

# 第四章：不等式

## ■ 不等式

- 有些量很难计算，不等式可以对这些量给出一个界
- 不等式也是下一章讨论收敛理论的基础

## ■ 关于概率的不等式

- Markov不等式
- Chebyshev不等式
- Hoeffding不等式

## ■ 关于期望的不等式

- Cauchy-Schwarze不等式
- Jensen不等式

# Markov不等式

- 4.1 定理（Markov不等式）：令 $X$ 为非负随机变量且假设 $\mathbb{E}(X)$ 存在，则对任意 $t > 0$ ，有

$$P(X > t) \leq \frac{\mathbb{E}(X)}{t} \quad (4.1)$$

- 当 $t = k\mu, \mu = \mathbb{E}(X)$ ， $\mathbb{P}(X > k\mu) \leq \frac{1}{k}$ 
  - 当 $k > 1$ 时，表示随机变量的取值离不会期望不会太远（离期望较远的概率很小，小于 $1/k$ ）
    - $\mathbb{P}(X > 2\mu) \leq 0.5, \mathbb{P}(X > 3\mu) \leq 0.33$
- 当 $0 < k \leq 1$ 时， $1/k \geq 1$ ，上式总是成立表示（ $\mathbb{P}(A) \leq 1$ ）

Markov 不等式证明:  $X$  为非负随机变量, 对任何  $t>0$ , 有  $\mathbb{P}(X > t) \leq \frac{\mathbb{E}(X)}{t}$

证明:  $\because X \geq 0$

$$\therefore \mathbb{E}(X) = \int_0^{\infty} xf(x)dx \quad (x \text{ 非负, 所以积分区间从 } 0 \text{ 开始})$$

$$= \int_0^t xf(x)dx + \int_t^{\infty} xf(x)dx \quad (\text{分成两段积分})$$

$$\geq \int_t^{\infty} xf(x)dx \quad (x \geq 0, t > 0, \therefore \int_0^t xf(x)dx > 0, \text{ 求和部分去掉正的一部分,}$$

不等号成立)

$$\geq t \int_t^{\infty} f(x)dx \quad (\text{将 } t \text{ 放到积分符号外面, 相当于令} \quad \text{由于} \quad \text{不等号成立})$$

$$= t\mathbb{P}(X > t)$$

$$\therefore \mathbb{P}(X > t) \leq \frac{\mathbb{E}(X)}{t}$$

# Markov不等式

- 将 $X$ 换成满足条件的 $r(X)$ , 上述结论也成立!

$$\mathbb{P}(r(X) > t) \leq \frac{\mathbb{E}(r(X))}{t}$$

- 当  $r(X) = \frac{(X - \mu)^2}{\sigma^2}$  ?

- Chebyshev不等式: Markov不等式的应用

# Chebyshev不等式

- 4.2 定理（Chebyshev不等式）：令  $\mu = \mathbb{E}(X), \sigma^2 = \mathbb{V}(X)$
- 则  $\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad \mathbb{P}(|Z| \geq k) \leq \frac{1}{k^2}$
- 其中  $Z = (X - \mu)/\sigma$
- $\mathbb{P}(|Z| > 2) \leq 1/4, \mathbb{P}(|Z| > 3) \leq 1/9.$
- $X$ 在其期望附近（ $t$ 邻域）的概率与方差 $\sigma^2$ 有关
  - $\sigma^2$ 越大，随机变量远离期望的概率越大（方差用于度量随机变量围绕均值的散布程度）
  - $\sigma^2$ 越小，随机变量在期望附近，远离期望的概率越小
  - 可用来证明样本均值会在其期望附近（样本数越多越接近，因为样本方差随 $n$ 增大而减小）

Chebyshev 不等式证明：令  $\mathbb{E}(X) = \mu, \mathbb{V}(X) = \sigma^2$ ，则  $\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2}$

