# Intermediate Statistics

By Junru Shao <junrushao1994@gmail.com>

## Misc

- Moment generating function: $M_x(t) = \mathbb{E}\left[\exp(tx)\right]$
- $\sigma$-sub-Gaussian: $M_{x-\mu} \le \exp\left(\sigma^2 t^2/2\right)$, 0-mean $[a, b]$ bounded RVs are $\frac{(b-a)}{2}$-sub-Gaussian
- Jensen's inequality: for $g$ convex, $g\left(\mathbb{E}[X]\right) \le \mathbb{E}\left[g(X)\right]$
- Cauchy inequality: $\left(\mathbb{E}[XY]\right)^2 \le \mathbb{E}\left[X^2\right]\mathbb{E}\left[Y^2\right]$
- maximin $\le$ minimax
- $\mathrm{Var}(X) = \mathbb{E}\left[X^2\right] - \mathbb{E}[X]^2$, $\mathrm{Var}(aX) = a^2\mathrm{Var}(X)$, $\mathrm{Var}\left(\sum_{i=1}^n \alpha_i X_i\right) = \sum_{i=1}^n \sum_{i=1}^n \alpha_i \alpha_j \mathrm{Cov}(X_i, X_j)$
- $\mathrm{Var}(XY) = \mathbb{E}\left[X^2\right]\mathbb{E}\left[Y^2\right] - \mathbb{E}[X]^2\mathbb{E}[Y]^2$
- Gamma function $\Gamma(z) = \int_0^\infty x^{z-1}\exp(-x)\,dx$. $\Gamma(n) = (n-1)!$ for $n \in \mathbb{N}$. $\Gamma(x+1) = x\Gamma(x)$
- Law of total expectation: $\mathbb{E}[X] = \mathbb{E}\left[\mathbb{E}[X \mid Y]\right]$, total variance: $\mathrm{Var}(Y) = \mathbb{E}\left[\mathrm{Var}(Y \mid X)\right] + \mathrm{Var}\left[\mathbb{E}(Y \mid X)\right]$
- Conditional gaussian: $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$ ($\boldsymbol{\mu}_1 \in \mathbb{R}^q$, $\boldsymbol{\mu}_2 \in \mathbb{R}^{N-q}$), $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$. Then
  - $\overline{\boldsymbol{\mu}} = \boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^{-1}(\boldsymbol{a} - \boldsymbol{\mu}_2)$, $\overline{\Sigma} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$
- Posterior for $\mu$ in Gaussian, with known $\sigma$, and prior $\mu \sim \mathcal{N}(m, \tau^2)$
  - $\mathbb{E}[\mu \mid X^n] = \frac{\tau^2}{\tau^2 + \sigma^2/n}\overline{X}_n + \frac{\sigma^2/n}{\tau^2 + \sigma^2/n}m$
  - $\mathrm{Var}(\mu \mid X^n) = \tau^2 \cdot \frac{\sigma^2}{n} / \left(\tau^2 + \frac{\sigma^2}{n}\right)$
- Berry-Essen Theorem: let $F_n(x) = P\left(\frac{\sqrt{n}}{\sigma}(\hat{\mu} - \mu) \le x\right)$, then $\sup|F_n(x) - \Phi(x)| \le \frac{33}{4}\mathbb{E}|X_1 - \mu|^3 / \sigma^3\sqrt{n}$
- KL divergence: $\mathrm{KL}(P\|Q) = \sum_i P(i)\log\frac{P(i)}{Q(i)}$
- Total variance distance: $\delta(P, Q) = \frac{1}{2}\int |P(x) - Q(x)|\,dx$

## Inequalities

- Markov: for $X \ge 0$, $P[X \ge t] \le \mathbb{E}[X]/t$
- Chebyshev: $P[|X - \mathbb{E}X| \ge k\sigma] \le 1/k^2$
- Chernoff bound: $P[X - \mu \ge u] \le \inf_{0 \le t \le b}\frac{M_{x-\mu}(t)}{\exp(tu)}$, where $M_{x-\mu}$ is finite $\forall |t| \le b$
  - (sub-)Gaussian tail bound: $P(|X - \mu| \ge k\sigma) \le 2\exp\left(-2k^2\right)$
- Hoeffding bound (bounded RV): $P(|x - \mu| \ge k(b-a)) \le 2\exp\left(-2k^2\right)$
  - for $n$ RVs, $P(|x - \mu| \ge t) \le 2\exp\left(-2n^2t^2 / \sum_{i=1}^n (b_i - a_i)^2\right)$
- Bernstein ($[a, b]$ bounded RV, $\mu = 0$, small $\sigma$, i.i.d.)
  - $P(|\mu - \hat{\mu}| \ge t) \le 2\exp\left(-\frac{nt^2}{2(\sigma^2 + (b-a)t)}\right)$
  - $P\left(|\mu - \hat{\mu}| \le 4\sigma\sqrt{\frac{\ln(2/\delta)}{n}} + \frac{4(b-a)\ln(2/\delta)}{n}\right) \ge 1 - \delta$
- Azuma's bound (difference bounded: $\left|f(x_1, \ldots, x_k, \ldots, x_n) - f(x_1, \ldots, x_k', \ldots, x_n)\right| \le L_k$)
  - $P[|f(x_1, \ldots, f_n) - \mu| \ge t] \le 2\exp\left(-2t^2 / \sum_{k=1}^n L_k^2\right)$
- U-statistics: $U(x_1, \ldots, x_n) = \frac{1}{\binom{n}{2}}\sum_{j<k} g(x_j, x_k)$ and $g$ is symmetric and $|g| \le B$
  - Azuma gives: $P[|U(x_1, \ldots, x_n) - \mu| \ge t] \le 2\exp\left(-nt^2/8B^2\right)$
- $\chi^2$ tail bound: $Y \sim X_n^2$ where $X_n \sim \mathcal{N}(0, 1)$. $P\left[\left|\frac{1}{n}Y - 1\right| \ge t\right] \le 2\exp\left(-nt^2/8\right)$
- Johnson-Lindenstrauss Lemma: for $X^n \in \mathbb{R}^d$, $m \ge 16\frac{\log(n/\delta)}{\varepsilon^2}$, $Z \in \mathbb{R}^{m \times d}$, $Z_{i,j} \sim \mathcal{N}(0, 1)$, $F(X_i) = \frac{Z}{\sqrt{m}}X_i$
  - $(1 - \varepsilon)\|X_i - X_j\|_2^2 \le \|F(X_i) - F(X_j)\|_2^2 \le (1 + \varepsilon)\|X_i - X_j\|_2^2$, holds w.p. $1 - \delta$

