

高維統計學

干燥症

2025 年 1 月 19 日

目录

1 测度集中	2 3 度量熵	33
1.1 Chernoff 方法	3	
1.1.1 次高斯隨機變量	3	4 隨機矩陣
1.1.2 次指數隨機變量	6	34
1.2 鞅方法	9	
1.3 熵方法	11	A 預備知識
1.3.1 Herbst 方法	13	35
1.3.2 熵的張量化	15	A.1 Landau 記號
1.4 等周不等式	17	35
1.5 傳輸成本不等式	18	A.2 概率論
1.5.1 Wasserstein 距離	18	35
1.5.2 KL 散度與傳輸成本不等式	19	A.2.1 積分變換、Stieltjes 積分
1.5.3 傳輸成本的張量化	21	36
1.5.4 非對稱耦合成本	23	A.2.2 矩生成函數、累計生成函數
		37
		A.2.3 Randon-Nikodym 導數、密度
		37
		A.2.4 條件期望、鞅、鞅差
		38
		A.2.5 方差的表示
		38
		A.2.6 耦合
		38
		A.3 凸分析與凸優化
		39
		A.3.1 Rademacher 定理
		39
		A.3.2 Fenchel 共軛
		39
		A.3.3 Lagrange 乘數法
		40
		A.4 矩陣
		40
		A.4.1 矩陣範數
		40
2 一致大數定律	23	B 定理證明
2.1 經驗過程	24	41
2.2 經驗過程的尾部概率界	25	B.1 定理 1.9 的證明
2.3 函數類的 Rademacher 複雜度	26	41
2.4 多項式識別函數類	29	
2.5 Vapnik-Červonenkis 維數	31	

記號表

$\lceil x \rceil, \lfloor x \rfloor$	分別為不小于、不大于 x 的最小整數
$\mathbf{x}_j^k, \mathbf{X}_j^k$	$(x_j, x_{j+1}, \dots, x_k), (X_j, X_{j+1}, \dots, X_k)$
$\mathbf{X}^{\setminus k}, \mathbf{X}^{\setminus k}$	除去第 k 個分量後的向量 $(x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n), (X_1, \dots, X_{k-1}, X_{k+1}, \dots, X_n)$
$\mathbb{P}(f) = \mathbb{E}_{\mathbb{P}}[f]$	$\int f d\mathbb{P} = \mathbb{E}_{X \sim \mathbb{P}}[f(X)]$
\mathbb{E}_k	保持其他分量不變, 對第 k 個分量求逐坐標期望 $\mathbb{E}_k[f(\mathbf{X})] = \mathbb{E}[f(\mathbf{X}) \mathbf{X}^{\setminus k}]$
vol	體積
$\mathcal{C}(\mathbb{Q}, \mathbb{P})$	分布對 (\mathbb{Q}, \mathbb{P}) 的所有可能的耦合

1 测度集中

A random variable that depends (in a 'smooth' way) on the influence of many independent variables (but not too much on any of them) is essentially constant.

—Michel Talagrand (1996)

驯服随机!

在大尺度上, 无界随机变量的概率分布函数通常有着纤细、绵长的尾部, 这意味着集中现象: 一切随机变量几乎是一个有界随机变量. 例如正态分布的 3σ 原则告诉我们, “几乎所有”的值都在平均值正负三个标准差的范围内 (若 $X \sim \mathcal{N}(\mu, \sigma^2)$, 则 $\mathbb{P}(|X - \mu| \leq 3\sigma) \approx 0.9973$).

设函数 $f: \mathbb{R} \rightarrow \mathbb{R}_{>0}$ 单调增, 对任意 $t > 0$, 注意到

$$\mathbb{E}[f(X)] \geq \mathbb{E}[f(X); X \geq t] \geq \mathbb{E}[f(t)\mathbb{I}_{\{X \geq t\}}] = f(t) \cdot \mathbb{P}[X \geq t].$$

于是我们有经典的 (广义)Markov 不等式:

$$\mathbb{P}[X \geq t] \leq \frac{\mathbb{E}[f(X)]}{f(t)}. \quad (1)$$

通过发展 Markov 不等式, 我们可以得到单变量函数的集中不等式. 例如对于随机变量 $Y \in L^p$, 取 $X = |Y - \mathbb{E}[Y]|$, $f(u) = u^p$ ($u \geq 0$), 可以得到基于高阶矩的集中不等式

$$\mathbb{P}[|Y - \mathbb{E}[Y]| \geq t] \leq \frac{\mathbb{E}[|Y - \mathbb{E}[Y]|^p]}{t^p}.$$

特别地, $p = 2$ 时就是经典的 Chebyshev 不等式.

在统计学的背景下, 统计量是多个随机变量的函数 $f(X_1, \dots, X_n)$. 当函数 f 具有一些良好的性质时, 我们可以通过估计每个随机变量 X_k 对 f 的贡献, 再结合每个 X_k 的集中不等式来得到 f 的集中不等式, 这种方法被称为张量化. 例如从 Chebyshev 不等式出发, 由于独立随机变量之和的方差等于方差的和, 我们可以得到独立随机变量 $\{X_k\}_{k=1}^n$ 之和的集中不等式:

$$\mathbb{P}\left[\left|\sum_{k=1}^n (X_k - \mathbb{E}[X_k])\right| \geq t\right] \leq \frac{\text{Var}(\sum_k X_k)}{t^2} = \frac{\sum_k \text{Var}(X_k)}{t^2}.$$

相较于大数定律的渐进视角, 这提供了一种非渐进的方法, 对于小样本也可以有很好的估计.

1-范数-通常与维数 n 有关

1.1 Chernoff 方法

可以看到, Markov 不等式 (1) 中尾部概率由 f 的增长速度所控制——这意味着选取增长速度最快的函数, 可以得到更有效的尾部概率不等式. 自然地, 我们考虑指数函数.

若中心化随机变量 $X - \mathbb{E}[X]$ 在 0 的某个邻域 I 内有矩生成函数, 即在 $\lambda \in I$ 上有 $M(\lambda) := \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] < \infty$, 那么 $\lambda \in I \cap [0, \infty)$ 时, 取 $f(u) = e^{\lambda u}$ 可以得到下述的 Chernoff 不等式:

$$\mathbb{P}[X - \mathbb{E}[X] \geq t] \leq \frac{\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}]}{e^{\lambda t}}, \quad \forall \lambda \in I \cap [0, \infty).$$

通过选取最优的参数 λ , 我们可以得到 Chernoff 界¹

$$\mathbb{P}[X \geq \mathbb{E}[X] + t] \leq \exp \left[\inf_{\lambda \in I \cap [0, \infty)} \left\{ \log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] - \lambda t \right\} \right].$$

于是随机变量的集中不等式可以从其矩生成函数/累计生成函数的界得到.

1.1 示例 (Gauss 随机变量的上偏差 inequality). Gauss 随机变量 $X \sim N(\mu, \sigma^2)$ 的矩生成函数

$$\mathbb{E}[e^{\lambda X}] = e^{\frac{\sigma^2 \lambda^2}{2} + \mu \lambda} \quad (2)$$

在 $\lambda \in \mathbb{R}$ 总是存在. 通过简单的求导可以看到最优参数 $\lambda^* = \frac{t}{\sigma^2}$, 于是有

$$\mathbb{P}[X \geq \mu + t] \leq e^{-\frac{t^2}{2\sigma^2}}, \quad \forall t \geq 0. \quad (3)$$

1.1.1 次高斯随机变量

我们将 Gauss 随机变量作为“模版”来研究其他随机变量: 如果某个随机变量的中心矩生成函数能被中心 Gauss 随机变量的矩生成函数所控制, 利用 Chernoff 方法, 它的尾概率也会被中心 Gauss 随机变量的尾概率控制.

1.2 定义 (次高斯随机变量). 称期望为 μ 的随机变量 X 为参数为 σ 的次高斯随机变量, 如果存在 $\sigma > 0$, 使得

$$\mathbb{E}[e^{\lambda(X - \mu)}] \leq e^{\sigma^2 \lambda^2 / 2}, \quad \forall \lambda \in \mathbb{R}.$$

容易看到, σ -次高斯随机变量 X 总是满足上偏差 inequality (3). 此外, 由于 $-X$ 也是次高斯的, 可以得到下偏差 inequality: $\mathbb{P}[X \leq \mu - t] \leq e^{-\frac{t^2}{2\sigma^2}}$ 对任意 $t \geq 0$ 成立, 于是次高斯随机变量满足集中不等式

$$\mathbb{P}[|X - \mu| \geq t] \leq 2e^{-\frac{t^2}{2\sigma^2}}, \quad \forall t > 0. \quad (4)$$

¹在最优的 k 下, 基于高阶矩的 Markov 不等式不会比 Chernoff 界中基于矩生成函数获得的界更差. 但在实际应用中, 由于矩生成函数计算的便利性, Chernoff 界仍有着广泛的应用.

1.3 注 (輔助函數的構造). 定义 1.2 实际上是说, 中心化的次高斯随机变量的累积生成函数 $K(\lambda) := \log \mathbb{E}[e^{\lambda X}] = O(\lambda^2)$, 等价地, $K''(\lambda) = O(1)$; 或者说 $G(\lambda) := \lambda^{-1} K(\lambda) = O(\lambda)$, 等价地, $G'(\lambda) = O(1)$. 引理 1.7 和定理 1.31 的证明分别利用了这两种等价表述来构造辅助函数. 命题 1.19 的证明则用到了更为复杂的构造.

类似的方法还可以用来得到次高斯随机变量最大值的上界.

1.4 定理 (次高斯隨機變量最大值的上界). 设 $\{X_k\}_{k=1}^n$ 为均值为 0, 参数为 σ 的次高斯随机变量 (对独立性未做要求).

$$(1) \mathbb{E}[\max_k X_k] \leq \sigma \sqrt{2 \log n}, \forall n \geq 1;$$

$$(2) \mathbb{E}[\max_k |X_k|] \leq \sigma \sqrt{2 \log(2n)} \leq 2\sigma \sqrt{\log n}, \forall n \geq 2.$$

證明. (1) 由 Jensen 不等式, 对任意 $\lambda > 0$,

$$\begin{aligned} \exp\left(\lambda \mathbb{E}\left[\max_{1 \leq k \leq n} X_k\right]\right) &\leq E\left[\exp\left(\lambda \max_{1 \leq i \leq n} X_k\right)\right] \\ &= E\left[\max_{1 \leq k \leq n} e^{\lambda X_k}\right] \leq \sum_{k=1}^n E[e^{\lambda X_k}] \leq n e^{\frac{\lambda^2 \sigma^2}{2}}. \end{aligned}$$

于是 $\mathbb{E}[\max_k X_k] \leq \inf_{\lambda > 0} \left\{ \frac{\log n}{\lambda} + \lambda \frac{\sigma^2}{2} \right\} = \sqrt{2\sigma^2 \log n}, \forall n \geq 1.$

(2) 由于上述结果的建立对独立性未做要求, 我们考虑 $X_{n+k} := -X_k, i = k, \dots, n$, 有

$$\max_{1 \leq k \leq n} |X_k| = \max_{1 \leq k \leq 2n} X_k,$$

于是将结果中 n 替换为 $2n$ 可见结果成立. □

1.5 定理 (次高斯隨機變量的等價定義). 对任意均值为 0 的随机变量 X , 下述几条命题等价:

(I) 存在常数 $\sigma \geq 0$ 使得

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}, \forall \lambda \in \mathbb{R};$$

(II) 存在常数 $c \geq 0$ 和 $Z \sim \mathcal{N}(0, \tau^2)$ 使得

$$\mathbb{P}[|X| \geq s] \leq c \mathbb{P}[|Z| \geq s], \forall s \geq 0;$$

(III) 存在常数 $\theta \geq 0$ 使得

$$\mathbb{E}[X^{2k}] \leq \frac{(2k)!}{2^k k!} \theta^{2k}, \forall k = 1, 2, \dots;$$

(IV) 存在常数 $\sigma \geq 0$ 使得

$$\mathbb{E} \left[\exp \left(\frac{\lambda X^2}{2\sigma^2} \right) \right] \leq \frac{1}{\sqrt{1-\lambda}}, \quad \forall \lambda \in [0, 1).$$

利用独立性和 Chernoff 方法, 容易看到偏差不等式关于次高斯参数具有可加性:

1.6 命题 (Hoeffding 界). 设随机变量序列 $\{X_k\}_{k=1}^n$ 中 X_k 为均值为 μ_k , 参数为 σ_k 的独立次高斯随机变量, 那么 $\sum_k X_k$ 是 $\sqrt{\sum_k \sigma_k^2}$ -次高斯的, 于是我们有上偏差不等式

$$\mathbb{P} \left[\sum_{k=1}^n (X_k - \mu_k) \geq t \right] \leq \exp \left(-\frac{t^2}{2 \sum_{k=1}^n \sigma_k^2} \right), \quad \forall t \geq 0.$$

直观上, 一个有界随机变量没有无限的尾部, 因此它应当是次高斯随机变量. 事实上, 我们可以得到有界随机变量的紧次高斯参数.

1.7 引理 (有界随机变量). 若随机变量 $X \in [a, b]$ a.s., 那么它是 $\frac{b-a}{2}$ -次高斯的.

证明. 定义随机变量 X_λ , 它的分布关于 X 的分布的 Radon-Nikodym 导数为 $\frac{e^{\lambda X}}{\mathbb{E}[e^{\lambda X}]}$. 于是 $X_\lambda \in [a, b]$ a.s.¹, 从而 $|X_\lambda - \frac{a+b}{2}| \leq \frac{b-a}{2}$ a.s..

记 $\mathbb{E}[X] = \mu$, 累计生成函数 $K(\lambda) := \log \mathbb{E}[e^{\lambda X}]$ 满足 $K(0) = 0, K'(0) = \frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \Big|_{\lambda=0} = \mu$, 且对任意 λ ,

$$\begin{aligned} K''(\lambda) &= \frac{\mathbb{E}[X^2 e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} - \left(\frac{\mathbb{E}[X e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} \right)^2 = \mathbb{E}[X_\lambda^2] - (\mathbb{E}[X_\lambda])^2 = \text{Var}(X_\lambda) \\ &= \text{Var} \left(X_\lambda - \frac{a+b}{2} \right) \leq \mathbb{E} \left[X_\lambda - \frac{a+b}{2} \right]^2 \leq \left(\frac{b-a}{2} \right)^2. \end{aligned}$$

于是 $K(\lambda)$ 的在原点的 Taylor 展开为 $K(\lambda) = \lambda\mu + \frac{\lambda^2}{2} K''(\xi) \leq \lambda\mu + \frac{\lambda^2}{2} \left(\frac{b-a}{2} \right)^2$. 从而 $\mathbb{E}[e^{\lambda(X-\mu)}] \leq \exp \left(\frac{\lambda^2}{2} \cdot \left(\frac{b-a}{2} \right)^2 \right)$, 即 X 是 $\frac{b-a}{2}$ -次高斯的. \square

1.8 推论. 经典的 Hoeffding 不等式由命题 1.6 结合引理 1.7 得到.

一个基本的事实是, Gauss 随机向量的 Lipschitz 函数的集中度与维数无关.

1.9 定理. 设 $\mathbf{X} \in \mathbb{R}^n$ 为 Gauss 随机向量, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 关于 2-范数是 L -Lipschitz 的, 那么 $f(\mathbf{X})$ 为参数不大于 L 的次高斯随机变量.

它的证明请参照附录, 这里我们对它的应用进行举例说明.

¹这是由于 X 的分布作为概率测度的零测集一定是 X_λ 的分布的零测集.

1.10 示例 (χ^2 尾部概率界). 设 $Y = \|\mathbf{Z}\|_2^2 \sim \chi_n^2$, 其中 $\mathbf{Z} \subseteq \mathbb{R}^n$ 为标准 Gauss 向量. 考虑随机变量 $V := \sqrt{Y}/\sqrt{n} = \|\mathbf{Z}\|_2/\sqrt{n}$, 它作为 \mathbf{Z} 的函数是 $1/\sqrt{n}$ -Lipschitz 的, 从而由定理 1.9, V 是 $1/\sqrt{n}$ -次高斯的且有

$$\mathbb{P}[V \geq \mathbb{E}[V] + \delta] \leq e^{-\frac{n\delta^2}{2}}, \quad \forall \delta \geq 0.$$

由 Jensen 不等式, $\mathbb{E}[V] \leq \sqrt{\mathbb{E}[V^2]} = \left(\frac{1}{n} \sum_k \mathbb{E}[Z_k]^2\right)^{\frac{1}{2}} = 1$, 代入有

$$e^{-\frac{n\delta^2}{2}} \geq \mathbb{P}\left[\sqrt{\frac{Y}{n}} \geq 1 + \delta\right] = \mathbb{P}[Y \geq n(1 + \delta)^2].$$

注意到对任意的 $\delta \in [0, 1]$, 都有 $(1 + \delta)^2 \leq 1 + 3\delta$, 于是做替换 $t = 3\delta$ 有

$$\mathbb{P}[Y \geq n(1 + t)] \leq e^{-\frac{nt^2}{18}}, \quad \forall t \in [0, 3].$$

1.11 示例 (点集的 Gauss 复杂度). 设 $\mathbf{X} \subseteq \mathbb{R}^n$ 为标准 Gauss 向量, 给定 \mathbb{R}^n 的子集 \mathcal{A} , 随机变量

$$Z(\mathcal{A}) := \sup_{\mathbf{a} \in \mathcal{A}} \left[\sum_{k=1}^n a_k X_k \right] = \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \mathbf{X} \rangle$$

在某种意义下度量了集合 \mathcal{A} 的大小, 它的期望 $\mathbb{E}_{\mathbf{X}}[Z(\mathcal{A})] =: \mathcal{G}(\mathcal{A})$ 称为点集 \mathcal{A} 的 Gauss 复杂度. 我们记 $f(\mathbf{x}) := \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \mathbf{x} \rangle$, 对任意向量 $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$, $\mathbf{a} \in \mathcal{A}$, 由 Cauchy-Schwarz 不等式,

$$\langle \mathbf{a}, \mathbf{x} \rangle - f(\mathbf{x}') \leq \langle \mathbf{a}, \mathbf{x} - \mathbf{x}' \rangle \leq \mathcal{W}(\mathcal{A}) \|\mathbf{x} - \mathbf{x}'\|_2,$$

其中 $\mathcal{W}(\mathcal{A}) = \sup_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a}\|_2$ 为 \mathcal{A} 的欧氏宽度. 对左侧 \mathbf{a} 在集合 \mathcal{A} 中取上确界, 再交换 \mathbf{x} 和 \mathbf{x}' 的位置我们有 $|f(\mathbf{x}) - f(\mathbf{x}')| \leq \mathcal{W}(\mathcal{A}) \|\mathbf{x} - \mathbf{x}'\|_2$, 即 f 是 $\mathcal{W}(\mathcal{A})$ -Lipschitz 的. 从而由定理 1.9, $Z = f(\mathbf{X})$ 是 $\mathcal{W}(\mathcal{A})$ -次高斯的.

1.1.2 次指数随机变量

很多时候随机变量的中心矩生成函数只会在 0 的某个邻域内存在, 相应地, 我们将次高斯随机变量的条件放宽.

