

<p>DEFINITION</p> <p><i>How do you compute the MLE estimator for parameter <math>\theta</math> of random variable with pdf <math>p_\theta(x)</math>?</i></p>	<p>DEFINITION</p> <p><i>How do you compute the method of moments estimator for parameter of interest <math>\gamma</math> of a random variable with probability density function <math>p_\theta(x)</math>?</i></p>
<p><i>How do you compute the <math>n^{th}</math> moment of a random variable <math>X</math>?</i></p>	<p><i>Let <math>X_1, \dots, X_n</math> be random variables with pdf</i></p> $p_\theta = \begin{cases} \frac{1}{\theta} & , \text{ if } 0 \leq x \leq \theta \\ 0 & \text{ otherwise} \end{cases}$ <p><i>What can you say about its likelihood function <math>L_X(\theta)</math>? What is the MLE of <math>\theta</math>?</i></p>
<p><i>What is the likelihood function <math>L_X(\theta)</math> for i.i.d. random variables <math>X_1, \dots, X_n</math> ?</i></p>	<p><i>Suppose <math>X_1, \dots, X_n</math> are independent random variables with pdf <math>p_\theta(x)</math>. What is the likelihood <math>L_X(\theta)</math> ?</i></p>
<p>THEOREM</p> <p><i>Law of total probability</i></p>	<p>DEFINITION</p> <p><i><math>\chi_k^2</math> distribution with <math>k</math> degrees of freedom</i></p>
<p>DEFINITION</p> <p><i>What is the expectation <math>\mathbb{E}[\mathbb{E}g(X, Y) \mid Y]</math> ?</i></p>	<p><i>If <math>X \sim \text{Poisson}(\lambda)</math> and <math>Y \sim \text{Poisson}(\mu)</math>, what can you say about the distribution of <math>X + Y</math> ?</i></p>

<p>Let <math>k</math> be the number of parameters to estimate/dimensions of <math>\gamma \in \Gamma</math>. Find a mapping <math>m : \Gamma \rightarrow \mathbb{R}^k</math> that maps possible values of <math>\gamma</math> to the <math>k</math> moments of <math>X</math>. If <math>m</math> is invertible, then plug in the vector of <math>k</math> empirical moments <math>\hat{\mu}</math> into it. Then the method of moments estimate is <math>\hat{\gamma} = m^{-1}(\hat{\mu})</math>.</p>	<p>Need to find the <math>\hat{\theta}</math> which solves the optimization problem <math>\hat{\theta} = \arg \max_{\phi} p_{\phi}(x)</math>.</p>
<p>The likelihood <math>L_X(\theta)</math> is 0 for all <math>\theta &lt; \max \{X_1, \dots, X_n\}</math>, and otherwise it's decreasing. So, MLE of <math>\theta</math> is <math>\max \{X_1, \dots, X_n\}</math></p>	<p>Integrate it as <math>\mathbb{E}(X^k) = \int X^k dP</math></p>
$L_X(\theta) = \prod_{i=1}^n p_{\theta}(x_i)$	<p>It's the probability <math>X_1, \dots, X_n</math> have a given value if the parameter is <math>\theta</math>. Formally <math>L_X(\theta) = \prod_{i=1}^n p_{\theta}(x_i)</math></p>
<p>Let <math>Y_1, \dots, Y_k</math> be i.i.d. <math>\mathcal{N}(0, 1)</math>. Then <math>\sum_{j=1}^k Y_j^2 \sim \chi_k^2</math>.</p>	<p>Let <math>(B_j)_{j \geq 1}</math> be a partition of <math>\Omega</math>, i.e. with</p> $\bigcup_{j \geq 1} B_j = \Omega$ <p>Then</p> $\begin{aligned} \mathbb{P}(A) &= \sum_j \mathbb{P}(A \cap B_j) \\ &= \sum_j \mathbb{P}(A   B_j) \mathbb{P}(B_j) \end{aligned}$
$X + Y \sim \text{Poisson}(\lambda + \mu)$	<p>By the iterated expectation lemma, it's <math>\mathbb{E}[\mathbb{E}g(X, Y)   Y] = \mathbb{E}g(X, Y)</math>.</p>

<p>Let <math>X_1, \dots, X_n</math> be Poisson (<math>\lambda_i</math>) distributed random variables, with a sum <math>Z = \sum_{i=1}^n X_i</math>. What can you say about the distribution <math>(X_1, \dots, X_n \mid Z)</math> ?</p>	<p>Given i.i.d. random variables <math>X_1, X_2</math> with distribution function <math>F_X</math>, what can you say about the distribution of <math>Z = \max \{X_1, X_2\}</math> ? What is its distribution function and density?</p>
<p>DEFINITION</p> <p>Sufficient statistics for distribution with random variable <math>X</math>.</p>	<p>THEOREM</p> <p>Factorization theorem of Neyman</p>
<p>How can you show that the MLE <math>\hat{\theta}_{MLE}</math> only depends on the sufficient statistic <math>S</math>?</p>	<p>DEFINITION</p> <p>Exponential family <math>\{P_\theta : \theta \in \Theta\}</math>.</p>
<p>How can you show that a random variable <math>X \sim \text{Poisson}(\theta), \theta &gt; 0</math> belongs to an exponential family? The pmf of Poisson is <math>\mathbb{P}(X = x) = e^{-\theta} \frac{\theta^x}{x!}</math>.</p>	<p>Show that for independent and identically distributed random variables <math>X_1, \dots, X_n</math> sampled according to some distribution from the exponential family</p> $p_\theta(x) = \exp[\underbrace{c(\theta)^T}_{\in \mathbb{R}^k} \underbrace{T(x)}_{\in \mathbb{R}^k} - d(\theta)]h(x),$ <p>the sum <math>\sum_{i=1}^n T(X_i)</math> is a sufficient statistic for <math>\theta</math>.</p>
<p>How can you show that a normal random variable <math>\mathcal{N}(\mu, \sigma^2)</math> with probability density function</p> $p_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$ <p>for <math>x \in \mathbb{R}</math> belongs to the exponential family?</p>	<p>DEFINITION</p> <p>What is the canonical form of a <math>k</math>-dimensional exponential family?</p>

<p><math>Z</math> has distribution function <math>F(z) = F_X^2(z)</math>, and its density is <math>f(z) = 2F_X(z)f_X(z)</math></p>	<p><math>(X_1, \dots, X_n \mid Z) \sim \text{Multinomial}(Z, p_1, \dots, p_n)</math> with</p> $p_i = \frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$
<p>A statistic <math>S</math> is sufficient if and only if the density function <math>p_\theta(x)</math> can be written as</p> $p_\theta(x) = g_\theta(S(x))h(x)$ <p>or all <math>x, \theta</math> and some functions <math>g_\theta(\cdot) \geq 0</math> and <math>h(\cdot) \geq 0</math>.</p>	<p>The statistic <math>S = S(X)</math> is called sufficient if the distribution of <math>X</math> given <math>S = s</math> does not depend on <math>\theta</math>.</p>
<p>A <math>k</math>-dimensional exponential family is a family of distributions <math>\{P_\theta : \theta \in \Theta\}</math> if the density functions of distributions in the family can be written in the form</p> $p_\theta(x) = \exp \left[ \sum_{j=1}^k c_j(\theta) T_j(x) - d(\theta) \right] h(x).$	<p>You can use the factorization theorem of Neyman:</p> $\begin{aligned} \hat{\theta}_{\text{MLE}} &= \arg \max_{\varphi \in \Theta} p_\varphi(x) \\ &= \arg \max_{\varphi \in \Theta} g_\varphi(S) h(x) \\ &= \arg \max_{\varphi \in \Theta} g_\varphi(S) \end{aligned}$
$\prod_{i=1}^n p_\theta(x_i) = \exp \left[ \underbrace{c(\theta) \sum_{i=1}^n T(x_i) - nd(\theta)}_{g_\theta(S(x))} \right] \underbrace{\prod_{i=1}^n h(x_i)}_{h(x)}$	<p>The pdf can be written as</p> $\begin{aligned} p_\theta(x) &= \mathbb{P}(X = x) = e^{-\theta} \frac{\theta^x}{x!} \\ &= \exp \left[ \underbrace{\log(\theta)}_{c(\theta)} \underbrace{x}_{T(x)} - \underbrace{\theta}_{d(\theta)} \right] \underbrace{\frac{1}{x!}}_{h(x)} \end{aligned}$
<p>It's obtained by setting <math>c(\theta)</math> to be an identity function</p> $p_\theta(x) = \exp \left[ \sum_{j=1}^k \theta_j T_j(x) - d(\theta) \right] h(x)$ <p>Typically this can be achieved by reparametrization.</p>	<p>Expand the square in the numerator and move <math>\sigma^2</math> into the exponential</p> $\begin{aligned} p_\theta(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(x-\mu)^2}{2\sigma^2} \right] \\ &= \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{x^2 - 2x\mu + \mu^2}{2\sigma^2} \right] \exp \left[ -\log \frac{1}{2} (\sigma^2) \right] \\ &= \exp \left[ -\frac{x^2}{2\sigma^2} + x \frac{\mu}{\sigma^2} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}} \\ &= \exp \left[ \underbrace{\left[ -\frac{1}{2\sigma^2} \right]}_{c(\theta)} \underbrace{\left[ \frac{x^2}{x} \right]}_{T(x)} - \underbrace{\left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right)}_{d(\theta)} \right] \underbrace{\frac{1}{\sqrt{2\pi}}}_{h(x)} \end{aligned}$

<p>DEFINITION</p> <p>For a collection of probability density functions <math>\{p_\theta : \theta \in \Theta\}</math> with <math>\Theta \subset \mathbb{R}</math>, what is a score function?</p>	<p>DEFINITION</p> <p>What is the Fisher information for a parameter <math>\theta</math>?</p>
<p>What can you say about the expected value of the score function <math>\mathbb{E}[s_\theta(x)]</math> ?</p>	<p>What can you say about the expected value of Fisher information with respect to the score function?</p>
<p>If <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math>, how can you write <math>\mathbb{E}_\theta T(x)</math> and <math>\text{Var}_\theta(T(x))</math> in terms of <math>c, d</math>, and/or <math>h</math> ?</p>	<p>If <math>p_\theta(x)</math> is a probability density function from the exponential family in the canonical form</p> $p_\theta(x) = \exp \left[ \sum_{j=1}^k \theta_j T_j(x) - d(\theta) \right] h(x),$ <p>what can you say about the expectation <math>\mathbb{E}_\theta[T(x)]</math> and <math>\text{Var}_\theta(T(x))</math> ?</p>
<p>Given <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math>, derive <math>\mathbb{E}_\theta[T(x)]</math> in terms of <math>\dot{c}(\theta)</math> and <math>\dot{d}(\theta)</math>.</p>	<p>How is the maximum likelihood estimate <math>\hat{\theta}_{\text{MLE}}</math> of a parameter <math>\theta</math> related to the score function?</p>
<p>Given <math>X \sim \text{Bernoulli}(\theta)</math> for <math>\theta \in (0, 1)</math>, and knowing that for <math>x \in \{0, 1\}</math> we can write</p> $\mathbb{P}(X = x) = \exp \left[ \left( \log \frac{\theta}{1 - \theta} \right) x - (-\log(1 - \theta)) \right]$ <p>how can you reparametrize <math>X</math> into canonical exponential family form?</p>	<p>Take a collection of i.i.d. random variables <math>X_1, \dots, X_n \sim \mathcal{N}(0, \sigma^2)</math>. What is a minimal sufficient statistic for <math>\sigma^2</math> ? Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \ln(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$

<p>It's defined as</p> $I(\theta) = \text{Var}_{\theta}(s_{\theta}(x))$	<p>The score function <math>s_{\theta}</math> is</p> $s_{\theta}(x) = \frac{d}{d\theta} \log p_{\theta}(x)$
<p>Under regularity assumptions <math>I(\theta) = -\mathbb{E}_{\theta}[\dot{s}_{\theta}(x)]</math></p>	<p>Under regularity conditions, <math>\mathbb{E}_{\theta}[s_{\theta}(x)] = 0</math>.</p>
$\mathbb{E}_{\theta}T(x) = \dot{d}(\theta)$ $\text{Var}_{\theta}(T(x)) = \ddot{d}(\theta)$	$\mathbb{E}_{\theta}T(x) = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}$ $\text{Var}_{\theta}(T(x)) = \frac{\left[\ddot{d}(\theta) - \ddot{c}(\theta)\frac{\dot{d}(\theta)}{\dot{c}(\theta)}\right]}{\dot{c}(\theta)^2}$
<p><math>\hat{\theta}_{\text{MLE}}</math> is the solution of</p> $\frac{1}{n} \sum_{i=1}^n s_{\theta}(x_i) = 0.$	<p>We will use the score function, and <math>\mathbb{E}_{\theta}[s_{\theta}(x)] = 0</math></p> $\log p_{\theta}(x) = c(\theta)T(x) - d(\theta) + \log h(x)$ $s_{\theta}(x) = \dot{c}(\theta)T(x) - \dot{d}(\theta)$ $\mathbb{E}_{\theta}s_{\theta}(x) = 0$ $\implies \mathbb{E}_{\theta}T(x) = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}$
<ol style="list-style-type: none"> <li>From the formula we see that <math>x_i^2</math> is minimal sufficient for any particular <math>X_i</math> since <math>\mu = 0</math>.</li> <li>We know that for any i.i.d. collection of exponential family random variables, the sum of <math>T(x)</math> is sufficient.</li> </ol> <p>So, the minimal sufficient statistic in this case is the sum (or e.g. mean) of squares <math>s(x) = \sum_{i=1}^n x_i^2</math>.</p>	<ol style="list-style-type: none"> <li>Define <math>\gamma = \left(\log \frac{\theta}{1-\theta}\right)</math>.</li> <li>Rewrite <math>(-\log(1-\theta))</math> using <math>\gamma</math> <math display="block">\theta = \frac{e^{\gamma}}{1+e^{\gamma}} \implies -\log(1-\theta) = \log(1+e^{\gamma})</math> </li> <li>Plug into the equation. <math display="block">\mathbb{P}(X=x) = \exp[\gamma x - \log(1+e^{\gamma})]</math> </li> </ol>

<p>Consider <math>X \sim P_\theta, \theta \in \Theta</math>, some parameter of interest <math>\gamma = g(\theta) \in \mathbb{R}</math>, and an estimator <math>T(x)</math> of <math>\gamma</math>. What is the bias of <math>T</math>? What is an unbiased estimator?</p>	<p>Consider <math>X \sim \text{Binomial}(n, \theta)</math>. For what class of functions <math>g(\theta)</math> is it possible to construct unbiased estimators? Why? Recall that the pmf of <math>X</math> is</p> $\mathbb{P}(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$
<p>Consider <math>X \sim \text{Binomial}(n, \theta)</math>. How could you show that <math>T(X) = \frac{X(X-1)}{n(n-1)}</math> is an unbiased estimator of <math>\theta^2</math>? Recall that:</p> <ol style="list-style-type: none"> <li>1. Pmf of <math>X</math> is <math>\mathbb{P}(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}</math>.</li> <li>2. <math>\mathbb{E}X = n\theta</math></li> <li>3. <math>\mathbb{E}X^2 = n\theta(1 - \theta) + (n\theta)^2</math>.</li> </ol>	<p>Let <math>X_1, \dots, X_n</math> be i.i.d. copies of <math>X</math>. What is an unbiased estimator of <math>\text{Var}(X)</math>?</p>
<p>If <math>X_1, \dots, X_n</math> are i.i.d. random variables with mean <math>\mu</math> and variance <math>\sigma^2</math>, what</p> $\text{Var} \left( \frac{X_1 + \dots + X_n}{n} \right)?$	<p>Let <math>X_1, \dots, X_n</math> i.i.d. copies of <math>X</math>. Define <math>\mu = \mathbb{E}(X)</math> and <math>\sigma^2 = \text{Var}(X)</math>. What are the steps to show that</p> $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ <p>is an unbiased estimator of <math>\sigma^2</math>?</p>
<p>Consider <math>X_1, \dots, X_n</math> i.i.d. copies of <math>X</math> with variance <math>\sigma^2</math>. Knowing that</p> $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ <p>is an unbiased estimator of <math>\sigma^2</math>, show that <math>S = \sqrt{S^2}</math> is not an unbiased estimator of <math>\sigma</math>.</p>	<p>What is the mean square error of an estimator <math>T(X)</math> trying to estimate <math>g(\theta)</math>?</p>
<p>What is <math>\text{MSE}_\theta(T)</math> in terms of <math>\text{Bias}_\theta(T)</math> and <math>\text{Var}_\theta(T)</math>?</p>	<p>Show that for an estimator <math>T(X)</math> trying to estimate <math>g(\theta)</math></p> $\text{MSE}_\theta(T) = \text{Bias}_\theta^2(T) + \text{Var}_\theta(T)$

