

Probability Inequalities

Thm 1 (Gaussian Tail Inequality): Let $X \sim \mathcal{N}(0, 1)$. Then

$$\mathbb{P}(|X| > \epsilon) \leq \frac{2}{\epsilon} e^{-\epsilon^2/2} \quad (1)$$

Additionally:

$$\mathbb{P}(|\bar{X}_n| > \epsilon) \leq \frac{1}{\sqrt{n}\epsilon} e^{-n\epsilon^2/2} \quad (2)$$

Thm 2 (Markov Inequality): Let X be a non-negative random variable s.t. $\mathbb{E}(X)$ exists. Then $\forall t > 0$

$$\mathbb{P}(X > t) \leq \frac{\mathbb{E}(X)}{t} \quad (3)$$

Thm 3 (Chebyshev's Inequality): Let $\mu = \mathbb{E}(X)$ and $\sigma^2 = \text{Var}(X)$. Then

$$\begin{aligned} \mathbb{P}(|X - \mu| \geq t) &\leq \frac{\sigma^2}{t^2} \\ \mathbb{P}(|(X - \mu)/\sigma| \geq t) &\leq \frac{1}{t^2} \end{aligned} \quad (4) \quad (5)$$

Lemma 4: Let $\mathbb{E}(X) = 0$ and $a \leq X \leq b$. Then

$$\mathbb{E}(e^{tX}) \leq e^{t^2(b-a)^2/8} \quad (6)$$

Lemma 5: Let X be any random variable. Then

$$\mathbb{P}(X > \epsilon) \leq \inf_{t \geq 0} e^{-t\epsilon} \mathbb{E}(e^{tX}) \quad (7)$$

Thm 6 (Hoeffding's Inequality): X_1, \dots, X_n iid, $\mathbb{E}(X_i) = \mu$, $a \leq X_i \leq b$. Then $\forall \epsilon > 0$

$$\mathbb{P}(|\bar{X} - \mu| \geq \epsilon) \leq 2e^{-2n\epsilon^2/(b-a)^2} \quad (8)$$

Thm 9 (McDiarmid):

Thm 12 (Cauchy-Schwartz inequality):

Thm 13 (Jensen's inequality):

Ex 15 (Kullback Leibler distance):

Thm 18?:

O_p and o_p :

Shattering

Note: remember uniform bounds and union bound.

\mathcal{A} picks out $G \subset F$.

$S(\mathcal{A}, F)$.

F shattered by \mathcal{A} if $S(\mathcal{A}, F) = 2^{|F|}$ (ie if \mathcal{A} picks out all $G \subset F$).

The shatter coefficient $s_n(\mathcal{A}) = \sup_{F \in \mathcal{F}_n} s(\mathcal{A}, F)$. Note $n = |F|$ and $s_n(\mathcal{A}) \leq 2^n$.

Thm 5:

The VC dimension $d(\mathcal{A}) = \text{largest } n \text{ s.t. } s_n(\mathcal{A}) = 2^n$.

Random Samples

For $X_1, \dots, X_n \sim F$ a statistic is any $T = g(X_1, \dots, X_n)$.

E.g. $\bar{X}_n, S_n = \sum_i (X_i - \bar{X}_n)^2 / (n-1), (X_{(1)}, \dots, X_{(n)})$

Note: $\mathbb{E}(\bar{X}_n) = \mathbb{E}(X_i)$, $\text{Var}(\bar{X}_n) = \text{Var}(X_i)/n$, $\mathbb{E}(S_n)^2 = \text{Var}(X_i)$.

Note: sum of bernoulli is binomial(n,p), sum of exp(beta) is gamma(n,beta), sum of standard normal is chi-squared(n dof).

Thm. 1: $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2) \implies \bar{X}_n \sim \mathcal{N}(\mu, \sigma^2/n)$.

Convergence

X, X_1, X_2, \dots random variables.

(1) X_n converges **almost surely** $X_n \xrightarrow{a.s.} X$ if $\forall \epsilon > 0$

$$\mathbb{P}(\lim_{n \rightarrow \infty} |X_n - X| < \epsilon) = 1 \quad (9)$$

(2) X_n converges **in probability** $X_n \xrightarrow{p} X$ if $\forall \epsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \epsilon) = 0 \quad (10)$$

(3) X_n converges **in quadratic mean** $X_n \xrightarrow{qm} X$ if

$$\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - X)^2] = 0 \quad (11)$$

(4) X_n converges **in distribution** $X_n \rightsquigarrow X$ if

$$\lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t) \quad (12)$$

$\forall t$ on which F_X is continuous.

Thm 7: Conv. a.s. and in q.m. imply conv. in prob. All three imply conv. in distribution. Conv. in distribution to a point-mass also implies conv. in prob.