1.12 定义 (次指数随机变量). 称期望为 μ 的随机变量 X 为 (ν, α) -次指数的, 如果存在非负参数对 (ν, α) 使得

$$\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{\nu^2 \lambda^2}{2}}, \quad \forall |\lambda| < \frac{1}{\alpha}.$$

如果记 $+\infty = \frac{1}{0}$, 那么参数为 σ 的次高斯随机变量是 $(\sigma, 0)$ -次指数的——参数 α 衡量了次指数随机变量与次高斯随机变量相差“多大”.

和次高斯随机变量类似, 我们利用 Chernoff 方法可以得到它的尾部不等式

1.13 命题 (次指数随机变量的上偏差 inequality). 设随机变量 X 是 (ν, α) -次指数的, 那么

$$\mathbb{P}[X \geq \mathbb{E}[X] + t] \leq \begin{cases} e^{-\frac{t^2}{2\nu^2}}, & 0 \leq t \leq \frac{\nu^2}{\alpha}, \\ e^{-\frac{t}{\alpha}}, & t > \frac{\nu^2}{\alpha}. \end{cases}$$

和次指数随机变量类似, 次指数随机变量也具有可加性, 只需改变相应的参数即可.

1.14 命题. 设随机变量序列 $\{X_k\}_{k=1}^n$ 中 X_k 为均值为 μ_k , 参数为 (ν_k, σ_k) 的独立次指数随机变量, 那么 $\sum_k X_k$ 是 (ν_*, α_*) -次指数的, 其中

$$\nu_* := \sqrt{\sum_{k=1}^n \nu_k^2}, \quad \alpha_* := \max_{1 \leq k \leq n} \alpha_k.$$

1.15 示例 (χ^2 集中度). 设随机变量 $Y = \sum_{k=1}^n Z_k^2 \sim \chi_n^2$, 其中 $Z_k \sim \mathcal{N}(0, 1)$ 相互独立, 且有 $E[Z_k^2] = 1$, 注意到中心化矩生成函数

$$\mathbb{E}[e^{\lambda(Z_k^2-1)}] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{\lambda(z^2-1)} \cdot e^{-z^2/2} dz = \frac{e^{-\lambda}}{\sqrt{1-2\lambda}}$$

在 $\lambda > \frac{1}{2}$ 时不存在, 在 $|\lambda| < \frac{1}{4}$ 时有 $\mathbb{E}[e^{\lambda(Z_k^2-1)}] \leq e^{2\lambda^2}$, 从而 Z_k^2 为 $(2, 4)$ -次指数的. 再由可加性, Y 是 $(2\sqrt{n}, 4)$ -次指数的, 从而有比示例 1.10 更紧的尾部概率界

$$\mathbb{P}[|Y - n| \geq nt] \leq 2e^{-\frac{nt^2}{8}}, \quad \forall t \in (0, 1). \quad (5)$$

1.16 示例 (Johnson-Lindenstrauss 嵌入). 我们考虑对 $N \geq 2$ 个向量 $\{\mathbf{u}_1, \dots, \mathbf{u}_N\} \subseteq \mathbb{R}^d$ 进行降维, 即寻求映射 $F: \mathbb{R}^d \rightarrow \mathbb{R}^m$ (其中 m 可能远小于 d), 能够保留向量集中的一些重要特征, 例如距离、内积或者范数. 例如我们寻求满足

$$1 - \delta \leq \frac{\|F(\mathbf{u}_i) - F(\mathbf{u}_j)\|_2^2}{\|\mathbf{u}_i - \mathbf{u}_j\|_2^2} \leq 1 + \delta, \quad \forall \mathbf{u}_i \neq \mathbf{u}_j, \quad (6)$$

的映射 F 其中 $\delta \in (0, 1)$ 为我们容许的误差. 下面我们验证 $F: \mathbf{u} \mapsto \mathbf{X}\mathbf{u}/\sqrt{m}$ 满足上述条件, 其中随机矩阵 $\mathbf{X} \in \mathbb{R}^{m \times d}$ 各元素为独立的标准正态分布随机变量. 记 $\mathbf{x}_i \in \mathbb{R}^d$ 为 \mathbf{X} 的第 i 行, 对于非零向量 $\mathbf{u} \in \mathbb{R}^d$, $\langle \mathbf{x}_i, \mathbf{u}/\|\mathbf{u}\|_2 \rangle \sim \mathcal{N}(0, 1)$,

$$Y := \frac{\|\mathbf{X}\mathbf{u}\|_2^2}{\|\mathbf{u}\|_2^2} = \sum_{i=1}^m \left\langle \mathbf{x}_i, \frac{\mathbf{u}}{\|\mathbf{u}\|_2} \right\rangle^2 \sim \chi_m^2$$

于是由(5), 对任意 $\delta \in (0, 1)$,

$$2e^{-\frac{m\delta^2}{8}} \geq \mathbb{P}\left[\left|\frac{\|\mathbf{X}\mathbf{u}\|_2^2}{m\|\mathbf{u}\|_2^2} - 1\right| \geq \delta\right] = \mathbb{P}\left[\frac{\|F(\mathbf{u})\|_2^2}{\|\mathbf{u}\|_2^2} \notin [1 - \delta, 1 + \delta]\right].$$

由于在至多相差一个正负号下, 非零的 $\mathbf{u}_i - \mathbf{u}_j \neq 0$ 至多有 $\binom{N}{2}$ 种取法, F 为线性映射, 我们有

$$\mathbb{P}\left[\frac{\|F(\mathbf{u}_i) - F(\mathbf{u}_j)\|_2^2}{\|\mathbf{u}_i - \mathbf{u}_j\|_2^2} \notin [1 - \delta, 1 + \delta], \mathbf{u}_i \neq \mathbf{u}_j\right] \leq 2\binom{N}{2}e^{-\frac{m\delta^2}{8}}.$$

于是对任意的 $\epsilon \in (0, 1)$, 取 $m \geq \frac{16}{\delta^2} \log \frac{N}{\epsilon}$ 可以使 F 满足(6)的概率大于 $1 - \epsilon$.

有时无法得到随机变量的矩生成函数, 我们可以通过控制随机变量的矩来控制矩生成函数. 称期望为 μ , 方差为 σ^2 的随机变量 X 满足参数为 b 的 **Bernstein 条件**, 如果

$$\left| \mathbb{E}[(X - \mu)^k] \right| \leq \frac{k!}{2} \sigma^2 b^{k-2}, \quad \forall k = 2, 3, 4, \dots \quad (7)$$

此时对任意 $|\lambda| < 1/b$, 我们有

$$\begin{aligned} \mathbb{E}[e^{\lambda(X-\mu)}] &= 1 + \frac{\sigma^2 \lambda^2}{2} + \sum_{k=3}^{\infty} \frac{\mathbb{E}[(X - \mu)^k]}{k!} \lambda^k \leq 1 + \frac{\sigma^2 \lambda^2}{2} \sum_{k=0}^{\infty} (|\lambda|b)^k \\ &= 1 + \frac{\sigma^2 \lambda^2}{2(1 - |\lambda|b)} \leq \exp\left(\frac{\sigma^2 \lambda^2}{2(1 - |\lambda|b)}\right). \end{aligned} \quad (8)$$

进一步地, 当 $|\lambda| < 1/2b$ 时, 我们有 $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{(\sqrt{2}\sigma)^2 \lambda^2}{2}}$, 即 X 为参数为 $(\sqrt{2}\sigma, 2b)$ 的次指数随机变量. 此时应用命题 1.13 可以得到满足 Bernstein 条件的随机变量的尾部概率界, 但是如果直接从不等式 (8) 出发, 在 Chernoff 不等式中取 $\lambda = \frac{t^2}{bt + \sigma^2} < \frac{1}{b}$, 至少可以得到常数项更紧的尾部概率界:

1.17 命题 (Bernstein 界). 若随机变量 X 满足 Bernstein 条件 (7), 那么

$$\mathbb{P}[|X - \mu| \geq t] \leq 2e^{-\frac{t^2}{2(\sigma^2 + bt)}}, \quad \forall t \geq 0. \quad (9)$$

1.18 注. Bernstein 条件有着广泛的应用: 有很多满足 Bernstein 条件无界随机变量, 即使对于有界随机变量, Bernstein 条件也可能会得到比 Hoeffding 界更紧的尾部概率界.

例如, 有界随机变量 $|X - \mu| \leq b$ 是 b -次高斯的, 我们可以得到它的 Hoeffding 界; 此外它还满足条件 (7), 从而我们还可以得到它的 Bernstein 界 (9). 对于充分小的 t , 随机变量的表现会和参数为 $\sigma = \sqrt{\mathbb{E}[(X - \mu)^2]} \leq b$ 的次高斯随机变量一致, 于是 Bernstein 界不会比原先的结果差——这是因为我们利用了方差的信息. 特别对于方差小但是可能取到较大值的随机变量, 此时 $\sigma \ll b$, 我们可以得到更紧的尾部概率界.

在随机变量只有单侧界时, 我们也可以得到单侧的偏差 inequality.

1.19 命题 (单侧 Bernstein 不等式). 若 $X \leq b$, a.s., 那么它的中心化矩生成函数

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \exp\left(\frac{\lambda^2 \mathbb{E}[X^2]/2}{1 - b\lambda/3}\right), \quad \forall \lambda \in \left[0, \frac{3}{b}\right).$$

相应的, 给定 n 个满足 $X_k \leq b$ a.s. 的独立随机变量, 我们有

$$\mathbb{P}\left[\sum_{k=1}^n (X_k - \mathbb{E}[X_k]) \geq n\delta\right] \leq \exp\left(-\frac{n\delta^2}{2\left(\frac{1}{n} \sum_k \mathbb{E}[X_k^2] + \frac{b\delta}{3}\right)}\right).$$

证明. 定义单增函数 $h(u) := \frac{2}{u^2}(e^u - u - 1)$, 于是当 $\lambda \geq 0$ 时, 中心化矩生成函数

$$\begin{aligned} \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] &= \left(1 + \lambda \mathbb{E}[X] + \frac{\lambda^2}{2} \mathbb{E}[X^2] h(\lambda X)\right) e^{-\lambda \mathbb{E}[X]} \leq \left(1 + \lambda \mathbb{E}[X] + \frac{\lambda^2}{2} \mathbb{E}[X^2] h(\lambda b)\right) e^{-\lambda \mathbb{E}[X]} \\ &\leq \exp\left(\lambda \mathbb{E}[X] + \frac{\lambda^2}{2} \mathbb{E}[X^2] h(\lambda b)\right) \cdot e^{-\lambda \mathbb{E}[X]} \leq \exp\left(\frac{\lambda^2}{2} \mathbb{E}[X^2] h(\lambda b)\right). \end{aligned}$$

而在 $k \geq 2$ 时, 有 $k! \geq 2 \cdot 3^{k-2}$, 由此可得在 $\lambda b < 3$ 时成立

$$h(\lambda b) = 2 \sum_{k=2}^{\infty} \frac{(\lambda b)^{k-2}}{k!} \leq \sum_{k=2}^{\infty} \left(\frac{\lambda b}{3} \right)^{k-2} = \frac{1}{1 - \frac{\lambda b}{3}}.$$

□

1.20 注. 对于有下界的随机变量 X , 我们只需把上述结果应用至 $-X$ 即可

最后, 我们给出一些次高斯随机变量的等价描述.

1.21 定理 (次高斯随机变量的等价定义). 对任意均值为 0 的随机变量 X , 下述几条命题等价:

(I) 存在非负参数 (ν, α) 使得

$$\mathbb{E}[e^{\lambda X}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}, \quad \forall |\lambda| < \frac{1}{\alpha}.$$

(II) 存在常数 $c_0 > 0$, 使得对于任意 $|\lambda| \leq c_0$, 都有 $\mathbb{E}[e^{\lambda X}] < \infty$.

(III) 存在常数 $c_1, c_2 > 0$, 使得

$$\mathbb{P}[|X| \geq t] \leq c_1 e^{-c_2 t}, \quad \forall t > 0.$$

(IV) 量 $\gamma := \sup_{k \geq 2} \left[\frac{\mathbb{E}[X^k]}{k!} \right]^{1/k}$ 是有限的.

1.2 鞅方法

在本章的开始, 我们讨论了独立随机变量和 $f(X_1, \dots, X_n) = \sum_k X_k$ 的一些尾概率界. 对于更一般的函数 f , 建立尾概率界的经典方法方法是利用鞅差分解, 再逐个做估计.

设随机向量 $\mathbf{X} = (X_1, \dots, X_n)$ 各分量独立, 函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 满足 $\mathbb{E}[f(\mathbf{X})] < \infty$. 考虑随机变量序列 $Y_0 = \mathbb{E}[f(\mathbf{X})]$, $Y_k = \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^k]$, $k = 1, \dots, n$. 由条件期望的性质易见 $Y_n = f(\mathbf{X})$ a.s. 且 $\{Y_k\}_{k=1}^n$ 是适应于 $\{\mathbf{X}_1^k\}_{k=1}^n$ 的鞅:

- 由 Jensen 不等式和重期望公式,

$$\mathbb{E}[|Y_k|] = \mathbb{E}[\mathbb{E}[|f(\mathbf{X})| | \mathbf{X}_1^k]] \leq \mathbb{E}[\mathbb{E}[|f(\mathbf{X})| | \mathbf{X}_1^k]] = \mathbb{E}[|f(\mathbf{X})|] < \infty$$

- $\mathbb{E}[Y_{k+1} | \mathbf{X}_1^k] = \mathbb{E}[\mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k+1}] | \mathbf{X}_1^k] = \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^k] = Y_k$.

从而 $f(\mathbf{X})$ 和 $\mathbb{E}[f(\mathbf{X})]$ 的偏差可以表示为鞅差分解

$$f(\mathbf{X}) - \mathbb{E}[f(\mathbf{X})] = Y_n - Y_0 = \sum_{k=1}^n (Y_k - Y_{k-1}) =: \sum_{k=1}^n D_k.$$

我们先来证明一个一般的鞅差序列的 Bernstein 界.

1.22 引理 (Azuma). 设鞅差序列 $\{(D_k, \mathcal{F}_k)\}_{k=1}^\infty$ 满足次指数条件

$$\mathbb{E}[e^{\lambda D_k} | \mathcal{F}_{k-1}] \leq e^{\frac{\lambda^2 \nu_k^2}{2}} \text{ a.e. }, \forall |\lambda| < \frac{1}{\alpha_k},$$

那么和式 $\sum_k D_k$ 是 $(\sqrt{\sum_k \nu_k^2}, \alpha_*)$ -次指数随机变量, 其中 $\alpha_* := \max_k \alpha_k$. 再沿用 Chernoff 方法, 可以得到集中不等式

$$\mathbb{P}\left[\left|\sum_{k=1}^n D_k\right| \geq t\right] \leq \begin{cases} 2 \exp\left(-\frac{t^2}{2 \sum_k \nu_k^2}\right), & 0 \leq t \leq \alpha_*^{-1} \sum_k \nu_k^2, \\ 2 \exp\left(-\frac{t}{2 \alpha_*}\right), & t > \alpha_*^{-1} \sum_k \nu_k^2. \end{cases}$$

证明. 我们使用控制鞅差和的标准方法, 对于 $|\lambda| < \alpha_*^{-1}$,

$$\begin{aligned} \mathbb{E}[e^{\lambda \sum_k D_k}] &= \mathbb{E}\left[\mathbb{E}[e^{\lambda \sum_k D_k} | \mathcal{F}_{n-1}]\right] = \mathbb{E}\left[e^{\lambda \sum_{k=1}^{n-1} D_k} \mathbb{E}[e^{\lambda D_n} | \mathcal{F}_{n-1}]\right] \\ &\leq \mathbb{E}\left[e^{\lambda \sum_{k=1}^{n-1} D_k}\right] e^{\frac{\lambda^2 \nu_n^2}{2}} \leq \dots \leq \exp\left(\frac{\lambda^2}{2} \cdot \sum_{k=1}^n \nu_k^2\right). \end{aligned}$$

□

1.23 定理 (Azuma-Hoeffding 不等式). 若鞅差序列 $\{(D_k, \mathcal{F}_k)\}_{k=1}^\infty$ 满足 $D_k \in [a_k, b_k]$ a.s., 则 $\sum_k D_k$ 是 $\frac{1}{2}\sqrt{\sum_k (b_k - a_k)^2}$ -次高斯的, 且有集中不等式

$$\mathbb{P}\left[\left|\sum_{k=1}^n D_k\right| \geq t\right] \leq 2 \exp\left(-\frac{2t^2}{\sum_k (b_k - a_k)^2}\right), \quad \forall t \geq 0.$$

证明. 由引理 1.7, 条件随机变量 $D_k | \mathcal{F}_{k-1}$ 是 $\frac{b_k - a_k}{2}$ -次高斯的. 再根据上一证明中控制鞅差和的方法, 不难得到 Hoeffding 型集中不等式. □

称函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 满足参数为 (L_1, \dots, L_n) 的有界差不等式, 如果对指标 $k = 1, \dots, n$, 总有

$$|f(\mathbf{x}_1^{k-1}, x_k, \mathbf{x}_{k-2}^n) - f(\mathbf{x}_1^{k-1}, x'_k, \mathbf{x}_{k-2}^n)| \leq L_k.$$

可以证明, 满足有界差不等式的函数一定有界, 而有界函数显然满足有界差不等式——从而二者等价, 但是具体的参数可以让我们做出更为精细的估计. 在有界差不等式的条件下, 每个随机变量对函数的贡献是容易估计的.

1.24 推论 (有界差不等式). 设函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 满足参数为 (L_1, \dots, L_n) 的有界差不等式, 随机向量 $\mathbf{X} = (X_1, \dots, X_n)$ 各分量独立, 那么 $f(\mathbf{X})$ 是 $\sqrt{\sum_k L_k^2}/2$ -次高斯的, 从而有集中不等式

$$\mathbb{P}[|f(\mathbf{X}) - \mathbb{E}[f(\mathbf{X})]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_k L_k^2}\right), \quad \forall t \geq 0. \quad (10)$$

證明. 定义随机变量

$$A_k := \inf_x \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k-1}, X_k = x] - \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k-1}],$$

$$B_k := \sup_x \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k-1}, X_k = x] - \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k-1}].$$

于是鞅差 $D_k = \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^k] - \mathbb{E}[f(\mathbf{X}) | \mathbf{X}_1^{k-1}] \in [A_k, B_k]$ a.s. 且区间长度

$$B_k - A_k = \sup_{x,y} \mathbb{E}[f(\mathbf{X}_1^k, x, \mathbf{X}_{k+1}^n) - f(\mathbf{X}_1^k, y, \mathbf{X}_{k+1}^n) | \mathbf{X}_1^{k-1}] \leq L_k.$$

于是由定理1.23可以得到集中不等式 □

1.25 注 (次高斯函数). 注意我们对随机向量 \mathbf{X} 仅做了各分量相互独立的要求, 这里的次高斯性更多的是函数 f 本身的性质——于是我们称函数 f 是 $\sqrt{\sum_k L_k^2}/2$ -次高斯的.

