

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

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad (4)$$

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

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): X_1, \dots, X_n indep't. If $\sup_{x_1, \dots, x_n, x'_i} |g(x_1, \dots, x_n) - g_i^*(x_1, \dots, x_n)| \leq c_i \quad \forall i, \implies$

$$\mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}(g(X_1, \dots, X_n)) \geq \epsilon) \leq e^{-2\epsilon^2/\sum_i c_i^2} \quad (9)$$

where $g_i^* = g$ with x_i replaced by x'_i .

Thm 12 (Cauchy-Schwartz inequality):

Thm 13 (Jensen's inequality):

Ex 15 (Kullback Leibler distance):

Thm 18:

O_p and o_p : $X_n = o_p(1)$ if $\forall \epsilon > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|X_n| > \epsilon) = 0$.

$X_n = O_p(1)$ if $\forall \epsilon > 0, \exists C > 0$ s.t. $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n| > C) \leq \epsilon$.

$X_n = o_p(a_n)$ if $X_n/a_n = o_p(1)$ and $X_n = O_p(a_n)$ if $X_n/a_n = O_p(1)$.

Shattering

Note: remember uniform bounds and union bound.

F a finite set, $|F| = n$, and $G \subset F$. \mathcal{A} is a class of sets.

\mathcal{A} picks out G if $\exists A \in \mathcal{A}$ s.t. $A \cap F = G$.

Let $S(\mathcal{A}, F) = |\{G \subset F \text{ picked out by } \mathcal{A}\}| \leq 2^n$.

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

Let \mathcal{F}_n be all finite sets with n elements.

The **shatter coefficient** $s_n(\mathcal{A}) = \sup_{F \in \mathcal{F}_n} S(\mathcal{A}, F) \leq 2^n$.

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

Thm 5: $\forall \epsilon > 0, \mathbb{P}(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \epsilon) \leq 8s_n(\mathcal{A})e^{-n\epsilon^2/32}$

Random Samples

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

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

Notes: $\mathbb{E}(\bar{X}_n) = \mathbb{E}(X_i), \text{Var}(\bar{X}_n) = \text{Var}(X_i)/n, \mathbb{E}(S_n)^2 = \text{Var}(X_i), X_{1, \dots, n} \sim \text{Bern}(p) \implies \sum_i X_i \sim \text{Bin}(n, p), X_{1, \dots, n} \sim \text{Exp}(\beta) \implies \sum_i X_i \sim \Gamma(n, \beta), X_{1, \dots, n} \sim \mathcal{N}(0, 1) \implies \sum_i X_i^2 \sim \chi_n.$

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 $\forall \epsilon > 0$

$$\mathbb{P}(\lim_{n \rightarrow \infty} |X_n - X| < \epsilon) = 1 \quad (10)$$

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

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| < \epsilon) = 1 \quad (11)$$

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

$$\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - X)^2] = 0 \quad (12)$$

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

$$\lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t) \quad (13)$$

$\forall t$ 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.

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

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

$$(b) \quad X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \implies X_n Y_n \xrightarrow{p} XY \quad (15)$$

$$(c) \quad X_n \xrightarrow{qm} X, Y_n \xrightarrow{qm} Y \implies X_n + Y_n \xrightarrow{qm} X + Y \quad (16)$$

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

$$(a) \quad X_n \rightsquigarrow X, Y_n \rightsquigarrow c \implies X_n + Y_n \rightsquigarrow X + c \quad (17)$$

$$(b) \quad X_n \rightsquigarrow X, Y_n \rightsquigarrow c \implies X_n Y_n \rightsquigarrow cX \quad (18)$$

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):

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.

$$p(x^n; \theta) = h(x^n) \times g(t; \theta) \quad (19)$$

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

Thm 15: T is a MSS if:

$$\frac{p(y^n; \theta)}{p(x^n; \theta)} \text{ constant in } \theta \iff T(y^n) = T(x^n) \quad (20)$$

Parametric Point Estimation

Method of Moments: Define equations

- (a) $(\sum_i X_i)/n = \mathbb{E}_{\hat{\theta}}(X_i)$
- (b) $(\sum_i X_i^2)/n = \mathbb{E}_{\hat{\theta}}(X_i^2)$
- (c) ...

And solve for $\hat{\theta}$.

Maximum Likelihood (MLE): The MLE is

$$\hat{\theta} = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} l(\theta) \quad (21)$$

Often suffices to solve for θ in $\frac{\partial l(\theta)}{\partial \theta} = 0$. The MLE is **equivariant**
 \implies if $\eta = g(\theta)$ then $\hat{\eta} = g(\hat{\theta})$.

Bayes Estimation: For prior $\pi(\theta)$, choose

$$\hat{\theta} = \mathbb{E}(\theta|x^n) = \int \theta \pi(\theta|x^n) d\pi \quad (22)$$

Mean Squared Error (MSE): The MSE is

$$\text{MSE} = \mathbb{E}(\hat{\theta} - \theta)^2 = \int (\hat{\theta} - \theta)^2 p(x^n; \theta) dx^n = \text{bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta}) \quad (23)$$

Defs: **bias**($\hat{\theta}$) = $\mathbb{E}(\hat{\theta}) - \theta$. We say $\hat{\theta}$ is **consistent** if $\hat{\theta} = \hat{\theta}_n \xrightarrow{P} \theta$.
 The **standard error** of $\hat{\theta}$, $\text{se}(\hat{\theta})$, is the standard deviation of $\hat{\theta}$.