<p>Only polynomials. This is because:</p> <ol style="list-style-type: none"> <li>1. An estimator is unbiased when <math>\mathbb{E}_\theta T(X) - g(\theta) = 0</math> for all <math>\theta \in \Theta</math>.</li> <li>2. <math>\mathbb{E}_\theta T(X) = \sum_{x=0}^n \binom{n}{x} \theta^x (1-\theta)^{n-x} T(x) = q(\theta)</math> where <math>q(\theta)</math> is a polynomial of degree <math>\leq n</math>.</li> <li>3. <math>\mathbb{E}_\theta T(X) - g(\theta) = 0</math> for all <math>\theta \in \Theta</math> means <math>g(\theta)</math> must be a polynomial.</li> </ol>	$\text{Bias}_\theta(T) := \mathbb{E}_\theta T(X) - g(\theta)$ <p>An estimator <math>T</math> is unbiased if <math>\text{Bias}_\theta(T) = 0</math> for all <math>\theta \in \Theta</math></p>
<p>An unbiased estimator of <math>\text{Var}(X)</math> is</p> $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	$\begin{aligned} \mathbb{E}_\theta T &= \frac{\mathbb{E}_\theta X^2 - \mathbb{E}_\theta X}{n(n-1)} \\ &= \frac{n\theta(1-\theta) + (n\theta)^2 - n\theta}{n(n-1)} \\ &= \frac{-n\theta^2 + n^2\theta^2}{n(n-1)} \\ &= \frac{n(n-1)\theta^2}{n(n-1)} = \theta^2 \end{aligned}$ <p>so <math>\forall \theta \in (0,1) \mathbb{E}_\theta T - \theta^2 = 0</math></p>
<p>First note the identity</p> $\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2 = \sum_{i=1}^n X_i^2 - n\bar{X}^2.$ $\begin{aligned} \mathbb{E}(S^2) &= \frac{1}{n-1} \mathbb{E} \left[ \sum_{i=1}^n (X_i - \bar{X})^2 \right] \\ &= \frac{1}{n-1} \mathbb{E} \left[ \sum_{i=1}^n ([X_i - \mu] + [\mu - \bar{X}])^2 \right] \quad (\text{Insert } \mu - \mu) \\ &= \frac{1}{n-1} \mathbb{E} \left[ \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2 \right] \quad (\text{Above identity}) \\ &= \frac{1}{n-1} \left[ n\sigma^2 - n \frac{\sigma^2}{n} \right] = \sigma^2 \end{aligned}$	<p>Note that <math>\text{Var}\left(\frac{X_i}{n}\right) = \frac{\sigma^2}{n^2}</math>. So it's</p> $\begin{aligned} \text{Var}\left(\frac{X_1 + \dots + X_n}{n}\right) &= \text{Var}\left(\frac{X_1}{n}\right) + \dots + \text{Var}\left(\frac{X_n}{n}\right) \\ &= \frac{1}{n^2} (\text{Var}(X_1) + \dots + \text{Var}(X_n)) \\ &= \frac{n}{n^2} \sigma^2 \\ &= \frac{1}{n} \sigma^2 \end{aligned}$
$\text{MSE}_\theta(T) = \mathbb{E}_\theta (T(X) - g(\theta))^2$	<p>Observe that</p> $\mathbb{E}(S) = \mathbb{E}\left(\sqrt{S^2}\right) \leq \sqrt{\mathbb{E}(S^2)} = \sigma$ <p>Since the middle inequality (holds for concave Jensen) is an equality if and only if <math>\text{Var}(S) = 0</math>, <math>S</math> is not an unbiased estimator of <math>\sigma</math>.</p>
<p>Let <math>q(\theta) = \mathbb{E}_\theta T(X)</math>. Then</p> $\begin{aligned} \mathbb{E}_\theta (T - g(\theta))^2 &= \mathbb{E}_\theta (T - q(\theta) + q(\theta) - g(\theta))^2 \\ &= \underbrace{\mathbb{E}_\theta (T - q(\theta))^2}_{\text{Var}_\theta(T)} + \underbrace{(q(\theta) - g(\theta))^2}_{\text{Bias}_\theta(T)^2} \\ &\quad + 2(q(\theta) - g(\theta)) \underbrace{\mathbb{E}_\theta (T - q(\theta))}_0 \\ &= \text{Bias}_\theta^2(T) + \text{Var}_\theta(T). \end{aligned}$	$\text{MSE}_\theta(T) = \text{Bias}_\theta^2(T) + \text{Var}_\theta(T)$



<p>When estimating a parameter <math>\gamma</math>, is it possible for a biased estimator to get a better MSE than an unbiased estimator?</p>	<p>DEFINITION</p> <p>What is a Uniform Minimum Variance Unbiased (UMVU) estimator?</p>
<p>LEMMA</p> <p>State the iterated variance lemma.</p>	<p>Show that for any random variables <math>Y</math> and <math>Z</math>,</p> $\text{Var}(Y) = \text{Var}(\mathbb{E}(Y \mid Z)) + \mathbb{E}(\text{Var}(Y \mid Z))$ <p>(part-1)</p>
<p>Show that for any random variables <math>Y</math> and <math>Z</math>,</p> $\text{Var}(Y) = \text{Var}(\mathbb{E}(Y \mid Z)) + \mathbb{E}(\text{Var}(Y \mid Z))$ <p>(part-2)</p>	<p>If <math>T</math> is an unbiased estimator and <math>S</math> is a sufficient statistic, what can you say about the variance of an estimator <math>T^* = \mathbb{E}(T \mid S)</math> in relation to the variance of <math>T</math> ?</p>
<p>Show that if <math>T</math> is an unbiased estimator for <math>g(\theta)</math> and <math>S</math> is a sufficient statistic, then the estimator defined by <math>T^* = \mathbb{E}(T \mid S)</math> is unbiased and fulfils <math>\text{Var}_\theta(T^*) \leq \text{Var}_\theta(T)</math> for all <math>\theta</math></p>	<p>What does it mean for a sufficient statistic <math>S</math> to be called complete?</p>
<p>LEMMA</p> <p>State the Lehmann-Scheffe lemma.</p>	<p>Prove that if <math>T</math> is an unbiased estimator of <math>g(\theta)</math> with finite variance for all <math>\theta</math>, and <math>S</math> is a sufficient and complete statistic, then <math>T^* = \mathbb{E}(T \mid S)</math> is UMVU.</p>

<p>An unbiased estimator <math>T^*</math> is called Uniform Minimum Variance Unbiased (UMVU) estimator if for any other unbiased estimator <math>T</math>,</p> $\text{Var}_\theta(T^*) \leq \text{Var}_\theta(T), \forall_\theta$	<p>Yes, for instance when estimating <math>\sigma^2</math> of a normal distribution. The estimator</p> $T_{\text{opt}} = \frac{1}{n+1} \sum_{i=1}^n (X_i - \bar{X})^2$ <p>is biased, but obtains a better MSE than the unbiased sample variance</p> $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
<p>Just repeatedly use <math>\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2</math>, linearity of <math>\mathbb{E}</math>, and iterated <math>\mathbb{E}</math></p> <p>1.</p> $\begin{aligned} \text{Var}(\mathbb{E}(Y   Z)) &= \mathbb{E}(\mathbb{E}(Y   Z))^2 - (\mathbb{E}(\mathbb{E}(Y   Z)))^2 \\ &= \mathbb{E}(\mathbb{E}(Y   Z))^2 - (\mathbb{E}(Y))^2 \end{aligned}$ <p>( Iterated <math>\mathbb{E}</math>)</p>	<p>Let <math>Y</math> and <math>Z</math> be two random variables. Then</p> $\text{Var}(Y) = \text{Var}(\mathbb{E}(Y   Z)) + \mathbb{E}(\text{Var}(Y   Z))$
<p>For all <math>\theta</math>, <math>\text{Var}_\theta(T^*) \leq \text{Var}_\theta(T)</math></p>	<p>2. (Linearity, iterated <math>\mathbb{E}</math>)</p> $\begin{aligned} \mathbb{E}(\text{Var}(Y   Z)) &= \mathbb{E}(\mathbb{E}(Y^2   Z) - (\mathbb{E}(Y   Z))^2) \\ &= \mathbb{E}(Y^2) - \mathbb{E}(\mathbb{E}(Y   Z))^2 \end{aligned}$ <p>3.</p> $\begin{aligned} &\text{Var}(\mathbb{E}(Y   Z)) + \mathbb{E}(\text{Var}(Y   Z)) \\ &= \mathbb{E}(\mathbb{E}(Y   Z))^2 - (\mathbb{E}(Y))^2 + \mathbb{E}(Y^2) - \mathbb{E}(\mathbb{E}(Y   Z))^2 \\ &= \mathbb{E}(Y^2) - (\mathbb{E}(Y))^2 = \text{Var}(Y) \end{aligned}$
<p>A sufficient statistic <math>S</math> is called complete if the following implication holds:</p> $\forall_\theta \mathbb{E}_\theta h(S) = 0 \implies \forall_\theta h(S) = 0, \mathbb{P}_\theta\text{-almost surely}$ <p>where <math>h</math> is a function of <math>S</math> not depending on <math>\theta</math>.</p>	<p><math>T^*</math> is an unbiased estimator for <math>g(\theta)</math> because</p> $\mathbb{E}_\theta T^* = \mathbb{E}_\theta \mathbb{E}(T   S) = \mathbb{E}_\theta T = g(\theta).$ <p>Now use iterated variance lemma</p> $\begin{aligned} \text{Var}_\theta(T) &= \text{Var}_\theta(\overbrace{\mathbb{E}(T   S)}^{T^*}) + \overbrace{\mathbb{E}_\theta \text{Var}(T   S)}^{\geq 0} \\ &\geq \text{Var}_\theta(T^*). \end{aligned}$
<p>Let <math>\tilde{T}</math> be unbiased. Define <math>T' = \mathbb{E}(\tilde{T}   S) = T'(S)</math>. We know now that <math>\text{Var}_\theta(T') \leq \text{Var}_\theta(\tilde{T})</math>. Now <math>T^*(S) = \mathbb{E}(T   S)</math> and <math>T'(S) = \mathbb{E}(\tilde{T}   S)</math> are unbiased. Since they're unbiased, it's true that for all <math>\theta</math></p> $\mathbb{E}_\theta [T^*(S) - T'(S)] = 0$ <p>Since <math>S</math> is complete, <math>T^*(S) - T'(S) = 0</math> <math>\mathbb{P}_\theta</math>-almost surely for all <math>\theta</math></p>	<p>Let <math>T</math> be an unbiased estimator of <math>g(\theta)</math> with finite variance for all <math>\theta</math>. Also let <math>S</math> be a sufficient and complete statistic. Then <math>T^* = \mathbb{E}(T   S)</math> is UMVU.</p>

<p>Consider <math>X_1, \dots, X_n</math> i.i.d. <math>\text{Poisson}(\theta)</math> random variables with <math>\theta &gt; 0</math>. Show that the statistic <math>S = \sum_{i=1}^n X_i</math> is complete. (part-1)</p>	<p>Consider <math>X_1, \dots, X_n</math> i.i.d. <math>\text{Poisson}(\theta)</math> random variables with <math>\theta &gt; 0</math>. Show that the statistic <math>S = \sum_{i=1}^n X_i</math> is complete. (part-2)</p>
<p>If <math>S</math> is a sufficient and complete statistic, and <math>T'(S) \in \mathbb{R}</math> is a function which only depends on <math>S</math>, what is an easy UMVU estimator of <math>\mathbb{E}_\theta T'</math> ?</p>	<p>LEMMA</p> <p>State the lemma about sufficiency and completeness of statistics for exponential family random variables.</p>
<p>Consider i.i.d. <math>X_1, \dots, X_n \sim \text{Gamma}(k\lambda)</math> random variables where <math>\theta = (k, \lambda) \in \mathbb{R}_+^2</math> with density</p> $p_\theta(x) = \frac{e^{-\lambda x} x^{k-1} \lambda^k}{\Gamma(k)}$ $= \exp[-\lambda x + (k-1) \log x + k \log \lambda - \log \Gamma(k)]$ <p>Sufficient and complete statistic for <math>\theta</math> ? Why?</p>	<p>Consider a collection of distributions on <math>\mathcal{X}</math> parametrized by one-dimensional <math>\theta \in \Theta</math>. What is the support condition of the Cramér-Rao lower bound?</p>
<p>Consider a collection of distributions on <math>\mathcal{X}</math> parametrized by one-dimensional <math>\theta \in \Theta</math>. What is the differentiability in <math>L^2</math> condition of the Cramér-Rao lower bound?</p>	<p>Consider a collection of distributions on <math>\mathcal{X}</math> parametrized by one-dimensional <math>\theta \in \Theta</math>. What do the following conditions imply about <math>s_\theta</math> ?</p> <ol style="list-style-type: none"> <li>1. The support <math>\{x : p_\theta(x) &gt; 0\}</math> does not depend on parameter <math>\theta</math>.</li> <li>2. There exists a function <math>s_\theta : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_\theta s_\theta^2(x) &lt; \infty</math> and</li> </ol> $\lim_{h \rightarrow 0} \mathbb{E}_\theta \left[ \frac{p_{\theta+h}(x) - p_\theta(x)}{h p_\theta(x)} - s_\theta(x) \right]^2 = 0$ <p>for all <math>\theta</math>.</p>
<p>Let <math>T</math> be an estimator of <math>g(\theta)</math> where <math>\theta</math> is one-dimensional. Assume the support of <math>P_\theta</math> does not depend on <math>\theta</math>, and the differentiability in <math>L^2</math> condition holds. If <math>\text{Var}_\theta(T) &lt; \infty</math> for all <math>\theta</math>, then what is a simple formula for <math>\dot{g}(\theta)</math> ?</p>	<p>State the Cauchy-Schwartz inequality for covariances/variances.</p>

<p>Since the monomials <math>\{1, x, x^2, \dots\}</math> form an independent system, we know that</p> $\sum_{k=0}^{\infty} \frac{x^k}{k!} h(k) = 0 \text{ for all } x > 0 \iff h(k) = 0 \text{ for all } k \geq 0$ <p>which is the desired implication to prove completeness of <math>S</math>.</p>	<p>Note that the sum <math>S</math> is distributed according to <math>\text{Poisson}(n\theta)</math>. Consider now a function <math>h</math> which has the property that</p> $\mathbb{E}_{\theta} h(X) = \sum_{k=0}^{\infty} e^{-n\theta} \frac{(n\theta)^k}{k!} h(k) = 0 \text{ for all } \theta > 0.$ <p>Rewrite expansion in an obviously independent basis</p> $f(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!} h(k) = 0 \text{ for all } x > 0.$
$p_{\theta}(x) = \exp[\underbrace{c(\theta)}_{\in \mathbb{R}^k} \underbrace{T(x)}_{\in \mathbb{R}^k} - d(\theta)] h(x)$ $\mathcal{C} := \left\{ c(\theta) = \begin{pmatrix} c_1(\theta) \\ \vdots \\ c_k(\theta) \end{pmatrix} \in \mathbb{R}^k; \theta \in \Theta \right\}$ <p>contains a <math>k</math>-dimensional open ball. Then <math>T = \begin{pmatrix} T_1 \\ \vdots \\ T_k \end{pmatrix}</math> is sufficient and complete.</p>	<p><math>T'</math> itself is an UMVU estimator of <math>\mathbb{E}_{\theta} T'</math></p>
<p>The support <math>\{x : p_{\theta}(x) &gt; 0\}</math> does not depend on parameter <math>\theta</math></p>	<p>The sufficient and complete statistic is <math>S = (\sum_{i=1}^n X_i, \sum_{i=1}^n \log X_i)</math>. This is because <math>c_1(\theta) = -\lambda, c_2(\theta) = k - 1</math>, so the space of <math>\mathcal{C} = (-\infty, 0) \times (-1, \infty)</math> is two dimensional, and we know from a lemma that this means that <math>T</math> is sufficient and complete.</p>
<p>For all <math>\theta</math> it's true that <math>\mathbb{E}_{\theta} s_{\theta}(X) = 0</math></p>	<p>There exists a function <math>s_{\theta} : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_{\theta} s_{\theta}^2(x) &lt; \infty</math> and</p> $\lim_{h \rightarrow 0} \mathbb{E}_{\theta} \left[ \frac{p_{\theta+h}(x) - p_{\theta}(x)}{h p_{\theta}(x)} - s_{\theta}(x) \right]^2 = 0$ <p>for all <math>\theta</math>.</p>
$\text{Cov}(X, Y)^2 \leq \text{Var}(X) \text{Var}(Y)$	$\dot{g}(\theta) = \text{Cov}(T, s_{\theta}(X))$