1.26 示例. 若函数 f 关于 Hamming 距离 $d_H(\mathbf{X}, \mathbf{Y}) = \sum_k \mathbb{I}_{\{x_k \neq y_k\}}$ 是 L -Lipschitz 的, 那么 f 满足参数为 (L, \dots, L) 的有界差不等式. 于是 f 是 $\frac{\sqrt{n}L}{2}$ -次高斯的.

1.27 示例 (点集的 Rademacher 复杂度). Rademacher 随机变量 ϵ 是指等概率地取 $\{-1, 1\}$ 的随机变量, Rademacher 随机向量 $\epsilon \subseteq \{-1, 1\}^n$ 的各分量是相互独立的 Rademacher 随机变量. 给定 \mathbb{R}^n 的子集 \mathcal{A} , 随机变量

$$Z(\mathcal{A}) := \sup_{\mathbf{a} \in \mathcal{A}} \left[\sum_{k=1}^n a_k \epsilon_k \right] = \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \epsilon \rangle$$

在某种意义下度量了集合 \mathcal{A} 的大小, 它的期望 $\mathbb{E}_\epsilon[Z(\mathcal{A})] =: \mathcal{R}(\mathcal{A})$ 称为点集 \mathcal{A} 的 Rademacher 复杂度. 这里我们说明在集合 \mathcal{A} 满足 $a_k^* = \sup_{\mathbf{a} \in \mathcal{A}} |a_k| < \infty, k = 1, \dots, n$ 的情况下, $Z(\mathcal{A})$ 是次高斯随机变量.

记 $f(\epsilon) := \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \epsilon \rangle$, 我们对第 k 个坐标分量进行扰动: 考虑 $\epsilon' \in \{-1, 1\}^n$ 满足 $\epsilon'_i = \epsilon_i, i \neq k$, 于是对任意 $\mathbf{a} \in \mathcal{A}$,

$$\langle \mathbf{a}, \epsilon \rangle - f(\epsilon') = \langle \mathbf{a}, \epsilon \rangle - \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \epsilon' \rangle \leq \langle \mathbf{a}, \epsilon - \epsilon' \rangle = a_k(\epsilon_k - \epsilon'_k) \leq 2|a_k| \leq 2a_k^*,$$

对左侧中 \mathbf{a} 在集合 \mathcal{A} 上取上确界有 $f(\epsilon) - f(\epsilon') \leq 2a_k^*$, 交换 ϵ 和 ϵ' 的位置可以得到 f 满足参数为 $(2a_1^*, \dots, 2a_n^*)$ 的有界差不等式. 再由定理1.23, $Z = f(\epsilon)$ 是 $\sqrt{\sum_k (a_k^*)^2}$ -次高斯的.

1.3 熵方法

设随机变量 $X \sim \mathbb{P}$, 给定凸函数 $\Phi: \mathbb{R} \rightarrow \mathbb{R}$, 若 X 和 $\Phi(X)$ 的期望存在, 则称概率分布空间上的泛函

$$\mathcal{H}_\Phi(X) := \mathbb{E}[\Phi(X)] - \Phi(\mathbb{E}[X])$$

为 X 的 Φ -熵. 由 Jensen 不等式易见 $\mathcal{H} \geq 0$, 它可以用来衡量随机变量随机性的大小:

- 若 $X = C$ a.e., 那么 $\mathcal{H}_\Phi(X) = 0$;

- 如果取 $\Phi(u) = u^2$, 此时熵 $\mathcal{H}_\Phi(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$ 就是方差. Chebyshev 不等式陈述的就是方差熵可以控制尾部概率这一事实;
- 如果取 $\Phi(u) = -\log u$, 随机变量 $Z = e^{\lambda X}$ 时,

$$\mathcal{H}_\Phi(Z) = -\lambda \mathbb{E}[X] + \log \mathbb{E}[e^{\lambda X}] = \log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}],$$

这样的熵对应的是中心化的累计生成函数, 可以用来计算尾部概率的 Chernoff 界.

本节我们总是考虑非负随机变量 $e^{\lambda X}$ 由凸函数 $\Phi: [0, \infty) \rightarrow \mathbb{R}$:

$$\Phi(u) = \begin{cases} u \log u, & u > 0; \\ 0, & u = 0. \end{cases}$$

诱导的熵:

$$\mathcal{H}(e^{\lambda X}) = \lambda \mathbb{E}[X e^{\lambda X}] - \mathbb{E}[e^{\lambda X}] \log \mathbb{E}[e^{\lambda X}] = \lambda M'_X(\lambda) - M_X(\lambda) \log M_X(\lambda), \quad (11)$$

其中 $M_X(\lambda) = \mathbb{E}[e^{\lambda X}]$ 为矩生成函数.

1.28 示例 (Gauss 随机变量的熵). 随机变量 $X \sim \mathcal{N}(\mu, \sigma^2)$ 的矩生成函数为 $M_X(\lambda) = e^{\lambda^2 \sigma^2 / 2 + \lambda \mu}$, 于是

$$\mathcal{H}(e^{\lambda X}) = (\lambda^2 \sigma^2 + \lambda \mu) M_X(\lambda) - \left(\frac{\lambda^2 \sigma^2}{2} + \lambda \mu \right) M_X(\lambda) = \frac{\lambda^2 \sigma^2}{2} \cdot M_X(\lambda). \quad (12)$$

1.29 定理 (熵的变分表示).

$$\mathcal{H}(e^{\lambda f(X)}) = \sup \left\{ \mathbb{E}[g(X) e^{\lambda f(X)}]; \mathbb{E}[e^{g(X)}] \leq 1 \right\} \quad (13)$$

证明. 考虑测度 $\mathbb{P}^g[A] = \mathbb{E}^g[\mathbb{I}_A] := \mathbb{E}[e^{g(X)}; A]$, 此时期望的 Jensen 不等式依然成立. 一方面, 我们有

$$\begin{aligned} \mathcal{H}(e^{\lambda f(X)}) - \mathbb{E}[g(X) e^{\lambda f(X)}] &= \mathbb{E} \left[(\lambda f(X) - g(X)) e^{\lambda f(X)} \right] - \mathbb{E}[e^{\lambda f(X)}] \log \mathbb{E}[e^{\lambda f(X)}] \\ &= \mathbb{E}^g \left[(\lambda f(X) - g(X)) e^{\lambda f(X) - g(X)} \right] - \mathbb{E}^g[e^{\lambda f(X) - g(X)}] \log \mathbb{E}^g[e^{\lambda f(X) - g(X)}] \\ &= \mathcal{H}^g(e^{\lambda f(X) - g(X)}) \geq 0. \end{aligned}$$

另一方面, 容易验证函数 $g(x) = \lambda f(x) - \log \mathbb{E}[e^{\lambda f(X)}]$ 恰好使不等式取到等号. \square

使用独立复制的方法可以对单变量函数熵的界做如下估计.

1.30 引理 (单变量函数熵的界). 设独立随机变量 $X, Y \sim \mathbb{P}$, 对任意函数 $g: \mathbb{R} \rightarrow \mathbb{R}$, 我们有

$$\mathcal{H}(e^{\lambda g(X)}) \leq \lambda^2 \mathbb{E}[(g(X) - g(Y))^2 e^{\lambda g(X)}; g(X) \geq g(Y)], \quad \forall \lambda > 0. \quad (14)$$

特别地, 当 X 支撑集为 $[a, b]$, g 为凸的 Lipschitz 函数时, 我们有

$$\mathcal{H}(e^{\lambda g(X)}) \leq \lambda^2 (b-a)^2 \mathbb{E}[(g'(X))^2 e^{\lambda g(X)}], \quad \forall \lambda > 0.$$

证明. 由 Jensen 不等式和 X 和 Y 的对称性

$$\begin{aligned} \mathcal{H}(e^{\lambda g(X)}) &= \mathbb{E}[\lambda g(X) e^{\lambda g(X)}] - \mathbb{E}[e^{\lambda g(Y)}] \log \mathbb{E}[e^{\lambda g(X)}] \leq \mathbb{E}[\lambda g(X) e^{\lambda g(X)}] - \mathbb{E}[\lambda g(X) e^{\lambda g(Y)}] \\ &= \lambda \mathbb{E}[g(X)(e^{\lambda g(X)} - e^{\lambda g(Y)})] = -\lambda \mathbb{E}[g(Y)(e^{\lambda g(X)} - e^{\lambda g(Y)})] \\ &= \frac{\lambda}{2} \mathbb{E}[(g(X) - g(Y))(e^{\lambda g(X)} - e^{\lambda g(Y)})] \\ &= \lambda \mathbb{E}[(g(X) - g(Y))(e^{\lambda g(X)} - e^{\lambda g(Y)}); g(X) \geq g(Y)]. \end{aligned}$$

再根据指数函数的凸性, 在 $\{g(X) \geq g(Y)\}$, $\lambda > 0$ 时, 有

$$e^{\lambda g(X)} - e^{\lambda g(Y)} \leq \lambda(g(X) - g(Y))e^{\lambda g(X)},$$

从而(14)成立. 进一步地, 当 g 为凸的 Lipschitz 函数时, 由 Rademacher 定理 A.10, g' 几乎处处存在. 于是在 $\{g(X) \geq g(Y)\}$, $\lambda > 0$, X 支撑集为 $[a, b]$ 的条件下, 我们有估计

$$(g(X) - g(Y))^2 \leq (g'(X))^2 (X - Y)^2 \leq (b-a)^2 (g'(X))^2 \quad \text{a.s.}$$

代入(14)中可见引理成立. □

1.3.1 Herbst 方法

和次高斯随机变量的想法类似, 如果 $e^{\lambda X}$ 的熵能像 Gauss 随机变量的熵那样被控制, 它的矩生成函数、尾部概率也会被控制.

1.31 定理 (Herbst 方法). 若对任意的 $\lambda \in I$ (这里 I 取 $[0, \infty)$ 或者 \mathbb{R}) 总有熵 $\mathcal{H}(e^{\lambda X}) \leq \frac{\lambda^2 \sigma^2}{2} \cdot \mathbb{E}[e^{\lambda X}]$, 那么

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right), \quad \forall \lambda \in I.$$

进一步地, 由 Chernoff 方法, $I = [0, \infty)$ 时, X 满足次高斯随机变量的上偏差不等式; 当 $I = \mathbb{R}$ 时, X 是 σ -次高斯的, 从而满足集中不等式 (4).

证明. 对于 $\lambda \in I \setminus \{0\}$, 令 $G(\lambda) := \frac{\log M_X(\lambda)}{\lambda}$, 于是

$$G'(\lambda) = \frac{\lambda M'_X(\lambda) - M_X(\lambda) \log M_X(\lambda)}{\lambda^2 M_X(\lambda)} = \frac{\mathcal{H}(e^{\lambda X})}{\lambda^2 M_X(\lambda)} \leq \frac{\sigma^2}{2}, \quad \forall \lambda \in I \setminus \{0\}.$$

当 $\lambda > 0$ 时, 对任意的 $0 < \lambda_0 < \lambda$, 在区间 $[\lambda_0, \lambda]$ 上积分有 $G(\lambda) - G(\lambda_0) \leq \frac{\sigma^2(\lambda - \lambda_0)}{2}$. 由洛必达法则, $\lim_{\lambda \downarrow 0} G(\lambda) = \lim_{\lambda \downarrow 0} \frac{M'_X(\lambda)}{M_X(\lambda)} = \mathbb{E}[X]$. 于是令 $\lambda_0 \rightarrow 0^+$ 有

$$\log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] = \lambda(G(\lambda) - \mathbb{E}[X]) \leq \frac{\lambda^2 \sigma^2}{2}, \quad \forall \lambda \geq 0.$$

类似的, 我们可以证明目标不等式在 $\lambda \leq 0$ 时成立. □

1.32 注. 反之, 若 X 为参数是 $\sigma/2$ -的次高斯随机变量, 那么有 $\mathcal{H}(e^{\lambda X}) \leq \frac{\lambda^2 \sigma^2}{2} \cdot \mathbb{E}[e^{\lambda X}]$.

證明. 考虑随机变量 $Z := e^{\lambda X} / \mathbb{E}[e^{\lambda X}]$, 那么

$$\mathbb{E}[Z \log Z] = \frac{1}{\mathbb{E}[e^{\lambda X}]} \mathbb{E}[e^{\lambda X} (\lambda X - \log \mathbb{E}[e^{\lambda X}])] = \frac{\mathcal{H}(e^{\lambda X})}{\mathbb{E}[e^{\lambda X}]}.$$

考虑引理 1.7 证明中的指数加权期望, 由 Jensen 不等式和次高斯假设, 我们还有

$$\begin{aligned} \mathbb{E}[Z \log Z] &= \frac{\mathbb{E}[\log Z \cdot e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} = \mathbb{E}_\lambda[\log Z] \leq \log \mathbb{E}_\lambda[Z] = \log \mathbb{E}[e^{2\lambda X}] - 2 \log \mathbb{E}[e^{\lambda X}] \\ &\leq 2\lambda \mathbb{E}[X] - 2 \log \mathbb{E}[e^{\lambda X}] + \frac{\lambda^2 \sigma^2}{2} \leq \frac{\lambda^2 \sigma^2}{2}. \end{aligned}$$

于是对任意 $\lambda \in \mathbb{R}$, 都有

$$\mathcal{H}(e^{\lambda X}) = \mathbb{E}[Z \log Z] \cdot \mathbb{E}[e^{\lambda X}] \leq \frac{\lambda^2 \sigma^2}{2} \cdot \mathbb{E}[e^{\lambda X}].$$

□

自然地, 我们可以将上述方法由次高斯随机变量推广至次指数随机变量.

1.33 命题 (Bernstein 熵的界). 若存在正常数 b, σ 使得

$$\mathcal{H}(e^{\lambda X}) \leq \lambda^2 [bM'_X(\lambda) + M_X(\lambda)(\sigma^2 - b\mathbb{E}[X])], \quad \forall \lambda \in [0, 1/b),$$

那么 X 满足上界

$$\log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \frac{\lambda^2 \sigma^2}{1 - b\lambda}, \quad \forall \lambda \in [0, 1/b).$$

进一步地, 由 Chernoff 方法, X 满足上偏差不等式

$$\mathbb{P}[X \geq \mathbb{E}[X] + t] \leq \exp\left(-\frac{t^2}{4\sigma^2 + 2bt}\right), \quad \forall t \geq 0.$$

證明. 类似地, 命题中的条件意味着

$$G'(\lambda) = \frac{\mathcal{H}(e^{\lambda X})}{\lambda^2 M_X(\lambda)} \leq \sigma^2 - b\mathbb{E}[X] + b \cdot \frac{M'_X(\lambda)}{M_X(\lambda)}.$$

在任意的区间 $[\lambda_0, \lambda] \subseteq (0, 1/b)$ 上积分有

$$G(\lambda) - G(\lambda_0) \leq (\sigma^2 - b\mathbb{E}[X])(\lambda - \lambda_0) + b(\log M_X(\lambda) - \log M_X(\lambda_0)).$$

再令 $\lambda_0 \rightarrow 0^+$ 有 $G(\lambda) - \mathbb{E}[X] \leq \lambda\sigma^2 + b\lambda(G(\lambda) - \mathbb{E}[X])$, 于是

$$\log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] = \lambda(G(\lambda) - \mathbb{E}[X]) \leq \frac{\lambda^2 \sigma^2}{1 - b\lambda}, \quad \forall \lambda \in [0, 1/b).$$

□

1.3.2 熵的张量化

鞅方法虽然可以得到随机变量函数的集中不等式, 但是它要求随机变量有界, 或者要求函数满足有界差不等式, 并没有利用充分各分量 X_k 的分布信息. 就像我们在注1.18中提到的那样, 即使有界的条件下, 进一步地利用分布的信息可以得到更好的估计. 例如定理1.23中, 有界随机变量 X_k 的次高斯参数可能远小于 $\frac{b_k - a_k}{2}$.

由于独立随机变量的联合分布是边缘分布的张量积, 我们可以计算每个随机变量对函数的贡献得到随机变量的函数熵的界¹——这样我们就利用了随机变量分布的信息. 这种方法被称为熵的张量化, 或者说, 熵具有次可加性.

具体地, 设 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 随机向量 $\mathbf{X} = (X_1, \dots, X_n)$ 各分量独立. 条件随机变量 $e^{\lambda f(\mathbf{X})} | \mathbf{X}^{\setminus k}$ 可以看作其它分量 $\mathbf{X}^{\setminus k}$ 已经给定, 关于 X_k 的函数. 于是我们可以在保持 $\mathbf{X}^{\setminus k}$ 不变的同时, 计算关于第 k 个分量的逐坐标条件熵

$$\mathcal{H}_k(e^{\lambda f(\mathbf{X})}) := \mathcal{H}(e^{\lambda f(\mathbf{X})} | \mathbf{X}^{\setminus k}).$$

1.34 引理 (熵的张量化). 在上述的假设和记号下,

$$\mathcal{H}(e^{\lambda f(\mathbf{X})}) \leq \mathbb{E} \left[\sum_{k=1}^n \mathcal{H}_k(e^{\lambda f(\mathbf{X})}) \right], \quad \forall \lambda > 0.$$

证明. 对任意满足 $\mathbb{E}[e^{g(\mathbf{X})}] \leq 1$ 的函数 g , 定义函数序列 $\{g^1, \dots, g^n\}$ 为

$$g^k(\mathbf{X}_k^n) := \log \mathbb{E}[e^{g(\mathbf{X})} | \mathbf{X}_k^n] - \log \mathbb{E}[e^{g(\mathbf{X})} | \mathbf{X}_{k+1}^n], \quad k = 1, \dots, n.$$

于是 $\sum_{k=1}^n g^k(\mathbf{X}_k^n) = g(\mathbf{X}) - \log \mathbb{E}[e^{g(\mathbf{X})}] \geq g(\mathbf{X})$, 利用这一分解可得

$$\mathbb{E}[g(\mathbf{X}) e^{\lambda f(\mathbf{X})}] \leq \sum_{k=1}^n \mathbb{E}[g^k(\mathbf{X}_k^n) e^{\lambda f(\mathbf{X})}] = \sum_{k=1}^n \mathbb{E} \left[\mathbb{E}[g^k(\mathbf{X}_k^n) e^{\lambda f(\mathbf{X})} | \mathbf{X}^{\setminus k}] \right]. \quad (15)$$

注意到条件随机变量 $g^k(\mathbf{X}_k^n) | \mathbf{X}^{\setminus k}$ 作为 X_k 的 (单变量) 函数满足

$$\mathbb{E}[e^{g^k(\mathbf{X}_k^n)} | \mathbf{X}^{\setminus k}] = \mathbb{E} \left[\frac{\mathbb{E}[e^{g(\mathbf{X})} | \mathbf{X}_k^n]}{\mathbb{E}[e^{g(\mathbf{X})} | \mathbf{X}_{k+1}^n]} \middle| \mathbf{X}^{\setminus k} \right] = 1,$$

另一方面, $e^{\lambda f(\mathbf{X})} | \mathbf{X}^{\setminus k}$ 也可以看作 X_k 的函数, 于是由定理1.29, 我们有 $\mathbb{E}[g^k(\mathbf{X}_k^n) e^{\lambda f(\mathbf{X})} | \mathbf{X}^{\setminus k}] \leq \mathcal{H}_k(e^{\lambda f(\mathbf{X})})$. 再次利用定理1.29, 对(15)中满足条件的 g 取上确界可知引理成立. \square

特别地, 熵的张量化在处理可分凸函数时, 可以得到很好的结论. 称函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 是可分凸的, 如果对任意指标 $k \in \{1, \dots, n\}$, 给定向量 $\mathbf{x}^{\setminus k} := (x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n) \in \mathbb{R}^{n-1}$, 单变量函数

$$y_k \mapsto f(x_1, \dots, x_{k-1}, y_k, x_{k+1}, \dots, x_n)$$

总是凸函数. 凸函数总是可分凸的.

