

Contents

1

1	Introduction	1	2
1.1	Basic Inequalities	1	3
1.2	Why bother?	2	4
1.2.1	Coin Tossing	2	5
1.2.2	Central Limit Theorem	2	6
2	Exponential Inequalities	4	7
2.1	Chevyshev vs Chernoff vs Hoeffding	4	8
2.2	Chernoff-Okamoto Inequalities	4	9
2.3	Hoeffding-Bernstein inequalities	5	10
3	Application to Estimation of Data Dimension	6	11

1 Introduction

1.1 Basic Inequalities

Theorem 1.1.1 (Markov's inequality). *For a random variable X with $\mathbf{P}\{X < 0\} = 0$ and $t > 0$, we have*

$$\mathbf{P}\{X \geq t\} \leq \frac{\mathbf{E} X}{t}.$$

It follows that for a non-decreasing function φ which only takes non-negative values,

$$\mathbf{P}\{X \geq t\} = \mathbf{P}\{\varphi(X) \geq \varphi(t)\} \leq \frac{\varphi(X)}{\varphi(t)}.$$

Proof. In the first place, note that

$$\begin{aligned} X &= X \cdot \mathbf{1}_{X \geq t} + X \cdot \mathbf{1}_{X < t} \\ &\geq t \cdot \mathbf{1}_{X \geq t} + 0, \end{aligned}$$

and thus,

$$\mathbf{E} X \geq t \cdot \mathbf{E} \mathbf{1}_{X \geq t} = t \cdot \mathbf{P}\{X \geq t\}.$$

For the second statement, apply the same argument on the random variable $Y := \varphi(X)$ and the constant $s := \varphi(t)$. \square

Theorem 1.1.2 (Chebyshev's inequality). *For $t > 0$ and a random variable X with mean $\mu = \mathbf{E} X$ and variance $\sigma^2 = \mathbf{Var} X$, then*

$$\mathbf{P}\{|X - \mu| \geq t\} \leq \sigma^2 t^{-2}.$$

Proof. Applying Markov's inequality with $\varphi : x \mapsto x^2$ we obtain,

$$\mathbf{P}\{|X - \mu| \geq t\} = \mathbf{P}\{|X - \mu|^2 \geq t^2\} \leq \frac{\mathbf{E} [(X - \mu)^2]}{t^2} = \sigma^2 t^{-2}.$$

\square

Theorem 1.1.3 (Jensen's inequality). *For any real valued random variable X and convex function φ*

$$\varphi(\mathbf{E} X) \leq \mathbf{E} \varphi(X)$$

1.2 Why bother?

The concentration inequalities are used to obtain information on how a random variable is distributed at some specific places of its domain. In the most common scenarios, these inequalities will be used to quantify how concentrated a random variable at its tails, for example,

$$\mathbf{P}\{|X - \mu| \geq t\} < f(t) \ll 1.$$

A concentration inequality is specially useful when this probability cannot be calculated at a low computational cost or estimated with high precision. The following will illustrate a case where using concentration inequalities achieves the best results.

1.2.1 Coin Tossing

A coin tossing game is fair if the chances of winning are equal to the chances of losing. We can verify from a sample of N games that the game is not rigged if the number of heads in the sample is not very distant from the average $N/2$. However, there's a chance that one may classify the coin as rigged, even when the coin is fair. By the *Law of Large Numbers*, we know that the larger the sample, the less likely it is to obtain a false positive. But let's ask ourselves how fast this probability converges to 0.

Let $S_N \sim \text{Bi}(N, 1/2)$ denote the number of heads in a fair coin tossing game. Then,

$$\mu = \mathbf{E} S_N = \frac{N}{2}, \quad \sigma^2 = \mathbf{Var} S_N = \frac{N}{4}.$$

For a fixed $\varepsilon > 0$, we may classify a coin tossing game as rigged if, after N trials, the ratio of heads vs tails in the sample is greater than $[1 + \varepsilon : 1 - \varepsilon]$, or similarly,

$$S_N \geq \mu + \frac{\varepsilon}{2}N = \frac{1 + \varepsilon}{2}N.$$

Using the Chebyshev inequality 1.1.2, we assert that

$$\mathbf{P}\left\{S_N \geq \mu + \frac{\varepsilon}{2}N\right\} \leq \mathbf{P}\left\{|S_N - \mu| \geq \frac{\varepsilon}{2}N\right\} \leq \sigma^2 \frac{4}{\varepsilon^2 N^2} = \frac{1}{\varepsilon^2 N}.$$

Therefore, the probability of bad events tends to 0 at least linearly with the number of games.