证明：  $\mathbb{P}(|X - \mu| \geq t) = \mathbb{P}((X - \mu)^2 \geq t^2)$  （两边同时平方）

$$\leq \frac{\mathbb{E}((X - \mu)^2)}{t^2} \quad (\text{Markov 不等式})$$

$$\leq \frac{\sigma^2}{t^2}$$

当  $t = k\sigma$ ，则  $\mathbb{P}(|X - \mu| \geq k\sigma) = \mathbb{P}\left(\frac{|X - \mu|}{\sigma} \geq k\right) \leq \frac{\sigma^2}{(k\sigma)^2} = \frac{1}{k^2}$

# Chebyshev不等式

- $X$ 在其期望附近 ( $t$ 邻域) 的概率与方差  $\sigma^2$  有关
- 另外一个变形:

$$\mathbb{P}(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

- $k=2?$   $\mathbb{P}(|X - \mu| \geq 2\sigma) \leq 0.25$
- $k=3?$   $\mathbb{P}(|X - \mu| \geq 3\sigma) \leq 0.11$ 
  - 高斯分布为0.9997

- 这个界很松, 因为Chebyshev不等式没有限定分布的形式, 所以应用广泛

- 对某些具体的分布来说, 可以得到更紧致的界, 如高斯分布  $Z \sim N(0,1)$

$$\mathbb{P}(Z \geq t) = \frac{1}{\sqrt{2\pi}} \int_t^\infty x e^{-x^2/2} dx \leq \frac{1}{\sqrt{2\pi}} \int_t^\infty \frac{x}{t} e^{-x^2/2} dx = \frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t}$$

Mill's inequality

$$\mathbb{P}(|Z| \geq t) = 2\mathbb{P}(Z \geq t) \leq \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}$$

# Chebyshev不等式

- 4.3例：假设我们在一个有 $n$ 个测试样本的测试集上测试一个预测方法（以神经网络为例）。若预测错误置 $X_i = 1$ ，预测正确则置 $X_i = 0$ 。则 $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ 为观测到的错误率。每个 $X_i$ 可视为有未知均值 $p$ 的Bernoulli分布。我们想知道真正的错误率 $p$ 。

- 直观地，我们希望 $\bar{X}_n$ 接近 $p$ 。但 $\bar{X}_n$ 有多大可能不在 $p$ 的 $\varepsilon$ 邻域内？

- $\mathbb{V}(\bar{X}_n) = \mathbb{V}(X_1)/n = p(1-p)/n,$


$$\mathbb{P}(|\bar{X}_n - p| \geq \varepsilon) \leq \frac{\mathbb{V}(\bar{X}_n)}{\varepsilon^2} = \frac{p(1-p)}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}.$$

- 由于对任意 $p$ 有 $p(1-p) \leq \frac{1}{4}$ ，所以当 $\varepsilon = 0.2$ ,  $n = 100$ 时，边界为0.0625。



# Hoeffding不等式

- 作用与Chebyshev不等式类似，但区间更紧致（增加了独立性约束）
- 4.4 定理（Hoeffding不等式）：设  $Y_1, \dots, Y_n$  相互独立，且  $\mathbb{E}(Y_i) = 0$ ，且  $a_i \leq Y_i \leq b_i$ 。令  $\varepsilon > 0$ ，则对任意  $t > 0$

$$Y_i = \frac{1}{n}(X_i - p) \quad \mathbb{P}\left(\sum_{i=1}^n Y_i \geq \varepsilon\right) \leq e^{-t\varepsilon} \prod_{i=1}^n e^{t^2(b_i - a_i)^2/8}$$


- 4.5 定理（Hoeffding不等式）：令  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$  则对任意  $\varepsilon > 0$ ，有

$$\mathbb{P}\left(|\bar{X}_n - p| > \varepsilon\right) \leq 2e^{-2n\varepsilon^2}$$