¹最为经典的例子是独立随机变量和的方差熵是随机变量的方差和.

1.35 命题. 令 $\{X_k\}_{k=1}^n$ 为区间 $[a, b]$ 上的独立随机变量, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 为可分凸函数且关于 2-范数为 L -Lipschitz 的, 那么对任意 $t \geq 0$, 成立

$$\mathbb{P}[f(\mathbf{X}) \geq \mathbb{E}f(\mathbf{X}) + t] \leq \exp\left(-\frac{t^2}{4L^2(b-a)^2}\right).$$

证明. 对于 $\lambda > 0$, 由引理 1.34、1.30,

$$\begin{aligned} \mathcal{H}(e^{\lambda f(\mathbf{X})}) &\leq \mathbb{E} \left[\sum_{k=1}^n \mathcal{H}_k(e^{\lambda f(\mathbf{X})}) \right] \leq \lambda^2(b-a)^2 \mathbb{E} \left[\sum_{k=1}^n \mathbb{E} \left[\left(\frac{\partial f(\mathbf{X})}{\partial X_k} \right)^2 e^{\lambda f(\mathbf{X})} \middle| \mathbf{X}^{\setminus k} \right] \right] \\ &= \lambda^2(b-a)^2 \mathbb{E} \left[\sum_{k=1}^n \left(\frac{\partial f(\mathbf{X})}{\partial X_k} \right)^2 e^{\lambda f(\mathbf{X})} \right] \leq \lambda^2(b-a)^2 L^2 \mathbb{E}[e^{\lambda f(\mathbf{X})}], \end{aligned}$$

其中 $\sum_k \left(\frac{\partial f(\mathbf{X})}{\partial x_k} \right)^2 = \|\nabla f\|_2^2 \leq L^2$ 由 Lipschitz 条件得到. 再由 Herbst 方法 (定理 1.31) 可见命题成立. \square

1.36 示例 (點集的 Rademacher 複雜度). 我们给出 $Z = f(\epsilon)$ 的比示例 1.27 中更紧的概率界. 容易看到 $f(\mathbf{x}) = \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \mathbf{x} \rangle$ 为凸函数: 对任意 $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\alpha \in (0, 1)$,

$$f(\alpha \mathbf{x} + (1-\alpha)\mathbf{y}) = \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \alpha \mathbf{x} + (1-\alpha)\mathbf{y} \rangle \leq \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, \alpha \mathbf{x} \rangle + \sup_{\mathbf{a} \in \mathcal{A}} \langle \mathbf{a}, (1-\alpha)\mathbf{y} \rangle = \alpha f(\mathbf{x}) + (1-\alpha)f(\mathbf{y}).$$

另一方面, 示例 1.11 告诉我们 f 是 $\mathcal{W}(\mathcal{A})$ -Lipschitz 的, 于是由命题 1.35,

$$\mathbb{P}[Z \geq \mathbb{E}[Z] + t] \leq \exp\left(-\frac{t^2}{16\mathcal{W}^2(\mathcal{A})}\right),$$

其中 $\mathcal{W}^2(\mathcal{A}) = \sup_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a}\|_2^2$ 一般会比示例 1.27 中的 $\sum_k \sup_{\mathbf{a} \in \mathcal{A}} a_k^2$ 小很多——甚至是 n 的数量级的差距.

1.37 示例 (隨機矩陣的譜範數). 随机矩阵 $\mathbf{X} \in \mathbb{R}^{n \times d}$ 的谱范数

$$\|\mathbf{X}\|_2 = \max_{\|\mathbf{v}\|_2=1} \max_{\|\mathbf{u}\|_2=1} \mathbf{u}^T \mathbf{X} \mathbf{v},$$

作为线性函数和的上确界是一个凸函数. 另一方面, 由三角不等式

$$|\|\mathbf{X}\|_2 - \|\mathbf{X}'\|_2| \leq \|\mathbf{X} - \mathbf{X}'\|_2 \leq \|\mathbf{X} - \mathbf{X}'\|_F.$$

其中 Frobenius 范数就是所有元素的 2-范数, 于是谱范数是 1-Lipschitz 的. 若随机矩阵 \mathbf{X} 中的独立同分布于 $[-1, 1]$ 上的某个零均值分布, 由命题 1.35, 我们有

$$\mathbb{P}[\|\mathbf{X}\|_2 \geq \mathbb{E}[\|\mathbf{X}\|_2] + \delta] \leq e^{-\frac{\delta^2}{16}}.$$

1.4 等周不等式

经典的等周不等式断言, 欧氏空间 (\mathbb{R}^n, ρ) 中相同体积的子集中, 球的表面积最小. 一种等价表述是, 使得给定体积的子集 \mathcal{A} 的 (一致) ϵ -扩张

$$\mathcal{A}^\epsilon := \{x \in \mathbb{R}^n : \rho(x, \mathcal{A}) < \epsilon\}$$

的体积 (作为 ϵ 的函数) 最小的集合 A 一定是球体. 这种表述绕过了表面积的概念, 并且可以推广到任意度量空间上. 它的经典证明基于数学家 Minkowski 的天才洞见.

Minkowski 将空间的向量加法这一代数结构同凸体的体积这一几何结构联系在一起, 提出了 **Minkowski 和** 的概念—— \mathbb{R}^n 中的凸集 \mathcal{C} 和 \mathcal{D} 的 Minkowski 和定义为

$$\mathcal{C} + \mathcal{D} = \{c + d : c \in \mathcal{C}, d \in \mathcal{D}\}$$

并证明了混合体积的 **Brunn-Minkowski 不等式**

$$(\text{vol}(\lambda\mathcal{C} + (1-\lambda)\mathcal{D}))^{\frac{1}{n}} \geq \lambda(\text{vol}(\mathcal{C}))^{\frac{1}{n}} + (1-\lambda)(\text{vol}(\mathcal{D}))^{\frac{1}{n}}, \forall \lambda \in [0, 1].$$

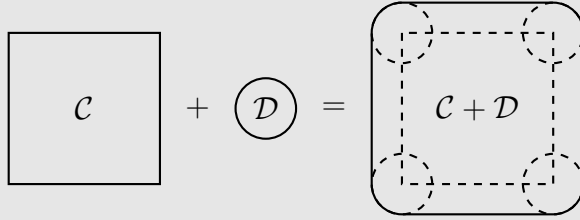


图 1: 正方形和圆的 Minkowski 和是一个具有圆角的正方形.

从这一定理出发, 很容易证明经典等周不等式: 记 \mathbb{B}^n 为 \mathbb{R}^n 中的单位球, 对任意 $\mathcal{A} \subseteq \mathbb{R}^n$ 满足 $\text{vol}(\mathcal{A}) = \text{vol}(\mathbb{B}^n)$,

$$\begin{aligned} (\text{vol}(\mathcal{A}^\epsilon))^{\frac{1}{n}} &= (\text{vol}(\mathcal{A} + \epsilon\mathbb{B}^n))^{\frac{1}{n}} = (1 + \epsilon) \left(\text{vol} \left(\frac{1}{1+\epsilon} \mathcal{A} + \frac{\epsilon}{1+\epsilon} \mathbb{B}^n \right) \right)^{\frac{1}{n}} \\ &\geq (\text{vol}(\mathcal{A}))^{\frac{1}{n}} + \epsilon(\text{vol}(\mathbb{B}^n))^{\frac{1}{n}} = (1 + \epsilon)(\text{vol}(\mathbb{B}^n))^{\frac{1}{n}} = (\text{vol}((\mathbb{B}^n)^\epsilon))^{\frac{1}{n}}. \end{aligned}$$

本节和下一节我们考虑度量空间 (\mathcal{X}, ρ) 中的集中不等式, 这需要对度量空间赋予一个概率测度 \mathbb{P} , 我们称三元组 $(\mathcal{X}, \rho, \mathbb{P})$ 为**度量测度空间**.

等周不等式可以表述为确定满足 $\mathbb{P}[A] \geq 1/2$ 、使得测度 $\mathbb{P}[A^\epsilon]$ 最小的集合 $A \subseteq \mathcal{X}$. 我们引入 $(\mathcal{X}, \rho, \mathbb{P})$ 上的集中度函数 $\alpha: \mathbb{R}_{\geq 0} \rightarrow [0, 1/2]$

$$\alpha_{\mathbb{P}}(\epsilon) := \sup_{A \subseteq \mathcal{X} : \mathbb{P}[A] \geq 1/2} \{1 - \mathbb{P}[A^\epsilon]\}.$$

于是等周不等式相当于确定 $\alpha_{\mathbb{P}}$ 的上界.

下面的定理说明, 集中度函数可以控制 Lipschitz 函数的尾部. 回忆 $f(X)$ 的**中位数**是指满足 $\mathbb{P}[f(X) \geq m_f] \geq 1/2, \mathbb{P}[f(X) \leq m_f] \geq 1/2$ 的某个常数 m_f .

1.38 定理 (Lévy 不等式). 设 $f: \mathcal{X} \rightarrow \mathbb{R}$ 关于 ρ 是 L -Lipschitz 连续的函数, $X \sim \mathbb{P}$, 有 $\mathbb{P}[|f(X) - m_f| \geq \epsilon] \leq 2\alpha(\epsilon/L)$. 特别地, 当 f 是 Lipschitz 连续函数时, 我们有

$$\mathbb{P}[|f(X) - m_f| \geq \epsilon] \leq 2\alpha(\epsilon/L).$$

證明. 令 $A := \{x \in \mathcal{X}: f(x) \leq m_f\}$, 于是 $\mathbb{P}[A] \geq 1/2$. 由扩张的定义, 对任意 $x \in A^{\epsilon/L}$, 存在 $y \in A$ 使得 $\rho(x, y) < \epsilon/L$. 于是 $f(x) < f(y) + |f(x) - f(y)| < m_f + \epsilon$, 进一步地, 我们有 $\mathbb{P}[A^{\epsilon/L}] \leq \mathbb{P}[f(X) < m_f + \epsilon]$. 取余集可以得到

$$\mathbb{P}[f(X) \geq m_f + \epsilon] \leq 1 - \mathbb{P}[A^{\epsilon/L}] \leq \alpha_{\mathbb{P}}(\epsilon/L).$$

对 $-f$ 运用相同的方法可以得到下偏差不等式, 结合起来可得集中不等式. \square

反之, Lipschitz 函数的集中不等式也蕴含着等周不等式, 于是两种 t 对尾部的控制是等价的.

1.39 定理. 若存在函数 $\beta: \mathbb{R}_{\geq 0} \rightarrow [0, 1]$ 使得对任意的 (\mathcal{X}, ρ) 上的 Lipschitz 函数都有

$$\mathbb{P}[f(X) \geq \mathbb{E}f(X) + \epsilon] \leq \beta(\epsilon), \quad \forall \epsilon \geq 0,$$

那么 $\alpha_{\mathbb{P}}(\epsilon) \leq \beta(\epsilon/2)$.

證明. 对任意 \mathcal{X} 中满足 $\mathbb{P}[A] \geq 1/2$ 的可测集 A , 构造 $f_A(x) := \rho(x, A) \wedge \epsilon$. 注意到在 A 上有 $f_A = 0$, 在 A 外有 $f_A \leq \epsilon$, 所以 $\mathbb{E}f_A(X) \leq \epsilon(1 - \mathbb{P}[A]) \leq \epsilon/2$. 于是我们有

$$1 - \mathbb{P}[A^{\epsilon}] = \mathbb{P}[X \in \bar{A}^{\epsilon}] = \mathbb{P}[f_A(X) \geq \epsilon] \leq \mathbb{P}\left[f_A(X) \geq \mathbb{E}f_A(X) + \frac{\epsilon}{2}\right] \leq \beta\left(\frac{\epsilon}{2}\right),$$

再对满足条件 $\mathbb{P}[A] \geq 1/2$ 的 A 取上确界即可. 其中最后一个不等式是由于 f_A 是一个 Lipschitz 函数:

- 若 $x, y \in A^{\epsilon}$, 则 $|f_A(x) - f_A(y)| = |\rho(x, A) - \rho(y, A)| \leq \rho(x, y)$;
- 若 $x, y \in \bar{A}^{\epsilon}$, 则 $|f_A(x) - f_A(y)| = |\epsilon - \epsilon| = 0 \leq \rho(x, y)$;
- 若 $x \in A^{\epsilon}, y \in \bar{A}^{\epsilon}$, 此时 $\rho(x, A) \geq \rho(y, A) - \rho(x, y) \geq \epsilon - \rho(x, y)$, 则 $|f_A(x) - f_A(y)| = \epsilon - \rho(x, A) \leq \epsilon - (\epsilon - \rho(x, y)) = \rho(x, y)$.

\square

1.5 傳輸成本不等式

1.5.1 Wasserstein 距離

给定 (\mathcal{X}, ρ) 上的两个概率分布 \mathbb{Q} 和 \mathbb{P} , 它们之间由 ρ 诱导的 **Wasserstein 距离**为

$$W_{\rho}(\mathbb{Q}, \mathbb{P}) := \sup_{\|f\|_{\text{Lip}} \leq 1} \{\mathbb{E}_{\mathbb{Q}}[f] - \mathbb{E}_{\mathbb{P}}[f]\} = \sup_{\|f\|_{\text{Lip}} \leq 1} \int f(d\mathbb{Q} - d\mathbb{P}).$$

可以证明这样的 W_ρ 构成了一个度量. 对任意耦合 $\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})$, $f: \mathcal{X} \rightarrow \mathbb{R}$ 为 Lipschitz 函数, 根据 Fubini 定理我们可以得到

$$\int \rho(x, x') d\mathbb{M}(x, x') \geq \int (f(x) - f(x')) d\mathbb{M}(x, x') = \int f(d\mathbb{P} - d\mathbb{Q}).$$

Kantorovich–Rubinstein 对偶告诉我们, 左侧关于耦合的下确界和右侧关于 Lipschitz 函数的上界是相等的, 即有

$$W_\rho(\mathbb{Q}, \mathbb{P}) = \inf_{\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})} \mathbb{E}_{\mathbb{M}}[\rho] = \inf_{\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})} \int_{\mathcal{X} \times \mathcal{X}} \rho(x, x') d\mathbb{M}(x, x').$$

于是 Wasserstein 距离也被称为**推土机距离**: 两个土堆的形状 (分布 \mathbb{P}, \mathbb{Q}) 确定, 传输成本 (距离 ρ) 确定, Wasserstein 距离就是最优的搬运方案 (最优的耦合 $\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})$) 下的传输成本, 等式右侧关于耦合的优化被称为最优传输问题.

1.40 示例 (Hamming 度量和全变差距离). 关于 Hamming 度量的 Wasserstein 距离 $W_{\text{Ham}}(\mathbb{Q}, \mathbb{P})$ 等价于全变差距离 $\|\mathbb{Q} - \mathbb{P}\|_{TV} := \sup_{A \subseteq \mathcal{X}} |\mathbb{Q}(A) - \mathbb{P}(A)|$, 再由 Kantorovich–Rubinstein 对偶, 我们有

$$\|\mathbb{Q} - \mathbb{P}\|_{TV} = \inf_{\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})} \int_{\mathcal{X} \times \mathcal{X}} \mathbb{I}_{\{x \neq x'\}} d\mathbb{M}(x, x') = \inf_{\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})} \mathbb{M}[Y \neq X],$$

其中 $X \sim \mathbb{P}$, $Y \sim \mathbb{Q}$.

证明. 函数 f 关于 Hamming 度量 Lipschitz 连续等价于 f 值域在某个区间 $[c, c+1]$ 中, 不失一般性地, 我们假定 $c = 0$. 记 \mathbb{Q} 和 \mathbb{P} 关于测度 ν 的密度¹分别为 q, p , 集合 $A = \{x \in \mathcal{X}: q(x) \geq p(x)\}$, 于是有

$$W_{\text{Ham}}(\mathbb{Q}, \mathbb{P}) = \sup_{f: \mathcal{X} \rightarrow [0,1]} \int_{\mathcal{X}} f(q-p) d\nu \leq \int_A (d\mathbb{Q} - d\mathbb{P}) \leq \|\mathbb{Q} - \mathbb{P}\|_{TV}.$$

另一方面, 对任意可测集 $B \subseteq \mathcal{X}$, 注意到示性函数 \mathbb{I}_B 是 Lipschitz 连续的, 于是

$$\mathbb{Q}(B) - \mathbb{P}(B) = \int \mathbb{I}_B(d\mathbb{Q} - d\mathbb{P}) \leq W_{\text{Ham}}(\mathbb{Q}, \mathbb{P}).$$

于是有 $\|\mathbb{Q} - \mathbb{P}\|_{TV} \leq W_{\text{Ham}}(\mathbb{Q}, \mathbb{P})$, 从而二者等价. □

1.5.2 KL 散度與傳輸成本不等式

\mathbb{Q} 和 \mathbb{P} 之间的 **Kullback–Leibler 散度** (亦称相对熵) 定义为

$$D(\mathbb{Q} \parallel \mathbb{P}) := \begin{cases} \mathbb{E}_{\mathbb{Q}} \left[\log \frac{d\mathbb{Q}}{d\mathbb{P}} \right] = \int_{\mathcal{X}} q(x) \log \frac{q(x)}{p(x)} \nu(dx), & \mathbb{Q} \ll \mathbb{P}, \\ +\infty, & \text{其它.} \end{cases} \quad (16)$$

其中 q, p 为 \mathbb{Q}, \mathbb{P} 的密度, ν 为 \mathcal{X} 上的 Lebesgue 测度. 尽管 KL 散度可以描述分布间的差异, 它实际上并不是一个度量: 它不满足对称性和三角不等式.

¹这样的测度 ν 是存在的, 例如 \mathbb{P} 和 \mathbb{Q} 都关于 $(\mathbb{P} + \mathbb{Q})/2$ 绝对连续.

称 (\mathcal{X}, ρ) 上的概率测度 \mathbb{P} 满足参数为 $\gamma > 0$ 的 ρ -传输成本不等式, 如果对任意的概率测度 \mathbb{Q} 总有

$$W_\rho(\mathbb{Q}, \mathbb{P}) \leq \sqrt{2\gamma D(\mathbb{Q} \parallel \mathbb{P})}. \quad (17)$$

这样的不等式在信息论中也被称为信息不等式, 例如经典的 Pinsker–Csiszár–Kullback 不等式: 对任意分布 \mathbb{Q}, \mathbb{P} ,

$$\|\mathbb{Q} - \mathbb{P}\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D(\mathbb{Q} \parallel \mathbb{P})}, \quad (18)$$

结合示例 1.40, 这实际上是说任意分布 \mathbb{P} 都满足参数 $\gamma = \frac{1}{4}$ 的 Hamming 度量传输成本不等式. 下述定理告诉我们, 在控制 Lipschitz 函数的矩生成函数时, 传输成本不等式成立.