## Convergence: estimator $\hat{\theta}_n$ is consistent iff $\hat{\theta}_n \xrightarrow{p} \theta$

- a.s. $P\left(\lim_{n \to \infty} X_n = X\right) = 1$; i.p. $\lim_{n \to \infty} P(|X_n - X| \ge \varepsilon) = 0$; q.m. $\lim_{n \to \infty} \mathbb{E}(X_n - X)^2 = 0$; d $\lim_{n \to \infty} F_{X_n}(t) = F_X(t)$
  - a.s. $\Rightarrow$ p; q.m. $\Rightarrow$ p; q.m. $\Leftarrow$ p if $|X_n|$ bounded; p $\Rightarrow$ d; p $\Leftarrow$ d if $X = c$; p $\not\Rightarrow$ $\mathbb{E}[X_n] = c$
- Continuous mapping: if $X_n \xrightarrow{p} X$, then for any continuous function $f$, $f(X_n) \xrightarrow{p} f(X)$; if $X_n \xrightarrow{d} X$, then for any continuous function $f$, $f(X_n) \xrightarrow{d} f(X)$
- Slutsky's theorem: if $X_n \xrightarrow{p} X$, $Y_n \xrightarrow{p} Y$, then $X_n + Y_n \xrightarrow{p} X + Y$, and $X_n Y_n \xrightarrow{p} XY$; if $X_n \xrightarrow{d} X$, $Y_n \xrightarrow{d} c$, then $X_n + Y_n \xrightarrow{d} X + c$, and $X_n Y_n \xrightarrow{d} cX$
- WLLN: $\frac{1}{n}\sum_{i=1}^n X_i \xrightarrow{p} \mu$ for i.i.d. $X_i$s with $\mathbb{E}(|X|) < \infty$ and $\mathrm{Var}(X) < \infty$
- CLT: $X^n$ independent RVs, mean and variance finite, then $\frac{\sqrt{n}}{\sigma}(\hat{\mu} - \mu) \xrightarrow{d} \mathcal{N}(0, 1)$, works also with estimated variance
  - Lyapunov CLT: $X^n$ are independent but not i.i.d. Define $s_n = \frac{1}{n}\sqrt{\sum_{i=1}^n \sigma_n^2}$. If $\lim_{n \to \infty} \frac{1}{s_n^3}\sum_{i=1}^n \mathbb{E}|X_i - $

$\mu_i|^3 < \infty$, we have $\frac{1}{s_n}\sum_{i=1}^n (X_i - \mu_i) \xrightarrow{d} \mathcal{N}(0, 1)$
  - multivariate CLT: $\sqrt{n}(\mu - \hat{\mu}) \xrightarrow{d} \mathcal{N}(0, \Sigma)$
- Delta method (CLT for function of RV): if $\frac{\sqrt{n}}{\sigma}(\hat{\mu} - \mu) \to \mathcal{N}(0, 1)$, $g$ continuous differentiable, $g'(\mu) \ne 0$
  - $\frac{\sqrt{n}}{\sigma}(g(\hat{\mu}) - g(\mu)) \to \mathcal{N}\left(0, [g'(\mu)]^2\right)$
  - multivariate delta method: $\sqrt{n}(g(\hat{\mu}) - g(\mu)) \to \mathcal{N}\left(0, \nabla_g^T(\mu)\Sigma\nabla_g(\mu)\right)$

## Empirical process

- Vapnik-Chervonenkis Theory
  - empirical CDF: $\hat{F}_n(x) = \frac{1}{n}\sum_{i=1}^n \mathbb{I}(X_i \le x)$
    * Glivenko-Cantelli Theorem: $\Delta = \sup_{x \in \mathbb{R}}\left|\hat{F}_n(x) - F_X(x)\right| \xrightarrow{p} 0$
  - empirical probability of set: $P_n(A) = \frac{1}{n}\sum_{i=1}^n \mathbb{I}(X_i \in A)$, and $\Delta(A) = \sup_{A \in \mathcal{A}}|P_n(A) - P(A)|$
  - empirical process: $\Delta(\mathcal{F}) = \sup_{f \in \mathcal{F}}\left|\frac{1}{n}\sum_{i=1}^n f(X_i) - \mathbb{E}[f]\right|$
  - Glivenko-Cantelli class: $\Delta(\mathcal{F}) \xrightarrow{p} 0$ where $X^n$ i.i.d.
- Shattering: $n$-th shattering coef of $\mathcal{A}$: $s(\mathcal{A}, n) = \max_{z^n} N_{\mathcal{A}}(z^n)$, where $N_{\mathcal{A}} = \#\{z^n \cap A\}$ (valid coloring)
- VC dimension: max $d$ that $s(\mathcal{A}, d) = 2^d$. There exists $\{z_i\}_{i=1}^d$ each of whose colorings are valid, but for any $(d+1)$ points, there exists an invalid coloring
  - you select the best $x_1, \ldots, x_n$, adversary assigns label $y_1, \ldots, y_n$
  - if $\mathrm{VC}_{\mathcal{A}} \ge n$, you can find $f \in \mathcal{A}$ that is consistent with the labels
- VC dimension example:
  - $\mathcal{A} = \{A_1, \ldots, A_N\}$, $V_{\mathcal{A}} \le \log_2 N$
  - intervals $[a, b]$ on the real line: 2
  - discs in $\mathbb{R}^2$: 3
  - closed balls in $\mathbb{R}^d$: $V_{\mathcal{A}} \le d + 2$
  - rectangles in $\mathbb{R}^d$: $2d$
  - half-space in $\mathbb{R}^d$: $d + 1$
  - convex polygons in $\mathbb{R}^2$: $\infty$
  - convex polygon with $d$ vertices: $2d + 1$
- Empirical risk minimization: $\hat{f} = \arg\min_{f \in \mathcal{F}} \hat{R}_n(f)$, optimal $f^* = \arg\min_{f \in \mathcal{F}} R(f)$
  - minimize $\Delta(\mathcal{F}) = R(\hat{f}) - R(f^*)$, consider minimize $\Delta(\mathcal{A}) = \sup_{A \in \mathcal{A}}|P_n(A) - P(A)|$
  - if $|A|$ is finite: $P(|P_n(A) - P(A)| \ge t) \le 2\exp\left(-2nt^2\right)$, $P(\Delta(\mathcal{A}) \ge t) \le 2|A|\exp\left(-2nt^2\right)$
  - if $|A|$ is not finite: $P(\Delta(\mathcal{A}) \ge t) \le 8s(\mathcal{A}, n)\exp\left(-nt^2/32\right)$
- Sauer's Lemma: if $\mathrm{VC}(\mathcal{A}) = d < \infty$, then for $n > d$, $s(\mathcal{A}, n) \le (n+1)^d$
  - another form: $S_{\mathcal{F}}(n) \le \sum_{k=1}^{\mathrm{VC}_{\mathcal{F}}}\binom{n}{k}$
- Rademacher complexity: $\mathcal{R}(\mathcal{F}) = \mathbb{E}_\epsilon \mathbb{E}_X\left[\sup_{f \in \mathcal{F}}\left|\frac{1}{n}\sum_{i=1}^n \epsilon_i f(X_i)\right|\right]$
  - $\Delta(\mathcal{F}) = \sup_{f \in \mathcal{F}}\left|\frac{1}{n}\sum_{i=1}^n f(X_i) - \mathbb{E}[f]\right|$
  - Rademacher Theorem: $\mathbb{E}[\Delta(\mathcal{F})] \le 2\mathcal{R}(\mathcal{F})$
  - finite class bound: if $|\mathcal{F}| = N$, $\forall f_i \in \mathcal{F}$ we have $\|f_i\|_\infty \le b$, then $\mathcal{R}(\mathcal{F}) \le 2b\sqrt{\frac{\log(2N)}{n}}$