<p>Let <math>X_1, \dots, X_n</math> be i.i.d. copies of <math>X</math>. Suppose <math>s_\theta = \frac{d}{d\theta} \log p_\theta</math>. How does the Fisher information change as the sample size <math>n</math> increases?</p>	<p>Let <math>X \sim \text{Geometric}(\theta), 0 &lt; \theta &lt; 1</math>. What is the score function <math>s_\theta</math> in this case? Recall that for <math>X \sim \text{Geometric}(\theta)</math>,</p> <ol style="list-style-type: none"> <li>1. <math>p_\theta(x) = \mathbb{P}_\theta(X = x) = \theta(1 - \theta)^{x-1}, x = 1, 2, \dots,</math></li> <li>2. <math>\mathbb{E}_\theta X = \frac{1}{\theta}</math></li> <li>3. <math>\text{Var}_\theta(X) = \frac{1-\theta}{\theta^2},</math></li> <li>4. <math>\log p_\theta(x) = \underbrace{(x-1)}_{T(x)} \underbrace{\log(1-\theta)}_{c(\theta)} - \underbrace{(-\log \theta)}_{d(\theta)}.</math></li> </ol>
<p>Consider a collection of distributions on <math>\mathcal{X}</math> parametrized by one-dimensional <math>\theta \in \Theta</math>. What is the Cramér-Rao lower bound for an unbiased estimator of <math>g(\theta)</math>?</p>	<p>Let <math>X \sim \text{Geometric}(\theta), 0 &lt; \theta &lt; 1</math>. What are the possible ways to find the Fisher information <math>I(\theta)</math> in this case? Recall that for <math>X \sim \text{Geometric}(\theta)</math>,</p> <ol style="list-style-type: none"> <li>1. <math>p_\theta(x) = \mathbb{P}_\theta(X = x) = \theta(1 - \theta)^{x-1}, x = 1, 2, \dots,</math></li> <li>2. <math>\mathbb{E}_\theta X = \frac{1}{\theta}</math></li> <li>3. <math>\text{Var}_\theta(X) = \frac{1-\theta}{\theta^2}</math></li> <li>4. <math>\log p_\theta(x) = \underbrace{(x-1)}_{T(x)} \underbrace{\log(1-\theta)}_{c(\theta)} - \underbrace{(-\log \theta)}_{d(\theta)}</math></li> <li>5. <math>s_\theta(x) = \frac{d}{d\theta} \log p_\theta(x) = -\frac{x-1}{1-\theta} + \frac{1}{\theta}.</math></li> </ol>
<p>Let <math>X_1, \dots, X_n \sim \text{Geometric}(\theta)</math> be i.i.d with <math>0 &lt; \theta &lt; 1</math>. What is the Cramér Rao lower bound for <math>g(\theta) = \frac{1}{\theta}</math> ? Recall that for <math>X \sim \text{Geometric}(\theta)</math>,</p> <ol style="list-style-type: none"> <li>1. <math>p_\theta(x) = \mathbb{P}_\theta(X = x) = \theta(1 - \theta)^{x-1}, x = 1, 2, \dots,</math></li> <li>2. <math>\mathbb{E}_\theta X = \frac{1}{\theta}</math></li> <li>3. <math>\text{Var}_\theta(X) = \frac{1-\theta}{\theta^2}</math></li> <li>4. <math>\log p_\theta(x) = \underbrace{(x-1)}_{T(x)} \underbrace{\log(1-\theta)}_{c(\theta)} - \underbrace{(-\log \theta)}_{d(\theta)},</math></li> <li>5. <math>s_\theta(x) = \frac{d}{d\theta} \log p_\theta(x) = -\frac{x-1}{1-\theta} + \frac{1}{\theta}</math></li> <li>6. <math>I(\theta) = \frac{1}{(1-\theta)\theta^2}.</math></li> </ol>	<p>Let <math>X_1, \dots, X_n \sim \text{Geometric}(\theta)</math> be i.i.d with <math>0 &lt; \theta &lt; 1</math>. What is the Cramér Rao lower bound for <math>g(\theta) = \theta</math> ? Recall that for <math>X \sim \text{Geometric}(\theta)</math>,</p> <ol style="list-style-type: none"> <li>1. <math>p_\theta(x) = \mathbb{P}_\theta(X = x) = \theta(1 - \theta)^{x-1}, x = 1, 2, \dots,</math></li> <li>2. <math>\mathbb{E}_\theta X = \frac{1}{\theta}</math></li> <li>3. <math>\text{Var}_\theta(X) = \frac{1-\theta}{\theta^2}</math></li> <li>4. <math>\log p_\theta(x) = \underbrace{(x-1)}_{T(x)} \underbrace{\log(1-\theta)}_{c(\theta)} - \underbrace{(-\log \theta)}_{d(\theta)}</math></li> <li>5. <math>s_\theta(x) = \frac{d}{d\theta} \log p_\theta(x) = -\frac{x-1}{1-\theta} + \frac{1}{\theta}</math></li> <li>6. <math>I(\theta) = \frac{1}{(1-\theta)\theta^2}.</math></li> </ol>
<p>What is the correlation between random variables <math>X</math> and <math>Y</math> ?</p>	<p>When is correlation <math>\rho(X, Y)</math> between <math>X</math> and <math>Y</math> equal to 1?</p>
<p>Which estimators for which class of distributions reach the Cramér-Rao lower bound?</p>	<p>What are the ways to find whether an estimator <math>T</math> reaches the Cramer-Rao lower bound?</p>

$s_{\theta}(x) = \frac{d}{d\theta} \log p_{\theta}(x) = -\frac{x-1}{1-\theta} + \frac{1}{\theta}$	$\underbrace{I^{(n)}(\theta)}_{\text{for sample}} = \text{Var} \left( s_{\theta}^{(n)}(x_1, \dots, x_n) \right)$ $= \text{Var} \left( \frac{d}{d\theta} \log p_{\theta}^{(n)}(x_1, \dots, x_n) \right)$ $= \text{Var} \left( \frac{d}{d\theta} \sum_{i=1}^n \log p_{\theta}(x_i) \right)$ $= \text{Var} \left( \sum_{i=1}^n s_{\theta}(x_i) \right) = \sum_{i=1}^n \text{Var}(s_{\theta}(x_i)) = n \underbrace{I(\theta)}_{\text{single ob}}$
<p>1. Directly, by <math>I(\theta) = \text{Var}_{\theta}(s_{\theta}(X))</math> :</p> $I(\theta) = \frac{\text{Var}_{\theta}(X)}{(1-\theta)^2} = \frac{1}{(1-\theta)^2} \cdot \frac{1-\theta}{\theta^2} = \frac{1}{(1-\theta)\theta^2}$ <p>2. By <math>\mathbb{E}_{\theta} \dot{s}_{\theta}(x) = -I(\theta)</math> :</p> $\dot{s}_{\theta}(x) = \frac{(x-1)}{(1-\theta)^2} - \frac{1}{\theta^2}$ $\mathbb{E}_{\theta} \dot{s}_{\theta}(X) = \frac{\frac{1}{\theta} - 1}{(1-\theta)^2} - \frac{1}{\theta^2} = -\frac{1}{(1-\theta)\theta^2} = -I(\theta)$	<p>Assume the support of <math>P_{\theta}</math> does not depend on <math>\theta</math>, and the differentiability in <math>L^2</math> condition holds. If <math>\text{Var}_{\theta}(T) &lt; \infty</math> for all <math>\theta</math>, then the Cramér-Rao lower bound is</p> $\text{Var}(T) \geq \frac{[\dot{g}(\theta)]^2}{I(\theta)}$
<p>We have</p> $\text{CRLB} = \frac{\dot{g}(\theta)^2}{I(\theta)} = \frac{1}{\frac{n}{(1-\theta)\theta^2}} = \frac{(1-\theta)\theta^2}{n}.$ <p>Remember that the Fisher information grows linearly in the sample size <math>n</math>.</p>	<p>We have</p> $\text{CRLB} = \frac{\dot{g}(\theta)^2}{I(\theta)} = \frac{\left[-\frac{1}{\theta^2}\right]^2}{\frac{n}{(1-\theta)\theta^2}} = \frac{(1-\theta)}{n\theta^2}.$ <p>Remember that the Fisher information grows linearly in the sample size <math>n</math>.</p>
<p>When there exist constants <math>a, b</math> such that <math>Y = aX + b</math>.</p>	<p>Correlation is defined as</p> $\rho(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$
<p>Two ways:</p> <p>1. Compute the bound by</p> $\text{CRLB} = \frac{(\dot{g}(\theta))^2}{I(\theta)}$ <p>and compare it to <math>\text{Var}(T)</math>.</p> <p>2. Check if the distribution is from the exponential family, and if <math>g(\theta) = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}</math></p>	<p>Assume the conditions of the Cramér-Rao lower bound hold with <math>s_{\theta} = \frac{d}{d\theta} \log p_{\theta}</math>. Let <math>T</math> be an unbiased estimator of <math>g(\theta)</math>, and suppose <math>\text{Var}_{\theta}(T) = \frac{(\dot{g}(\theta))^2}{I(\theta)}</math> for all <math>\theta</math>. Then for all <math>x</math> and <math>\theta</math>,</p> $p_{\theta}(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)$ <p>and <math>g(\theta) = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}</math>.</p>

<p>For a random vector taking values in <math>Z \in \mathbb{R}^k</math>, what is the covariance <math>\text{Cov}(Z)</math> ?</p>	<p>What does it mean for a symmetric matrix <math>V \in \mathbb{R}^{k \times k}</math> to be positive-definite?</p>
<p>PROOF</p> <p>Show that a covariance matrix <math>\Sigma := \text{Cov}(Z)</math> is positive semi-definite.</p>	<p>Let <math>V \succ 0</math> be a positive definite matrix. How is <math>\sqrt{V}</math> defined? Given <math>\sqrt{V}</math>, how can you compute <math>V</math> ? (part-1)</p>
<p>Let <math>V \succ 0</math> be a positive definite matrix. How is <math>\sqrt{V}</math> defined? Given <math>\sqrt{V}</math>, how can you compute <math>V</math> ? (part-2)</p>	<p>Let <math>V \succ 0</math> be a positive definite matrix. What is <math>V^{-\frac{1}{2}}</math> ?</p>
<p>Write <math>\text{Var}(X)</math> in terms of covariance.</p>	<p>If <math>X</math> and <math>Y</math> are independent, what is <math>\text{Cov}(X, Y)</math> ?</p>
<p>Is it true that <math>\text{Cov}(X, Y) = \text{Cov}(Y, X)</math> ?</p>	<p>Expand these statements into simpler forms:</p> <ol style="list-style-type: none"> <li>1. <math>\text{Cov}(aX, Y)</math></li> <li>2. <math>\text{Cov}(X + c, Y)</math></li> <li>3. <math>\text{Cov}(X + Y, Z)</math></li> </ol>

<p><math>V</math> is positive-definite, written as <math>V \succ 0</math>, if <math>a^T V a &gt; 0</math> for all <math>a \neq 0</math>.</p>	<p>Covariance is the matrix defined by</p> $\begin{aligned}\text{Cov}(Z) &= \mathbb{E}(ZZ^T) - \mathbb{E}(Z)\mathbb{E}(Z)^T \\ &= \mathbb{E}((Z - \mu)(Z - \mu)^T)\end{aligned}$
<p>We can write <math>V = Q\Lambda Q^T</math> where</p> <ul style="list-style-type: none"> <li>• <math>Q^T Q = I</math></li> <li>• <math>Q = (q_1, \dots, q_k)</math> is a matrix of eigenvectors <math>q_i</math>, and</li> <li>• <math>\Lambda = \begin{pmatrix} x_1 &amp; \cdots &amp; 0 \\ \cdots &amp; \cdots &amp; \cdots \\ 0 &amp; \cdots &amp; x_k \end{pmatrix}</math> is a diagonal matrix of eigenvalues <math>x_i</math>.</li> </ul>	<p>Note that variance is non-negative, and we have</p> $\begin{aligned}\text{Var}(a^T Z) &= \mathbb{E}\left((a^T Z - a^T \mu)(a^T Z - a^T \mu)^T\right) \\ &= \mathbb{E}(a^T (Z - \mu)(Z - \mu)^T a) \\ &= a^T \mathbb{E}((Z - \mu)(Z - \mu)^T) a \\ &= a^T \Sigma a\end{aligned}$ <p>for all <math>a \in \mathbb{R}^k</math>.</p>
$V^{-\frac{1}{2}} := (V^{-1})^{\frac{1}{2}} = \left(V^{\frac{1}{2}}\right)^{-1}$	<p>Then <math>\sqrt{V} := Q\sqrt{\Lambda}Q^T</math> where <math>\sqrt{\Lambda} = \begin{pmatrix} \sqrt{x_1} &amp; \cdots &amp; 0 \\ \cdots &amp; \cdots &amp; \cdots \\ 0 &amp; \cdots &amp; \sqrt{x_k} \end{pmatrix}</math> We can find <math>V</math> by <math>V^{\frac{1}{2}}V^{\frac{1}{2}} = V</math></p>
$\text{Cov}(X, Y) = 0$	$\text{Var}(X) = \text{Cov}(X, X)$
<ol style="list-style-type: none"> <li>1. <math>\text{Cov}(aX, Y) = a \text{Cov}(X, Y)</math></li> <li>2. <math>\text{Cov}(X + c, Y) = \text{Cov}(X, Y)</math></li> <li>3. <math>\text{Cov}(X + Y, Z) = \text{Cov}(X, Z) + \text{Cov}(Y, Z)</math></li> </ol>	<p>Yes</p>



<p>Consider a multidimensional parameter space <math>\theta \in \Theta \subset \mathbb{R}^k</math>, and a function <math>g(\theta)</math> for which we want to construct an estimator. How do you define</p> <ol style="list-style-type: none"> <li>1. partial derivative <math>\dot{g}(\theta)</math></li> <li>2. score vector <math>s_\theta(\cdot)</math></li> <li>3. Fisher information matrix <math>I(\theta)</math> ?</li> </ol>	<p>Consider a multidimensional parameter space <math>\theta \in \Theta \subset \mathbb{R}^k</math>, and a function <math>g(\theta)</math>. What is the Cramér-Rao lower bound on unbiased estimators of <math>g(\theta)</math> ?</p>
<p>Consider a parameter space <math>\Theta = \{\theta_0, \theta_1\}</math>, an action space <math>\mathcal{A} = [0, 1]</math>, a risk</p> $R(\theta, \phi) = \begin{cases} E_{\theta_0} \phi(X) & \text{if } \theta = \theta_0 \\ 1 - E_{\theta_1} \phi(X) & \text{if } \theta = \theta_1, \end{cases}$ <p>(continues next card)</p>	<p>(continued) and a Neyman-Pearson test</p> $\phi_{NP}(X) = \begin{cases} 1 & p_1(X)/p_0(X) > c \\ q & p_1(X)/p_0(X) = c \\ 0 & p_1(X)/p_0(X) < c \end{cases}$ <p>for <math>c &gt; 0, q \in [0, 1]</math>. Show that if <math>R(\theta, \phi_{NP}) \neq 0</math>, then <math>\phi_{NP}</math> is admissible.</p>
<p>PROOF</p> <p>Prove for the discrete case that if a statistic <math>S(X)</math> is sufficient then we can factorize the mass function as</p> $p_\theta(x) = g_\theta(S(x))h(x)$ <p>for all <math>x</math> and <math>\theta</math></p>	<p>PROOF</p> <p>Prove for the discrete case that if we can factorize the mass function as</p> $p_\theta(x) = g_\theta(S(x))h(x)$ <p>for all <math>x</math> and <math>\theta</math> then <math>S(X)</math> is a sufficient statistic.</p>
<p>Show that for <math>X_1, \dots, X_n</math> i.i.d. Uniform <math>[0, \theta]</math> random variables with <math>\theta &gt; 0</math> <math>\max\{X_1, \dots, X_n\}</math> is a sufficient statistic.</p>	<p>If <math>p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x)</math> is a probability density, then what is an equation for <math>d(\theta)</math>?</p>
<p>Show that if <math>p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x)</math> then <math>\mathbb{E}_\theta T(X) = \dot{d}(\theta)</math></p>	<p>Show that if <math>p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x)</math> then <math>\text{Var}_\theta T(X) = \ddot{d}(\theta)</math>. (part-1)</p>