1.41 定理 (Bobkov–Götze). 设度量空间 \mathcal{X}, ρ 上的随机变量 $X \sim \mathbb{P}$, 下列命题等价:

1. 对任意 Lipschitz 连续的函数 f , $f(X)$ 是 σ -次高斯的;
2. 概率测度 \mathbb{P} 满足参数为 σ^2 的 ρ -传输不等式, 即对任意概率测度 \mathbb{Q} 均成立

$$W_\rho(\mathbb{Q}, \mathbb{P}) \leq \sqrt{2\sigma^2 D(\mathbb{Q} \parallel \mathbb{P})}.$$

证明的关键在于注意到相对熵和矩生成函数间的密切关系, 下述的结果可以追溯到统计力学最早的历史.

1.42 引理 (Gibbs 变分原理).

$$\log \mathbb{E}_{\mathbb{P}}[e^f] = \sup_{\mathbb{Q}} \{\mathbb{E}_{\mathbb{Q}}[f] - D(\mathbb{Q} \parallel \mathbb{P})\}$$

证明. 这里我们假定 f 有界来避免可积性的问题 (否则, 考虑 $f \wedge M$ 的结果, 再对 M 取上确界即可). 考虑测度 $\tilde{\mathbb{P}}$, 它关于 \mathbb{P} 的 Radon–Nikodym 导数为 $\frac{e^f}{\mathbb{E}_{\mathbb{P}}[e^f]}$, 当 $\mathbb{Q} \ll \mathbb{P}$ 时,

$$\begin{aligned} \log \mathbb{E}_{\mathbb{P}}[e^f] - D(\mathbb{Q} \parallel \tilde{\mathbb{P}}) &= \log \mathbb{E}_{\mathbb{P}}[e^f] - \int \log \frac{d\mathbb{Q}}{d\tilde{\mathbb{P}}} d\mathbb{Q} \\ &= \log \mathbb{E}_{\mathbb{P}}[e^f] - \int \log \frac{d\mathbb{Q}}{d\mathbb{P}} d\mathbb{Q} + \int \log \frac{d\mathbb{Q}}{d\tilde{\mathbb{P}}} d\mathbb{Q} \\ &= \mathbb{E}_{\mathbb{Q}}[f] - D(\mathbb{Q} \parallel \mathbb{P}). \end{aligned}$$

于是两侧关于 \mathbb{Q} 的上确界相同, 特别是左侧中 $-D(\mathbb{Q} \parallel \tilde{\mathbb{P}})$ 关于 \mathbb{Q} 的上确界为 0. □

定理 1.41 的证明. 由定义, 条件 (1) 等价于

$$\log \mathbb{E}_{\mathbb{P}}[\exp(\lambda(f - \mathbb{E}_{\mathbb{P}}[f]))] \leq \frac{\lambda^2 \sigma^2}{2}, \quad \forall \lambda \in \mathbb{R}, \|f\|_{\text{Lip}} \leq 1.$$

结合引理 1.42, 这等价于

$$\sup_{\lambda \in \mathbb{R}} \sup_{\|f\|_{\text{Lip}} \leq 1} \sup_{\mathbb{Q}} \left\{ \lambda(\mathbb{E}_{\mathbb{Q}}[f] - \mathbb{E}_{\mathbb{P}}[f]) - D(\mathbb{Q} \parallel \mathbb{P}) - \frac{\lambda^2 \sigma^2}{2} \right\} \leq 0.$$

交换上确界的顺序并求出关于 f 和 λ 的上确界, 我们可以得到 (2) 的等价表述:

$$\sup_{\mathbb{Q}} \left\{ \frac{W_{\rho}(\mathbb{Q}, \mathbb{P})^2}{2\sigma^2} - D(\mathbb{Q} \parallel \mathbb{P}) \right\} \leq 0$$

□

1.5.3 傳輸成本的張量化

定理 1.41 描述了任意度量空间 (\mathcal{X}, ρ) 上 Lipschitz 函数的次高斯性质, 但这并不是一个高维的结果, 而推广至高维的关键思想是张量化. 在此之前, 我们先介绍 KL 散度的链法则.

1.43 引理 (Kullback–Leibler 散度的鏈法則). 给定两个 n 维分布 $\mathbb{Q} \ll \mathbb{P}$, 它们之间的 KL 散度可以分解为

$$D(\mathbb{Q} \parallel \mathbb{P}) = D(\mathbb{Q}_1 \parallel \mathbb{P}_1) + \sum_{j=2}^n \mathbb{E}_{\mathbb{Q}_1^{j-1}} \left[D \left(\mathbb{Q}_{j|j-1}(\cdot | \mathbf{X}_1^{j-1}) \parallel \mathbb{P}_{j|j-1}(\cdot | \mathbf{X}_1^{j-1}) \right) \right],$$

其中 $\mathbb{Q}_j, \mathbb{P}_j$ 分别为 \mathbb{Q}, \mathbb{P} 关于第 j 个分量的边缘分布, $\mathbb{Q}_{j|j-1}, \mathbb{P}_{j|j-1}$ 分别为给定前 $j-1$ 维随机变量下, \mathbb{Q}, \mathbb{P} 关于第 j 维的条件分布.

證明. 我们只说明 $n = 2$ 的情形, 更高维的结果可以通过归纳法得到. 记 \mathbb{Q} 和 \mathbb{P} 关于测度 ν 的密度分别为 q, p ,

$$\begin{aligned} D(\mathbb{Q} \parallel \mathbb{P}) &= \int q(x_1, x_2) \log \frac{q(x_1, x_2)}{p(x_1, x_2)} d\nu(x_1, x_2) \\ &= \int q_1(x_1) q_{2|1}(x_2|x_1) \log \frac{q_1(x_1) q_{2|1}(x_2|x_1)}{p_1(x_1) p_{2|1}(x_2|x_1)} d\nu(x_1, x_2) \\ &= \int q_1(x_1) \log \frac{q_1(x_1)}{p_1(x_1)} d\nu_1(x_1) + \int q_1(x_1) \int q_{2|1}(x_2|x_1) \log \frac{q_{2|1}(x_2|x_1)}{p_{2|1}(x_2|x_1)} d\nu(x_1, x_2) \\ &= D(\mathbb{Q}_1 \parallel \mathbb{P}_1) + \mathbb{E}_{\mathbb{Q}_1} \left[D \left(\mathbb{Q}_{2|1}(\cdot | \mathbf{X}_1) \parallel \mathbb{P}_{2|1}(\cdot | \mathbf{X}_1) \right) \right]. \end{aligned}$$

□

1.44 命題 (傳輸成本不等式的張量化). 若对于每个 $k = 1, \dots, n$, 度量空间 \mathcal{X} 上的概率测度 \mathbb{P}_k 满足参数为 γ_k 的 ρ_k -传输成本不等式, 那么乘积空间 \mathcal{X}^n 上的乘积测度 $\mathbb{P} := \otimes_{k=1}^n \mathbb{P}_k$ 满足传输成本不等式

$$W_{\rho}(\mathbb{Q}, \mathbb{P}) \leq \sqrt{2 \left(\sum_{k=1}^n \gamma_k \right) D(\mathbb{Q} \parallel \mathbb{P})}, \quad \text{对任意 } \mathcal{X}^n \text{ 上的分布 } \mathbb{Q}, \quad (19)$$

其中 $\rho(x, y) := \sum_{k=1}^n \rho_k(x_k, y_k)$.

證明. 由 Kantorovich–Rubinstein 对偶, 只需说明对任意 \mathcal{X}^n 上的分布 \mathbb{Q} , 存在某个耦合 $\mathbb{M} \in \mathcal{C}(\mathbb{Q}, \mathbb{P})$ 使得不等式 $\mathbb{E}_{\mathbb{M}}[\rho] \leq \sqrt{2(\sum_k \gamma_k) D(\mathbb{Q} \parallel \mathbb{P})}$ 成立. 我们用归纳法证明这样的分布 \mathbb{M} 的存在性.

考虑随机变量 $(\mathbf{X}_1^n, \mathbf{Y}_1^n) \sim \mathbb{M}$, 对于 $j = 1, \dots, n$, 记 \mathbb{M}_j 为 (X_j, Y_j) 的边缘分布. 特别地, $\mathbb{M}_1 \in \mathcal{C}(\mathbb{Q}_1, \mathbb{P}_1)$ 为最佳耦合, 即有

$$\mathbb{E}_{\mathbb{M}_1}[\rho_1] = W_{\rho_1}(\mathbb{Q}_1, \mathbb{P}_1) \leq \sqrt{2\gamma_1 D(\mathbb{Q}_1 \parallel \mathbb{P}_1)}.$$

记 \mathbb{M}_1^j 为 $(\mathbf{X}_1^j, \mathbf{Y}_1^j)$ 的联合分布, $\mathbb{M}_{j|j-1}$ 为给定 $(\mathbf{X}_1^{j-1}, \mathbf{Y}_1^{j-1})$ 下 (X_j, Y_j) 的条件分布. 现在假定我们得到了满足不等式

$$\mathbb{E}_{\mathbb{M}_1^{j-1}} \left[\sum_{k=1}^{j-1} \rho_k \right] \leq \sqrt{2 \left(\sum_{k=1}^{j-1} \gamma_k \right) D(\mathbb{Q}_1^{j-1} \parallel \mathbb{P}_1^{j-1})}$$

的耦合 \mathbb{M}_1^{j-1} , 由条件, 存在 $(\mathbb{Q}_{j|j-1}(\cdot | \mathbf{x}_1^{j-1}), \mathbb{P}_j)$ 的耦合 $\mathbb{M}_{j|j-1}(\cdot | \mathbf{x}_1^{j-1}, \mathbf{y}_1^{j-1})$ 使得

$$\mathbb{E}_{\mathbb{M}_{j|j-1}}[\rho_j] \leq \sqrt{2\gamma_j D(\mathbb{Q}_{j|j-1}(\cdot | \mathbf{x}_1^{j-1}) \parallel \mathbb{P}_j)}, \quad \forall \mathbf{x}_1^{j-1} \in \mathcal{X}^{j-1},$$

于是利用条件期望的塔形性质, 由平方根函数的凹性和 Jensen 不等式、Cauchy-Schwarz 不等式、KL 散度的链法则, 可以得到

$$\begin{aligned} \mathbb{E}_{\mathbb{M}}[\rho] &= \sum_{k=1}^n \mathbb{E}_{\mathbb{M}}[\rho_k] = \mathbb{E}_{\mathbb{M}_1}[\rho_1] + \sum_{j=2}^n \mathbb{E}_{\mathbb{M}_1^{j-1}}[\mathbb{E}_{\mathbb{M}_{j|j-1}}[\rho_j]] \\ &\leq \sqrt{2\gamma_1 D(\mathbb{Q}_1 \parallel \mathbb{P}_1)} + \sum_{j=2}^n \sqrt{2\gamma_j \mathbb{E}_{\mathbb{M}_1^{j-1}} \left[D(\mathbb{Q}_{j|j-1}(\cdot | \mathbf{X}_1^{j-1}) \parallel \mathbb{P}_j) \right]} \\ &\leq \sqrt{2 \left(\sum_{k=1}^n \gamma_k \right)} \sqrt{D(\mathbb{Q}_1 \parallel \mathbb{P}_1) + \sum_{j=2}^n \mathbb{E}_{\mathbb{M}_1^{j-1}} \left[D(\mathbb{Q}_{j|j-1}(\cdot | \mathbf{X}_1^{j-1}) \parallel \mathbb{P}_j) \right]} \\ &= \sqrt{2 \left(\sum_{k=1}^n \gamma_k \right) D(\mathbb{Q} \parallel \mathbb{P})}. \end{aligned}$$

□

1.45 示例 (有界差不等式). 下面我们利用张量化的传输成本不等式重新得到推论 1.24. 设函数 $f: \mathcal{X}^n \rightarrow \mathbb{R}$ 满足参数为 (L_1, \dots, L_n) 的有界差不等式, 利用三角不等式可以得到, f 关于重尺度化的 Hamming 度量 $\rho(x, y) := \sum_k L_k \mathbb{I}_{\{x_k \neq y_k\}}$ 是 Lipschitz 连续的. 于是由不等式 18, 每个边缘分布 \mathbb{P}_k 满足参数为 $\gamma_k = L_k^2/4$ 的 $L_k \mathbb{I}_{\{x_k \neq y_k\}}$ -传输成本不等式. 结合定理 1.44 和定理 1.41, Lipschitz 函数 f 是 $\sqrt{\sum_k L_k^2}/2$ -次高斯的.

1.46 定理. 若度量测度空间 $(\mathcal{X}, \rho, \mathbb{P})$ 中的概率测度满足 ρ -传输成本不等式(17), 那

么它的集中度满足

$$\alpha_{\mathbb{P}}(\epsilon) \leq 2 \exp\left(-\frac{\epsilon^2}{2\gamma}\right). \quad (20)$$

證明. 考虑任意满足 $\mathbb{P}[A] \geq \frac{1}{2}$ 的集合 A 和 $\epsilon > 0$, 只需证明 $B := \bar{A}^\epsilon$ 的测度总是小于不等式(20)的右侧. 若 $\mathbb{P}[B] = 0$, 则不等式显然成立, 下面我们总假设 $\mathbb{P}[B] > 0$.

考虑 $\mathbb{P}_A, \mathbb{P}_B$ 为在 A 和 B 上的条件分布, \mathbb{M} 为它们的任意耦合, 于是

$$\begin{aligned} \int_{\mathcal{X} \times \mathcal{X}} \rho(x, x') \, d\mathbb{M}(x, x') &= \int_{A \times B} \rho(x, x') \, d\mathbb{M}(x, x') \\ &\geq \rho(A, B) \int_{A \times B} d\mathbb{M} = \rho(A, B) \geq \epsilon. \end{aligned}$$

对所有可能的耦合取下确界, 可得 $W_\rho(A, B) \geq \epsilon$. 再由三角不等式和 ρ -传输成本不等式, 我们有 (这里根号下似乎少了个 2)

$$\begin{aligned} \epsilon &\leq W_\rho(\mathbb{P}, \mathbb{P}_A) + W_\rho(\mathbb{P}, \mathbb{P}_B) \leq \sqrt{\gamma D(\mathbb{P}_A \| \mathbb{P})} + \sqrt{\gamma D(\mathbb{P}_B \| \mathbb{P})} \\ &\leq \sqrt{2\gamma} [D(\mathbb{P}_A \| \mathbb{P}) + D(\mathbb{P}_B \| \mathbb{P})]^{1/2}. \end{aligned}$$

另一方面, \mathbb{P}_A 的密度为 $p_A(x) = \frac{\mathbb{P}(x) \mathbb{I}_A(x)}{\mathbb{P}[A]}$, 于是 $D(\mathbb{P}_A \| \mathbb{P}) = -\log \mathbb{P}[A]$, $D(\mathbb{P}_B \| \mathbb{P}) = -\log \mathbb{P}[B]$, 从而有 $\epsilon^2 \leq -2\gamma \log(\mathbb{P}[A]\mathbb{P}[B])$, 等价地

$$\mathbb{P}[B] \leq (\mathbb{P}[A])^{-1} \exp\left(-\frac{\epsilon^2}{2\gamma}\right) \leq 2 \exp\left(-\frac{\epsilon^2}{2\gamma}\right).$$

□

1.5.4 非對稱耦合成本

定义

$$\mathcal{C}(\mathbb{Q}, \mathbb{P}) = \sqrt{\int \left(1 - \frac{d\mathbb{Q}}{d\mathbb{P}}\right)_+^2 d\mathbb{P}}$$

2 一致大數定律

2. Talagrand 不等式 3. 假设类的引入 (0-下水平集...

设 $\{X_i\}_{i=1}^n$ 是来自分布函数 F 的独立同分布样本, F 经典的无偏估计是经验分布函数

$$\hat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{(-\infty, t]}(X_i).$$

经典的 **Glivenko-Cantelli** 定理告诉我们, 经验分布函数 \hat{F}_n 是 F 在一致范数下的强相合估计, 即 \hat{F}_n 和 F 之间的 **Kolmogorov** 距离几乎处处收敛到 0:

$$\|\hat{F}_n - F\|_\infty := \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \xrightarrow{a.s.} 0$$

2.1 經驗過程

设 $\{X_i\}_{i=1}^n$ 是来自分布 \mathbb{P} 的 n 个独立同分布样本, 经验分布为 $\mathbb{P}_n(A) := \frac{1}{n} \sum_{i=1}^n \mathbb{I}_A(X_i)$, 其中集合 $A \subseteq \mathcal{X}$. 考虑区域 \mathcal{X} 上的 \mathbb{P} -可积实值函数类 \mathcal{F} , 函数 $f \in \mathcal{F}$ 关于初始测度 \mathbb{P} 和经验测度 \mathbb{P}_n 的积分分别为

$$\mathbb{P}(f) = \int f d\mathbb{P} = \mathbb{E}_{X \sim \mathbb{P}}[f(X)], \quad \mathbb{P}_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i).$$

于是经验过程 $X_f = \mathbb{P}(f) - \mathbb{P}_n(f)$, $f \in \mathcal{F}$ 衡量了经验期望和总体期望的偏差. 经验过程理论的研究对象是 \mathbb{P} 在函数类 \mathcal{F} 上被 \mathbb{P}_n 一致逼近的性质, 更为具体的, 研究随机量

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} |\mathbb{P}_n(f) - \mathbb{P}(f)| = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}[f(X)] \right|$$

的测度集中现象和随机过程 $\{(\mathbb{P} - \mathbb{P}_n)(f) : f \in \mathcal{F}\}$ 的概率极限理论.

第一个问题涉及 $|\mathbb{P}_n - \mathbb{P}|_F$ 在 F 一致有界的情况下围绕其均值的集中性. 问题的关键在于: 变量 $|\mathbb{P}_n - \mathbb{P}|_F$ 在其均值附近有多集中? 是否可以得到关于差值 $|\mathbb{P}_n - \mathbb{P}|_F - \mathbb{E}|\mathbb{P}_n - \mathbb{P}|_F$ 的指数不等式, 且达到与经典不等式对 $\sum_{i=1}^n \xi_i$ (其中 ξ_i 是中心化且有界的) 同样的精度? 还是需要为我们同时考虑无穷多个独立随机变量之和而付出额外的代价?

对此问题的令人惊讶的答案是: 在适当定义参数 (规模和方差) 的情况下, 经典的指数不等式在经验过程上依然成立. 这是经验过程理论中最重要、最强大的结果之一, 称为 Talagrand 不等式. 在本节稍后部分, 我们将回顾实随机变量的经典指数不等式, 为后续的经验过程部分提供背景.