## Sufficiency

- Sufficient statistics: $T(X_1, \ldots, X_n)$ is sufficient for $\theta$ if $p(X^n \mid T = t; \theta)$ does not depend on $\theta$
  - factorization: $T(X_1, \ldots, X_n)$ is sufficient for $\theta$ iff $p(X_1, \ldots, X_n; \theta) = h(X_1, \ldots, X_n)g(T; \theta)$
- Minimal sufficient statistics $\Rightarrow$ for RVs $\{x_n\}$ and $\{y_n\}$, $T(X^n) = T(Y^n) \Leftrightarrow \frac{p(Y^n; \theta)}{p(X^n; \theta)}$ does not depend on $\theta$
- Rao-Blackwell Theorem: $\hat{\theta}$ an estimator, $T$ SS, $\tilde{\theta} = \mathbb{E}\left[\hat{\theta} \mid T\right]$, then $R(\tilde{\theta}, \theta) \le R(\hat{\theta}, \theta)$ under squared loss

## Parameter Estimation

- Method of moments (MOM) : solve for $\theta$ such that $\mathbb{E}\left[X^j\right] = \frac{1}{n}\sum_{i=1}^n X_i^j$
- MLE: $\hat{\theta}$ maximize $\mathcal{L}(X^n; \theta)$
  - equivariance: if we replace $\theta$ by $\eta = g(\theta)$, the MLE estimators satisfy $\hat{\eta} = g(\hat{\theta})$ and $\mathcal{L}(\hat{\eta}) = \mathcal{L}(\hat{\theta})$
    * profile likelihood: if $g$ not invertible, let $\mathcal{L}^*(\eta) = \sup_{\theta: g(\theta) = \eta} \mathcal{L}(\theta)$, then $\hat{\eta} = g(\hat{\theta})$ and $\mathcal{L}^*(\eta) = \mathcal{L}(\hat{\theta})$
  - MLE, ERM and KL: $R_n(\theta, \hat{\theta}) = \frac{1}{n}\sum_{i=1}^n \log\frac{p(X_i; \theta)}{p(X_i; \hat{\theta})} = $ constant $- $ MLE
    * population risk is $R(\theta, \hat{\theta}) = \mathbb{E}_{X; \theta}\log\frac{p(X; \theta)}{p(X; \hat{\theta})} = \mathrm{KL}\left(p(X; \theta)\|p(X; \hat{\theta})\right)$
  - conditions for MLE to be consistent
    * strong identifiability: $\inf_{\tilde{\theta}: |\tilde{\theta} - \theta| \ge \varepsilon} \mathrm{KL}\left(p(X; \theta)\|p(X; \hat{\theta})\right) > 0$
    * uniform LLN: $\sup_{\tilde{\theta}}\left|R_n(\tilde{\theta}, \theta) - R(\tilde{\theta}, \theta)\right| \xrightarrow{p} 0$

- MLE asymptotics: $\sqrt{n}\left(\hat{\theta}-\theta\right) \xrightarrow{d} \mathcal{N}\left(0, I_1^{-1}(\theta)\right)$, or $\sqrt{n}\left(\hat{\tau}-\tau\right) \rightsquigarrow \mathcal{N}\left(0, (g')^T I_1^{-1} g'\right)$ where $\tau = g(\theta)$, conditions:
  * 1) $\theta$ is identifiable. 2) $p(X; \theta)$ is thrice differentiable function of $\theta$; 3) The range of $X$ does not depend on $\theta \Leftarrow$ interchange diff w.r.t $\theta$ and int over $X$. 4) $\theta$ is in the interior of $\Theta$. 5) dimension space does not change with $n$
- Influence functions: $\psi(x) = \frac{\nabla_\theta \log p(x; \theta)}{I(\theta)}$, $\hat{\theta} = \theta + \frac{1}{n}\sum_{i=1}^n \psi(X_i) + $ Remainder. Robust if $\psi$ is bounded
  - asymptotically linear estimators: satisfy $\hat{\theta} \approx \theta + \frac{1}{n}\sum_{i=1}^n \psi(X_i)$
  - any sufficiently well-behaved (regular) estimator is asymptotically linear
- Asymptotic Relative Efficiency (ARE): two estimators $W_n$ and $V_n$ estimating $\tau(\theta)$ where $\sqrt{n}(W_n - \tau(\theta)) \rightsquigarrow \mathcal{N}(0, \sigma_w^2)$, $\sqrt{n}(V_n - \tau(\theta)) \rightsquigarrow \mathcal{N}(0, \sigma_v^2)$. Then $\text{ARE}(V_n, W_n) = \sigma_w^2/\sigma_v^2$. In general, MLE estimator $\hat{\theta}$ satisfies $\text{ARE}\left(\tilde{\theta}, \hat{\theta}\right) \leq 1$ for any other estimator $\tilde{\theta}$.