1.2.2 Central Limit Theorem

The proof of the following theorems can be found in (ref)

Theorem 1.2.1. *Let X_i be a i.i.d. sample. Let $S_N = \sum_{i=1}^N X_i$, with mean $\mu = \mathbf{E} S_N$ and variance $\sigma^2 = \mathbf{Var} S_N$. If*

$$Z_N = \frac{S_N - N \cdot \mathbf{E} X_i}{\sqrt{N \cdot \mathbf{Var} X_i}} = \frac{S_N - \mu}{\sqrt{N} \sigma},$$

1 Introduction

64 then,

$$65 \quad Z_N \rightarrow Z \sim \mathcal{N}(0, 1), \text{ in distribution.}$$

66

□

67 **Theorem 1.2.2** (Tails of the Normal Distribution). *Let $Z \sim \mathcal{N}(0, 1)$, for $t > 0$ we have*

$$68 \quad \left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-t^2}{2}\right) \leq \mathbf{P}\{Z \geq t\} \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-t^2}{2}\right).$$

69

□

70 With that in mind, we might naively assume that better bounds can be obtained by
71 using the previous theorem. For a large enough N we can say that for the coin tossing,

$$72 \quad Z_N = \frac{S_N - N/2}{\sqrt{N/4}}$$

73

$$74 \quad \Rightarrow \mathbf{P}\left\{S_N \geq \frac{1+\varepsilon}{2}N\right\} = \mathbf{P}\left\{Z_N \geq \varepsilon\sqrt{N}\right\} \sim \mathbf{P}\left\{Z \geq \varepsilon\sqrt{N}\right\}.$$

75 However, this raises the question of whether we can draw the following conclusion from
76 Theorem 1.2.2:

$$77 \quad \mathbf{P}\left\{S_N \geq \frac{1+\varepsilon}{2}N\right\} \leq \frac{1}{\varepsilon\sqrt{N}} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-\varepsilon^2 \cdot N}{2}\right).$$

78 Unfortunately, the answer is no. The following theorem will show why.

79 **Theorem 1.2.3** (Convergence Rate for Central Limit Theorem). *For Z_N, Z in Theo-*
80 *rem 1.2.1, we have:*

$$81 \quad |\mathbf{P}\{Z_N \geq t\} - \mathbf{P}\{Z \geq t\}| \in O\left(\frac{1}{\sqrt{N}}\right).$$

82

□

83 Since the approximation error is greater than the bound, the previous results cannot
84 be taken into account.

85 In the context of coin tossing, this may not matter at all because the linear bound
86 obtained using Chebyshev's inequality indicates that the probability of wrongly classi-
87 fying a fair coin as a rigged coin converges at least linearly to zero. Even the Central
88 Limit Theorem shows in a less precise way this convergence. However, for some specific
89 problems in statistics, these basic tools are not precise enough to solve them. In the fol-
90 lowing chapters, we will show some examples where better crafted strategies are needed
91 in order to get bounds to the tails of the random variables.

2 Exponential Inequalities

92

Even if we are satisfied with the linear convergence rate provided by Chebyshev's inequality, there is a simple way to improve this bound with the following result.

93

94

Theorem 2.0.1 (Hoeffding's inequality). *Let X_1, \dots, X_N be independent random variables, where $X_i \sim \text{Be}(p_i)$. Then, for every $t > 0$, we have*

95

96

$$P \left\{ \sum (X_i - \mathbf{E} X_i) \geq t \right\} \leq \exp \left(\frac{-2t^2}{N} \right)$$

97

Proof. TODO

□

98

Another exponential bound we can derive with a similar idea is the following

99

Theorem 2.0.2 (Chernoff's inequality). *Let X_1, \dots, X_N be independent random variables, where for every $i = 1, \dots, N$, $X_i \in [a_i, b_i]$. Define $S_N = \sum X_i$ and let $\mu = \mathbf{E} S_N$. Then, for every $t > 0$, we have*

100

101

102

$$P \{ S_N \geq \mu + t \} \leq \exp \left(\frac{-2t^2}{\sum (a_i - b_i)^2} \right)$$

103

2.1 Chevyshev vs Chernoff vs Hoeffding

104

2.2 Chernoff-Okamoto Inequalities

105

Applying Markov's Inequality to $Y = e^{uX}$, we can assert that

106

$$\mathbf{P}\{X \geq \lambda + t\} \leq e^{-u(\lambda+t)} \mathbf{E} e^{uX} = e^{-u(\lambda+t)} (1 - p + pe^u)^n.$$

107

The right hand equation is minimized when,

108

$$e^u = \frac{\lambda + t}{(n - \lambda - t)} \cdot \frac{1 - p}{p}.$$

109

Therefore, for $0 \leq t \leq n - \lambda$,

110

$$\mathbf{P}\{X \geq \lambda + t\} \leq \left(\frac{\lambda}{\lambda + t} \right)^{\lambda+t} \left(\frac{n - \lambda}{n - \lambda - t} \right)^{n-\lambda-t} \quad (2.1)$$