可测函数类的上确界未必可测, 恐有不测之忧

- 如果 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \xrightarrow{\mathbb{P}} 0$, 则称函数类 \mathcal{F} 为分布 \mathbb{P} 上的一个 **Glivenko-Cantelli 类**, 或者 \mathcal{F} 满足 **Glivenko-Cantelli 律**;
- 如果 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \xrightarrow{a.s.} 0$, 则称函数类 \mathcal{F} 满足 **强 Glivenko-Cantelli 律**.

经典的 Glivenko-Cantelli 定理实际上是在示性函数类 $\mathcal{F} = \{\mathbb{I}_{(-\infty, t]} : t \in \mathbb{R}\}$ 上的强一致定律. 对更为广义的函数类的研究则始于 Vapnik-Červonenkis(1971) 和 Dudley(1978) 的工作, 这在统计中有着十分重要的作用.

量 $\mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}]$ 实际上是经验测度 \mathbb{P}_n 和 \mathbb{P} 之间的 Wasserstein 距离 $W(\mathbb{P}, \mathbb{P}_n)$.

2.1 示例 (M 估计 [GN15] P110). 若分布 \mathbb{P}_{θ^*} 中的 $\theta^* \in \Theta$ 未知, 其中空间 Θ 可能是 \mathbb{R}^d , 对应参数估计问题; 或者是某个函数类 \mathcal{G} , 对应非参数问题.

θ^* 的估计总是基于最小化损失函数 $\mathcal{L}_{\theta}(X)$: 最优的 θ 应当使得总体风险 $R(\theta, \theta^*) := \mathbb{E}_{\mathbb{P}_{\theta^*}}[\mathcal{L}_{\theta}(X)]$ 达到最小. 然而在实践中, 我们通常无法获得总体数据, 只能根据有限个样本 $\{X_i\}_{i=1}^n$, 在 Θ 的某个子集 Θ_0 上最小化经验风险得到估计

$$\hat{\theta} = \arg \min_{\theta \in \Theta_0} \hat{R}_n(\theta, \theta^*) = \arg \min_{\theta \in \Theta_0} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\theta}(X_i).$$

我们希望经验风险和总体风险足够接近, 即控制过度风险 $\mathbb{E}(\hat{\theta}, \theta^*) := R(\hat{\theta}, \theta^*) - \inf_{\theta \in \Theta_0} R(\theta, \theta^*)$.

为了方便起见, 我们假设存在某个 $\theta_0 \in \Theta_0$ 满足 $R(\theta_0, \theta^*) = \inf_{\theta \in \Theta_0} R(\theta, \theta^*)$. 于是过度风险可以做如下估计

$$\mathbb{E}(\hat{\theta}, \theta^*) = \underbrace{\left[R(\hat{\theta}, \theta^*) - \hat{R}_n(\hat{\theta}, \theta^*) \right]}_{T_1} + \underbrace{\left[\hat{R}_n(\hat{\theta}, \theta^*) - \hat{R}_n(\theta_0, \theta^*) \right]}_{T_2} + \underbrace{\left[\hat{R}_n(\theta_0, \theta^*) - R(\theta_0, \theta^*) \right]}_{T_3}.$$

其中 $|T_1| = \left| \frac{1}{n} \sum_i \mathcal{L}_\theta(X_i) - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\mathcal{L}_\theta(X)] \right| \leq \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{L}_{\Theta_0}}$, 这需要需要考虑损失函数类 $\mathcal{L}_{\Theta_0} := \{\mathcal{L}_\theta(\cdot) : \theta \in \Theta_0\}$ 的一致大数定律; 而 $\hat{\theta}$ 最小化了经验风险, 可以看到 $T_2 \leq 0$; T_3 对应的则是控制随机变量 $\frac{1}{n} \sum_i \mathcal{L}_{\theta_0}(X_i)$ 及其期望 $\mathbb{E}_{\mathbb{P}_{\theta_0}}[\mathcal{L}_{\theta_0}(X)]$ 之间的偏差, 这里 θ_0 是一个未知但非随机的值, 因此可以用测度集中的方法来得到. 当然 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{L}_{\Theta_0}}$ 也是 T_3 的上界, 于是过度风险至多是 $2\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{L}_{\Theta_0}}$.

2.2 經驗過程的尾部概率界

设 $\mathbf{X}_1^n = (X_1, \dots, X_n)$ 来自乘积分布 $\mathbb{P} = \otimes_{i=1}^n \mathbb{P}_i$, 其中 \mathbb{P}_i 的支撑集 $\mathcal{X}_i \subseteq \mathcal{X}$. 对于定义域为 \mathcal{X} 函数类 \mathcal{F} , 考虑随机变量

$$Z := \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n f(X_i) \right\}.$$

对于 $\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_i f(X_i) \right|$, 只需在函数类 $\tilde{\mathcal{F}} := \mathcal{F} \cup (-\mathcal{F})$ 上考虑上确界即可:

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) \right| = \sup_{f \in \tilde{\mathcal{F}}} \left\{ \max \left\{ \frac{1}{n} \sum_{i=1}^n f(X_i), -\frac{1}{n} \sum_{i=1}^n f(X_i) \right\} \right\} = \sup_{f \in \tilde{\mathcal{F}}} \left\{ \frac{1}{n} \sum_{i=1}^n f(X_i) \right\}.$$

我们把 Hoeffding

2.2 定理 (泛函 Hoeffding 不等式). 若对每个 $f \in \mathcal{F}$, 都有 $f(\mathcal{X}_i) \subseteq [a_{i,f}, b_{i,f}]$, $i = 1, \dots, n$, 那么对任意 $\delta \geq 0$, 成立

$$\mathbb{P}[Z \geq \mathbb{E}[Z] + \delta] \leq \exp \left(-\frac{n\delta^2}{4L^2} \right), \quad (21)$$

其中 $L^2 = \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_i (b_{i,f} - a_{i,f})^2 \right\}$.

證明. 为了简便, 我们使用非重尺度化的 $Z = \sup_{f \in \mathcal{F}} \left\{ \sum_i f(X_i) \right\}$, 它是 \mathbf{X}_1^n 的泛函.

定义 $Z_j: x_j \mapsto Z(X_1, \dots, X_{j-1}, x_j, X_{j+1}, \dots, X_n)$

对于 $\lambda > 0$, 由引理 1.34、1.30,

对 $f \in \mathcal{F}$, 定义 $\mathcal{A}(f) := \{(x_1, \dots, x_n) : Z = \sum_i f(x_i)\}$ □

2.3 定理 (經驗過程的 *Talagrand* 集中度). 若可数函数类 \mathcal{F} 是 b -一致有界的, 那么对任意 $\delta > 0$, 成立

$$\mathbb{P}[Z \geq \mathbb{E}[Z] + \delta] \leq 2 \exp \left(-\frac{n\delta^2}{8e\mathbb{E}\Sigma^2 + 4b\delta} \right),$$

其中 $\Sigma^2 = \sup_{f \in \mathcal{F}} \frac{1}{n} f^2(X_i)$.

2.3 函数类的 *Rademacher* 复杂度

Peter L. Bartlett 与 Shahar Mendelson (此人是 Empirical Process 的专家) 提出了用 Rademacher / Gaussian Complexity 来研究 Risk Bounds 的方法

一致大数定律的一个度量是函数类 \mathcal{F} 的 **Rademacher 复杂度**. 设 \mathcal{F} 为区域 \mathcal{X} 上的函数类, 对于点集 $\mathbf{x}_1^n \in \mathcal{X}^n$, 记 $\mathcal{F}(\mathbf{x}_1^n) := \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\}$. \mathcal{F} 关于点集 \mathbf{x}_1^n 的 Rademacher 复杂度为

$$\mathcal{R}_n(\mathcal{F}) := \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \langle f(\mathbf{x}_1^n), \epsilon \rangle \right] = \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| \right].$$

进一步地, 函数类 \mathcal{F} 关于经验分布 \mathbb{P}_n 的 Rademacher 复杂度为

$$\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) := \mathbb{E}_{\mathbf{X}, \epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{k=1}^n \epsilon_k f(X_k) \right| \right].$$

这可以看作随机向量 $(f(X_1), \dots, f(X_n))_{f \in \mathcal{F}}$ 和噪声向量 ϵ 之间最大相关关系的平均值.

对于一致有界函数类 \mathcal{F} , 我们将看到“ $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \approx \mathcal{R}_{\mathbb{P}_n}(\mathcal{F})$ ”, 而 Rademacher 复杂度的界是更为容易得到的.

2.4 引理. 若函数类 \mathcal{F} 是 b -一致有界的, 那么随机变量 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 是 $\frac{b}{\sqrt{n}}$ -次高斯的.

證明. 引入函数的中心化记号 $\bar{f}(x) := f(x) - \mathbb{E}[f(X)]$, 则 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 可以简记为 $\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(X_i) \right|$. 考虑函数 $G(x_1, \dots, x_n) = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right|$, 它与各坐标的顺序置换无关. 于是只需对第一个坐标分量进行扰动, 就可以说明它满足参数为 $(\frac{2b}{n}, \dots, \frac{2b}{n})$ 的有界差不等式: 设向量 $\mathbf{X} = (x_1, \dots, x_n)$, $\mathbf{Y} = (y_1, \dots, y_n)$ 满足 $x_i = y_i, i > 1$, 说明 $|G(\mathbf{X}) - G(\mathbf{Y})| < \frac{2b}{n}$ 即可. 对任意 $f \in \mathcal{F}$, 由于 $\|f\|_{\infty} \leq b$,

$$\left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \sup_{h \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \bar{h}(y_i) \right| \leq \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(x_i) \right| - \left| \frac{1}{n} \sum_{i=1}^n \bar{f}(y_i) \right| \leq \frac{1}{n} |\bar{f}(x_1) - \bar{f}(y_1)| \leq \frac{2b}{n}.$$

结合推论 1.24, 可以看到 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 是 $\frac{b}{\sqrt{n}}$ -次高斯的. \square

2.5 定理. 设函数类 \mathcal{F} 是 b -一致有界的, 对任意 $n \in \mathbb{Z}_+$, $\delta \geq 0$, 我们有

$$\mathbb{P}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) + \delta] \geq 1 - \exp\left(-\frac{n\delta^2}{2b^2}\right).$$

于是当函数类 \mathcal{F} 满足 $\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) = o(1)$ 时, $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 以指数速度几乎确定收敛到 0, 即 \mathcal{F} 为 \mathbb{P} 上的 *Glivenko-Cantelli* 类.

在证明之前, 我们给出一个函数类上确界期望的不等式, 它和 Fatou 不等式或者 Jensen 不等式类似: 对于可积函数类 \mathcal{G} , 有 $\mathbb{E}[g(X)] \leq \mathbb{E}[\sup_{g \in \mathcal{G}} |g(X)|]$, 于是再对左侧取上确界可以得到

$$\sup_{g \in \mathcal{G}} \mathbb{E}[g(X)] \leq \mathbb{E}\left[\sup_{g \in \mathcal{G}} |g(X)|\right]. \quad (22)$$

进一步地, 对于凸的非减函数 Φ , 结合 Jensen 不等式, 我们有

$$\sup_{g \in \mathcal{G}} \Phi(\mathbb{E}[|g(X)|]) \leq \Phi\left(\mathbb{E}\left[\sup_{g \in \mathcal{G}} |g(X)|\right]\right) \leq \mathbb{E}\left[\Phi\left(\sup_{g \in \mathcal{G}} |g(X)|\right)\right] \quad (23)$$

证明. 我们首先证明 $2\mathcal{R}_{\mathbb{P}_n}(\mathcal{F})$ 是 $\mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}]$ 的上界, 这可以使用对称化技巧来得到: 设 $\mathbf{Y} = (Y_1, \dots, Y_n)$ 与 \mathbf{X} 独立同分布, $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ 为独立的 Rademacher 向量, 于是 $\epsilon_i(f(X_i) - f(Y_i))$ 和 $f(X_i) - f(Y_i)$ 有相同的分布:

$$\begin{aligned} \mathbb{P}[f(X_i) - f(Y_i) \leq t] &= \frac{1}{2}\mathbb{P}[f(X_i) - f(Y_i) \leq t] + \frac{1}{2}\mathbb{P}[f(Y_i) - f(X_i) \leq t] \\ &= \mathbb{P}[\epsilon_i = 1] \cdot \mathbb{P}[\epsilon_i(f(X_i) - f(Y_i)) \leq t | \epsilon_i = 1] + \mathbb{P}[\epsilon_i = -1] \cdot \mathbb{P}[\epsilon_i(f(X_i) - f(Y_i)) \leq t | \epsilon_i = -1] \\ &= \mathbb{P}[\epsilon_i(f(X_i) - f(Y_i)) \leq t]. \end{aligned}$$

结合不等式(22)、三角不等式

$$\begin{aligned} \mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] &= \mathbb{E}_{\mathbf{X}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}_{Y_1}[f(Y_1)] \right| \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}_{Y_i}[f(Y_i)]) \right| \right] = \mathbb{E}_{\mathbf{X}} \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathbf{Y}} \left[\frac{1}{n} \sum_{i=1}^n (f(X_i) - f(Y_i)) \right] \right| \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(Y_i)) \right| \right] = \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - f(Y_i)) \right| \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] + \mathbb{E}_{\mathbf{Y}, \epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Y_i) \right| \right] = 2\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}). \end{aligned}$$

于是根据 $\frac{b}{\sqrt{n}}$ -次高斯随机变量 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 的上偏差不等式, 对任意 $\delta \geq 0$, 我们有

$$\mathbb{P}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) + \delta] \geq \mathbb{P}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq \mathbb{E}[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] + \delta] \geq 1 - \exp\left(-\frac{nt^2}{2b^2}\right).$$

□

对称化技巧实际上考虑了随机变量 $\|\mathbb{S}_n\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_i \epsilon_i f(X_i) \right|$ ——它的期望就是 Rademacher 复杂度, 我们有更强的结论。

2.6 命题. 对于任意非减的凸函数 $\Phi: \mathbb{R} \rightarrow \mathbb{R}$, 我们有

$$\mathbb{E}_{\mathbf{X}, \epsilon} \left[\Phi \left(\frac{1}{2} \|\mathbb{S}\|_{\mathcal{F}} \right) \right] \leq \mathbb{E}_{\mathbf{X}} [\Phi (\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}})] \leq \mathbb{E}_{\mathbf{X}, \epsilon} [\Phi (2\|\mathbb{S}\|_{\mathcal{F}})]$$

2.7 注. 特别地, 取 $\Phi(t) = t$ 可以得到

$$\frac{1}{2} \mathbb{E}_{\mathbf{X}, \epsilon} [\|\mathbb{S}\|_{\mathcal{F}}] \leq \mathbb{E}_{\mathbf{X}} [\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] \leq 2 \mathbb{E}_{\mathbf{X}, \epsilon} [\|\mathbb{S}\|_{\mathcal{F}}] \quad (24)$$

证明. 右侧不等式可以看作是上一定理的证明的简单推广: 结合不等式(23)、三角不等式, 利用 Φ 的凸性, 我们有

$$\begin{aligned} \mathbb{E}_{\mathbf{X}} [\Phi (\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}})] &= \mathbb{E}_{\mathbf{X}} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathbf{Y}} \left[\frac{1}{n} \sum_{i=1}^n (f(X_i) - f(Y_i)) \right] \right| \right) \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(Y_i)) \right| \right) \right] = \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \epsilon} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - f(Y_i)) \right| \right) \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \epsilon} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| + \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Y_i) \right| \right) \right] \\ &\leq \frac{1}{2} \mathbb{E}_{\mathbf{X}, \epsilon} \left[\Phi \left(2 \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right) \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}, \epsilon} \left[\Phi \left(2 \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(Y_i) \right| \right) \right] \\ &= \mathbb{E}_{\mathbf{X}, \epsilon} \left[\Phi \left(2 \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right) \right] = \mathbb{E}_{\mathbf{X}, \epsilon} [\Phi (2\|\mathbb{S}\|_{\mathcal{F}})]. \end{aligned}$$

下面我们证明左侧不等式, 由不等式(23)、三角不等式和 Φ 的非减性、 Φ 的凸性, 我们有

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \epsilon} \left[\Phi \left(\frac{1}{2} \|\mathbb{S}\|_{\mathcal{F}} \right) \right] &= \mathbb{E}_{\mathbf{X}, \epsilon} \left[\Phi \left(\frac{1}{2} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - \mathbb{E}_{Y_i}[f(Y_i)]) \right| \right) \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \epsilon} \left[\Phi \left(\frac{1}{2} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(X_i) - f(Y_i)) \right| \right) \right] = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\Phi \left(\frac{1}{2} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(Y_i)) \right| \right) \right] \\ &\leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[\Phi \left(\frac{1}{2} \left\{ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}_{\mathbb{P}}[f]) \right| + \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(Y_i) - \mathbb{E}_{\mathbb{P}}[f]) \right| \right\} \right) \right] \\ &\leq \frac{1}{2} \mathbb{E}_{\mathbf{X}} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}_{\mathbb{P}}[f]) \right| \right) \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}} \left[\Phi \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(Y_i) - \mathbb{E}_{\mathbb{P}}[f]) \right| \right) \right] \\ &= \mathbb{E}_{\mathbf{X}} [\Phi (\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}})]. \end{aligned}$$

□

2.8 命题. 设函数类 \mathcal{F} 是 b -一致有界的, 对任意 $n \in \mathbb{Z}_+$, $\delta \geq 0$, 我们有

$$\mathbb{P} \left[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \geq \frac{1}{2} \mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) - \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbb{P}}[f]|}{2\sqrt{n}} - \delta \right] \geq 1 - \exp \left(-\frac{n\delta^2}{2b^2} \right).$$

于是当函数类 \mathcal{F} 的 Rademacher 复杂度 $\mathcal{R}_{\mathbb{P}_n}(\mathcal{F})$ 存在远离 0 的下界时, $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 不可能以概率收敛于零, 即 \mathcal{F} 不可能为 \mathbb{P} 上的 Glivenko-Cantelli 类.

证明. 根据三角不等式容易得到估计

$$\|\mathbb{S}_n\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) - \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{E}_{\mathbb{P}}[f] \right| \geq \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| - \sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbb{P}}[f]| \cdot \frac{|\sum_{i=1}^n \epsilon_i|}{n}.$$

再由 Cauchy-Schwarz 不等式可得 $\mathbb{E} \left[\left| \sum_{i=1}^n \epsilon_i \right| \right] \leq \sqrt{\mathbb{E} \left[\left(\sum_{i=1}^n \epsilon_i \right)^2 \right]} = \sqrt{\mathbb{E} \left[\sum_{i=1}^n \epsilon_i^2 \right]} = \sqrt{n}$. 于是由不等式(24), $\frac{b}{\sqrt{n}}$ -次高斯随机变量 $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ 的下偏差不等式, 可见

$$\begin{aligned} 1 - \exp \left(-\frac{n\delta^2}{2b^2} \right) &\leq \mathbb{P} [\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \geq \mathbb{E} [\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] - \delta] \\ &\leq \mathbb{P} \left[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \geq \frac{1}{2} \mathbb{E}_{\mathbf{X}, \epsilon} [\|\mathbb{S}\|_{\mathcal{F}}] - \delta \right] \leq \mathbb{P} \left[\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \geq \frac{1}{2} \mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) - \frac{\sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbb{P}}[f]|}{2\sqrt{n}} - \delta \right] \end{aligned}$$

□

于是 Rademacher 复杂度是 1 的高阶无穷小为满足 Glivenko-Cantelli 性质提供了一个充分必要条件. 下面两节我们寻求控制 Rademacher 复杂度的方法.