## Decision Theory
- Fisher information
  - score function: $s(\theta) = \sum_{i=1}^n \nabla_\theta \log p(X_i; \theta)$. data dependent, but $\mathbb{E}_{X^n; \theta}[s(\theta)] = 0$
  - fisher information: $I(\theta) = \mathbb{E}\left[s(\theta) s^T(\theta)\right] = \text{cov}(s(\theta))$. data independent
  - single sample: $I_1(\theta) = \mathbb{E}\left[-\nabla_\theta^2 \log p(X; \theta)\right]$; $n$ samples: $I(\theta) = nI_1(\theta)$
  - Cramér-Rao bound: $\text{Var}\left(\hat{\theta}\right) \geq 1/nI_1(\theta)$. Multivariate: $\text{Var}\left(\hat{\theta}\right) \succeq \frac{1}{n} I_1^{-1}(\theta)$
  - efficient estimator: unbiased & achieve Cramér-Rao bound
- Risk w.r.t loss $L\left(\theta, \hat{\theta}\right)$ is $R\left(\theta, \hat{\theta}\right) = \mathbb{E}_{X^n; \theta}\left[L\left(\theta, \hat{\theta}\right)\right]$
  - MSE: $L$ is squared error. $\text{MSE} = \mathbb{E}_\theta\left[\left(\hat{\theta}-\theta\right)^2\right] = \left(\mathbb{E}_\theta\left[\hat{\theta}\right] - \theta\right)^2 + \text{Var}_\theta\left(\hat{\theta}\right)$
  - minimax risk: $R_n = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R\left(\theta, \hat{\theta}\right) = \inf_{\hat{\theta}} \sup_\theta \mathbb{E}_{X^n; \theta}\left[L\left(\theta, \hat{\theta}\right)\right]$
- Minimax estimator: $\hat{\theta}$ that minimize the worst case $R\left(\hat{\theta}, \cdot\right)$: $\sup_\theta R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_\theta R(\hat{\theta}, \theta)$
  - maximum risk: $\overline{R}\left(\hat{\theta}\right) = \sup_{\theta \in \Theta} R\left(\theta, \hat{\theta}\right) = \sup_{\theta \in \Theta} \mathbb{E}_{X^n; \theta}\left[L\left(\theta, \hat{\theta}\right)\right]$
  - if Bayes estimator $\hat{\theta}$ for prior distribution $\pi$ satisfies $\overline{R}\left(\hat{\theta}\right) \leq B_\pi\left(\hat{\theta}\right)$, then $\hat{\theta}$ is minimax. $\pi$ is least favorable prior
    * Corollary: if risk $R\left(\theta, \hat{\theta}\right)$ is constant for a Bayes estimator $\hat{\theta}$, then $\hat{\theta}$ is also minimax
  - if $L\left(\theta, \hat{\theta}\right) = l\left(\theta - \hat{\theta}\right)$ where $l$ is convex and bounded, symmetric about the origin, $X \sim \mathcal{N}(\theta, \Sigma)$, then $X$ is the unique minimax estimator of $\theta$
  - MLE estimator is approximately minimax as $n \to \infty$
- Bayes estimator (natural estimator): $\hat{\theta}$ that minimizes Bayes risk
  - Bayes risk: $B_\pi\left(\hat{\theta}\right) = \mathbb{E}_\pi\left[R\left(\theta, \hat{\theta}\right)\right] = \mathbb{E}_\pi \mathbb{E}_{X^n; \theta}\left[L\left(\theta, \hat{\theta}\right)\right] \leq \overline{R}\left(\hat{\theta}\right)$
  - marginal: $m(X^n) = \int p(X^n | \theta) \pi(\theta) \, d\theta = \frac{p(\theta, X^n)}{\pi(\theta | X^n)}$ (normalizer in Bayes)
  - Bayes estimator $\hat{\theta}_B$ minimizes posterior risk: $r\left(\hat{\theta} | X^n\right) = \mathbb{E}_{\theta | X^n}\left[L\left(\theta, \hat{\theta}\right)\right]$
    * squared loss: posterior mean $\mathbb{E}[\theta | X^n = x^n]$; absolute loss: median; 0/1 loss: mode
    * under squared loss, $r\left(\hat{\theta} | X^n\right) = \text{Var}(\theta | X^n)$
  - Bayes risk equivalence: $B_\pi\left(\hat{\theta}\right) = \int r\left(\hat{\theta} | X^n\right) m(X^n) \, dX^n$

## Exponential family: $p(X; \theta) = \exp\left[\sum_{i=1}^s \eta_i(\theta) T_i(x) - A(\theta)\right] h(x)$
- Canonical parametrization: $\eta_i(\theta) = \theta_i$. $A(\theta) = \log\left[\int_X \exp\left[\sum_{i=1}^s \theta_i T_i(x)\right] h(x) \, dx\right]$
- Log partition: $\frac{\partial A(\theta)}{\partial \theta_i} = \mathbb{E}[T_i(X)]$, $\frac{\partial^2 A(\theta)}{\partial \theta_i \partial \theta_j} = \text{cov}(T_i(X), T_j(X))$. Therefore, $A$ is convex
- Sufficient statistics: $T(X^n) = \left(\sum_{i=1}^n T_1(X_i), \ldots, \sum_{i=1}^n T_s(X_i)\right)$
- Concave log-likelihood: $\mathcal{LL}(\theta; X^n) \propto \sum_{i=1}^s \theta_i \sum_{j=1}^n T_i(x_j) - nA(\theta)$, $\mathcal{LL} \propto \langle \theta, T \rangle - nA(\theta)$
- Minimal representation: the $T_i$s are linearly independent
  - over-complete exponential families are not statistically identifiable

## Hypothesis testing: Null hypothesis: $H_0 : \theta \in \Theta_0$. Alternative hypothesis: $H_1 : \theta \in \Theta_1$
- TN = retain True; FP = reject True (type 1); FN = retain False (type 2); TP = reject False
  - type 1: the incorrect rejection of a true $H_0$
  - type 2: the failure to reject a false $H_0$
- Statistical significance (from Wikipedia)
  - definition: unlikely to have occurred under $H_0$, when $p < \alpha$
  - significance level $\alpha$: $\Pr[\text{reject } H_0 | H_0 \text{ is true}] = \Pr[\text{type 1}]$

- Power (prob of rejection $\theta$): should be large for $\theta \in \Theta_1$, small for $\theta \in \Theta_0$
  - probability of rejection when true parameter is $\theta$ (Type 1 error): $\beta(\theta) = P_\theta(X^n \in R)$
  - always favor the null hypothesis and only consider tests that control the Type-I error.
  - size $s = \sup_{\theta \in \Theta_0} \beta(\theta)$. level $\alpha$: $s \leq \alpha$
  - power of a test: $\Pr(\text{reject } H_0 | H_1 \text{ is true}) = 1 - \Pr[\text{type 2}]$
- $p$-value: defined for data $x$, $H_0 : \theta \in \Theta_0$, smallest $\alpha$ at which we would reject $H_0$
  - definition: $p = \inf\{\alpha : T(x) \in R_\alpha\} = \sup_{\theta \in \Theta_0} P_\theta(T(X) \geq T(x))$
  - $p \sim \text{Unif}(0, 1)$ under $\Theta_0$
  - a small p-value indicates strong evidence against $H_0$
- Neyman-Pearson Test ($H_0 : \theta = \theta_0$ vs $H_1 : \theta = \theta_1$)
  - $T = \frac{\mathcal{L}(\theta_1)}{\mathcal{L}(\theta_0)}$, $R_\alpha = \{T > k_\alpha\}$; set $k \in \{t : P_{\theta_0}(T > t) \leq \alpha\}$ to obtain tests of level $\alpha$.
- Wald Test ($H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$): MLE $\hat{\theta}$, $T = \frac{\hat{\theta}-\theta}{\hat{\sigma_0}}$, reject $|T_n| \geq z_{\alpha/2}$
  - Assumption: asymptotically normal estimator $\hat{\theta}$ that satisfies $\hat{\theta} \xrightarrow{d} \mathcal{N}(\theta_0, \sigma_0^2)$ under $H_0$. Can replace $\sigma_0$ with its plug-in estimation $\text{se}\left(\hat{\theta}\right)$
  - For MLE $\hat{\theta}$: $T = \sqrt{nI_1(\theta_0)}\left(\hat{\theta} - \theta_0\right)$, $R_\alpha = \{|T| \geq z_{\alpha/2}\}$
  - Power: let $\Delta = \sqrt{nI_1(\theta_0)}(\theta - \theta_0)$, $\beta(\theta) = 1 - \Phi(\Delta + z_{\alpha/2}) + \Phi(\Delta - z_{\alpha/2})$. Large if $|\theta - \theta_0|$ or $n$ large
- Likelihood Ratio Test ($H_0 : \theta \in \Theta_0$, $H_1 : \theta \notin \Theta_0$), $p$-value $= P\left(\chi_{\#\text{param}}^2 > \lambda\right)$
  - $\lambda(X^n) = 2\log \frac{\sup_{\theta \in \Theta} \mathcal{L}(\theta)}{\sup_{\theta \in \Theta_0} \mathcal{L}(\theta)} \sim \chi_d^2$, $R_\alpha = \left\{\chi_{d,\alpha}^2 > \lambda\right\}$, where $d = \dim(\Theta) - \dim(\Theta_0)$
- Goodness-of-fit test:
  - $\chi^2$ test for multinomial distributions
    * one sample testing: $H_0 : \mathbf{p} = \mathbf{p}_0 = (p_{01}, \ldots, p_{0k})$ vs $H_1 : \mathbf{p} \neq \mathbf{p}_0$. $Z^k$ is bin counts
      · $T = \sum_{i=1}^k \frac{(Z_i - np_{0i})^2}{np_{0i}} \sim \chi_{k-1}^2$, $R_\alpha = \{\chi_{k-1,\alpha}^2 > T\}$, $p$-value $= \chi_{k-1,T}^2$
    * two sample testing: $H_0 : \mathbf{p}_x = \mathbf{p}_y$ vs $H_1 : \mathbf{p}_x \neq \mathbf{p}_y$. $Z_x, Z_y$ are bin counts for $x, y$
      · Let $\hat{p}_i = \frac{Z_{xi} + Z_{yi}}{n_x + n_y}$, then $T = \sum_{i=1}^k \frac{(Z_{xi} - n_x\hat{p}_i)^2}{n_x\hat{p}_i} + \frac{(Z_{yi} - n_y\hat{p}_i)^2}{n_y\hat{p}_i} \sim \chi_{k-1}^2$, $R_\alpha = \{\chi_{k-1,\alpha}^2 > T\}$
  - permutation test: observe $\{X_i\}_{i=1}^n \sim P$, $\{Y_j\}_{j=1}^m \sim Q$. We have $H_0 : P = Q$ vs $H_1 : P \neq Q$
    * $T = \frac{1}{N!}\sum_L \mathbb{I}(g(X, Y, L) > g(X, Y, L^*)) \sim \text{Unif}(0, 1)$, $R_\alpha = \{T < \alpha\}$, $p$-value $= \frac{1}{N!}\sum_{i=1}^{N!} \mathbb{I}(T_i > T_{\text{obs}})$