$\text{Var}_\theta(T) \geq \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta)$	$1. \dot{g}(\theta) := \begin{bmatrix} \delta g(\theta)/\delta \theta_1 \\ \vdots \\ \delta g(\theta)/\delta \theta_k \end{bmatrix}$ $2. s_\theta(\cdot) := \begin{bmatrix} \delta \log p_\theta/\delta \theta_1 \\ \vdots \\ \delta \log p_\theta/\delta \theta_k \end{bmatrix}$ $3. I(\theta) = \mathbb{E}_\theta s_\theta(X) s_\theta(X)^T = \text{Cov}(s_\theta(X))$
<p>1. If <math>R(\theta_0, \phi) \leq R(\theta_0, \phi_{NP})</math>. Then</p> $\begin{aligned} R(\theta_1, \phi) &= R(\theta_1, \phi) - R(\theta_1, \phi_{NP}) + R(\theta_1, \phi_{NP}) \\ &\geq c \underbrace{[R(\theta_0, \phi_{NP}) - R(\theta_0, \phi)]}_{\geq 0} + R(\theta_1, \phi_{NP}) \end{aligned}$ <p>2. If <math>R(\theta_1, \phi) \leq R(\theta_1, \phi_{NP})</math>. Then</p> $\begin{aligned} R(\theta_0, \phi) &= R(\theta_0, \phi) - R(\theta_0, \phi_{NP}) + R(\theta_0, \phi_{NP}) \\ &\geq \frac{1}{c} \underbrace{[R(\theta_1, \phi_{NP}) - R(\theta_1, \phi)]}_{\geq 0} + R(\theta_0, \phi_{NP}) \end{aligned}$	<p>We simply need to use the Neyman-Pearson lemma in form</p> <ul style="list-style-type: none"> <li><math>\bullet \frac{[R(\theta_1, \phi) - R(\theta_1, \phi_{NP})]}{c[R(\theta_0, \phi_{NP}) - R(\theta_0, \phi)]}, \text{ or} \geq</math></li> <li><math>\bullet \frac{[R(\theta_1, \phi_{NP}) - R(\theta_1, \phi)]}{c[R(\theta_0, \phi) - R(\theta_0, \phi_{NP})]} \leq</math></li> </ul>
<p>Suppose <math>S(x) = s</math>, and <math>p_\theta(x) = g_\theta(S(x))h(x)</math>. Then</p> $\begin{aligned} P_\theta(X = x \mid S = s) &= \frac{P_\theta(X = x)}{P_\theta(S = s)} \\ &= \frac{g_\theta(S)h(x)}{\sum_{\tilde{x}: S(\tilde{x})=s} g_\theta(s)h(\tilde{x})} \\ &= \frac{h(x)}{\sum_{\tilde{x}: S(\tilde{x})=s} h(\tilde{x})} \end{aligned}$ <p>does not depend on <math>\theta</math>.</p>	<p>Suppose <math>S</math> is sufficient, and that <math>S(x) = s</math>.</p> $\begin{aligned} p_\theta(x) &= P_\theta(X = x) = P(X = x \mid S = s)P_\theta(S = s) \\ &:= h(x)g_\theta(s) \end{aligned}$
<p><math>d(\theta)</math> is a normalizing constant, so we can write</p> $d(\theta) = \log \left( \int \exp[\theta T(x)] h(x) dx \right)$	<p>If we write the density function using indicator functions, we get</p> $\begin{aligned} \prod_{i=1}^n p_\theta(x_i) &= \prod_{i=1}^n \frac{1_{[0, \theta]}(x)}{\theta} = \frac{1 \{0 \leq x_i \leq \theta \forall i\}}{\theta^n} \\ &= \underbrace{1 \{\min\{x_1, \dots, x_n\} \geq 0\}}_{h(x_1, \dots, x_n)} \underbrace{\frac{1 \{\max\{x_1, \dots, x_n\} \leq \theta\}}{\theta^n}}_{g_\theta(\max\{x_1, \dots, x_n\})} \end{aligned}$ <p>Since we can factorize the density this way, by the Neyman Factorization theorem <math>\max\{X_1, \dots, X_n\}</math> is a sufficient statistic.</p>
<p>Since <math>d(\theta)</math> is a normalizing constant, we can write</p> $\begin{aligned} d(\theta) &= \log \left( \int \exp[\theta T(x)] h(x) dx \right) \\ &= \log \left( \int \exp[\theta T] h \right) \\ \dot{d}(\theta) &= \frac{1}{\int \exp[\theta T] h} \int \exp[\theta T] T h. \end{aligned}$	<p>Since <math>d(\theta)</math> is a normalizing constant, we can write</p> $\begin{aligned} d(\theta) &= \log \left( \int \exp[\theta T(x)] h(x) dx \right) \\ &= \log \left( \int \exp[\theta T] h \right) \\ \dot{d}(\theta) &= \frac{1}{\int \exp[\theta T] h} \int \exp[\theta T] T h \\ &= \int \exp[\theta T - d(\theta)] T h \\ &= \int p_\theta T = \mathbb{E}_\theta T(X). \end{aligned}$

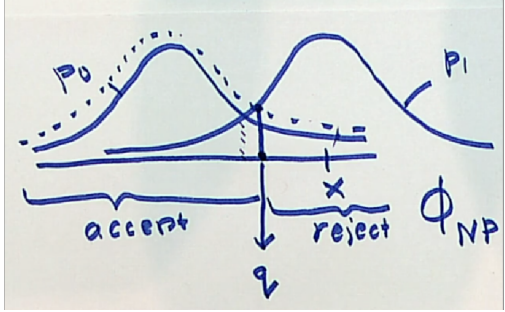
<p>Show that if <math>p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x)</math> then <math>\text{Var}_\theta T(X) = \dot{d}(\theta)</math>. (part-2)</p>	<p>PROOF</p> <p>Prove that if the support of a density <math>p_\theta</math> does not depend on <math>\theta</math>, and if there exists a function <math>s_\theta : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_\theta s_\theta^2(X) &lt; \infty</math> and</p> $\lim_{h \rightarrow 0} \mathbb{E}_\theta \left[ \frac{p_{\theta+h}(x) - p_\theta(x)}{hp\theta(x)} - s_\theta(x) \right]^2 = 0$ <p>for all <math>\theta</math>, and <math>T</math> is an unbiased estimator of <math>g(\theta)</math> with finite variance, then <math>\dot{g}(\theta) = \text{Cov}_\theta(T, s_\theta(x))</math> (part-1)</p>
<p>PROOF</p> <p>Prove that if the support of a density <math>p_\theta</math> does not depend on <math>\theta</math>, and if there exists a function <math>s_\theta : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_\theta s_\theta^2(X) &lt; \infty</math> and</p> $\lim_{h \rightarrow 0} \mathbb{E}_\theta \left[ \frac{p_{\theta+h}(x) - p_\theta(x)}{hp\theta(x)} - s_\theta(x) \right]^2 = 0$ <p>for all <math>\theta</math>, and <math>T</math> is an unbiased estimator of <math>g(\theta)</math> with finite variance, then <math>\dot{g}(\theta) = \text{Cov}_\theta(T, s_\theta(x))</math> (part-2)</p>	<p>PROOF</p> <p>Prove that if the support of a density <math>p_\theta</math> does not depend on <math>\theta</math>, and if there exists a function <math>s_\theta : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_\theta s_\theta^2(X) &lt; \infty</math> and</p> $\lim_{h \rightarrow 0} \mathbb{E}_\theta \left[ \frac{p_{\theta+h}(x) - p_\theta(x)}{hp\theta(x)} - s_\theta(x) \right]^2 = 0$ <p>for all <math>\theta</math>, and <math>T</math> is an unbiased estimator of <math>g(\theta)</math> with finite variance, then <math>\dot{g}(\theta) = \text{Cov}_\theta(T, s_\theta(x))</math> (part-3)</p>
<p>PROOF</p> <p>Sketch proof that if the support of a density <math>p_\theta</math> does not depend on <math>\theta</math>, and if there exists a function <math>s_\theta : \mathcal{X} \rightarrow \mathbb{R}</math> with <math>\mathbb{E}_\theta s_\theta^2(X) &lt; \infty</math> and</p> $\lim_{h \rightarrow 0} \mathbb{E}_\theta \left[ \frac{p_{\theta+h}(x) - p_\theta(x)}{hp\theta(x)} - s_\theta(x) \right]^2 = 0$ <p>for all <math>\theta</math>, and <math>T</math> is an unbiased estimator of <math>g(\theta)</math> with finite variance, then</p> $\text{Var}_\theta(T) \geq \frac{[\dot{g}(\theta)]^2}{I(\theta)}$	<p>Suppose the support of <math>P_\theta</math> does not depend on <math>\theta</math>, and the differentiability in <math>L^2</math> condition holds. Further assume that <math>T</math> is an unbiased estimator of <math>g(\theta)</math> and it reaches the Cramér-Rao lower bound. What can you say about <math>p_\theta(x)</math> ? What about <math>g(\theta)</math> ?</p>
<p>DEFINITION</p> <p>What is a pivot for a test?</p>	<p>How is a pivot <math>Z(X, \gamma)</math> used to construct a test?</p>
<p>Let <math>X</math> be distributed according to a known symmetric distribution function <math>F_0</math> shifted by a parameter <math>\theta = \mu</math>. If <math>\hat{\mu}</math> is an equivariant estimator of <math>\mu</math>, what is a pivot function you could take to test if <math>\mu = \mu_0</math> ?</p>	<p>Let <math>X_1, \dots, X_n</math> be distributed i.i.d. according to a distribution function from family <math>\mathcal{F}_0 = \{F_0(\cdot) = \Phi(\cdot/\sigma) : \sigma &gt; 0\}</math> shifted by a parameter <math>\theta = \mu</math>. What is a pivot function you could take to test if <math>\mu = \mu_0</math> ?</p>

$ \begin{aligned} \frac{g(\theta+h)-g(\theta)}{h} &= \frac{\mathbb{E}_{\theta+h}T - \mathbb{E}_{\theta}T}{h} = \int T \frac{p_{\theta+h} - p_{\theta}}{h} \\ &= \int T \left[ \frac{(p_{\theta+h} - p_{\theta}) p_{\theta}}{h p_{\theta}} \right] \\ &= \int T \left[ \underbrace{\frac{(p_{\theta+h} - p_{\theta}) p_{\theta}}{h p_{\theta}} - s_{\theta} p_{\theta}}_{(1)} \right] \\ &\quad + \underbrace{\int T s_{\theta} p_{\theta} (-\mathbb{E}_{\theta}T s_{\theta} + \mathbb{E}_{\theta}T s_{\theta})}_{(2)} \end{aligned} $	$ \begin{aligned} \ddot{d}(\theta) &= \frac{1}{\underbrace{\int \exp[\theta T] h}_{e^{d(\theta)}}} \int \exp[\theta T] T^2 h \\ &\quad - \frac{1}{\underbrace{(\int \exp[\theta T] h)^2}_{(e^{d(\theta)})^2}} \left( \int \exp[\theta T] T h \right)^2 \\ &= \int \exp[\theta T - d(\theta)] h T^2 - \left( \int p_{\theta}(T) \right)^2 \\ &= \mathbb{E}_{\theta} T^2 - (\mathbb{E}_{\theta} T)^2 = \text{Var}_{\theta}(T) \end{aligned} $
<p>2. This term is the covariance because</p> $ \begin{aligned} \text{Cov}_{\theta}(T, s_{\theta}(X)) &= \mathbb{E}_{\theta} T s_{\theta} - (\mathbb{E}_{\theta} T) \underbrace{(\mathbb{E}_{\theta} s_{\theta})}_{=0} \\ &= \mathbb{E}_{\theta} T s_{\theta} = \int T s_{\theta} p_{\theta}. \end{aligned} $ $ \begin{aligned} \dot{g}(\theta) &= \lim_{h \rightarrow 0} \frac{g(\theta+h) - g(\theta)}{h} \\ &= \mathbb{E}_{\theta} T s_{\theta} = \text{Cov}_{\theta}(T, s_{\theta}(X)). \end{aligned} $	<p>Now we have</p> <p>1. By squaring (1) we can use Cauchy-Schwarz bound on it.</p> $ \begin{aligned} &\left( \int T \left[ \frac{(p_{\theta+h} - p_{\theta}) p_{\theta}}{h p_{\theta}} - s_{\theta} p_{\theta} \right] \right)^2 \leq \\ &\underbrace{\left( \int T^2 p_{\theta} \right)}_{< \infty \text{ by finite var.}} \underbrace{\left( \int \left[ \frac{p_{\theta+h} - p_{\theta}}{h p_{\theta}} - s_{\theta} \right]^2 p_{\theta} \right)}_{\xrightarrow{h \rightarrow 0} 0 \text{ by assum. on } s_{\theta}} \end{aligned} $
<p>By a lemma we know that</p> <ol style="list-style-type: none"> <li>1. <math>p_{\theta}(x) = \exp[c(\theta)T(X) - d(\theta)]h(x)</math></li> <li>2. that <math>c(\theta)</math> and <math>d(\theta)</math> are differentiable, and</li> <li>3. that <math>g(\theta) = \dot{d}(\theta)/\dot{c}(\theta)</math> for all <math>\theta</math></li> </ol>	<p>It is possible for us to show that <math>\dot{g}(\theta) = \text{Cov}_{\theta}(T, s_{\theta}(X))</math>. From this point we can just use Cauchy-Schwarz</p> $\text{Cov}_{\theta}(T, s_{\theta}(X))^2 \leq \text{Var}_{\theta}(T) \text{Var}_{\theta}(s_{\theta}(X))$ <p>Hence</p> $[\dot{g}(\theta)]^2 \leq \text{Var}_{\theta}(T) I(\theta)$ <p>so we get the desired</p> $\text{Var}_{\theta}(T) \geq \frac{[\dot{g}(\theta)]^2}{I(\theta)}$
<p>Let <math>G(\cdot)</math> be the distribution</p> $G(\cdot) := P_{\theta}(Z(X, g(\theta)) \leq \cdot).$ <p>By definition of pivot, it does not depend on <math>X</math>. Now for a test <math>g(\theta) \neq \gamma_0</math> of level <math>\alpha</math> we compute the critical values</p> $q_L := q_{\sup}^G\left(\frac{\alpha}{2}\right), q_R := q_{\inf}^G\left(1 - \frac{\alpha}{2}\right)$ <p>and define the test</p> $\phi(X, \gamma_0) := \begin{cases} 1 & \text{if } Z(X, \gamma_0) \notin [q_L, q_R] \\ 0 & \text{otherwise.} \end{cases}$	<p>A pivot is a function <math>Z(X, \gamma)</math> depending on data <math>X</math> and on the parameter <math>\gamma</math>, such that for all <math>\theta \in \Theta</math>, the distribution</p> $P_{\theta}(Z(X, g(\theta)) \leq \cdot) =: G(\cdot)$ <p>does not depend on <math>\theta</math>.</p>
<p>Recall that <math>\bar{X}_n</math> is equivariant, so the distribution of <math>\bar{X}_n - \mu</math> does not depend on <math>\mu</math>. To make it a pivot, we need to make it not depend on <math>\sigma</math> as well. So, we take</p> $Z(X, \mu) := \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n}$ <p>where <math>S_n^2</math> is the sample variance. Then <math>Z(X, \mu_0)</math> is distributed according to the Student distribution with <math>n - 1</math> degrees of freedom.</p>	<p>We could take <math>Z(X, \mu) := \hat{\mu} - \mu</math> as the pivot. By equivariance, this function has a distribution <math>G</math> depending only on <math>F_0</math>.</p>