2.4 多项式识别函数类

2.9 定义 (多项式识别). 称区域 \mathcal{X} 上的函数类 \mathcal{F} 有阶为 $\nu \geq 1$ 的多项式识别, 如果对任意正整数 n 和点集 $\mathbf{x}_1^n = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$, 集合

$$\mathcal{F}(\mathbf{x}_1^n) := \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\} \subseteq \mathbb{R}^n$$

的基数 $\text{card}(\mathcal{F}(\mathbf{x}_1^n)) \leq (n+1)^\nu$.

很多时候函数类的基数 $\text{card}(\mathcal{F}) = \infty$, 但是

2.10 引理. 若区域 \mathcal{X} 上的函数类 \mathcal{F} 有阶为 $\nu \geq 1$ 的多项式识别, 那么对任意正整数 n , 函数类 \mathcal{F} 关于点集 $\mathbf{x}_1^n = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$ 的 Rademacher 复杂度

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right| \right] \leq 2D(\mathbf{x}_1^n) \sqrt{\frac{\nu \log(n+1)}{n}},$$

其中 $D(\mathbf{x}_1^n) = \sup_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum_i f(x_i)^2}$ 为向量集合 $\mathcal{F}(\mathbf{x}_1^n)/\sqrt{n}$ 的 ℓ^2 半径。

證明. 考虑集合 $A = \{\frac{1}{n}(f(x_1), \dots, f(x_n)) : f \in \mathcal{F} \cup (-\mathcal{F})\} \subseteq \mathbb{R}^n$, 由于 \mathcal{F} 有阶为 ν 的多项式识别, 我们有 $\text{card}(A) \leq 2(n+1)^\nu$. 回忆示例 1.27, 不等式左侧等价于点集 A 的 Rademacher 复杂度, 即

$$\mathcal{R}(A) = \mathbb{E}_\epsilon \left[\sup_{\mathbf{a} \in A} \langle \mathbf{a}, \boldsymbol{\epsilon} \rangle \right] = \mathbb{E}_\epsilon \left[\max_{\mathbf{a} \in A} \langle \mathbf{a}, \boldsymbol{\epsilon} \rangle \right],$$

其中 $\langle \mathbf{a}, \boldsymbol{\epsilon} \rangle$ 总是 $D(\mathbf{x}_1^n)/\sqrt{n}$ -次高斯的, 于是由定理 1.4

$$\mathbb{E}_\epsilon \left[\max_{\mathbf{a} \in A} \langle \mathbf{a}, \boldsymbol{\epsilon} \rangle \right] \leq D(\mathbf{x}_1^n)/\sqrt{n} \cdot \sqrt{2 \log(\text{card}(A))} \leq 2D(\mathbf{x}_1^n) \sqrt{\frac{\nu \log(n+1)}{n}}.$$

□

2.11 注. 若函数类 \mathcal{F} 是 b -一致有界的, 那么它关于经验分布 \mathbb{P}_n 的 Rademacher 复杂度

$$\begin{aligned} \mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) &= \mathbb{E}_{\mathbf{X}} \left[\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \right] \right] = \mathbb{E}_{\mathbf{X}} \left[\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \middle| \mathbf{X} \right] \right] \\ &\leq \sup_{\mathbf{x}} \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| \middle| \mathbf{X} = \mathbf{x} \right] \leq 2 \sup_{\mathbf{x}} \sup_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum_{i=1}^n f(x_i)^2} \cdot \sqrt{\frac{\nu \log(n+1)}{n}} \\ &\leq 2b \sqrt{\frac{\nu \log(n+1)}{n}}. \end{aligned}$$

结合定理 2.5, 我们可以看到具有多项式识别的有界函数类总是 *Glivenko-Cantelli* 的. 例如, 经典的 *Glivenko-Cantelli* 定理考虑的函数类 $\mathcal{F} = \{\mathbb{I}_{(-\infty, t]} : t \in \mathbb{R}\}$ 是 1-一致有界的, 我们可以得到 *Glivenko-Cantelli* 定理的定量版本.

2.12 推論 (*Glivenko-Cantelli* 定理-定量版本). 对任意 $\delta \geq 0$,

$$\mathbb{P} \left[\|\hat{F}_n - F\|_\infty \geq 8 \sqrt{\frac{\log(n+1)}{n}} + \delta \right] \leq \exp \left(-\frac{n\delta^2}{2} \right).$$

于是 $\|\hat{F}_n - F\|_\infty$ 以指数速度几乎确定收敛于 0.

證明. 对于任意样本 \mathbf{x}_1^n , 考虑次序样本 $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, 我们有

$$\mathcal{F}(\mathbf{x}_1^n) \subseteq \{(0, 0, \dots, 0), (1, 0, \dots, 0), (1, 1, \dots, 0), \dots, (1, 1, \dots, 1)\}.$$

于是 $\text{card}(\mathcal{F}(\mathbf{x}_1^n)) \leq n+1$, 即多项式识别阶数 $\nu = 1$, 于是 $\mathcal{R}_{\mathbb{P}_n}(\mathcal{F}) \leq 2 \sqrt{\frac{\log(n+1)}{n}}$. 结合定理 2.5 可见不等式成立. □

2.5 Vapnik-Červonenkis 維數

对于布尔值函数类, 即值域为 $\{0, 1\}$ 的函数构成的类, 我们常用 **Vapnik-Chervonenkis 维数** (简称 VC 维数) 来衡量它的复杂度.

hypothesis class/ concept class

例如集合类 \mathcal{S} 的示性函数类 $\mathbb{I}_{\mathcal{S}} := \{\mathbb{I}_S : S \in \mathcal{S}\}$, 为了记号的方便, 我们将集合类 \mathcal{S} 等价于函数类 $\mathbb{I}_{\mathcal{S}}$.

集合 Λ 被函数类打散是指无论我们对每个点如何赋予布尔值标签, 都有函数类中的一个函数将它实现. 例如对于点集 $\mathbf{x}_1^n = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$, 共有 2^n 种赋予布尔值标签的方式, 换言之, 集合 $\mathcal{F}(\mathbf{x}_1^n)$ 至多有 2^n 个元素, 于是 \mathbf{x}_1^n 被打散等价于 $\text{card}(\mathcal{F}(\mathbf{x}_1^n)) = 2^n$. 函数类的 VC 维定义为可以被其成员打散的点的最大数目.

2.13 定义 (VC 维数). 我们称集合 $\Lambda \subseteq \mathcal{X}$ 被 \mathcal{F} 打散, 如果对任意映射 $g: \Lambda \rightarrow \{0, 1\}$, 都存在某个 $f \in \mathcal{F}$ 使得 $f|_{\Lambda} = g$. 函数类 \mathcal{F} 的 VC 维数为能被 \mathcal{F} 打散的集合的最大基数: 如果

$$\text{vc}(\mathcal{F}) = \max\{\text{card}(\Lambda) : \text{card}(\mathcal{F}(\Lambda)) = 2^{\text{card}(\Lambda)}\} < \infty$$

则称函数类 \mathcal{F} 为 **VC 类**, 否则, 记 $\text{vc}(\mathcal{F}) = \infty$.

2.14 示例 (\mathbb{R} 上的区间). 推论 2.12 的证明本质上考虑了 \mathbb{R} 上左侧半区间类 $\mathcal{S}_{\text{left}} := \{(-\infty, t] : t \in \mathbb{R}\}$ 的示性函数类, 它的多项式识别的阶为 1. 对于 $x_1 < x_2$, $\mathcal{S}_{\text{left}}(x_1, x_2) = \{(0, 0), (0, 1), (1, 1)\}$, 于是 $\text{vc}(\mathcal{S}_{\text{left}}) = 1$.

进一步地, 双侧区间类 $\mathcal{S}_{\text{two}} : \{(b, a] : a, b \in \mathbb{R}, b < a\}$ 的示性函数类可以打散任意的两点集, 但是对于三个不同的点 $x_1 < x_2 < x_3$, 它不能选出集合 $\{x_1, x_3\}$: 即 $(1, 0, 1) \notin \mathcal{S}_{\text{two}}(\mathbf{x}_1^3)$, 于是 $\text{vc}(\mathcal{S}_{\text{two}}) = 2$. 此外, $\text{card}(\mathcal{S}_{\text{two}}(\mathbf{x}_1^n)) \leq (n+1)^2$, 于是 \mathcal{S}_{two} 多项式识别的阶为 2.

根据定义, 若 $\text{vc}(\mathcal{F}) < n$, 我们只能得到指数增长的结果 $\text{card}(\mathcal{F}(\mathbf{x}_1^n)) \leq 2^n - 1$. 但是在上述示例中, 我们看到 VC 维数和多项式识别的阶数似乎存在着一定的联系. 事实上, 利用组合的方法, 我们可以得到如下结论.

2.15 引理 (Sauer-Shelah). 设 \mathcal{S} 是 VC 类, 对任意点集 \mathbf{x}_1^n , 其中 $n > \text{vc}(\mathcal{S})$, 我们有

$$\text{card}(\mathcal{S}(\mathbf{x}_1^n)) \leq \sum_{i=1}^{\text{vc}(\mathcal{S})} \binom{n}{i} \leq (n+1)^{\text{vc}(\mathcal{S})}.$$

证明. □

VC 类在有限次集合运算下保持不变, 这被称为 VC 稳定性.

2.16 命题 (VC 稳定性). 若 \mathcal{S} 和 \mathcal{T} 是 VC 类, 那么下述集合类也是 VC 类:

- (1) $\mathcal{S}^c := \{S^c : S \in \mathcal{S}\};$
- (2) $\mathcal{S} \cap \mathcal{T} := \{S \cap T : S \in \mathcal{S}, T \in \mathcal{T}\};$
- (3) $\mathcal{S} \cup \mathcal{T} := \{S \cup T : S \in \mathcal{S}, T \in \mathcal{T}\}.$

证明. (1) 若点集 \mathbf{x}_1^n 可以被 \mathcal{S} 打散, 对任意 $S \in \mathcal{S}$, 由于 S^c 会给 \mathbf{x}_1^n 和 S 完全相反的布尔值标签, 于是 $\text{card}(\mathcal{S}^c(\mathbf{x}_1^n)) = 2^n = \text{card}(\mathcal{S}(\mathbf{x}_1^n))$, 即 $\text{vc}(\mathcal{S}^c) = \text{vc}(\mathcal{S})$.

(2) 注意到对任意的 $S \in \mathcal{S}, T \in \mathcal{T}$, 我们有 $\mathbb{I}_{S \cap T} = \mathbb{I}_S \cdot \mathbb{I}_T$, 结合引理 2.15,

$$\text{card}(\mathcal{S} \cap \mathcal{T}(\mathbf{x}_1^n)) \leq \text{card}(\mathcal{S}(\mathbf{x}_1^n)) \cdot \text{card}(\mathcal{T}(\mathbf{x}_1^n)) \leq (n+1)^{\text{vc}(\mathcal{S})+\text{vc}(\mathcal{T})}.$$

(3) 由 $S \cup T = (S^c \cap T^c)^c$ 可见成立. □

实值函数类 \mathcal{F} 可以通过取 0-下水平集来定义相伴集合类 $\mathcal{S}(\mathcal{F}) := \{S_f : f \in \mathcal{F}\}$, 其中 $S_f := \{x \in \mathcal{X} : f(x) \leq 0\}$ 称为函数 f 的 0-下水平集. 很多重要的集合类, 例如半平面、椭球体, 都可以用这种方式来表达.

2.17 命题. 设函数类 \mathcal{G} 为 \mathbb{R}^d 上实值函数的线性空间, 其中 $\dim(\mathcal{G}) < \infty$, 那么 $\mathcal{S}(\mathcal{G})$ 的 VC 维数至多为 $\dim(\mathcal{G})$.

证明. 采用反证法, 假设存在点集 $\mathbf{x}_1^n = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$ 可以被 $\mathcal{S}(\mathcal{G})$ 打散, 其中 $n = \dim(\mathcal{G}) + 1$. 由于 \mathcal{G} 构成了线性空间, 集合 $\mathcal{G}(\mathbf{x}_1^n) = \{(g(x_1), \dots, g(x_n)) : g \in \mathcal{G}\}$ 构成了 \mathbb{R}^n 的子空间, 并且它的维数至多为 $\dim(\mathcal{G}) = n - 1 < n$. 因此, 存在某个非零向量 $\gamma \in \mathbb{R}^n$ 满足 $\langle \gamma, g(\mathbf{x}_1^n) \rangle = \sum_i \gamma_i g(x_i) = 0, \forall g \in \mathcal{G}$. 不失一般性地, 假定 γ 至少存在一个正的分量, 于是

$$\sum_{i: \gamma_i \leq 0} (-\gamma_i) g(x_i) = \sum_{i: \gamma_i > 0} \gamma_i g(x_i), \quad \forall g \in \mathcal{G}.$$

由于 $\mathcal{S}(\mathcal{G})$ 可以将 \mathbf{x}_1^n 打散, 于是存在 $g \in \mathcal{G}$ 使得 $S_g \cap \mathbf{x}_1^n = \{x_i : \gamma_i \leq 0\}$, 此时等式右侧严格正而左侧非正, 推出矛盾. □

2.18 示例 (\mathbb{R}^d 中的半平面). \mathbb{R}^n 中的半平面 $H_{\mathbf{a},b} := \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{a}, \mathbf{x} \rangle \leq b\}$ 可以看作线性函数 $l_{\mathbf{a},b}(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle - b$ 的 0-下水平集. 全体线性函数构成类 $\mathcal{L}^d := \{l_{\mathbf{a},b} : (\mathbf{a}, b) \in \mathbb{R}^d \times \mathbb{R}\}$, 由线性代数的知识, 不难看出它是 $d+1$ 维的线性空间, 于是 $\mathcal{S}(\mathcal{L}^d)$ 的 VC 维数不会大于 $d+1$.

2.19 示例 (\mathbb{R}^d 中的球). 考虑 \mathbb{R}^d 中全体的欧氏球 $\mathcal{S}_{\text{euc}}^d := \{S_{\mathbf{a},b} : (\mathbf{a}, b) \in \mathbb{R}^d \times \mathbb{R}_+\}$, 其中 $S_{\mathbf{a},b} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{a}\|_2 \leq b\}$ 为以 \mathbf{a} 为圆心、 b 为半径的欧氏球, 可以看作是 $s_{\mathbf{a},b} := \|\mathbf{x} - \mathbf{a}\|_2 - b$ 的 0-下水平集. 考虑特征映射 $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d+2}, \mathbf{x} \mapsto (1, x_1, \dots, x_n, \|\mathbf{x}\|_2^2)$, 函

数 $g_c(\boldsymbol{x}) := \langle \boldsymbol{c}, \boldsymbol{x} \rangle$ 构成的函数类 $\mathcal{G} = \{g_c: \boldsymbol{c} \in \mathbb{R}^{d+2}\}$ 是一个 $d+2$ 维的线性空间, 并且 $s_{a,b}$ 就属于此类. 于是 $\text{vc}(\mathcal{S}_{euc}^d) \leq d+2$.¹

8.3.5 Empirical processes via VC dimension (HDP)

2.20 定理 (VC 类的一致大数定律). 设 \mathcal{F} 为 VC 类,

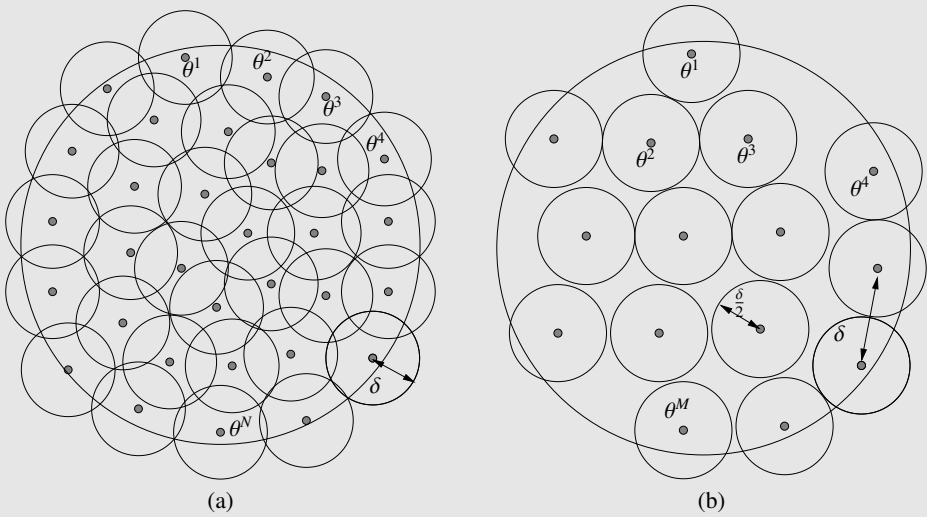
$$\mathbb{E} [\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}] \leq C \sqrt{\frac{\text{vc}(\mathcal{F})}{n}}.$$

3 度量熵

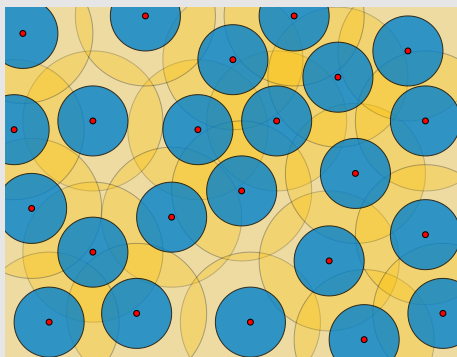
对于次高斯

在度量空间理论中, δ -网, δ -覆盖, δ -填充是几个紧密相关的定义, 可以用于描述点集的良好分布.

集合 $\mathbb{T} \subseteq \mathcal{X}$ 关于度量 ρ 的一个 δ -覆盖是指
不要相聚太远、也不要过于拥挤——稀疏有致的点



¹一个更加细致的分析可以得到它的 VC 维数实际上是 $d+1$.



4 隨機矩陣

A 預備知識

A.1 Landau 記號

A.2 概率論

A.1 定理 (Borel-Cantelli 引理). 设 $(A_n)_{n \in \mathbb{N}}$ 为事件序列.

1. 若 $\sum_n \mathbb{P}(A_n) < \infty$, 则 $\mathbb{P}(A_n \text{ i.o.}) = 0$.
2. 若 A_n 相互独立且 $\sum_n \mathbb{P}(A_n) = \infty$, 则 $\mathbb{P}(A_n \text{ i.o.}) = 1$.