## Multiple testing
- Family-Wise Error Rate (FWER): $P[\text{exist false rej}]$. False Discovery Rate (FDR): $\mathbb{E}[\#\text{false rej}/\#\text{rej}]$
  - FWER $\geq$ FDR, control FWER is controlling FDR; under global null, FDR = FWER
- Multiple Testing: $d$ tests in total
  - Sidak method: reject any test if its $p$-value $\leq 1 - (1-\alpha)^{1/d} = \alpha_t$ ($p$-value independent), then FWER $\leq \alpha$
  - Bonferroni method: $p$-value $\leq \alpha/d$, then FWER $\leq \alpha$
  - Holm's procedure: $i^* = \min\left\{i : p_i > \frac{\alpha}{d-i+1}\right\}$, reject all $H_i$ for $i < i^*$, then FWER $\leq \alpha$
  - BH procedure: $t_i = \frac{i\alpha}{d}$, $i^* = \max\{i : p_i < t_i\}$, reject all $H_i$ for $i \leq i_{\max}$, then FDR $\leq \alpha$

## Confidence set: a random set $C(X, \alpha)$ for $\alpha \in (0, 1)$ that satisfies $P(\theta \in C(X)) \geq 1 - \alpha$
- Probability inequalities: if a bounded $\hat{\theta}$ is an unbiased estimator of $\theta$, we can apply Hoeffding's bound
- Inverting a test: $C(X, \alpha) = \{\theta : X \in A(\theta, \alpha)\}$, where $A$ is the acceptance region of the test with "center" $\theta$
  - Wald Interval: for MLE $\hat{\theta}$ we have $\text{se}\left(\hat{\theta}\right) = \frac{1}{\sqrt{nI_1(\hat{\theta})}}$. Then $C(X, \alpha) = \left(\hat{\theta} - z_{\alpha/2}\text{se}\left(\hat{\theta}\right), \hat{\theta} + z_{\alpha/2}\text{se}\left(\hat{\theta}\right)\right)$
    * Note: can also be applied to another asymptotically normal estimators
    * delta method: $\tau\left(\hat{\theta}_n\right) \pm z_{\alpha/2}\text{se}\left|\tau'\left(\hat{\theta}_n\right)\right|$
  - Likelihood Interval: for MLE $\hat{\theta}$ we have $C(X, \alpha) = \left\{\theta : \frac{\mathcal{L}(\theta)}{\mathcal{L}(\hat{\theta})} > \exp\left(-\frac{1}{2}\chi_{d,\alpha}^2\right)\right\}$, where $d = \dim \Theta$
- Pivots: a function $Q(X, \theta)$ whose distribution does not depend on $\theta$. Therefore, a interval for $\theta$ is $\{\theta : Q(X, \theta) \in C(Q, 1 - \alpha)\}$

## Causal inference:
- Average treatment effect: $\tau = \mathbb{E}(Y(1) - Y(0))$. Association: $\alpha = \mathbb{E}[Y(1) | W = 1] - \mathbb{E}[Y(0) | W = 0]$
  - $\alpha$ is easy to estimate: we have $\hat{\alpha} = \frac{1}{m}\sum_{i:W_i=1} Y^{obs} - \frac{1}{n-m}\sum_{i:W_i=0} Y^{obs} = \sum_{i=1}^n \frac{W_i Y_i(1)}{m} - \frac{(1-W_i)Y_i(0)}{n-m}$
  - If $W \perp (Y(0), Y(1))$, then $\alpha = \tau$ and $\hat{\alpha} = \hat{\tau}$ is an unbiased estimator for $\tau$. Otherwise, we suffer from selection bias
- Exact p-values: we test hypothesis $H_0 : \mathbb{E}[Y(1)] = \mathbb{E}[Y(0)]$ vs $H_1 : \mathbb{E}[Y(1)] \neq \mathbb{E}[Y(0)]$
  - define $L^*$ as the original treatment assignment and $\mathcal{L}$ as the set of possible assignments

– randomly assign treatments. For each assignment $L$ we calculate $\hat{\tau}_L = \frac{1}{m}\sum_{i:W_{Li}=1} Y^{obs} - \frac{1}{n-m}\sum_{i:W_{Li}=0} Y^{obs}$

– we then have $p = \frac{1}{|\mathcal{L}|}\sum_L \mathbb{I}_{\{|\hat{\tau}_L|>|\hat{\tau}_{L^*}|\}}$