- 其中  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$

Hoeffding 不等式证明:

$$\text{证明: } \mathbb{P}\left(\sum_{i=1}^n Y_i \geq \varepsilon\right) = \mathbb{P}\left(t \sum_{i=1}^n Y_i \geq t\varepsilon\right) = \mathbb{P}\left(e^{t \sum_{i=1}^n Y_i} \geq e^{t\varepsilon}\right) \leq e^{-t\varepsilon} \mathbb{E}\left(e^{t \sum_{i=1}^n Y_i}\right) \quad (\text{Markov 不等式})$$

$$= e^{-t\varepsilon} \prod_{i=1}^n \mathbb{E}\left(e^{tY_i}\right) \quad (1)$$

$$\because a_i \leq Y_i \leq b_i, \therefore Y_i = \alpha b_i + (1-\alpha)a_i, \alpha = (Y_i - a_i)/(b_i - a_i)$$

$$\text{由于 } e^{ty} \text{ 为凸函数, 所以 } e^{tY_i} \leq \frac{(Y_i - a_i)}{(b_i - a_i)} e^{tb_i} + \frac{(b_i - Y_i)}{(b_i - a_i)} e^{ta_i}$$

$$\mathbb{E}\left(e^{tY_i}\right) \leq -\frac{a_i}{(b_i - a_i)} e^{tb_i} + \frac{b_i}{(b_i - a_i)} e^{ta_i} = e^{g(u)} \quad (\mathbb{E}(Y_i) = 0)$$

$$\text{其中 } u = t(b_i - a_i), g(u) = -\gamma a_i + \log(1 - \gamma + \gamma e^u), \gamma = -a_i / b_i - a_i$$

$$g(0) = g'(0) = 0, g''(u) \leq 1/4, \text{ for } u > 0$$

$$\text{根据 Talayor 展开, } g(u) = g(0) + ug'(0) + u^2 g''(\xi) = \frac{u^2}{2} g''(\xi) \leq \frac{1}{8} u^2 = \frac{1}{8} t^2 (b_i - a_i)^2$$

$$\text{所以 } \mathbb{E}\left(e^{tY_i}\right) \leq e^{g(u)} \leq e^{t^2(b_i - a_i)^2/8}, \text{ 代入 (1), 不等式得证。}$$

# Hoeffding不等式

- 4.6 例：令  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$   $n=100$ ,  $\varepsilon=0.2$
- 则根据Chebyshev不等式，有

$$\mathbb{P}\left(\left|\overline{X}_n - p\right| > \varepsilon\right) \leq 0.0625.$$

$$\mathbb{P}\left(\left|\overline{X}_n - p\right| \geq \varepsilon\right) \leq \mathbb{V}(\overline{X}_n) / \varepsilon^2$$

- 根据Hoeffding不等式，有

$$\mathbb{P}\left(\left|\overline{X}_n - p\right| > .2\right) \leq 2e^{-2(100)(.2)^2} = 0.00067.$$

$$\mathbb{P}\left(\left|\overline{X}_n - p\right| \geq \varepsilon\right) \leq 2e^{-2n\varepsilon^2}$$

- 结果远远小于0.0625。

# Hoeffding不等式

- 可用来计算二项分布中的参数 $p$ 的置信区间

- 对给定的  $\alpha > 0$  , 令

$$\varepsilon_n = \left\{ \frac{1}{2n} \log \left( \frac{2}{\alpha} \right) \right\}^{\frac{1}{2}}$$

- 则根据Hoeffding不等式

$$\mathbb{P} \left( \left| \overline{X}_n - p \right| > \varepsilon_n \right) \leq 2e^{-2n\varepsilon_n^2} = \alpha$$

- 令  $C = (\overline{X}_n - \varepsilon, \overline{X}_n + \varepsilon)$  , 则  $\mathbb{P}(C \notin p) = \mathbb{P}(|X_n - p| > \varepsilon) \leq \alpha$