證明. 1. 由 \mathbb{P} 上半连续、 σ -可加性和 Cauchy 收敛准则:

$$\mathbb{P} \left(\lim_{n \rightarrow \infty} \bigcup_{m \geq n} A_m \right) = \lim_{n \rightarrow \infty} \mathbb{P} \left(\bigcup_{m \geq n} A_m \right) \leq \lim_{n \rightarrow \infty} \sum_{m \geq n} \mathbb{P}(A_m) = 0.$$

2. 由 De Morgan 律和 \mathbb{P} 下半连续

$$\mathbb{P} \left(\left(\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} A_n \right)^c \right) = \mathbb{P} \left(\bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} A_n^c \right) = \lim_{m \rightarrow \infty} \mathbb{P} \left(\bigcap_{n=m}^{\infty} A_n^c \right).$$

而对任意 $m \in \mathbb{N}$, 由不等式 $\log(1-x) \leq -x$ 在 $x \in [0, 1]$ 成立, 总有

$$\begin{aligned} \mathbb{P} \left(\bigcap_{n=m}^{\infty} A_n^c \right) &= \lim_{N \rightarrow \infty} \mathbb{P} \left(\bigcap_{n=m}^N A_n^c \right) = \prod_{n=m}^{\infty} (1 - \mathbb{P}(A_n)) \\ &= \exp \left(\sum_{n=m}^{\infty} \log(1 - \mathbb{P}(A_n)) \right) \leq \exp \left(- \sum_{n=m}^{\infty} \mathbb{P}(A_n) \right) = 0. \end{aligned}$$

□

A.2.1 積分變換、Stieltjes 積分

A.2 定理 (積分變換). 设 $f: (\mathcal{X}, \mathcal{A}, \mu) \rightarrow (E, \mathcal{E})$ 为可测映射, g 为 (E, \mathcal{E}) 上的可测函数, 则

$$\int_{f^{-1}(B)} g \circ f \, d\mu = \int_B g \, df_* \mu.$$

證明. 只需证明对可测性成立即可. 对于 $g = \mathbb{I}_F$, $F \in \mathcal{E}$, 有

$$\begin{aligned} \int_B \mathbb{I}_F \, df_* \mu &= f_* \mu(B \cap F) = \mu(f^{-1}(B) \cap f^{-1}(F)) = \int_{f^{-1}(B)} \mathbb{I}_{f^{-1}(F)} \, d\mu \\ &= \int_{f^{-1}(B)} \mathbb{I}_F(f(x)) \mu(dx) = \int_{f^{-1}(B)} \mathbb{I}_F \circ f \, d\mu. \end{aligned}$$

□

A.2.2 矩生成函数、累计生成函数

对于随机变量 X , 若函数 $M_X(\lambda) := \mathbb{E}[e^{\lambda X}]$ 在 0 的开邻域 I 内存在, 则称 M_X 为 X 的矩生成函数 (moment generating function). 我们可以利用矩生成函数来导出 X 的各阶矩:

$$\frac{d^n}{d\lambda^n} M_X(\lambda) = \frac{d^n}{d\lambda^n} \mathbb{E} \left[1 + \lambda X + \frac{\lambda^2}{2!} X^2 + \dots \right] = \mathbb{E} X^n + \frac{\lambda}{n+1} \mathbb{E} X^{n+1} + \dots,$$

于是 $\mathbb{E} X^n = M_X^{(n)}(0)$. 随机变量 X 的累计生成函数 (cumulant generating function) 是它的矩生成函数的自然对数

$$K_X(\lambda) := \log \mathbb{E}[e^{\lambda X}],$$

它的一个优点是它是加性函数——对于互相独立的随机变量 X 和 Y ,

$$K_{X+Y}(\lambda) = \log (\mathbb{E}[e^{\lambda X}] \cdot \mathbb{E}[e^{\lambda Y}]) = K_X(\lambda) + K_Y(\lambda).$$

此外还可以按如下方式计算随机变量的各阶矩.

A.3 引理. 若非负随机变量 $X \in L^p$, $p > 0$, 则有

$$\mathbb{E} X^p = \int_0^\infty p x^{p-1} \mathbb{P}(X > x) \, dx. \quad (25)$$

特别的, 对于 $X \geq 0$, 有

$$\mathbb{E} X = \int_0^\infty \mathbb{P}(X > x) \, dx.$$

进一步地, 若 X 取值范围为 \mathbb{N} , 则有

$$\mathbb{E} X = \sum_{k=0}^\infty \mathbb{P}(X \geq k).$$

證明.

$$\begin{aligned} \mathbb{E} X^p &= \int_\Omega X^p \, d\mathbb{P} = \int_\Omega \int_0^Y p x^{p-1} \, dx \, d\mathbb{P} = \int_\Omega \int_0^\infty p x^{p-1} \mathbb{I}_{\{X > x\}} \, dx \, d\mathbb{P} \\ &= \int_0^\infty p x^{p-1} \int_\Omega \mathbb{I}_{\{X > x\}} \, d\mathbb{P} \, dx = \int_0^\infty p x^{p-1} \mathbb{P}(X > x) \, dx. \end{aligned}$$

□

A.2.3 Radon-Nikodym 導數、密度

Radon-Nikodym 导数是定义密度和条件期望的关键.

A.4 定理 (Radon-Nikodym 定理). 设 μ, ν 为可测空间 $(\mathcal{X}, \mathcal{A})$ 上的两个概率测度, ν 关于 μ 绝对连续, 即对于满足 $\mu(A) = 0$ 的 $A \in \mathcal{A}$, 一定有 $\nu(A) = 0$, 记做 $\nu \ll \mu$. 存在 \mathcal{X} 上的非负函数 f , 使得 $\nu(A) = \int_A f d\mu$, 且 f 在 μ -a.e. 意义下唯一, 记做 $f = \frac{d\nu}{d\mu}$.

A.5 示例 (分布的密度). 随机变量 $X: (\mathcal{X}, \mathcal{A}, \mathbb{P}) \rightarrow (\mathbb{R}, \mathcal{B}, \mu)$ 的分布是前推测度 $X_*\mathbb{P}(B) := \mathbb{P} \circ X^{-1}(B) = \mathbb{P}[X \in B]$, 它的关于 μ 的密度由 Radon-Nikodym 导数给出:

$$p_X = \frac{dX_*\mathbb{P}}{d\mu}.$$

从而由积分变换定理 A.2,

$$\begin{aligned} \mathbb{P}[X \in B] &= \int_{X^{-1}(B)} d\mathbb{P} = \int_B dX_*\mathbb{P} = \int_B p_X d\mu = \mathbb{E}[p_X; B] \\ \mathbb{E}[f(X); B] &= \int_B f dX_*\mathbb{P} = \int_B f(x) p_X(x) d\mu(x) \end{aligned}$$

A.6 示例 (指数加权). 若随机变量 X 的矩生成函数 $\mathbb{E}[e^{\lambda X}]$ 在某个开区间 I 存在, 考虑指数加权期望 \mathbb{E}_λ

$$\mathbb{E}_\lambda[f(X)] := \frac{\mathbb{E}[f(X)e^{\lambda X}]}{\mathbb{E}[e^{\lambda X}]} = \int f(X) \frac{e^{\lambda X}}{\mathbb{E}[e^{\lambda X}]} d\mathbb{P}$$

它借用了 Gibbs 分布的思想, 用指数权重调整对随机变量的关注, 并将矩生成函数作为配分函数. 也可以看作定义了新的概率测度 \mathbb{P}^λ , 它关于 \mathbb{P} 的 Radon-Nikodym 导数即为 $\frac{e^{\lambda X}}{\mathbb{E}[e^{\lambda X}]}$. 这被用在引理 1.7、定理 1.29、引理 1.42 等定理的证明.

A.2.4 條件期望、鞅、鞅差

给定概率空间 $(\Omega, \mathcal{F}_0, \mathbb{P})$, 子 σ -域 $\mathcal{F} \subset \mathcal{F}_0$, 随机变量 $X \in \mathcal{F}_0$ 可积. 称 Y 为 X 关于 \mathcal{F} 的条件期望, 如果

$$(1) Y \in \mathcal{F}; \quad (2) \text{ 对任意 } A \in \mathcal{F}, \mathbb{E}(Y; A) = \mathbb{E}(X; A).$$

可以证明这样的 Y 存在唯一 (a.s.), 且 $E|Y| < \infty$, 记做 $\mathbb{E}(X|\mathcal{F})$. 我们可以把 $X|\mathcal{F}$ 看作随机变量, 称为条件随机变量. 在这样的记号下, X 等价于 $X|\{\emptyset, \Omega\}$.

$$\text{条件概率 } \mathbb{P}(A|\mathcal{F}) = \mathbb{E}[\mathbb{I}_A|\mathcal{F}]$$

条件期望具有许多性质, 这里我们主要使用以下几个:

(i) 特别地, 如果 $X \in \mathcal{F}$, 则 $\mathbb{E}(X|\mathcal{F}) = X$ a.s.;

(ii) (全期望公式) $\mathbb{E}(\mathbb{E}(X|\mathcal{F})) = \mathbb{E}X$; (取 $A = \Omega \in \mathcal{F}$ 即可)

(iii) (Jensen 不等式) 若 M 为凸函数且 $\mathbb{E}X, \mathbb{E}M(X) < \infty$, 则 $\mathbb{E}(M(X)|\mathcal{F}) \geq M(\mathbb{E}(X|\mathcal{F}))$;

(iv) (塔性质) 若 $\mathcal{F}_1 \subset \mathcal{F}_2$, 则 $\mathbb{E}(\mathbb{E}(X|\mathcal{F}_1)|\mathcal{F}_2) = \mathbb{E}(X|\mathcal{F}_1) = \mathbb{E}(\mathbb{E}(X|\mathcal{F}_2)|\mathcal{F}_1)$.

随机变量序列 $\{X_k\}$ 是适应于 $\{\mathcal{F}_k\}$ 的鞅, 如果满足

$$(1) \mathbb{E}|X_k| < \infty; \quad (2) X_k \in \mathcal{F}_k; \quad (3) \mathbb{E}(X_{k+1}|\mathcal{F}_k) = X_k.$$

如果我们记 $D_k := X_k - X_{k-1}$, 容易验证 $\{D_k\}$ 期望为 0, 并且也是适应于 $\{\mathcal{F}_k\}$ 的鞅, 我们称其为鞅差.

A.2.5 方差的表示

方差的通常计算方式为 $\text{Var } X = \mathbb{E}[X - \mathbb{E}X]^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$, 这里我们介绍两种其他的表示方式.

A.7 引理 (方差的变分表示). 设随机变量 $X \in L^2$, 那么

$$\text{Var } X = \inf_{a \in \mathbb{R}} \mathbb{E}(X - a)^2.$$

证明. 记 $f(a) = \mathbb{E}(X - a)^2 = a^2 - 2\mathbb{E}X \cdot a + \mathbb{E}X^2$ 为二次函数, 不难看出 f 在 $\mathbb{E}X$ 有最小值 $-(\mathbb{E}X)^2 + \mathbb{E}X^2 = \text{Var } X$. \square

A.8 引理 (独立复制). 设随机变量 $X \in L^2$, X' 为 X 的独立复制, 那么

$$\text{Var } X = \frac{1}{2} \mathbb{E}(X - X')^2 = \mathbb{E}(X - X')_+^2 = \mathbb{E}(X - X')_-^2.$$

证明. 由独立性, $\mathbb{E}(X - X')^2 = \mathbb{E}X - 2\mathbb{E}X \cdot \mathbb{E}X + \mathbb{E}X^2 = 2 \text{Var } X$. 另一方面, $X - X'$ 和 $X' - X$ 有相同的分布, 于是 $\mathbb{E}(X - X')_+^2 = \mathbb{E}(X - X')_-^2$ 且两者之和即 $\mathbb{E}(X - X')^2$. \square

A.2.6 耦合

耦合是一种应用广泛的概率技术: 比较两个概率测度 \mathbb{Q}, \mathbb{P} , 我们可以考虑具有边缘分布 \mathbb{Q}, \mathbb{P} 的乘积概率空间.

为了比较概率空间 \mathcal{X} 上两个概率测度 \mathbb{Q}, \mathbb{P} , 我们可以

很多情况下, 构造乘积空间

的耦合, 是指 $\mathcal{X} \times \mathcal{X}$ 上的联合分布 \mathbb{M} , 其边缘分布满足

满足第一和第二坐标的边缘分布分别是 \mathbb{Q} 和 \mathbb{P} .

显然乘积测度 $\mathbb{Q} \otimes \mathbb{P}$ 是 (\mathbb{Q}, \mathbb{P}) 的耦合,

耦合并不唯一, 记为 $\mathcal{C}(\mathbb{Q}, \mathbb{P})$.

A.9 示例. 例如甲乙两人射击命中的概率分别为 0.75、0.5. 现在两人独立地重复射击 100 次, 记 N_A 和 N_B 分别为两人射中的次数, 试证明对任意 $1 \leq k \leq 100$, 都有

$$\mathbb{P}(N_A \geq k) \geq \mathbb{P}(N_B \geq k).$$

这一命题直观上是正确的, 但是如果我们显式地写出两者的概率

$$\mathbb{P}(N_A \geq k) = \sum_{n=k}^{100} \binom{100}{n} 0.75^n 0.25^{100-n}, \quad \mathbb{P}(N_B \geq k) = \sum_{n=k}^{100} \binom{100}{n} 0.5^{100}.$$

会发现两者是难以比较的. 回忆分布的定义, 两个随机变量同分布并不意味着它们的样本空间相同, 换言之, 相同的“输出”并不意味着有相同的“输入”. 借助这样的想法, 我们可以通过匹配适当的样本空间来实现这一目的: 设 $Z \sim U[0, 1]$, 引入随机变量

$$Y_{A,i} = \begin{cases} 1, & 0 \leq Z \leq 0.75 \\ 0, & 0.75 < Z \leq 1, \end{cases} \quad Y_{B,i} = \begin{cases} 1, & 0 \leq Z \leq 0.5 \\ 0, & 0.5 < Z \leq 1, \end{cases}$$

其中 $Y_{A,i}, Y_{B,i}$ 分别表示甲和乙第 i 次射中的次数. 显然有 $Y_{A,i} \geq Y_{B,i}, \forall 1 \leq i \leq 100$, 并且 N_A 和 $\sum_{n=1}^{100} Y_{A,i}$ 、 N_B 和 $\sum_{n=1}^{100} Y_{B,i}$ 同分布, 于是自然地

$$\mathbb{P}(N_A \geq k) = \mathbb{P}\left(\sum_{n=1}^{100} Y_{A,i} \geq k\right) \geq \mathbb{P}\left(\sum_{n=1}^{100} Y_{B,i} \geq k\right) = \mathbb{P}(N_B \geq k).$$

A.3 凸分析与凸优化

A.3.1 Rademacher 定理

A.10 定理 (Rademacher). 任意凸的 *Lipschitz* 函数几乎处处有导数

A.3.2 Fenchel 共轭

Fenchel 共轭是 Fourier 变换在凸分析中的对应. 对于实 Hilbert 空间 \mathcal{X} 上的正则函数 $g: \mathcal{X} \rightarrow (-\infty, +\infty]$, 即 $\text{dom } f := \{x \in \mathcal{X}: f(x) \in \mathbb{R} \neq \emptyset\}$, 其在 $u \in \mathcal{X}$ 的 **Fenchel 共轭** 为

$$f^*(u) = \sup_{x \in \mathcal{X}} \{\langle x, u \rangle - f(x)\}. \quad (26)$$

于是由定义, 对任意 $x \in \mathcal{X}$, $f^*(u) \geq \langle x, u \rangle - f(x)$, 从而有 **Fenchel-Young 不等式**

$$f(x) + f^*(u) \geq \langle x, u \rangle, \quad \forall x, u \in \mathcal{X} \quad (27)$$

此外, f^* 是凸的、下半连续的, 这是由于它是放射连续函数族 $(\langle x, \cdot \rangle - f(x))_{x \in \mathcal{X}}$ 的上确界. 对偶 $f = f^{**}$ 当且仅当 f 是凸的、下半连续函数

A.3.3 Lagrange 乘数法

Lagrange 乘数法是常用的约束优化方法 [Bur04]

我们寻求极小化函数 $f(\mathbf{x})$. 在单一等式约束 $c(\mathbf{x}) = 0$ 的情况下, 最优解 \mathbf{x}^* 处一定有 $\nabla f(\mathbf{x}^*) = \lambda \nabla c(\mathbf{x}^*)$, 其中 λ 称作 (未确定的)Lagrange 乘数——否则 $\nabla c(\mathbf{x})$ 不与 $\nabla f(\mathbf{x})$ 平行, 那么沿着 $c(\mathbf{x}) = 0$ 的运动一定有

与 $\nabla c(\mathbf{x})$ 垂直的方向

A.4 矩阵

A.4.1 矩阵范数

矩阵的 Frobenius 范数是 $\|M\|_F = \sqrt{\sum_{ij} m_{ij}^2}$.

矩阵的谱范数或者 ℓ_2 算子范数等价于最大奇异值

$$\|M\|_2 = \max_{\|v\|_2=1} \|Mv\|_2 = \max_{\|v\|_2=1} \max_{\|u\|_2=1} u^T M v.$$

B 定理證明

B.1 定理1.9的證明

Wasserstein 距离构成了度量

證明. [正定性]

[对称性] 对任意 $\epsilon > 0$, 存在满足 Lipschitz 函数 f 使得 $W_\rho(\mathbb{Q}, \mathbb{P}) - \epsilon \leq \mathbb{E}_\mathbb{Q}[f] - \mathbb{E}_\mathbb{P}[f]$, 于是

$$W_\rho(\mathbb{P}, \mathbb{Q}) \geq \mathbb{E}_\mathbb{P}[-f] - \mathbb{E}_\mathbb{Q}[-f] = \mathbb{E}_\mathbb{Q}[f] - \mathbb{E}_\mathbb{P}[f] \geq W_\rho(\mathbb{Q}, \mathbb{P}) - \epsilon$$

令 $\epsilon \rightarrow 0$, 可得 $W_\rho(\mathbb{P}, \mathbb{Q}) \geq W_\rho(\mathbb{Q}, \mathbb{P})$. 类似地, $W_\rho(\mathbb{Q}, \mathbb{P}) \geq W_\rho(\mathbb{P}, \mathbb{Q})$, 从而二者相等.

[三角不等式]

□

王家卫在《一代宗师》里寄出一句台词:

人生要是无憾, 那多无趣?

而我说: 算法要是无憾, 那应该是过拟合了。

參考文獻

- [BL12] H Bauschke and Yves Lucet, *What is a fenchel conjugate*, Notices of the AMS **59** (2012), no. 1, 44–46.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*, Oxford University Press, 02 2013.
- [Bur04] Christopher J. C. Burges, *Some notes on applied mathematics for machine learning*, pp. 21–40, Springer Berlin Heidelberg, 2004.
- [Çin11] Erhan Çinlar, *Probability and stochastics*, Graduate Texts in Mathematics, vol. 261, Springer New York, 2011.
- [GN15] Evarist Giné and Richard Nickl, *Mathematical foundations of infinite-dimensional statistical models*, Cambridge University Press, 2015.
- [Tro23] Joel A. Tropp, *Acm 217: Probability in high dimensions*, August 2023.
- [Ver18] Roman Vershynin, *High-dimensional probability: An introduction with applications in data science*, Cambridge University Press, September 2018.
- [vH16] Ramon van Handel, *Probability in high dimension*, APC 550 Lecture Notes Princeton University, December 2016.
- [Wai19] Martin J. Wainwright, *High-dimensional statistics: A non-asymptotic viewpoint*, Cambridge University Press, February 2019.