- No unmeasured confounding: we have $W \perp (Y(0), Y(1))|X$, where $X$ is the confounding variable
  – $\tau = \mathbb{E}_X[\mathbb{E}(Y^{obs}|X, W=1)] - \mathbb{E}_X[\mathbb{E}(Y^{obs}|X, W=0)] = \mathbb{E}[\mu_1(X)] - \mathbb{E}[\mu_0(X)]$
  – Suppose we can estimate $\mu_1$ and $\mu_0$ by $\hat{\mu}_1$ and $\hat{\mu}_0$. Then we directly have the plug-in estimator $\hat{\tau} = \frac{1}{n}\sum_{i=1}^n[\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)]$
  – Alternative estimator (Horvitz-Thompson): $\hat{\tau} = \frac{1}{n}\sum_{i=1}^n[\frac{Y_i^{obs}W_i}{\pi(X_i)} - \frac{Y_i^{obs}(1-W_i)}{1-\pi(X_i)}]$, where $\pi(x) = \mathbb{E}(W|X=x)$ is the propensity score

**Regression**: estimate $r(x) = \mathbb{E}[Y|X=x]$, risk of $\hat{r}$: $R(\hat{r}) = \mathbb{E}[\mathcal{L}(Y, \hat{r}(X))]$
- If joint distribution of $(X, Y)$ is known, under square, risk is minimized by $\hat{r}(x) = \mathbb{E}[Y|X=x]$
- Kernel regression: $r(x) = \sum_{i=1}^n w_i(x)Y_i$, where $w_i = K\left(\frac{x-x_i}{h}\right)/\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)$
  – bandwidth $h$: larger $h$, smoother regression function, more bias, less variance
  – Gaussian kernel: $K(x) = \frac{1}{\sqrt{2\pi}}\exp\left(-x^2/2\right)$
  – simple analysis
    * assumption: 1) $y_i = r(x_i) + \epsilon_i$; 2) $x_i$ is 1-dimensional, equally spaced in $[0, 1]$; 3) $r(x) = \mathbb{E}[Y|X=x]$ is $L$-Lipschitz $\left|\frac{d}{dx}r(x)\right| \leq L$; 4) noise i.i.d: $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$; 5) use spherical kernel $k(x) = \mathbb{I}(-1 \leq x \leq 1)$
    * if bandwidth $h \geq 1/n$, then bias $b(x) \leq Lh$, variance $v(x) \leq \frac{\sigma^2}{nh}$.
    * let bandwidth $h = \left(\frac{\sigma^2}{2nL^2}\right)^{1/3}$, then $\hat{R}(\hat{r}, r) \leq 2\left(\frac{L\sigma^2}{n}\right)^{2/3}$
  – Extension: if $r$ is $\beta$-smooth, then $b(x) \approx h^{2\beta}$, $v(x) \approx \frac{1}{nh^d}$. Consequently, $R(\hat{r}, r) \approx n^{-2\beta/(2\beta+d)}$
    * the curse of dimensionality: the rate gets exponentially slow as $d$ increases (we only get linearly slow in parametric regression).
    * get rid: smoothness / parametric (e.g. linear) / sparsity assumption
- Gaussian sequence model: observe $y^n$ where $y_i = \theta_i + \epsilon_i$, $\epsilon_i \sim \mathcal{N}\left(0, \frac{\sigma^2}{n}\right)$
  – minimax estimator under squared loss: $\hat{\theta} = [y_1, \ldots, y_d]^T$, $l_2$ risk: $R(\hat{\theta}, \theta) = \frac{\sigma^2 d}{n}$. Not consistent when $d \gg n$
  – hard thresholding $t$: $\hat{\theta}_i = y_i\mathbb{I}(|y_i| \geq t)$. Solution to: $\hat{\theta} = \arg\min_\theta \frac{1}{2}\|y - \theta\|_2^2 + \frac{t^2}{2}\sum_{i=1}^d \mathbb{I}(\theta_i \neq 0)$
  – soft thresholding $t$: $\hat{\theta}_i = \text{sign}(y_i)\max(|y_i| - t, 0)$. Solution to: $\hat{\theta} = \arg\min_\theta \frac{1}{2}\|y - \theta\|_2^2 + t\sum_{i=1}^d |\theta_i|$
  – analyzing hard thresholding
    * maximum of Gaussians: if $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, then w.p. $\geq 1 - \delta$, $\max_{i=1}^d |\epsilon_i| \leq \sigma\sqrt{2\log(2d/\delta)}$
    * let threshold $t = 2\sigma\sqrt{2\log(2d/\delta)/n}$, then w.p. $\geq 1 - \delta$, $\left\|\hat{\theta} - \theta\right\|_2^2 \leq 9\sum_{i=1}^d \min\left(\theta_i^2, \frac{t^2}{4}\right)$
    * risk upper bound: $R(\hat{\theta}, \theta) \lesssim \sum_{i=1}^d \min\left(\theta_i^2, \frac{\sigma^2\log d}{n}\right)$
      · worst case: $R(\hat{\theta}, \theta) \lesssim \frac{\sigma^2 d\log d}{n}$, similar to minimax estimator
      · $s$ sparse: $R(\hat{\theta}, \theta) \lesssim \frac{\sigma^2 s\log d}{n}$, consistent
      · $l_1$ sparse: if $\sum_{i=1}^d |\theta| \leq R$, then $R(\hat{\theta}, \theta) \lesssim 2R\sigma\sqrt{\frac{\log d}{n}}$
- Linear regression: observe $\{(x_i, y_i)\}_{i=1}^n$, model $y_i = \langle x_i, \beta^*\rangle + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, $\hat{\Sigma} = \frac{1}{n}\sum_{i=1}^n x_i x_i^T$
  – $\hat{\beta} = \arg\min_\beta \frac{1}{2}\sum_{i=1}^n (y_i - \langle x_i, \beta\rangle)^2 = \hat{\Sigma}^{-1}X^T y = (X^TX)^{-1}X^T y \sim \mathcal{N}\left(\beta^*, \sigma^2(X^TX)^{-1}\right)$
    * consider $X \in \mathbb{R}^{n\times d}$ as random variables: random design matrix
    * consider $X \in \mathbb{R}^{n\times d}$ as fixed: fixed design matrix
  – in-sample prediction error $\mathbb{E}\left[\left\|X\hat{\beta} - X\beta\right\|_2^2/n\right]$
    * $X\hat{\beta} \sim \mathcal{N}\left(X\beta^*, \sigma^2 X(X^TX)^{-1}X^T\right)$
    * $\mathbb{E}\left[\left\|X\hat{\beta} - X\beta\right\|_2^2/n\right] = \sigma^2\mathbb{E}\left[\text{tr}\left(X(X^TX)^{-1}X^T\right)\right]/n = \sigma^2 d/n$
  – $l_2$ error $\mathbb{E}\left[\left\|\hat{\beta} - \beta\right\|_2^2\right] = \frac{\sigma^2}{n}\text{tr}\left(\hat{\Sigma}^{-1}\right)$
    * if eigenvalues of $\hat{\Sigma}^{-1}$ are lower bounded by $c$, then $l_2$ error $\leq \frac{\sigma^2 d}{cn}$
    * $I_n(\beta) = n\mathbb{E}\left[\frac{X^TX}{n\sigma^2}\right] = \frac{n\Sigma}{\sigma^2}$
    * $\sqrt{n}(\hat{\beta} - \beta^*) \xrightarrow{d} \mathcal{N}(0, \sigma^2\Sigma^{-1})$