<p>Let <math>X_1, \dots, X_n</math> be distributed i.i.d. according to distribution function from family <math>\mathcal{F}_0 = \{F_0 \text{ symmetric and continuous at } x = 0\}</math> shifted by a parameter <math>\theta = \mu</math>. What is a pivot function you could take to test if <math>\mu = \mu_0</math> ?</p>	<p>Consider two samples  <math>X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)</math> and  <math>Y_1, \dots, Y_m \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu + \gamma, \sigma^2)</math>. Find a pivot for testing <math>\gamma = \gamma_0</math>. How is such a test called?  (part-1)</p>
<p>Consider two samples  <math>X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)</math> and  <math>Y_1, \dots, Y_m \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu + \gamma, \sigma^2)</math>. Find a pivot for testing <math>\gamma = \gamma_0</math>. How is such a test called?  (part-2)</p>	<p>Consider two samples <math>X_1, \dots, X_n</math> and <math>Y_1, \dots, Y_m</math> which are sampled i.i.d. from distributions with distribution functions <math>F_X(\cdot)</math> and <math>F_Y(\cdot) = F_X(\cdot - \gamma)</math> correspondingly. Find a pivot for testing <math>\gamma = \gamma_0</math>. How is such a test called?</p>
<p>Let <math>V \in \mathbb{R}^{k \times k}</math> be positive definite. Let <math>c \in \mathbb{R}^k</math>  Show that</p> $\max_{a \in \mathbb{R}^k} \frac{(a^T c)^2}{a^T V a} = c^T V^{-1} c$	<p>What is the Cramér-Rao lower bound for an unbiased estimator <math>T \in \mathbb{R}</math> if the parameter space <math>\theta \subset \mathbb{R}^k</math> is multidimensional?</p>
<p>What's a lower bound <math>A</math> such that for the covariance matrix <math>\Sigma_\theta</math> of an unbiased estimator <math>T</math> we have</p> $\Sigma_\theta - A \succeq 0?$	<p>How is the <math>i</math>th diagonal element of the Fisher information matrix <math>I^i(\theta)</math> related to the <math>i</math>th diagonal element of its inverse <math>I_{ii}(\theta)^{-1}</math> ?</p>
<p>What is the power of a test <math>\phi</math> ?</p>	<p>If <math>X_1, \dots, X_n</math> and <math>Y_1, \dots, Y_n</math> are i.i.d. sequences of random variables with distributions <math>F_X \sim \mathcal{N}(\mu, \sigma^2)</math> and <math>F_Y \sim \mathcal{N}(\mu + \gamma, \sigma^2)</math> respectively. How is <math>\bar{Y} - \bar{X}</math> distributed?</p>

<p>Note that <math>\bar{Y} - \bar{X} \sim \mathcal{N}\left(\gamma, \frac{\sigma^2}{m} + \frac{\sigma^2}{n}\right)</math></p> $\frac{\bar{Y} - \bar{X} - \gamma}{\sigma \sqrt{\frac{m+n}{mn}}} \sim \mathcal{N}(0, 1).$ <p>Now to find the pivot itself, we need to find a distribution not depending on <math>\sigma</math>. So, we replace <math>\sigma</math> by the pooled sample variance</p> $S^2 := \frac{1}{m+n-2} \left[ \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (Y_j - \bar{Y})^2 \right]$	<p>We could take</p> $Z(X, \mu) := \sum_{i=1}^n 1 \{X_i \geq \mu\}$ <p>which is equivariant for <math>\mu</math>. <math>G</math> is then distributed according to the Binomial(<math>n, p</math>) distribution with parameter <math>p = 1/2</math>, which does not depend on <math>\mu</math>, so it's a pivot.</p>
<p>Let <math>N = n + m</math>, and then <math>Z_1, \dots, Z_N = (X_1, \dots, X_n, Y_1, \dots, Y_m)</math>. Let <math>R_i</math> be the rank of <math>Z_i</math> if we were to sort all of the <math>Z_i</math>'s. Then the pivot is</p> $Z(X, Y, \gamma) = \sum_{i=1}^n R_i$ <p>This is the Wilcoxon test.</p>	<p>Now finally our pivot is</p> $Z(X, Y, \gamma) := \sqrt{\frac{mn}{m+n}} \left[ \frac{\bar{Y} - \bar{X} - \gamma}{S} \right].$ <p>Its distribution <math>G</math> is Student (<math>n + m - 2</math>). This is the so called Student's test.</p>
$\text{Var}(T) \geq \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta)$	<p>Let <math>b := V^{\frac{1}{2}}a, d := V^{-\frac{1}{2}}c</math>. Then <math>a^T V a = b^T b = \ b\ ^2</math> and <math>a^T c = b^T d</math>. Now we write</p> $\frac{(a^T c)^2}{a^T V a} = \frac{ b^T d ^2}{\underbrace{\ b\ ^2}_{\text{Cauchy-Schwarz}}} \leq \frac{\ b\ ^2 \ d\ ^2}{\ b\ ^2} = \ d\ ^2 = d^T d = c^T V^{-1} c$
$I^{ii}(\theta) \geq I_{ii}(\theta)^{-1}$	<p>It's the inverse of the Fisher information <math>I(\theta)^{-1}</math>, so we have</p> $\Sigma_\theta - I(\theta)^{-1} \succeq 0$
<p>It's</p> $\bar{Y} - \bar{X} \sim \mathcal{N}\left(\gamma, \frac{\sigma^2}{m} + \frac{\sigma^2}{n}\right)$	<p>The power of a test <math>\phi(x, \gamma_0)</math> is <math>\mathbb{P}_\theta(\phi(X, \gamma_0) = 1)</math> for <math>g(\theta) \neq \gamma_0</math>.</p>

<p>DEFINITION</p> <p><i>What is a randomized test?</i></p>	<p><i>How are non-randomized tests related to randomized tests?</i></p>
<p><i>What is the risk of a randomized test <math>\phi</math> when <math>\theta \in \{\theta_0, \theta_1\}</math> ?</i></p>	<p>DEFINITION</p> <p><i>What is a Neyman Pearson test?</i></p>
<p>DEFINITION</p> <p><i>What is a Neyman Pearson test? (plot)</i></p>	<p>LEMMA</p> <p><i>State the Neyman-Pearson lemma.</i></p>
<p>PROOF</p> <p><i>Sketch a proof of the Neyman-Pearson lemma (part-1)</i></p>	<p>PROOF</p> <p><i>Sketch a proof of the Neyman-Pearson lemma (part-2)</i></p>
<p><i>What is the level of a test?</i></p>	<p><i>What does it mean for a test to be uniformly most powerful?</i></p>

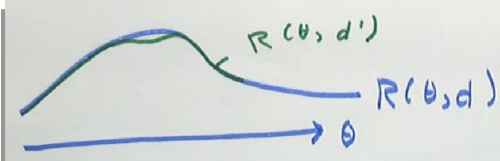
<p>A non-randomized test is a special case of randomized tests where <math>\phi \in \{0, 1\}</math>, and <math>\mathbb{E}(\phi(X)) = \mathbb{P}(\phi(X) = 1)</math></p>	<p>A randomized test at level <math>\alpha</math> is a statistic <math>\phi : x \rightarrow [0, 1]</math> such that <math>\mathbb{E}(\phi(X)) \leq \alpha</math> for <math>\phi</math> under the null hypothesis.</p>
<p>Consider the problem of testing whether <math>\theta = \theta_0</math> or <math>\theta = \theta_1</math>. Let <math>p_0(p_1)</math> be the density of <math>P_{\theta_0}(P_{\theta_1})</math> with respect to some dominating measure <math>\nu</math> (for example <math>\nu = P_{\theta_0} + P_{\theta_1}</math>). A Neyman Pearson test is then</p> $\phi_{\text{NP}}(x) = \begin{cases} 1 & \text{if } \frac{p_1(x)}{p_0(x)} > c \\ q & \text{if } \frac{p_1(x)}{p_0(x)} = c \\ 0 & \text{if } \frac{p_1(x)}{p_0(x)} < c \end{cases}$	$R(\theta, \phi) = \begin{cases} \mathbb{E}_{\theta_0} \phi(x) & \text{if } \theta = \theta_0 \\ 1 - \mathbb{E}_{\theta_1} \phi(x) & \text{if } \theta = \theta_1 \end{cases}$
<p>If <math>\theta \in \{\theta_0, \theta_1\}</math>, then for any test <math>\phi</math> and for a Neyman Pearson test <math>\phi_{\text{NP}}</math> with threshold <math>c</math> the risks fulfil the inequality</p> $[R(\theta_1, \phi) - R(\theta_1, \phi_{\text{NP}})] \geq c [R(\theta_0, \phi_{\text{NP}}) - R(\theta_0, \phi)]$	
$\begin{aligned} &\geq \int_{\frac{p_1}{p_0} > c} [\phi_{\text{NP}} - \phi] c p_{\theta_0} + \int_{\frac{p_1}{p_0} = c} [\phi_{\text{NP}} - \phi] c p_{\theta_0} \\ &+ \int_{\frac{p_1}{p_0} < c} [\phi_{\text{NP}} - \phi] c p_{\theta_0} \\ &= c [\mathbb{E}_{\theta_0} \phi_{\text{NP}}(x) - \mathbb{E}_{\theta_0} \phi(x)] \\ &= c [R(\theta_0, \phi_{\text{NP}}) - R(\theta_0, \phi)] . \end{aligned}$	$\begin{aligned} R(\theta_1, \phi) - R(\theta_1, \phi_{\text{NP}}) &= \mathbb{E}_{\theta_1} \phi_{\text{NP}}(x) - \mathbb{E}_{\theta_1} \phi(x) \\ &= \int [\phi_{\text{NP}} - \phi] p_{\theta_1} \\ &= \int_{\frac{p_1}{p_0} > c} \underbrace{[\phi_{\text{NP}} - \phi]}_{1 - \phi \geq 0} \underbrace{p_{\theta_1}}_{\geq c p_0} \\ &+ \int_{\frac{p_1}{p_0} = c} [\phi_{\text{NP}} - \phi] \underbrace{p_{\theta_1}}_{= c p_0} + \int_{\frac{p_1}{p_0} < c} \underbrace{[\phi_{\text{NP}} - \phi]}_{0 - \phi \leq 0} \underbrace{p_{\theta_1}}_{< c p_0} \end{aligned}$
<p>Consider <math>\theta \in \Theta</math> in a parameter space with a null hypothesis <math>H_0 : \theta \in \Theta_0 \subset \Theta</math> and alternative hypothesis <math>H_1 : \theta \in \Theta_1 \subset \Theta</math> which are disjoint. A test <math>\phi</math> is uniformly most powerful (UMP) at level <math>\alpha</math> if <math>\text{level}(\phi) \leq \alpha</math> and</p> $\mathbb{E}_{\theta} \phi(x) = \sup \{ \mathbb{E}_{\theta} \phi' : \text{level}(\phi') \leq \alpha \}, \forall \theta \in \Theta_1 .$ <p>So, it attains the supremum of power over all tests of this level.</p>	<p>Consider <math>\theta \in \Theta</math> in a parameter space with a null hypothesis <math>H_0 : \theta \in \Theta_0 \subset \Theta</math> and alternative hypothesis <math>H_1 : \theta \in \Theta_1 \subset \Theta</math> which are disjoint. Then the level of a test <math>\phi \in [0, 1]</math> is</p> $\text{level}(\phi) := \sup_{\theta \in \Theta_0} \underbrace{\mathbb{E}_{\theta} \phi(x)}_{\text{Probability of error of first kind}} .$



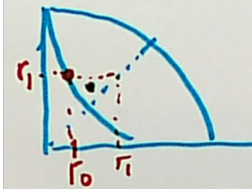
<p>What does it mean for a test <math>\phi</math> to be unbiased?</p>	<p>DEFINITION</p> <p>What is a uniformly most powerful unbiased (UMPU) test?</p>
<p>LEMMA</p> <p>State the Rao Blackwell lemma.</p>	<p>PROOF</p> <p>Let <math>d : \mathcal{X} \rightarrow \mathcal{A}</math> be a decision and <math>d^* := \mathbb{E}(d(X) \mid S)</math> where <math>S</math> is sufficient. Suppose <math>\mathcal{A}</math> is a convex subset of <math>\mathbb{R}^p</math>, and <math>a \mapsto L(\theta, a)</math> is convex for all <math>\theta</math>. Then for all <math>\theta</math></p> $R(\theta, d^*) \leq R(\theta, d)$
<p>DEFINITION</p> <p>What does it mean for an estimator <math>T</math> to be location equivariant?</p>	<p>What are some examples of location equivariant estimators?</p>
<p>What does it mean for a loss function <math>L</math> to be called location invariant?</p>	<p>Show that under a location invariant loss, the risk of a location equivariant estimator does not depend on its location parameter.</p>
<p>DEFINITION</p> <p>What is a uniform minimum risk equivariant (UMRE) estimator?</p>	<p>LEMMA</p> <p>State lemma about construction of UMRE estimators.</p>

<p>An UMPU test is uniformly most powerful among all <math>\phi'</math> unbiased with level <math>(\phi') \leq \alpha</math></p>	<p>A test <math>\phi</math> is called unbiased if its power is at least the level, in other words if for all <math>\theta_1 \in \Theta_1</math></p> $\mathbb{E}_{\theta_1} \phi(X) \geq \sup_{\theta \in \Theta_0} \mathbb{E}_{\theta} \phi(X) = \text{level}(\phi)$
$\begin{aligned} R(\theta, d) &= \mathbb{E}_{\theta} L(\theta, d(X)) \\ &= \mathbb{E}_{\theta} \mathbb{E}(L(\theta, d(X)) \mid S) \quad (\text{Iterated expectations}) \\ &\geq \mathbb{E}_{\theta} L(\theta, \mathbb{E}(d(X) \mid S)) \quad (\text{Jensen}) \\ &= \mathbb{E}_{\theta} L(\theta, d^*(S)) \\ &= R(\theta, d^*) \end{aligned}$	<p>Let <math>d : \mathcal{X} \rightarrow \mathcal{A}</math> be a decision and <math>d^* := \mathbb{E}(d(X) \mid S)</math> where <math>S</math> is sufficient. Suppose <math>\mathcal{A}</math> is a convex subset of <math>\mathbb{R}^p</math>, and <math>a \mapsto L(\theta, a)</math> is convex for all <math>\theta</math>. Then for all <math>\theta</math></p> $R(\theta, d^*) \leq R(\theta, d)$
<p>Mean, median.</p>	<p>An estimator <math>T</math> is called location equivariant if <math>T(x + c) = T(x) + c</math> for all <math>c \in \mathbb{R}, x \in \mathbb{R}^n</math></p>
$\begin{aligned} R(\theta, T) &= \mathbb{E}_{\theta} L(\theta, T(X)) \\ &= \mathbb{E}_{\theta} L_0(T(X) - \theta) \quad (\text{Loss invariance}) \\ &= \mathbb{E}_{\theta} L_0(T(X - \theta)) \quad (\text{Estimator equivariance}) \\ &= \mathbb{E}_{\theta} L_0(T(\varepsilon)) \\ &= \mathbb{E}_0 L_0(T(\varepsilon)) \\ &= R(0, T). \end{aligned}$	<p>A loss function <math>L : \underbrace{\mathbb{R}}_{\Theta} \times \underbrace{\mathbb{R}}_{\mathcal{A}}</math> is called location invariant if for all <math>(\theta, a) \in \mathbb{R}^2</math>,</p> $L(\theta, a) = L(\theta + c, a + c).$ <p>We can also write</p> $L(\theta, a) = L(0, a - \theta) =: L_0(a - \theta)$
<p>Let <math>T</math> be an equivariant estimator, <math>Y_i := X_i - X_n, i = 1, \dots, n, Y := (Y_1, \dots, Y_n)</math>, and define</p> $T^*(Y) = \arg \min_v \mathbb{E}(L_0(v + \varepsilon_n) \mid Y).$ <p>Moreover, let <math>T^*(X) := T^*(Y) + X_n</math>. Then <math>T^*</math> is UMRE.</p>	<p>An estimator <math>T</math> is called uniform minimum risk equivariant if it is equivariant and</p> $R(0, T) = \min \{ R(0, T') : T' \text{ equivariant} \}$

<p>LEMMA</p> <p><i>State Basu's lemma.</i></p>	<p>PROOF</p> <p><i>Consider a random variable <math>X \sim P_\theta, \theta \in \Theta</math>. Let <math>T = T(X)</math> be a sufficient and complete statistic, and let <math>Y = Y(X)</math> have a distribution not depending on <math>\theta</math>. Then <math>T</math> and <math>Y</math> are independent.</i></p>
<p><i>If an estimator <math>T</math> is unbiased, sufficient, complete, equivariant, what other property does it have?</i></p>	<p><i>Let <math>X_1, \dots, X_n</math> be i.i.d. variables distributed according to <math>\mathcal{N}(\mu, \sigma^2)</math>. Assume also that <math>\sigma^2 = \sigma_0^2</math> is known. Show that <math>T(X) = \bar{X}</math> is an UMRE estimator of <math>\mu</math></i></p>
<p><i>Let <math>X_1, \dots, X_n</math> be i.i.d. variables distributed according to <math>\mathcal{N}(\mu, \sigma^2)</math>. Assume also that <math>\sigma^2 = \sigma_0^2</math> is known. What is <math>\mathbb{E}_\theta[\bar{X} \mid X - \bar{X}]</math> ? Why?</i></p>	<p><i>What does it mean for an estimator <math>d'</math> to be strictly better than <math>d</math> ?</i></p>
<p><i>What does it mean for an estimator <math>d</math> to be admissible?</i></p>	<p><i>What does it mean that <math>\tilde{P} \gg P</math> ?</i></p>
<p><i>What does it mean for a decision <math>d</math> to be minimax?</i></p>	<p><i>Show that if for a Neyman-Pearson test <math>\phi_{NP}</math> it's true that</i></p> $R(\theta_0, \phi_{NP}) = R(\theta_1, \phi_{NP}),$ <p><i>then <math>\phi_{NP}</math> is minimax.</i></p>