111

2 Exponential Inequalities

112 **Theorem 2.2.1.** *Let X be random variable with the binomial distribution $\text{Bi}(n, p)$ with*
113 *$\lambda := np = \mathbf{E} X$, then for $t \geq 0$,*

114
$$\mathbf{P}\{X \geq \lambda + t\} \leq \exp\left(-\frac{t^2}{2(\lambda + t/3)}\right) \quad (2.2)$$

115
$$\mathbf{P}\{X \leq \lambda - t\} \leq \exp\left(-\frac{t^2}{2\lambda}\right) \quad (2.3)$$

117 **Used in:** Theorem 3.0.2

118 *Proof.* (TODO I've already written the proof on paper) □

119 2.3 Hoeffding-Bernstein inequalities

120 **Theorem 2.3.1.** *Let $\|f\|_\infty < c$, $\mathbf{E} f(X_1, \dots, X_m) = 0$ and $\sigma^2 = \mathbf{E} f^2(X_1, \dots, X_m)$.*
121 *Then for any $t > 0$,*

122
$$\mathbf{P}\{U_m^n(f, P) > t\} \leq \exp\left(\frac{\frac{n}{m}t^2}{2\sigma^2 + \frac{2}{3}ct}\right) \quad (2.1)$$

123 **Used in:** Theorem 3.0.4

124 *Proof.* Proposition 2.3(a) M.A. Arcones, E. Gine, Limit theorems for U-processes, Ann.
125 Probab. 21 (1993) 14941542

126 <https://sci-hub.se/https://www.jstor.org/stable/2244585> □

3 Application to Estimation of Data Dimension

The article (ref) explains how we can estimate the dimension d of a manifold M embedded on a Euclidean space of dimension m , say \mathbb{R}^m . First, we are going to introduce the method they used, and then, we will show how does the exponential inequalities can be used to prove two important results in the paper. The procedure starts with an example on a uniformly distributed sample on a d -sphere $\mathbb{S}^{d-1} \subset \mathbb{R}^d$, but will be later generalized for samples of any distribution on any manifold.

In the first place, let Z_1, \dots, Z_k be a i.i.d. sample uniformly distributed on \mathbb{S}^{d-1} . Then, we have the following formula for the variance of the angles between $Z_i, Z_j, i \neq j$:

$$\beta_d := \mathbf{Var}(\arccos \langle Z_i, Z_j \rangle) = \begin{cases} \frac{\pi^2}{4} - 2 \sum_{j=1}^k (2j-1)^{-2}, & \text{if } d = 2k+1 \text{ is odd,} \\ \frac{\pi^2}{12} - 2 \sum_{j=1}^k (2j)^{-2}, & \text{if } d = 2k+2 \text{ is even.} \end{cases} \quad (3.1)$$

The previous formula for the angle variance is proven in (ref) and will be skipped (TODO: Should it really be skipped?). In order to give more insight on how we will be choosing the estimator \hat{d} for the dimension of the sphere, consider the following theorem.

Theorem 3.0.1 (Bounds for β_d). *For every $d > 1$,*

$$\frac{1}{d} \leq \beta_d \leq \frac{1}{d-1}.$$

Proof. The even and the odd cases must be distinguished:

(1): When $d = 2k+2$ is even: In the first place, remember that,

$$\lim_{k \rightarrow \infty} \sum_{j=1}^k j^{-2} = \frac{\pi^2}{6}.$$

It follows that

$$\begin{aligned} \beta_d &= \frac{\pi^2}{12} - 2 \sum_{j=1}^k (2j)^{-2} = \frac{\pi^2}{12} - \frac{1}{2} \sum_{j=1}^k j^{-2} \\ &= \frac{1}{2} \sum_{j=k+1}^{\infty} j^{-2}. \end{aligned}$$

3 Application to Estimation of Data Dimension

149 Since $(j^{-2})_{j \in \mathbb{N}}$ is a monotonically decreasing sequence, it follows that (**TODO:**
150 Improve array syntax)

$$\begin{aligned} \frac{1}{d} &= \frac{1}{2k+2} = \frac{1}{2} \int_{k+1}^{\infty} x^{-2} dx \\ 151 \quad &\leq \beta_d \leq \frac{1}{2} \int_{k+1/2}^{\infty} x^{-2} dx \\ &= \frac{1}{2k+1} = \frac{1}{d-1}. \end{aligned}$$

