \left( \sum_{i=1}^t X_i + X_{t+1} - \mathbf{E} \left[ \sum_{i=1}^t X_i + X_{t+1} \right] \right)^2 \right] \quad (3.6a)$$

$$\begin{aligned} &= \mathbf{E} \left[ \left( \sum_{i=1}^t X_i - \mathbf{E} \left[ \sum_{i=1}^t X_i \right] \right)^2 + (X_{t+1} - \mathbf{E}[X_{t+1}])^2 + 2 \left( \sum_{i=1}^t X_i - \mathbf{E} \left[ \sum_{i=1}^t X_i \right] \right) (X_{t+1} - \mathbf{E}[X_{t+1}]) \right] \\ &= \sum_{i=1}^t \mathbf{V}[X_i] + \mathbf{V}[X_{t+1}] + 2\mathbf{E} \left[ X_{t+1} \sum_{i=1}^t X_i - \mathbf{E}[X_{t+1}] \sum_{i=1}^t X_i - \mathbf{E} \left[ \sum_{i=1}^t X_i \right] X_{t+1} + \mathbf{E} \left[ \sum_{i=1}^t X_i \right] \mathbf{E}[X_{t+1}] \right] \\ &= \sum_{i=1}^{t+1} \mathbf{V}[X_i] + 2 \left( \mathbf{E} \left[ X_{t+1} \sum_{i=1}^t X_i \right] - \mathbf{E} \left[ \sum_{i=1}^t X_i \right] \mathbf{E}[X_{t+1}] \right) \end{aligned} \quad (3.6b)$$

$$= \sum_{i=1}^{t+1} \mathbf{V}[X_i] + 2 \left( \mathbf{E}[X_{t+1}] \mathbf{E} \left[ \sum_{i=1}^t X_i \right] - \mathbf{E} \left[ \sum_{i=1}^t X_i \right] \mathbf{E}[X_{t+1}] \right) = \sum_{i=1}^{t+1} \mathbf{V}[X_i] \quad (3.6c)$$

Eq. (3.6b) because  $X_{t+1}$  and  $\sum_{i=1}^t X_i$  are independent random variables (after combination). (3)  
In conclusion, the equality is satisfied for all  $n \in \mathbb{Z}^+$ .  $\square$