E.g. from class: Showed conv. in prob $\not\Rightarrow$ conv. a.s.. Showed conv. in prob $\not\Rightarrow$ conv. in q.m.. Showed conv. in distro $\not\Rightarrow$ conv. in prob.

Thm 10a: X, X_n, Y, Y_n random variables. Then

$$(a) \quad X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \implies X_n + Y_n \xrightarrow{p} X + Y \quad (13)$$

$$(b) \quad X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \implies X_n Y_n \xrightarrow{p} XY \quad (14)$$

$$(c) \quad X_n \xrightarrow{qm} X, Y_n \xrightarrow{qm} Y \implies X_n + Y_n \xrightarrow{qm} X + Y \quad (15)$$

Thm 10b (Slutzky's Thm): X, X_n, Y_n random variables. Then

$$(a) \quad X_n \rightsquigarrow X, Y_n \rightsquigarrow c \implies X_n + Y_n \rightsquigarrow X + c \quad (16)$$

$$(b) \quad X_n \rightsquigarrow X, Y_n \rightsquigarrow c \implies X_n Y_n \rightsquigarrow cX \quad (17)$$

Thm 12 (Law of Large Numbers): X_1, \dots, X_n iid, $\mathbb{E}(X_i) = \mu \implies \bar{X}_n \xrightarrow{qm} \mu$.

Thm 14 (CLT): X_1, \dots, X_n iid, $\mathbb{E}(X_i) = \mu$, $\text{Var}(X_i) = \sigma^2$

$\implies \sqrt{n}(\bar{X}_n - \mu)/\sigma \rightsquigarrow \mathcal{N}(0, 1)$

$\implies \bar{X}_n \rightsquigarrow \mathcal{N}(\mu, \sigma^2/n)$

$\implies \sqrt{n}(\bar{X}_n - \mu)/S_n \rightsquigarrow \mathcal{N}(0, 1)$

Thm 18 (delta method): If $\sqrt{n}(Y_n - \mu)/\sigma \rightsquigarrow \mathcal{N}(0, 1)$, $g'(\mu) \neq 0$

$\implies \sqrt{n}(g(Y_n) - g(\mu))/|g'(\mu)|\sigma \rightsquigarrow \mathcal{N}(0, 1)$

ie $Y_n \approx \mathcal{N}(\mu, \sigma^2/n) \implies g(Y_n) \approx \mathcal{N}(g(\mu), g'(\mu)^2 \sigma^2/n)$

Thm 18b (2nd order delta method):?? Should I include this?

Sufficiency

If $X_1, \dots, X_n \sim p(x; \theta)$, T **sufficient** for θ if $p(x^n | t; \theta) = p(x^n | t)$.

Thm 9 (factorization): for $X^n \sim p(x; \theta)$, $T(X^n)$ sufficient for θ if the joint probability can be factorized as.

$$p(x^n; \theta) = h(x^n) \times g(t; \theta) \quad (18)$$

Sufficient stat T is a **minimal sufficient statistic (MSS)** if $T = g(U)$ for all other sufficient stats.

Thm 15: T is a MSS if:

$$\frac{p(y^n; \theta)}{p(x^n; \theta)} \text{ constant in } \theta \iff T(y^n) = T(x^n) \quad (19)$$

Parametric Point Estimation

make sure i've defined: $\mathbb{E}_\theta(\hat{\theta})$, bias, sampling distro, standard error, $\hat{\theta}_n$ consistent.

Method of Moments

Maximum Likelihood (MLE)

Bayes Estimation

Mean Squared Error (MSE)

ex: MSE for normal

Bias and Variance

Risks and Estimators

Let $x^n = x_1, \dots, x_n$.

The **risk** of an estimator $\hat{\theta} = \hat{\theta}(x^n)$ for $x^n \sim p(x^n; \theta)$ is

$$R(\theta, \hat{\theta}) = \mathbb{E}[L(\theta, \hat{\theta})] = \int L(\theta, \hat{\theta}) p(x^n; \theta) dx^n \quad (20)$$

The **posterior risk** of $\hat{\theta}$ given a prior $\pi(\theta)$ is

$$r(\hat{\theta}|x^n) = \int L(\theta, \hat{\theta}) \pi(\theta|x_1, \dots, x_n) d\theta \quad (21)$$

where $\pi(\theta|x^n) = \frac{\mathbb{P}(x^n; \theta) \pi(\theta)}{m(x^n)}$ is the posterior over θ .

The **Bayes risk** of $\hat{\theta}$ given a prior $\pi(\theta)$ is

$$B_\pi(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta \quad (22)$$

The **max risk** of $\hat{\theta}$ is

$$\overline{R}(\hat{\theta}) = \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \quad (23)$$

The **minimax risk** is

$$R_n = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \quad (24)$$