<ol style="list-style-type: none"> <li>1. Label <math>\mathcal{Y}</math> to be the space of values of <math>Y</math>, and take an arbitrary subset. <math>A \subset \mathcal{Y}</math>.</li> <li>2. Define <math>h(T) := \mathbb{P}(Y \in A   T) - \mathbb{P}(Y \in A)</math>.</li> <li>3. By iterated expectations, <math>\mathbb{E}_\theta h(T) = 0</math> for all <math>\theta</math>.</li> <li>4. By completeness of <math>T</math>, <math>h(T) = 0</math> almost surely, so <math display="block">\mathbb{P}(Y \in A   T) - \mathbb{P}(Y \in A) = 0</math> almost surely.</li> <li>5. Thus, <math>T</math> and <math>Y</math> are independent.</li> </ol>	<p>Consider a random variable <math>X \sim P_\theta, \theta \in \Theta</math>. Let <math>T = T(X)</math> be a sufficient and complete statistic, and let <math>Y = Y(X)</math> have a distribution not depending on <math>\theta</math>. Then <math>T</math> and <math>Y</math> are independent.</p>
<p>Note that <math>T</math> is unbiased and equivariant. We know it is also sufficient and complete from properties of exponential family distributions. From these four properties we can conclude it is UMRE.</p>	<p>It is also uniform minimum risk equivariant (UMRE).</p>
<p>An estimator <math>d'</math> is called strictly better than <math>d</math> if</p> <ol style="list-style-type: none"> <li>1. <math>R(\theta, d') \leq R(\theta, d)</math> for all <math>\theta</math>, and</li> <li>2. <math>R(\theta, d') &lt; R(\theta, d)</math> for at least one <math>\theta</math>.</li> </ol> 	<p>Since <math>\bar{X}</math> is sufficient and complete, we can use Basu's lemma to conclude</p> $\mathbb{E}_\theta[\bar{X}   X - \bar{X}] = \mathbb{E}_\theta \bar{X} = \mu$
<p><math>\tilde{P} \gg P</math> means "<math>\tilde{P}</math> dominates <math>P</math>", so</p> $\tilde{P}(A) = 0 \implies P(A) = 0$	<p>An estimator <math>d</math> is admissible if there is no strictly better estimator <math>d'</math> than it.</p>
<p>Suppose <math>R(\theta_0, \phi_{NP}) = R(\theta_1, \phi_{NP}) = r</math>. Consider another test <math>\phi'</math> with</p> $r' = \max \{ R(\theta_0, \phi'), R(\theta_1, \phi') \}.$ <p>By Neyman-Pearson lemma</p> $[r - R(\theta_1, \phi)] \leq c [R(\theta_0, \phi) - r]$ $r(1 + c) \leq cR(\theta_0, \phi) + R(\theta_1, \phi) \leq r'(1 + c)$ $r \leq r'$ <p>so <math>\phi_{NP}</math> is minimax.</p>	<p>A decision <math>d</math> is minimax if</p> $\sup_{\theta \in \Theta} R(\theta, d) = \inf_{d'} \sup_{\theta \in \Theta} R(\theta, d')$

<p>Suppose that a Neyman-Pearson test <math>\phi_{NP}</math> is minimax. Show that then</p> $R(\theta_0, \phi_{NP}) = R(\theta_1, \phi_{NP})$	<p>DEFINITION</p> <p>What is a Bayes risk?</p>
<p>DEFINITION</p> <p>What is a Bayes estimator?</p>	<p>Consider <math>\Theta = \{\theta_0, \theta_1\}</math> with prior probabilities <math>w(\theta_0) = w_0, w(\theta_1) = w_1 = 1 - w_0</math>. The risk of a decision <math>\phi</math> is then</p> $R(\theta, \phi) = \begin{cases} \mathbb{E}_{\theta_0} \phi(X) & \text{if } \theta = \theta_0 \\ 1 - \mathbb{E}_{\theta_1} \phi(X) & \text{if } \theta = \theta_1 \end{cases}$ <p>Derive a decision <math>\phi : \mathcal{X} \rightarrow [0, 1]</math> which minimizes the Bayes risk. (part-1)</p>
<p>Consider <math>\Theta = \{\theta_0, \theta_1\}</math> with prior probabilities <math>w(\theta_0) = w_0, w(\theta_1) = w_1 = 1 - w_0</math>. The risk of a decision <math>\phi</math> is then</p> $R(\theta, \phi) = \begin{cases} \mathbb{E}_{\theta_0} \phi(X) & \text{if } \theta = \theta_0 \\ 1 - \mathbb{E}_{\theta_1} \phi(X) & \text{if } \theta = \theta_1 \end{cases}$ <p>Derive a decision <math>\phi : \mathcal{X} \rightarrow [0, 1]</math> which minimizes the Bayes risk. (part-2)</p>	<p>LEMMA</p> <p>State the Bayes estimator construction lemma.</p>
<p>PROOF</p> <p>Given data <math>X = x</math>, consider <math>\theta</math> as a random variable with density <math>w(\vartheta   x)</math>. Let</p> $l(x, a) := \mathbb{E}(L(\theta, a)   X = x)$ $= \int_{\Theta} L(\vartheta, a) w(\vartheta   x) d\mu(\vartheta),$ <p><math>d(x) := \arg \min_a l(x, a)</math>. Then <math>d</math> is the Bayes decision <math>d_{\text{Bayes}}</math>.</p>	<p>What is a maximum a posteriori estimator?</p>
<p>What are the posterior, likelihood, and prior?</p>	<p>Consider the action space <math>\mathcal{A} = \mathbb{R}</math>, parameter space <math>\Theta \subset \mathbb{R}</math>, and loss function <math>L(\theta, a) = (\theta - a)^2</math>. What is <math>d_{\text{Bayes}}</math> for <math>\theta</math>? Why?</p>

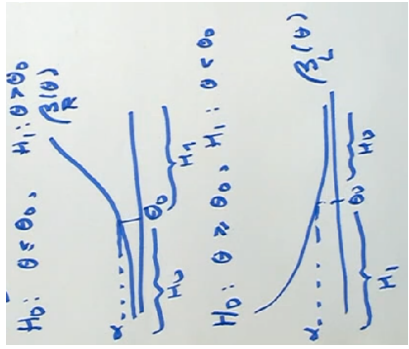
<p>Let <math>\Theta \subset \mathbb{R}^k</math>. Let <math>w(\theta)</math> be given weights such that <math>\int_{\Theta} w(\theta) d\theta = 1</math>. The Bayes risk of a decision <math>d</math> is then</p> $r_w(d) = \int_{\Theta} R(\theta, d) d\mu(\theta).$ <p>So, it's the expected risk when <math>\theta</math> is a random variable with density <math>w</math>.</p>	<p>Let <math>S = \{(R(\theta_0, \phi), R(\theta_1, \phi)) : \phi : \mathcal{X} \rightarrow [0, 1]\}</math>. Note that <math>S</math> is convex (see image). Thus, if <math>r_0 &lt; r_1</math>, we can find a test <math>\phi</math> with <math>r_0 &lt; r'_0 &lt; r_1</math> and <math>r'_1 &lt; r_1</math>. So then <math>\phi_{NP}</math> is not minimax. Similarly for <math>r_0 &gt; r_1</math>.</p> 
$\begin{aligned} r_w(\phi) &= w_0 \mathbb{E}_{\theta_0} \phi(X) + w_1 (1 - \mathbb{E}_{\theta_1} \phi(X)) \\ &= \int (w_0 p_0 - w_1 p_1) \phi + w_1 \\ &= \int_{w_0 p_0 - w_1 p_1 > 0} (w_0 p_0 - w_1 p_1) \phi \\ &\quad + \int_{w_0 p_0 - w_1 p_1 = 0} (w_0 p_0 - w_1 p_1) \phi \\ &\quad + \int_{w_0 p_0 - w_1 p_1 < 0} (w_0 p_0 - w_1 p_1) \phi + w_1 \end{aligned}$	<p>A Bayes estimator is the one which minimizes the Bayes risk for some prior</p> $d_{\text{Bayes}} = \arg \min_{d'} r_w(d')$
<p>Given data <math>X = x</math>, consider <math>\theta</math> as a random variable with density <math>w(\vartheta   x)</math>. Let</p> $l(x, a) := \mathbb{E}(L(\theta, a)   X = x) = \int_{\Theta} L(\vartheta, a) w(\vartheta   x) d\mu(\vartheta)$ <p>and</p> $d(x) := \arg \min_a l(x, a)$ <p>Then <math>d</math> is the Bayes decision <math>d_{\text{Bayes}}</math>.</p>	<p>So the optimal decision is</p> $\phi_{\text{Bayes}} = \begin{cases} 1 & \text{if } p_1/p_0 > w_0/w_1 \\ q & \text{if } p_1/p_0 = w_0/w_1 \\ 0 & \text{if } p_1/p_0 < w_0/w_1 \end{cases}$
<p>The maximum a posteriori estimator is</p> $\theta_{\text{MAP}}(x) = \arg \max_{\vartheta \in \Theta} w(\vartheta   x)$	<p>Let <math>d'</math> be some decision.</p> $\begin{aligned} r_w(d') &= \int_{\Theta} R(\vartheta, d') w(\vartheta) d\mu(\vartheta) \\ &= \mathbb{E}(R(\vartheta, d')) \\ &= \mathbb{E}(\mathbb{E}(L(\vartheta, d')   \vartheta)) \\ &= \mathbb{E}(L(\vartheta, d')) && \text{(Iterated } \mathbb{E}) \\ &= \mathbb{E}(\mathbb{E}(L(\vartheta, d')   X)) && \text{(Rev. iterated } \mathbb{E}) \\ &= \mathbb{E}(l(X, d')) \\ &\geq \mathbb{E}(l(X, d)) \\ &= r_w(d) \quad d = \arg \min_{d'} \dots \end{aligned}$
<p><math>d_{\text{Bayes}}(x) = \mathbb{E}(\theta   x)</math>. This is because we are minimizing</p> $l(x, a) = \mathbb{E}([\theta - a]^2   X = x)$ <p>and the minimum is simply</p> $\min_{a \in \mathbb{R}} l(x, a) = \mathbb{E}(\theta   X = x)$	<p>In the Bayesian context, we have</p> $\underbrace{w(\vartheta   x)}_{\text{posterior}} \propto \underbrace{p(x   \vartheta)}_{\text{likelihood}} \underbrace{w(\vartheta)}_{\text{prior}}$

<p>Show that for quadratic loss, the Bayes risk of any estimator can be written as</p> $r_w(T') = r_w(T_{\text{Bayes}}) + \mathbb{E} \left( [T_{\text{Bayes}}(X) - T'(X)]^2 \right)$	<p>Consider <math>X   \theta \sim \text{Poisson}(\theta), \theta \sim \text{Gamma}(k, \lambda)</math> with <math>k, \lambda</math> given. The mass functions are</p> <ol style="list-style-type: none"> <li>1. <math>p(x   \theta) = e^{-\theta} \frac{\theta^x}{x!}</math> for <math>x \in \{0, 1, \dots\}</math></li> <li>2. <math>w(\vartheta) = e^{-\lambda\vartheta} \vartheta^{k-1} \lambda^k / \Gamma(k)</math></li> </ol> <p>Also note that <math>\mathbb{E}(\theta) = \frac{k}{\lambda}</math>. What is the Bayes estimator for the mean squared error loss? What is the maximum a posteriori estimator for <math>\theta</math>? (part-1)</p>
<p>Consider <math>X   \theta \sim \text{Poisson}(\theta), \theta \sim \text{Gamma}(k, \lambda)</math> with <math>k, \lambda</math> given. The mass functions are</p> <ol style="list-style-type: none"> <li>1. <math>p(x   \theta) = e^{-\theta} \frac{\theta^x}{x!}</math> for <math>x \in \{0, 1, \dots\}</math></li> <li>2. <math>w(\vartheta) = e^{-\lambda\vartheta} \vartheta^{k-1} \lambda^k / \Gamma(k)</math></li> </ol> <p>Also note that <math>\mathbb{E}(\theta) = \frac{k}{\lambda}</math>. What is the Bayes estimator for the mean squared error loss? What is the maximum a posteriori estimator for <math>\theta</math>? (part-2)</p>	<p>Consider <math>X   \theta \sim \text{Binomial}(n, \theta), \theta \sim \text{Beta}(r, s)</math> with <math>r, s &gt; 0</math> given. The density/mass functions are</p> <ol style="list-style-type: none"> <li>1. <math>p(x   \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}</math> for <math>x = 0, 1, \dots, n</math>.</li> <li>2. <math>w(\theta) = \theta^{r-1} (1 - \theta)^{s-1} \frac{\Gamma(r+s)}{\Gamma(r)\Gamma(s)}, 0 &lt; \theta &lt; 1</math>.</li> </ol> <p>Also note that <math>\mathbb{E}(\theta) = \frac{r}{r+s}</math>. What is the Bayes estimator under the quadratic loss?</p>
<p>THEOREM</p> <p>State the theorem about UMP tests for <math>\theta</math> of exponential family distributions with increasing <math>c(\theta)</math>.</p>	<p>Suppose <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math> where <math>\theta \in \Theta</math> is an interval in <math>\mathbb{R}</math> and <math>c(\cdot)</math> is strictly increasing. Let</p> $\phi(T) = \begin{cases} 1 & \text{if } T > t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T < t_0 \end{cases}$ <p>where <math>t_0</math> and <math>q</math> are such that <math>\mathbb{E}_{\theta_0} \phi(T) = \alpha</math>. Let <math>\beta(\theta) = \mathbb{E}_\theta \phi(T)</math>. Show in the discrete case that <math>\theta \mapsto \beta(\theta)</math> is increasing.</p> <p>(part-1)</p>
<p>Suppose <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math> where <math>\theta \in \Theta</math> is an interval in <math>\mathbb{R}</math> and <math>c(\cdot)</math> is strictly increasing. Let</p> $\phi(T) = \begin{cases} 1 & \text{if } T > t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T < t_0 \end{cases}$ <p>where <math>t_0</math> and <math>q</math> are such that <math>\mathbb{E}_{\theta_0} \phi(T) = \alpha</math>. Let <math>\beta(\theta) = \mathbb{E}_\theta \phi(T)</math>. Show in the discrete case that <math>\theta \mapsto \beta(\theta)</math> is increasing.</p> <p>(part-2)</p>	<p>Suppose <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math> where <math>\theta \in \Theta</math> is an interval in <math>\mathbb{R}</math> and <math>c(\cdot)</math> is strictly increasing. Let</p> $\phi(T) = \begin{cases} 1 & \text{if } T > t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T < t_0 \end{cases}$ <p>where <math>t_0</math> and <math>q</math> are such that <math>\mathbb{E}_{\theta_0} \phi(T) = \alpha</math>. Let <math>\beta(\theta) = \mathbb{E}_\theta \phi(T)</math>. Show in the discrete case that <math>\theta \mapsto \beta(\theta)</math> is increasing.</p> <p>(part-3)</p>
<p>Suppose <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math> where <math>\theta \in \Theta</math> is an interval in <math>\mathbb{R}</math> and <math>c(\cdot)</math> is strictly increasing. Let</p> $\phi(T) = \begin{cases} 1 & \text{if } T > t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T < t_0 \end{cases}$ <p>where <math>t_0</math> and <math>q</math> are such that <math>\mathbb{E}_{\theta_0} \phi(T) = \alpha</math>. Let <math>\beta(\theta) = \mathbb{E}_\theta \phi(T)</math>. Knowing that <math>\theta \mapsto \beta(\theta)</math> is increasing, show that <math>\phi</math> is uniformly most powerful (UMP) for <math>\tilde{H}_0 : \theta = \theta_0, \tilde{H}_1 : \theta = \theta_1</math> for all <math>\theta_1 &gt; \theta_0</math>.</p>	<p>Suppose <math>X_1, \dots, X_n</math> are i.i.d. from Exponential(<math>\theta</math>), <math>\theta &gt; 0</math>. The probability density function for the one-dimensional exponential distribution is</p> $p_\theta(x) = \exp[-\theta x - (-\log \theta)] 1\{x > 0\}.$ <p>Consider a hypothesis <math>H_0 : \theta \leq \theta_0</math>, and <math>H_1 : \theta &gt; \theta_0</math>. What is a UMP test for <math>H_0</math>?</p> <p>(part-1)</p>