– issues when $d$ increases: 1) $\hat{\Sigma}$ not invertible; 2) error dominated by $d$; 3) too many solutions of $y = X\beta$ in high dimension
- High dimensional regression
  – hard thresholding type estimator
    * penalized form: $\hat{\beta} = \arg\min_\beta \frac{1}{2}\|y - X\beta\|_2^2 + \frac{t^2}{2}\sum_{i=1}^d \mathbb{I}(\beta_i \neq 0)$
    * constraint form (best subset regression): $\hat{\beta} = \arg\min_\beta \frac{1}{2}\|y - X\beta\|_2^2$, subject to $\sum_{i=1}^d \mathbb{I}(\beta_i \neq 0) \leq k$
  – soft thresholding type estimator (LASSO: least absolute selection and shrinkage operator)
    * penalized form: $\hat{\beta} = \arg\min_\beta \frac{1}{2}\|y - X\beta\|_2^2 + t\sum_{i=1}^d |\beta_i|$
    * constraint form: $\hat{\beta} = \arg\min_\beta \frac{1}{2}\|y - X\beta\|_2^2$, subject to $\sum_{i=1}^d |\beta_i| \leq k$
    * prediction error, w.p. $\geq 1 - \delta$, $\frac{1}{n}\left\|X\hat{\beta} - X\beta\right\|_2^2 \leq 4\sigma\|\beta^*\|_1\sqrt{\frac{2\log 2d/\delta}{n}}$

**Bayesian inference**
- Bayesian vs Frequentist
  – probability: subjective degree of belief; limiting frequency
  – goal: analyze belief; create procedures with frequency guarantee
  – $\theta$: random variable; fixed
  – $X$: random variable; random variable
  – use Bayes: Yes, to update beliefs; Yes if it leads to procedure with good frequentist behavior, otherwise no
- Setup: given a collection of distributions $\{P_\theta : \theta \in \Theta\}$; prior: $\theta \sim \pi$; likelihood: $p(X^n|\theta) \sim P_\theta$
  – compute posterior belief: $\pi(\theta|X) = p(X|\theta)\pi(\theta) / \int_\theta p(X|\theta)\pi(\theta) \propto p(X|\theta)\pi(\theta)$
  – Bayesian: success if you can calculate posterior
  – frequentist: success if the posterior concentrates on the true parameter $\theta^*$
- Credible set $C_\alpha$ for posterior distribution: $\int_{C_\alpha}\pi(\theta|X^n)d\theta = 1 - \alpha$
- Frequentist view
  – Consistency: $\forall\epsilon > 0$, $\pi\left(\{\theta : \|\theta - \theta^*\| \geq \epsilon\}\middle|X^n\right) \to 0$
  – Convergence rate (based on consistency): for $\forall\delta > 0$, we define its convergence rate $\epsilon = \sup\{\epsilon' : \pi_\theta(\|\theta - \theta^*\| \geq \epsilon'|X) \leq \delta\}$
- Bernstein-von Mises Theorem: if prior is continuous and strictly positive around $\theta^*$, then for a fixed dimension $d$, as $n \to \infty$, the posterior is close to Gaussian:
  – $\left\|\pi(\theta; X^n) - \mathcal{N}\left(\hat{\theta}_n, 1/nI_1(\hat{\theta}_n)\right)\right\|_{TV} \to 0$, where $\hat{\theta}_n$ is MLE, and TV is total variance
  – in high-dimensional or non-parametric setting, strictly positive around $\theta^*$ is hard to be true
  – if it is true, can use credible sets like Wald test
- Ideas in choosing prior: 1) by convenience (minimax); 2) doesn't matter (low-dimensional / parametric); 3) based on data: empirical Bayes; 4) non-informative prior; 5) Joffrey's prior: $\pi(\theta) \sim \sqrt{I(\theta)}$: invariant under transformations; 6) hierarchical prior
- Regularizer: 1) LASSO: posterior mode with Laplace prior; 2) ridge regression: Gaussian prio; 3) estimating Bernoulli prob with laplace smoothing: posterior mean with Beta prior

**Monte Carlo**
- MC Integration: compute $\mu = \mathbb{E}_{X\sim P}[f(X)]$
  – directly sampling from $P$: $\hat{\mu} = \frac{1}{n}\sum_{i=1}^n f(X_i)$. If cannot sample from $P$
  – importance weighting: $\mu = \mathbb{E}_{X\sim Q}\left[\frac{p(X)}{q(X)}f(X)\right] = \mathbb{E}_{X\sim Q}[w(X)f(X)]$, $\hat{\mu} = \frac{1}{n}\sum_{i=1}^n w(X_i)f(X_i)$
  – if doesn't work, use Markov Chain Monte Carlo
- Markov chain: transition probability $T(x_i, x_{i+1}) = P(x_{i+1}|x_i)$
  – limiting distribution always exists: $\pi(x) = P(\lim_{n\to\infty}x_n = x)$
  – detailed balance: condition for reaching limiting distribution: $\forall x, y, \pi(x)T(x, y) = \pi(y)T(y, x)$
  – Markov LLN: MC $\{X^n\}$, limiting distribution $\pi$, then $\frac{1}{n}\sum_{i=1}^n f(X_i) \to \mathbb{E}_\pi[f(x)]$ (weak dependence)
- Metropolis-Hastings: at step $i$, $y \sim q(y|X=x_i)$, accept w.p. $r = \min\left\{1, \frac{f(y)q(x|y)}{f(x)q(y|x)}\right\}$
  – proposal distribution $q$, often $q(y|x_i = x) \sim \mathcal{N}(x, \sigma^2)$
  – if $q(\cdot|\cdot)$ is symmetric, we sample more from high-prob regions, and accept w.p. 1 if prob going up
  – limiting distribution of this Markov chain is $f$: $T(y, x) = q(y|x)r$, verify $f(x)T(y, x) = f(y)T(x, y)$

**Bootstrap**: test the variability of a point estimate (variance, confidence set)
- Given $X^n \sim P$, estimate empirical distribution $P_n(A) = \frac{1}{n}\sum_{i=1}^n \mathbb{I}(X_i \in A)$, then draw from $P_n$
- Bootstrap variance estimate
  – draw bootstrap sample $X^{*n} \sim P_n$, compute $\hat{\theta}_n^* = g(X^{*n})$
  – repeat $B$ times, yielding $\hat{\theta}_{n, 1}^*, \hat{\theta}_{n, 2}^*, \ldots, \hat{\theta}_{n, B}^*$, $\overline{\theta} = \frac{1}{B}\sum_{j=1}^B \hat{\theta}_{n, j}^*$
  – bootstrap variance $\hat{s}^2 = \frac{1}{B}\sum_{j=1}^B \left(\hat{\theta}_{n, j}^* - \overline{\theta}\right)^2$
- Bootstrap confidence interval