Risks and Estimators

$L(\theta, \hat{\theta})$ is the **loss** of an estimator $\hat{\theta} = \hat{\theta}(x^n)$ for $x^n \sim p(x^n; \theta)$.
 The **risk** of this $\hat{\theta}$ is

$$R(\theta, \hat{\theta}) = \mathbb{E}[L(\theta, \hat{\theta})] = \int L(\theta, \hat{\theta}) p(x^n; \theta) dx^n \quad (24)$$

When $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$, the risk is the MSE.

The **max risk** of $\hat{\theta}$ over a set $\theta \in \Theta$ is

$$\bar{R}(\hat{\theta}) = \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \quad (25)$$

The **minimax estimator** is

$$\hat{\theta} = \arg \inf_{\hat{\theta}} \bar{R}(\hat{\theta}) \quad (26)$$

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

$$B_{\pi}(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta \quad (27)$$

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

$$r(\hat{\theta}|x^n) = \int L(\theta, \hat{\theta}) \pi(\theta|x^n) d\theta \quad (28)$$

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

The **Bayes estimator** is

$$\hat{\theta} = \arg \inf_{\hat{\theta}} B_{\pi}(\hat{\theta}) = \arg \inf_{\hat{\theta}} r(\hat{\theta}|x^n) \quad (29)$$

which equals the posterior mean $\mathbb{E}(\theta|x^n)$ when $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$,
 the posterior median when $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$, and the posterior mode
 when $L(\theta, \hat{\theta}) = \mathbb{I}[\theta \neq \hat{\theta}]$.

Thm 10: If $\hat{\theta}$ is a Bayes estimator for some prior π and $R(\theta, \hat{\theta})$ is constant, then $\hat{\theta}$ is a minimax estimator.

Note: The MLE is approximately minimax (as n increases, if dimension of the parameter is fixed).

Distributions

Discrete distributions:

$$(a) \text{ Bernoulli } f(x|p) = p^x (1-p)^{1-x}, \quad x \in \{0, 1\} \quad (30)$$

$$(b) \text{ Binomial } f(x|n, p) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x \in \{0, 1, \dots, n\} \quad (31)$$

$$(c) \text{ Poisson } f(x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x \in \{0, 1, 2, \dots\} \quad (32)$$

Continuous distributions:

$$(a) \text{ Uniform } f(x|a, b) = \frac{1}{b-a}, \quad x \in [a, b] \quad (33)$$

$$(b) \text{ Normal } f(x|\mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)}, \quad x \in \mathbb{R} \quad (34)$$

$$(c) \text{ Gamma } f(x|\alpha, \beta) = \frac{1}{\Gamma(\alpha) \beta^\alpha} x^{\alpha-1} e^{-x/\beta}, \quad x \in \mathbb{R}_+, \alpha, \beta > 0 \quad (35)$$

Expected Values

The **mean** or **expected value** of $g(X)$ is

$$\mathbb{E}(g(X)) = \int g(x) dF(x) = \int g(x) dP(x) \quad (36)$$

Related properties and definitions:

$$(a) \mu = \mathbb{E}(X) \quad (37)$$

$$(b) \mathbb{E}(\sum_i c_i g_i(X_i)) = \sum_i c_i \mathbb{E}(g_i(X_i)) \quad (38)$$

$$(c) \mathbb{E}\left(\prod_i X_i\right) = \prod_i \mathbb{E}(X_i), \quad X_1, \dots, X_n \text{ indep't} \quad (39)$$

$$(d) \text{Var}(X) = \sigma^2 = \mathbb{E}((X - \mu)^2) \quad \text{is the } \mathbf{variance} \text{ of } X \quad (40)$$

$$(e) \text{Var}(X) = \mathbb{E}(X^2) - \mu^2 \quad (41)$$

$$(f) \text{Var}\left(\sum_i a_i X_i\right) = \sum_i a_i^2 \text{Var}(X_i), \quad X_1, \dots, X_n \text{ indep't} \quad (42)$$

$$(g) \text{Cov}(X, Y) = \mathbb{E}((X - \mu_X)(Y - \mu_Y)) \quad \text{is the } \mathbf{covariance} \quad (43)$$

$$(h) \text{Cov}(X, Y) = \mathbb{E}(XY) - \mu_X \mu_Y \quad (44)$$

$$(i) \rho(X, Y) = \text{Cov}(X, Y) / \sigma_X \sigma_Y, \quad -1 \leq \rho(X, Y) \leq 1 \quad (45)$$

The **conditional expectation** of Y given X is the random variable $g(X) = \mathbb{E}(Y|X)$, where

$$\mathbb{E}(Y|X = x) = \int y f(y|x) dy \quad (46)$$

$$\text{and } f(y|x) = f_{X,Y}(x, y) / f_X(x) \quad (47)$$

The *Law of Total/Iterated Expectation* is

$$\mathbb{E}(Y) = \mathbb{E}[\mathbb{E}(Y|X)] \quad (48)$$

The *Law of Total Variance* is

$$\text{Var}(Y) = \text{Var}[\mathbb{E}(Y|X)] + \mathbb{E}[\text{Var}(Y|X)] \quad (49)$$

The *Law of Total Covariance* is

$$\text{Cov}(X, Y) = \mathbb{E}(\text{Cov}(X, Y|Z)) + \text{Cov}(\mathbb{E}(X|Z), \mathbb{E}(Y|Z)) \quad (50)$$