<p>The posterior is</p> $  \begin{aligned}  w(\vartheta   x) &\propto p(x   \vartheta)w(\vartheta) \\  &\propto e^{-\vartheta} \vartheta^x e^{-\lambda \vartheta} \vartheta^{k-1} \\  &= e^{-(1+\lambda)\vartheta} \vartheta^{x+k-1}  \end{aligned}  $ <p>so we can see that <math>\vartheta   X = x \sim \text{Gamma}(x+k, 1+\lambda)</math> with <math>\mathbb{E}(\theta   X) = \frac{X+k}{1+\lambda}</math>. This is the Bayes estimator for mean squared error loss.</p>	$  \begin{aligned}  r_w(T') &= \mathbb{E}L(\theta, T'(X)) \\  &= \mathbb{E}\mathbb{E}[L(\theta, T'(X))   X] \quad (\text{Iterated } \mathbb{E}) \\  &= \mathbb{E}\mathbb{E}[\theta - T'(X)]^2   X] \\  &= \underbrace{\mathbb{E} \text{Var}(\theta   X)}_{r_w(T_{\text{Bayes}})} + \mathbb{E} \left[ \underbrace{(\mathbb{E}(\theta   X) - T'(X))^2}_{T_{\text{Bayes}}(X)}   X \right] \\  &= r_w(T_{\text{Bayes}}) + \mathbb{E}([T_{\text{Bayes}}(X) - T'(X)]^2).  \end{aligned}  $ <p>4<sup>th</sup> equality: <math>MSE = Var + Bias^2</math></p>
<p>First compute the posterior</p> $  \begin{aligned}  w(\vartheta   x) &\propto \vartheta^x (1-\vartheta)^{n-x} \vartheta^{r-1} (1-\vartheta)^{s-1} \\  &= \vartheta^{x+r-1} (1-\vartheta)^{n-x+s-1}.  \end{aligned}  $ <p>We can see that it is <math>\theta   X = x \sim \text{Beta}(x+r, n-x+s)</math>. So, the Bayes estimator under quadratic loss is <math>\mathbb{E}(\theta   X) = \frac{X+r}{n+r+s}</math>.</p>	<p>Now we need to find the maximum of <math>w(\vartheta   x)</math>. We can do it under logarithm, so it's easier</p> $  \log w(\vartheta   x) = (x+k-1) \log \vartheta - (1+\lambda) \vartheta + \underbrace{c(x)}_{\text{normalizing constant to ignore}}  $ $  \frac{\delta}{\delta \vartheta} w(\vartheta   x) = \frac{x+k-1}{\vartheta} - (1+\lambda) \triangleq 0  $ $  \theta_{\text{MAP}}(x) = \frac{x+k-1}{1+\lambda}  $
<p>We have <math>P_\theta(T=t) = \exp[c(\theta)t - d(\theta)] \sum_{x:T(x)=t} h(x)</math>. For <math>\tilde{\theta} &gt; \theta</math> it's true that the ratio</p> $  \frac{P_{\tilde{\theta}}(T=t)}{P_\theta(T=t)} = \exp[\underbrace{\{c(\tilde{\theta}) - c(\theta)\}}_{>0} t - \{d(\tilde{\theta}) - d(\theta)\}]  $ <p>is increasing in <math>t</math>. To find the difference <math>\mathbb{E}_{\tilde{\theta}}\phi(T) - \mathbb{E}_\theta\phi(T)</math> we first need to consider whether <math>p_{\tilde{\theta}}(t)</math> is greater than, equal, or less than <math>p_\theta(t)</math>.</p>	<p>Suppose <math>\mathbb{P}</math> is a one-dimensional exponential family <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math>. Assume also that <math>c(\theta)</math> is a strictly increasing function of <math>\theta</math>. Then a UMP test is</p> $  \phi(T) = \begin{cases} 1 & \text{if } T > t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T < t_0, \end{cases}  $ <p>where <math>q</math> and <math>t_0</math> are chosen in such a way that <math>\mathbb{E}_\theta\phi(T) = \alpha</math>.</p>
$  \begin{aligned}  &= \sum_{t \leq s_0} \underbrace{\phi(t)}_{\leq \phi(s_0)} \underbrace{[p_{\tilde{\theta}}(t) - p_\theta(t)]}_{\leq 0} + \sum_{t > s_0} \underbrace{\phi(t)}_{\geq \phi(s_0)} \underbrace{[p_{\tilde{\theta}}(t) - p_\theta(t)]}_{\geq 0} \\  &\geq \phi(s_0) \sum_{t \leq s_0} [p_{\tilde{\theta}}(t) - p_\theta(t)] + \phi(s_0) \sum_{t > s_0} [p_{\tilde{\theta}}(t) - p_\theta(t)] \\  &= \phi(s_0) \underbrace{\sum_t [p_{\tilde{\theta}}(t) - p_\theta(t)]}_{=0} = 0.  \end{aligned}  $	<p>Note that if <math>p_{\tilde{\theta}}(t) &lt; p_\theta(t)</math> for all <math>t</math>, then we would also have</p> $  1 = \sum_t p_{\tilde{\theta}}(t) < \sum_t p_\theta(t) = 1  $ <p>which is a contradiction. The same happens if <math>p_{\tilde{\theta}}(t) &gt; p_\theta(t)</math> for all <math>t</math>. There must be a point <math>s_0</math> such that for all <math>t \leq s_0, p_{\tilde{\theta}}(t) \leq p_\theta(t)</math>, but for all <math>t &gt; s_0</math> we instead have <math>p_{\tilde{\theta}}(t) &gt; p_\theta(t)</math>.</p> $  \mathbb{E}_{\tilde{\theta}}\phi(T) - \mathbb{E}_\theta\phi(T) = \sum_t \phi(t) [p_{\tilde{\theta}}(t) - p_\theta(t)] = (cont)  $
<p>In the case of our sample we have</p> $  p_\theta(x) = \prod_{i=1}^n p_\theta(x_i) = \exp\left[\underbrace{-\theta}_{c(\theta)} \underbrace{\sum_{i=1}^n x_i}_{T(x_1, \dots, x_n)} - (-\log \theta)\right] 1\{x > 0\}  $ <p>So, using the lemma for exponential family UMP tests with strictly increasing <math>c(\theta)</math> we can construct a UMP test <math>\phi</math> such that</p>	<p>Since <math>\beta(\theta) = \mathbb{E}_\theta\phi(T)</math> is increasing, we have</p> $  \sup_{\theta \leq \theta_0} \mathbb{E}_\theta\phi(T) = \mathbb{E}_{\theta_0}\phi(T) = \alpha  $ <p>Since <math>t_0</math> and <math>q</math> are chosen such that <math>\mathbb{E}_{\theta_0}\phi(T) = \alpha</math>, we can see that they do not depend on the alternative hypothesis <math>\theta_1</math>. Thus, <math>\phi</math> is uniformly most powerful (UMP).</p>



<p>Suppose <math>X_1, \dots, X_n</math> are i.i.d. from Exponential <math>(\theta)</math>, <math>\theta &gt; 0</math>. The probability density function for the one-dimensional exponential distribution is</p> $p_\theta(x) = \exp[-\theta x - (-\log \theta)]1\{x > 0\}.$ <p>Consider a hypothesis <math>H_0 : \theta \leq \theta_0</math>, and <math>H_1 : \theta &gt; \theta_0</math>. What is a UMP test for <math>H_0</math>? (part-2)</p>	<p>Consider the problem of testing for parameter <math>\theta</math> in exponential family distributions with strictly increasing <math>c(\theta)</math>. How do you construct UMP tests for the right-sided and left-sided alternative hypotheses? (part-1)</p>
<p>Consider the problem of testing for parameter <math>\theta</math> in exponential family distributions with strictly increasing <math>c(\theta)</math>. How do you construct UMP tests for, the right-sided and left-sided alternative hypotheses? (part-2)</p>	<p>THEOREM</p> <p>State the theorem about exponential family UMPU tests.</p>
<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>\mathcal{N}(\mu, \sigma^2)</math> with <math>\sigma^2 = \sigma_0^2</math> known. Find a UMPU test for <math>H_0 : \mu = \mu_0, H_1 : \mu \neq \mu_0</math>. Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p_\theta(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$ <p>(part-1)</p>	<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>\mathcal{N}(\mu, \sigma^2)</math> with <math>\sigma^2 = \sigma_0^2</math> known. Find a UMPU test for <math>H_0 : \mu = \mu_0, H_1 : \mu \neq \mu_0</math>. Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p_\theta(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$ <p>(part-2)</p>
<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>\mathcal{N}(\mu, \sigma^2)</math> with <math>\sigma^2 = \sigma_0^2</math> known. Find a UMPU test for <math>H_0 : \mu = \mu_0, H_1 : \mu \neq \mu_0</math>. Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p_\theta(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$ <p>(part-3)</p>	<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>\mathcal{N}(\mu, \sigma^2)</math> with <math>\sigma^2 = \sigma_0^2</math> known. Find a UMPU test for <math>H_0 : \mu = \mu_0, H_1 : \mu \neq \mu_0</math>. Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p_\theta(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$ <p>(part-4)</p>
<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>\mathcal{N}(\mu, \sigma^2)</math> with <math>\sigma^2 = \sigma_0^2</math> known. Find a UMPU test for <math>H_0 : \mu = \mu_0, H_1 : \mu \neq \mu_0</math>. Keep in mind that for any normal random variable <math>Y \sim \mathcal{N}(\mu, \sigma^2)</math> you can write the probability density as</p> $p_\theta(y) = \exp \left[ \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix} \begin{bmatrix} y^2 \\ y \end{bmatrix} - \left( \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right) \right] \frac{1}{\sqrt{2\pi}}$ <p>(part-5)</p>	<p>Let <math>T</math> be an unbiased estimator of <math>g(\theta)</math> where <math>\theta \in \mathbb{R}^k</math>. Sketch a proof that then, under regularity conditions,</p> $\text{Var}_\theta(T) \geq \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta)$ <p>(part-1)</p>

<ul style="list-style-type: none"> <li>Right-sided alternative: <math>c(\cdot)</math> strictly increasing,  <math display="block">H_0 : \theta \leq \theta_0, H_1 : \theta &gt; \theta_0, \phi_R(T) = \begin{cases} 1 &amp; T &gt; t_0 \\ q &amp; t = t_0 \\ 0 &amp; T &lt; t_0 \end{cases}</math></li> <li>Left-sided alternative: <math>c(\cdot)</math> strictly decreasing,  <math display="block">H_0 : \theta \geq \theta_0, H_1 : \theta &lt; \theta_0, \phi_L(T) = \begin{cases} 1 &amp; T &lt; t_0 \\ q &amp; t = t_0 \\ 0 &amp; T &gt; t_0 \end{cases}</math></li> </ul>	$\phi(T) = \begin{cases} 1 & \text{if } t \leq t_0 \\ 0 & \text{if } t > t_0 \end{cases}$ <p>with a <math>t_0</math> for which <math>\mathbb{E}_{\theta_0}(T) = \alpha</math> holds. In particular,</p> $P_{\theta_0}(T \leq t_0) = P(\theta_0 T \leq \theta_0 t_0) \text{ (mult. by } \theta_0 \rightarrow \text{mean 1)} \\ = G(\theta_0, t_0) = \alpha. \quad (\text{Sum of Exp}(1) \text{ is Gamma}(n, 1))$ <p>So <math>\theta_0 t_0 = G^{-1}(\alpha) \Rightarrow t_0 = G^{-1}(\alpha)/\theta_0</math> and then <math>\phi</math> is UMP.</p>
<p>Suppose <math>p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x)</math> where <math>\theta \in \Theta</math> is an interval in <math>\mathbb{R}</math>, and <math>c(\cdot)</math> is strictly increasing. Let</p> $\phi(T) = \begin{cases} 1 & \text{if } T \notin [t_L, t_R] \\ q_L & \text{if } T = t_L \\ q_R & \text{if } T = t_R \\ 0 & \text{if } T \in (t_L, t_R) \end{cases}$ <p>where <math>q_L, q_R, t_L, t_R</math> are such that</p> <ol style="list-style-type: none"> <li><math>\mathbb{E}_{\theta_0} \phi(T) = \alpha</math></li> <li><math>\left. \frac{d}{d\theta} \mathbb{E}_\theta \phi(T) \right _{\theta=\theta_0} = 0</math>.</li> </ol> <p>Then <math>\phi</math> is UMP unbiased (UMPU) for <math>H_0 : \theta = \theta_0, H_1 : \theta \neq \theta_0</math>.</p>	
$\begin{aligned} \mathbb{E}_\mu \phi(T) &= P_\mu(T \notin (t_L, t_R)) = 1 - P_\mu(T \in (t_L, t_R)) \\ &= 1 - P_\mu(T \leq t_R) - P_\mu(T \leq t_L) \\ &= 1 - P_\mu\left(\frac{T - n\mu}{\sqrt{n\sigma_0^2}} \leq \frac{t_R - n\mu}{\sqrt{n\sigma_0^2}}\right) \\ &\quad - P_\mu\left(\frac{T - n\mu}{\sqrt{n\sigma_0^2}} \leq \frac{t_L - n\mu}{\sqrt{n\sigma_0^2}}\right) \\ &= 1 - \Phi\left(\frac{t_R - n\mu}{\sqrt{n\sigma_0^2}}\right) - \Phi\left(\frac{t_L - n\mu}{\sqrt{n\sigma_0^2}}\right) \end{aligned}$	<p>Looking at the density function of a single observation from a normal distribution, we can see that <math>c(\mu) = \frac{\mu}{\sigma_0^2}</math> with <math>T(x) = x</math>. So, for the entire sample, our statistic is</p> $T(x) = \sum_{i=1}^n x_i \sim \mathcal{N}(n\mu, n\sigma_0^2).$
<p>We take <math>(t_R - n\mu_0) = -(t_L - n\mu_0)</math>, because the other possibility (<math>t_R = t_L</math>) is a test which always returns 0. Knowing that <math>\Phi(-x) = 1 - \Phi(x)</math>, set</p> $P_{\mu_0}(T \notin (t_L, t_R)) = 2 \left( 1 - \Phi\left(\frac{t_R - n\mu_0}{\sqrt{n\sigma_0}}\right) \right) \triangleq \alpha.$ $\Rightarrow \Phi\left(\frac{t_R - n\mu_0}{\sqrt{n\sigma_0}}\right) = \frac{1 - \alpha}{2}$ <p>Hence we have <math>t_R = n\mu_0 + \sqrt{n\sigma_0} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)</math> and <math>t_L = n\mu_0 - \sqrt{n\sigma_0} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)</math></p>	<p>Now we want to find</p> $\begin{aligned} \left. \frac{d}{d\theta} \mathbb{E}_\theta \phi(T) \right _{\theta=\theta_0} &= \left. \frac{d}{d\theta} P_\mu(T \notin (t_L, t_R)) \right _{\theta=\theta_0} \\ &= -\frac{n}{\sqrt{n\sigma_0}} \varphi\left(\frac{t_R - n\mu}{\sqrt{n\sigma_0}}\right) \\ &\quad + \frac{n}{\sqrt{n\sigma_0}} \varphi\left(\frac{t_L - n\mu}{\sqrt{n\sigma_0}}\right) \Big _{\theta=\theta_0} \triangleq 0 \\ &\Leftrightarrow (t_R - n\mu_0)^2 = (t_L - n\mu)^2 \\ &\Leftrightarrow (t_R - n\mu_0) = -(t_L - n\mu_0) \text{ or } t_R = t_L \end{aligned}$
<p>As in the one-dimensional Cramer-Rao lower bound proof, we can show that for <math>j = 1, \dots, k</math>,</p> $\dot{g}_j(\theta) = \text{Cov}_\theta(T, s_{\theta,j}(X)).$ <p>Hence for all <math>a \in \mathbb{R}^k</math>,</p> $\begin{aligned}  a^T \dot{g}(\theta) ^2 &=  \text{Cov}_\theta(T, a^T s_\theta(X)) ^2 \\ &\leq \text{Var}_\theta(T) \text{Var}_\theta(a^T s_\theta(X)) \\ &= \text{Var}_\theta(T) a^T I(\theta) a \end{aligned}$	<p>Finally, our UMPU test is</p> $\phi(T) = \begin{cases} 1 & \text{if }  T - n\mu_0  > \sqrt{n\sigma_0} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \\ 0 & \text{otherwise} \end{cases}$