- draw bootstrap sample $X^{*n} \sim P_n$, compute $\hat{\theta}_n^* = g(X^{*n})$
- repeat $B$ times, yielding $\hat{\theta}_{n,1}^*, \hat{\theta}_{n,2}^*, ..., \hat{\theta}_{n,B}^*$
- calculate Bootstrap CDF: $\hat{G}(t) = \frac{1}{B} \sum_{j=1}^{B} \mathbb{I}\left(\sqrt{n}\left(\hat{\theta}_{n,j}^* - \hat{\theta}_n\right) \leq t\right)$
- $C_n = \left[\hat{\theta}_n - \frac{g_{1-\alpha/2}}{\sqrt{n}}, \hat{\theta}_n - \frac{g_{\alpha/2}}{\sqrt{n}}\right]$ where $g_{\alpha/2} = \hat{G}^{-1}(\alpha/2)$, $g_{1-\alpha/2} = \hat{G}^{-1}(1-\alpha/2)$

- Bootstrap theorem: $F_n(t) = P\left(\sqrt{n}\left(\hat{\theta}_n - \theta\right) \leq t\right)$, $\hat{F}_n(t) = P\left(\sqrt{n}\left(\hat{\theta}_n^* - \hat{\theta}_n\right) \leq t\right)$, where $\hat{\theta}_n$ is the estimator based on $X^n$.
  - Example: if $\mu_3 = E|X_i|^3 < \infty$, and $\hat{\theta}_n$ is sample mean, then $\sup_t \left|\hat{F}_n(t) - F_n(t)\right| = O_P\left(\frac{1}{\sqrt{n}}\right)$

**Model selection**
- Definition: A scheme $S$ is model selection consistent: as sample size $n \to \infty$, $P(S \text{ selects wrong model}) \to 0$
- Cross validation:
  - prediction consistent: as test size $n_{\text{te}} \to \infty$, risk estimation converges to its expectation.
  - not model selection consistent
- Akaike Information Criterion (AIC): define for model space $\Theta$
  - $\text{AIC}(\Theta) = 2l(\hat{\theta}) - 2\dim\Theta$, where parameter $\hat{\theta}$ is the MLE within $\Theta$ and $l\left(\hat{\theta}\right) = \log P\left(X \mid \hat{\theta}\right)$
  - not model selection consistent
- Bayesian Information Criterion (BIC): define for model space $\Theta$ with sample size $n$.
  - $\text{BIC}(\Theta) = 2l\left(\hat{\theta}\right) - \dim\Theta \log n$, where parameter $\hat{\theta}$ is the MLE within $\Theta$ and $l\left(\hat{\theta}\right) = \log P(X \mid \hat{\theta})$
  - compared to AIC, prefer sparser / simpler models
  - model selection consistent
- Application of AIC and BIC: compare among different model spaces and select the one that maximizes AIC / BIC.
- Main take-away
  - if goal is "prediction", and have reasonable sample size & computational budget $\Rightarrow$ cross validation
  - if goal is "prediction", less samples / computational budget $\Rightarrow$ AIC
  - if goal is "selecting true model" $\Rightarrow$ BIC

| | $\mathcal{N}(\mu, \Sigma)$ | Gamma$(\alpha, \beta)$ | Exp$(\lambda)$ | Poisson$(\lambda)$ |
|---|---|---|---|---|
| support | $\mu + \text{span}(\Sigma)$ | $(0, +\infty)$ | $(0, +\infty)$ | $k \in \mathbb{N}^+ \cup \{0\}$ |
| pdf | $\frac{\exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)}{(2\pi)^{d/2}|\det(\Sigma)|^{1/2}}$ | $\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$ | $\lambda e^{-\lambda x}$ | $\frac{\lambda^k \exp(-\lambda)}{k!}$ |
| cdf | - | $\frac{1}{\Gamma(\alpha)}\gamma(\alpha, \beta x)$ | $1 - e^{-\lambda x}$ | $e^{-\lambda}\sum_{i=0}^{\lfloor k \rfloor} \frac{\lambda^i}{i!}$ |
| mean | $\mu$ | $\alpha/\beta$ | $\lambda^{-1}$ | $\lambda$ |
| variance | $\Sigma$ | $\alpha/\beta^2$ | $\lambda^{-2}$ | $\lambda$ |
| mgf | $\exp\left(\mu^T t + \frac{1}{2} t^T \Sigma t\right)$ | $(1-t/\beta)^{-\alpha}$ for $t < \beta$ | $\frac{\lambda}{\lambda - t}$ for $t < \lambda$ | $\exp(\lambda(z-1))$ |
| conv | $\mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$ | Gamma$(\alpha_1 + \alpha_2, \beta)$ | - | Poisson$(\lambda_1 + \lambda_2)$ |
| fisher | $\begin{bmatrix} 1/\sigma^2 & \\ & 1/2\sigma^4 \end{bmatrix}$ | - | $\lambda^{-2}$ | $\lambda^{-1}$ |
| MoM | $\hat{\mu} = \overline{X}_n,\ \sigma^2 = \frac{1}{n}\sum_{i=1}^n (X - \overline{X}_n)^2$ | | | $\hat{\lambda} = \overline{X}_n$ |
| MLE | same | | | same |

| | Ber$(p)$ | Rademacher | Binom$(n, p)$ | Geometric$(p)$ | $\chi^2(k)$ |
|---|---|---|---|---|---|
| support | $\{0, 1\}$ | $\{-1, 1\}$ | $k \in \mathbb{N} \cup \{0\}$ | $\mathbb{N}^+$ | $x \in (0, \infty)$ |
| pdf | - | - | $\binom{n}{k} p^k (1-p)^{n-k}$ | $(1-p)^{k-1} p$ | $\frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-\frac{x}{2}}$ |
| cdf | - | - | - | $1 - (1-p)^k$ | - |
| mean | $p$ | $0$ | $np$ | $1/p$ | $k$ |
| variance | $p(1-p)$ | $1$ | $np(1-p)$ | $(1-p)/p^2$ | $2k$ |
| mgf | $q + pe^t$ | $\exp(t^2/2)$ | $(1 - p + pe^t)^n$ | $\frac{pe^t}{1-(1-p)\exp(t)}$ | $(1-2t)^{-k/2}$ for $t < \frac{1}{2}$ |
| conv | - | - | Binom$(n+m, p)$ | - | - |
| fisher | $1/p(1-p)$ | - | $n/p(1-p)$ | - | - |
| MoM | | | $\hat{p} = \overline{X}_n/n$ | $\hat{p} = 1/\overline{X}_n$ | |
| MLE | | | same | same | |

| | Beta$(\alpha, \beta)$ |
|---|---|
| support | $(0, 1)$ |
| pdf | $\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} x^{\alpha-1} (1-x)^{\beta-1}$ |
| cdf | - |
| mean | $\frac{\alpha}{\alpha+\beta}$ |
| variance | $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ |
| mgf | $1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r}\right) \frac{t^k}{k!}$ |
| conv | - |
| fisher | $\begin{bmatrix} \text{Var}[\ln X] & \text{Cov}[\ln X, \ln(1-X)] \\ \text{Cov}[\ln X, \ln(1-X)] & \text{Var}[\ln(1-X)] \end{bmatrix}$ |
| MoM | |
| MLE | |

Table 1: Probability distribution