- 则  $\mathbb{P}(C \in p) \geq 1 - \alpha$  。

- 称 $C$ 为  $1 - \alpha$  置信区间。

# Cauchy-Schwarze不等式

- 4.8 定理（Cauchy-Schwarze不等式）：若 $X$ 、 $Y$ 是有限方差，则

$$\mathbb{E}(|XY|) \leq \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$$

- 例：协方差不等式

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] \leq \sqrt{\mathbb{E}((X - \mu_X)^2)\mathbb{E}((Y - \mu_Y)^2)} = \sigma_X \sigma_Y$$

$$(\text{Cov}(X, Y))^2 \leq \sigma_X^2 \sigma_Y^2$$

$$-1 \leq \rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \leq 1$$

# Jensen不等式

- 4.9 定理（ Jensen不等式）： 如果 $g$ 是凸的， 则

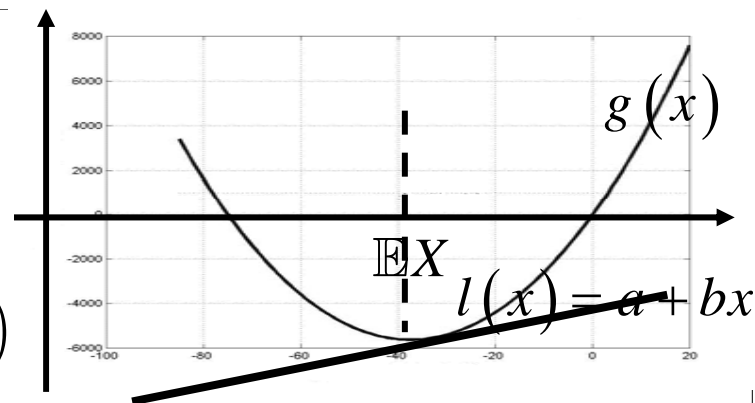
$$\mathbb{E}g(X) \geq g(\mathbb{E}(X))$$

- 如果 $g$ 是凹的， 则

$$\mathbb{E}g(X) \leq g(\mathbb{E}(X))$$

- $g(x) = x^2 \Rightarrow \mathbb{E}(X^2) \geq (\mathbb{E}X)^2$
- $g(x) = \frac{1}{x} \Rightarrow \mathbb{E}\left(\frac{1}{X}\right) \geq \frac{1}{\mathbb{E}X}$
- $g(x) = \log x \Rightarrow \mathbb{E}(\log X) \leq \log(\mathbb{E}X)$

证明:  $g(x)$  为凸函数, 则  $\mathbb{E}g(X) \geq g(\mathbb{E}X)$



证明: 令在点  $x = \mathbb{E}X$  处, 函数  $g(x)$  的切线为  $l(x) = a + bx$

由于  $g(x)$  为凸函数,  $g(x)$  位于其切线之上, 即

$$g(x) \geq l(x) = a + bx$$

所以  $\mathbb{E}g(X) \geq \mathbb{E}(a + bX)$  (两边同取期望, 不等式仍然成立)

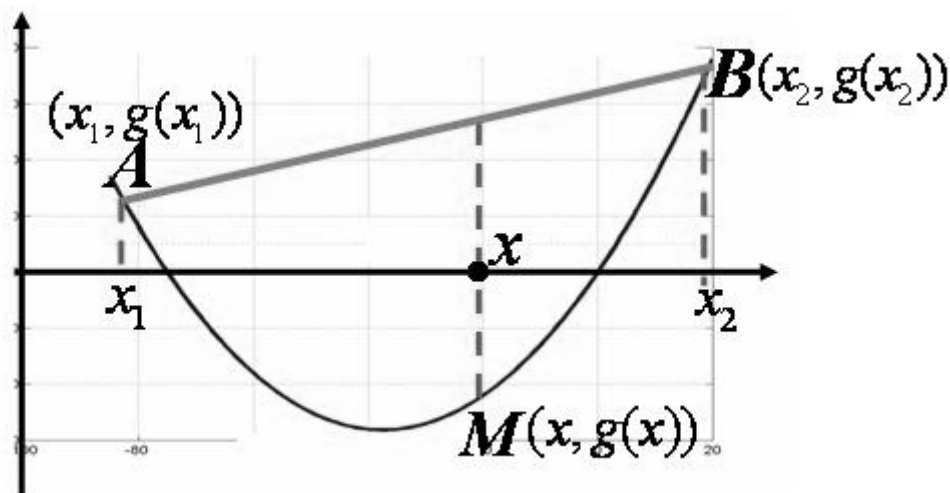
$$= a + b\mathbb{E}X \quad (\text{期望的线性性质})$$