152 (2): When $d = 2k + 3$ is odd: On the other hand, note that

$$\begin{aligned} \lim_{k \rightarrow \infty} \sum_{j=1}^k (2j-1)^{-2} &= \lim_{k \rightarrow \infty} \sum_{j=1}^{2k-1} j^{-2} - \sum_{j=1}^{k-1} (2j)^{-2} \\ 153 \quad &= \lim_{k \rightarrow \infty} \sum_{j=1}^{2k-1} j^{-2} - \frac{1}{4} \sum_{j=1}^{k-1} j^{-2} \\ &= \frac{\pi^2}{6} - \frac{\pi^2}{24} = \frac{\pi^2}{8} \end{aligned}$$

154 Then,

$$\begin{aligned} \beta_d &= \frac{\pi^2}{4} - 2 \sum_{j=1}^k (2j-1)^{-2} \\ 155 \quad &= 2 \sum_{j=k+1}^{\infty} (2j-1)^{-2}. \end{aligned}$$

156 Using a similar argument we conclude that (**TODO:** Improve array syntax)

$$\begin{aligned} \frac{1}{d} &= \frac{1}{2k+1} = 2 \int_{k+1}^{\infty} (2x-1)^{-2} dx \\ 157 \quad &\leq \beta_d \leq 2 \int_{k+1/2}^{\infty} (2x-1)^{-2} dx \\ &= \frac{1}{2k+2} = \frac{1}{d-1}. \end{aligned}$$

158 □

159 Knowing that for every $d > 1$, β_d is in the interval $[\frac{1}{d}, \frac{1}{d-1}]$, we are going to guess the
160 dimension of the sphere by estimating β_d , and then, taking d from the lower bound of
161 the interval where our estimator is. Since β_d is the variance of the angles in our sphere,
162 our best choice for an estimator is the angle's sample variance,

$$163 \quad U_k = \binom{k}{2}^{-1} \sum_{i < j \leq k} \left(\arccos \langle Z_i, Z_j \rangle - \frac{\pi^2}{2} \right)^2. \quad (3.2)$$

3 Application to Estimation of Data Dimension

In Proposition 1. of (ref) the authors prove that it's the Minimum Variance Unbiased Estimator for β_d on the unit sphere. However, the authors also prove that this result can be generalized for any manifold with samples of any distribution.

Let X_1, \dots, X_n be a i.i.d. sample from a random distribution P on a manifold $M \subset \mathbb{R}^m$, and let $p \in M$ a point. For $C > 0 \in \mathbb{R}$, let $k = \lceil C \ln(n) \rceil$ and define $R(n) = L_{k+1}(p)$ as the distance between p and its $(k+1)$ -nearest neighbor. W.L.O.G. assume that $p = 0$ and that X_1, \dots, X_k are the k -nearest neighbors of p

Theorem 3.0.2 (Bound k -neighbors). *For any sufficiently large $C > 0$, we have that, there exists n_0 such that, with probability 1, for every $n \geq n_0$,*

$$R(n) \leq f_{p,P,C}(n) = O(\sqrt[d]{\ln(n)/n}). \quad (3.3)$$

The function $f_{p,P,C}$ is a deterministic function which depends on p , P and C .

(TODO: I need help connecting the previous theorem with the following idea)

Let $\pi : \mathbb{R}^m \rightarrow T_p M$ be the orthogonal projection on the Tangent Space of M at p . Also, define $W_i := \pi(X_i)$ and then normalize,

$$Z_i := \frac{X_i}{\|X_i\|}, \quad \widehat{W}_i := \frac{W_i}{\|W_i\|}.$$

What follows from the previous and the following lemma is that if we know that W_1, \dots, W_k behave similar to a uniformly distributed sample on the sphere \mathbb{S}^d , then, Z_1, \dots, Z_k (the normalized k -nearest neighbors of p) also behave like they are uniformly distributed on \mathbb{S}^d .

Theorem 3.0.3 (Projection Distance Bounds). *For any $i < j \leq n$,*

$$(i) \quad \|X_i - \pi(X_i)\| = O(\|\pi(X_i)\|^2) \quad (3.4)$$

$$(ii) \quad \|Z_i - \widehat{W}_i\| = O(\|\pi(X_i)\|) \quad (3.5)$$

(iii) *The inner products (cosine of angles) can be bounded as it follows:*

$$|\langle Z_i, Z_j \rangle - \langle \widehat{W}_i, \widehat{W}_j \rangle| \leq Kr, \quad (3.6)$$

for a constant $K \in \mathbb{R}$, whenever $r \geq \max(\|\pi(X_i)\|, \|\pi(X_j)\|)$.

The last result in the article shows that if we estimate β_d as we did with $U_{k,n} = U_k$ in equation (3.2), and then, extract \widehat{d} from the interval where $U_{k,n}$ is located, it follows that,

Theorem 3.0.4 (Consistency). *When $n \rightarrow \infty$,*

$$\mathbf{P}\{\widehat{d} \neq d\} \rightarrow 0.$$

Proofs

Theorem 3.0.2. □