$$= l(\mathbb{E}X) \quad (l \text{ 的定义})$$

$$= g(\mathbb{E}X) \quad (\text{在切点 } x = \mathbb{E}X \text{ 处, } l \text{ 等于函数值})$$

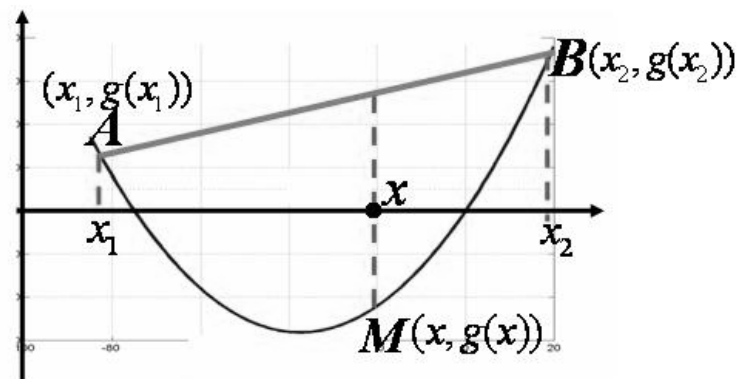
# 凸函数

- 如果对所有的  $x, y$ ,  $0 < \lambda < 1$  , 满足
$$g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y)$$
- 则函数  $g(x)$  为凸函数 (convex),  $-g(x)$  为凹函数 (concave)
  - 凸: 装水, 如  $g(x) = x^2$
  - 凹: 溢出水, 如  $g(x) = \log x$



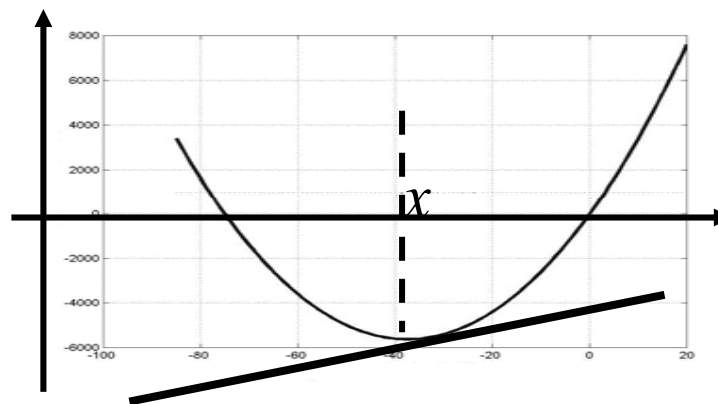


# 凸函数



## ■ 几何意义

- 连接  $(a, g(a)), (b, g(b))$  两点的弦，永远在  $y=g(x)$  之上
- 凸光滑函数上任一点的切线在曲线的下方



# 下节课内容：随机变量序列的收敛性

- 随机样本：IID样本  $X_1, X_2, \dots, X_n$  ,  $X_i \sim F$

- 统计量：对随机样本概述

$$Y = T(X_1, X_2, \dots, X_n)$$

- $Y$ 为随机变量， $Y$ 的分布称为统计量的采样分布
- 如：样本均值、样本方差、样本中值

- 收敛性：当样本数量 $n$ 趋向无穷大时，统计量的变化
  - 大样本理论、极限定理、渐近理论

# 作业

- 作业4: Chp4: 第1、2、4题