<p>Let <math>T</math> be an unbiased estimator of <math>g(\theta)</math> where <math>\theta \in \mathbb{R}^k</math>. Sketch a proof that then, under regularity conditions,</p> $\text{Var}_{\theta}(T) \geq \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta)$ <p>(part-2)</p>	<p>Let <math>X \sim \mathcal{N}(\theta, 1)</math> with <math>\sigma \in \mathbb{R}</math>, and <math>\theta \sim \mathcal{N}(c, \tau^2)</math> with <math>c, \tau^2</math> known. What is the Bayes estimator for quadratic loss? (part-1)</p>
<p>Let <math>X \sim \mathcal{N}(\theta, 1)</math> with <math>\sigma \in \mathbb{R}</math>, and <math>\theta \sim \mathcal{N}(c, \tau^2)</math> with <math>c, \tau^2</math> known. What is the Bayes estimator for quadratic loss? (part-2)</p>	<p>What does it mean for an estimator <math>T</math> to be extended Bayes?</p>
<p>THEOREM</p> <p>State the theorem about minimax estimators of constant risk</p>	<p>Suppose the risk <math>R(\theta, T) = R(T)</math> does not depend on <math>\theta</math>. Show that</p> <ol style="list-style-type: none"> <li>1. <math>T</math> is admissible <math>\implies T</math> is minimax,</li> <li>2. <math>T</math> is Bayes <math>\implies T</math> is minimax,</li> <li>3. <math>T</math> is extended Bayes <math>\implies T</math> is minimax.</li> </ol> <p>(part-1)</p>
<p>Suppose the risk <math>R(\theta, T) = R(T)</math> does not depend on <math>\theta</math>. Show that</p> <ol style="list-style-type: none"> <li>1. <math>T</math> is admissible <math>\implies T</math> is minimax,</li> <li>2. <math>T</math> is Bayes <math>\implies T</math> is minimax,</li> <li>3. <math>T</math> is extended Bayes <math>\implies T</math> is minimax.</li> </ol> <p>(part-2)</p>	<p>Suppose <math>X \sim \text{Binomial}(n, \theta)</math> with <math>\theta \in (0, 1)</math>. Suppose also the prior <math>\theta \sim \text{Beta}(r, s)</math> and quadratic loss <math>L(\vartheta, a) = (a - \vartheta)^2</math>. Knowing that the Bayes estimator minimizing the loss is</p> $T_{\text{Bayes}}(X) = \mathbb{E}(\theta \mid X) = \frac{x + r}{n + r + s}$ <p>find a minimax estimator for <math>\theta</math>. (part-1)</p>
<p>Suppose <math>X \sim \text{Binomial}(n, \theta)</math> with <math>\theta \in (0, 1)</math>. Suppose also the prior <math>\theta \sim \text{Beta}(r, s)</math> and quadratic loss <math>L(\vartheta, a) = (a - \vartheta)^2</math>. Knowing that the Bayes estimator minimizing the loss is</p> $T_{\text{Bayes}}(X) = \mathbb{E}(\theta \mid X) = \frac{x + r}{n + r + s}$ <p>find a minimax estimator for <math>\theta</math>. (part-2)</p>	<p>What does it mean for a Bayes estimator to be unique?</p>

$ \begin{aligned} w(\vartheta \mid x) &\propto p(x \mid \vartheta)w(\vartheta) \\ &= \varphi(x - \vartheta) \frac{1}{\tau} \varphi\left(\frac{\vartheta - c}{\tau}\right) \\ &\propto \exp\left[-\frac{1}{2}(x - \vartheta)^2 - \frac{1}{2\tau^2}(\vartheta - c)^2\right] \\ &\propto \exp\left[x\vartheta - \frac{1}{2}\vartheta^2 - \frac{1}{2\tau^2}\vartheta^2 + \frac{\vartheta c}{\tau^2}\right] \\ &= \exp\left[-\frac{1}{2}\left(-2\vartheta \underbrace{\left(x + \frac{c}{\tau^2}\right)}_a + \vartheta^2 \underbrace{\left(1 + \frac{1}{\tau^2}\right)}_b\right)\right] \end{aligned} $	<p>Then we can show</p> $\text{Var}_\theta(T) \geq \max_{a \in \mathbb{R}^k} \frac{ a^T \dot{g}(\theta) ^2}{a^T I(\theta) a} = \dot{g}^T(\theta) I(\theta)^{-1} \dot{g}(\theta)$ <p>(Last equality: auxiliary lemma, Cauchy-Schwarz, algebra)</p>
<p>An estimator <math>T</math> is called extended Bayes if there exists a sequence of priors <math>(w_m)_{m \geq 1}</math> such that for the Bayes estimator <math>T_m</math> for each prior <math>w_m</math> we have</p> $r_{w_m}(T) - r_{w_m}(T_m) \xrightarrow{m \rightarrow \infty} 0$	<p>To come up with an expression which only has one <math>\theta</math>, we will complete the square</p> $-2a\vartheta + b\vartheta^2 = b\left(-\frac{2a}{b}\vartheta + \vartheta^2\right) = b\left(\vartheta - \frac{a}{b}\right)^2 - \frac{a^2}{b}$ <p>hence from the squared expression <math>\left(\vartheta - \frac{a}{b}\right)^2</math> we can see that the mean is <math>\frac{a}{b}</math> (the distribution <math>w(\vartheta \mid x)</math> is symmetric). So plugging in the values for <math>a</math> and <math>b</math>, the Bayes estimator for quadratic loss is</p> $T_{\text{Bayes}} = \mathbb{E}(\theta \mid X) = \frac{\tau^2 X + c}{\tau^2 + 1}$
<ol style="list-style-type: none"> <li>1. Suppose <math>\sup_\theta R(\theta, T') \leq \sup_\theta R(\theta, T) = R(T)</math> for all <math>\theta</math>. Since <math>T</math> is admissible, <math>T'</math> is never strictly better than <math>T</math>, so <math>R(\theta, T') = R(T)</math> for all <math>\theta</math>.</li> <li>2. This is implied by (3).</li> <li>3. Since <math>T</math> is extended Bayes, for any <math>\varepsilon</math> there exists an <math>m</math> such that</li> </ol> $r_{w_m}(T) \leq r_{w_m}(T_m) + \varepsilon$	<p>Suppose the risk <math>R(\theta, T) = R(T)</math> does not depend on <math>\theta</math>. Then</p> <ol style="list-style-type: none"> <li>1. <math>T</math> is admissible <math>\implies T</math> is minimax,</li> <li>2. <math>T</math> is Bayes <math>\implies T</math> is minimax,</li> <li>3. <math>T</math> is extended Bayes <math>\implies T</math> is minimax.</li> </ol>
<p>To find the minimax estimator, we will find parameters <math>r</math> and <math>s</math> for the prior such that risk of <math>T_{\text{Bayes}}</math> is constant in <math>\theta</math>.</p> $ \begin{aligned} R(\theta, T_{\text{Bayes}}) &= \mathbb{E}_\theta(T_{\text{Bayes}} - \theta)^2 \\ &= \text{Var}_\theta(T_{\text{Bayes}}) + \text{Bias}_\theta^2(T_{\text{Bayes}}) \\ &= \frac{n\theta(1-\theta)}{(n+r+s)^2} + \left[\frac{n\theta+r}{n+r+s} - \frac{(n+r+s)\theta}{n+r+s}\right]^2 \\ &= \frac{[(r+s)^2 - n]\theta^2 + [n - 2r(r+s)]\theta + r^2}{(n+r+s)^2} \end{aligned} $	<p>3. Now</p> $ \begin{aligned} R(T) &= r_{w_m}(T) \\ &\leq r_{w_m}(T_m) + \varepsilon \\ &\leq r_{w_m}(T') + \varepsilon \quad (T_m \text{ is optimal for } r_{w_m}) \\ &\leq \sup_{\vartheta} R(\vartheta, T') + \varepsilon \end{aligned} $ <p>since the <math>\varepsilon</math> is arbitrary, we know the risk <math>R(T)</math> is less than the worst case risk <math>\sup_{\vartheta} R(\vartheta, T') + \varepsilon</math> of any other estimator <math>T'</math></p>
<p>A <math>w</math>-Bayes estimator <math>T</math> is called unique if <math>r_w(T') = r_w(T)</math> implies <math>P_\theta(T = T') = 1</math> for all <math>\theta</math></p>	<p>Since we don't want the risk to depend on <math>\theta</math>, we need to set the coefficients' in front of <math>\theta^2</math> and <math>\theta</math> to 0</p> $(r+s)^2 - n = 0, n - 2r(r+s) = 0.$ <p>Solving for <math>r</math> and <math>s</math> gives <math>r = s = \sqrt{n}/2</math>. Plugging these values into the estimator gives</p> $T = \frac{X + \sqrt{n}/2}{n + \sqrt{n}}.$ <p>Since <math>T</math> is Bayes and its risk does not depend on <math>\theta</math>, <math>T</math> is minimax.</p>

<p>LEMMA</p> <p>State the lemma about sufficient conditions for a Bayes estimator <math>T</math> to be admissible.</p>	<p>Show that if <math>T</math> is the unique Bayes estimator for the prior density <math>w</math>, then <math>T</math> is also admissible.</p>
<p>Suppose that for all <math>T'</math>, <math>R(\theta, T')</math> is continuous in <math>\theta</math>, and for all open <math>U \subset \Theta</math> the prior probability</p> $\Pi(U) := \int w(\vartheta) d\mu(\vartheta)$ <p>of <math>U</math> is strictly positive. Show that if <math>T</math> is a Bayes estimator for the prior <math>w</math>, then <math>T</math> is also admissible. (part-1)</p>	<p>Suppose that for all <math>T'</math>, <math>R(\theta, T')</math> is continuous in <math>\theta</math>, and for all open <math>U \subset \Theta</math> the prior probability</p> $\Pi(U) := \int w(\vartheta) d\mu(\vartheta)$ <p>of <math>U</math> is strictly positive. Show that if <math>T</math> is a Bayes estimator for the prior <math>w</math>, then <math>T</math> is also admissible. (part-2)</p>
<p>LEMMA</p> <p>State the lemma about admissibility of extended Bayes estimators.</p>	<p>Suppose that <math>T</math> is extended Bayes, and that for all <math>T'</math>, the risk <math>R(\theta, T')</math> is continuous in <math>\theta</math>. Furthermore, assume that for all open sets <math>U \subset \Theta</math>,</p> $\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0$ <p>where <math>\Pi_m(U) := \int_U w_m(\vartheta) d\mu_m(\vartheta)</math> is the probability of <math>U</math> under the prior <math>\Pi_m</math>. Show that then <math>T</math> is admissible. (part-1)</p>
<p>Suppose that <math>T</math> is extended Bayes, and that for all <math>T'</math>, the risk <math>R(\theta, T')</math> is continuous in <math>\theta</math>. Furthermore, assume that for all open sets <math>U \subset \Theta</math>,</p> $\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0$ <p>where <math>\Pi_m(U) := \int_U w_m(\vartheta) d\mu_m(\vartheta)</math> is the probability of <math>U</math> under the prior <math>\Pi_m</math>. Show that then <math>T</math> is admissible. (part-2)</p>	<p>Suppose that <math>T</math> is extended Bayes, and that for all <math>T'</math>, the risk <math>R(\theta, T')</math> is continuous in <math>\theta</math>. Furthermore, assume that for all open sets <math>U \subset \Theta</math>,</p> $\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0$ <p>where <math>\Pi_m(U) := \int_U w_m(\vartheta) d\mu_m(\vartheta)</math> is the probability of <math>U</math> under the prior <math>\Pi_m</math>. Show that then <math>T</math> is admissible. (part-3)</p>
<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Show that</p> $T \text{ is admissible} \implies \begin{cases} (i) & 0 < a < 1 \\ \text{or} \\ (ii) & a = 1, b = 0 \end{cases}$	<p>PROOF</p> <p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math> Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Sketch steps to prove that if <math>0 &lt; a &lt; 1</math> then <math>T</math> is admissible. (part-1)</p>

<p>Suppose some other estimator <math>T'</math> has a lower risk than <math>T</math> for all <math>\theta</math></p> $R(\theta, T') \leq R(\theta, T).$ <p>If that's true, then also</p> $r_w(T') \leq r_w(T).$ <p>Since <math>T</math> is <math>w</math>-Bayes, this must be an equality <math>r_w(T') = r_w(T)</math>. Since <math>T</math> is unique, <math>T = T'P_\theta</math>-almost surely for all <math>\theta</math>. So, <math>R(\theta, T') = R(\theta, T)</math> and thus <math>T</math> is admissible</p>	<p>Suppose that <math>T</math> is a Bayes estimator for the prior density <math>w</math>. Then either of these conditions is sufficient for <math>T</math> to be admissible:</p> <ol style="list-style-type: none"> <li>1. <math>T</math> is the unique Bayes decision,</li> <li>2. for all <math>T'</math>, <math>R(\theta, T')</math> is continuous in <math>\theta</math>, and for all open <math>U \subset \Theta</math> the prior probability <math>\Pi(U) := \int w(\vartheta)d\mu(\vartheta)</math> of <math>U</math> is strictly positive.</li> </ol>
$\begin{aligned} r_w(T') &= \int R(\vartheta, T') w(\vartheta)d\mu(\vartheta) \\ &= \int_U R(\vartheta, T') w(\vartheta)d\nu(\vartheta) + \int_{U^c} R(\vartheta, T') w(\vartheta)d\nu(\vartheta) \\ &= \int_U R(\vartheta, T) w(\vartheta)d\nu(\vartheta) - \varepsilon\Pi(U) \\ &\quad + \int_{U^c} R(\vartheta, T) w(\vartheta)d\nu(\vartheta) \\ &= r_w(T) - \varepsilon\Pi(U) < r_w(T). \end{aligned}$ <p>So we came at contradiction, since <math>T</math> is Bayes for prior <math>w</math>.</p>	<p>Suppose <math>T</math> is not admissible, so there exists another estimator <math>T'</math> which fulfils <math>R(\theta, T') \leq R(\theta, T)</math> for all <math>\theta</math>, and <math>R(\theta_0, T') &lt; R(\theta_0, T)</math> for some specific <math>\theta_0</math>. The assumptions imply that there exists an <math>\varepsilon &gt; 0</math> for which there is an open set <math>U</math> with <math>\theta_0 \in U</math> such that</p> $R(\theta, T') \leq R(\theta, T) - \varepsilon \text{ for all } \theta \in U.$ <p>Then the risk fulfils</p>
<p>Suppose <math>T</math> is not admissible, so there exists another estimator <math>T'</math> which fulfils <math>R(\theta, T') \leq R(\theta, T)</math> for all <math>\theta</math>, and <math>R(\theta_0, T') &lt; R(\theta_0, T)</math> for some specific <math>\theta_0</math>. So, there exists an <math>\varepsilon &gt; 0</math> and <math>U</math> open with <math>\theta_0 \in U</math> such that <math>R(\theta, T') \leq R(\theta, T)</math> for all <math>\theta \in U</math>. We can then say (continued)</p>	<p>Suppose that <math>T</math> is extended Bayes, and that for all <math>T'</math>, the risk <math>R(\theta, T')</math> is continuous in <math>\theta</math>. Furthermore, assume that for all open sets <math>U \subset \Theta</math>,</p> $\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0$ <p>where <math>\Pi_m(U) := \int_U w_m(\vartheta)d\mu_m(\vartheta)</math> is the probability of <math>U</math> under the prior <math>\Pi_m</math>. Then <math>T</math> is admissible.</p>
<p>which means that <math>\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)}</math> does not converge to 0, which contradicts our assumptions.</p>	$\begin{aligned} r_{w_m}(T') &\leq r_{w_m}(T) - \varepsilon\Pi_m(U) \\ \frac{r_{w_m}(T) - r_{w_m}(T')}{\Pi_m(U)} &\geq \varepsilon \\ \frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} &= \underbrace{\frac{r_{w_m}(T) - r_{w_m}(T')}{\Pi_m(U)}}_{\geq \varepsilon} \\ &\quad + \underbrace{\frac{r_{w_m}(T') - r_{w_m}(T_m)}{\Pi_m(U)}}_{\geq 0} \end{aligned}$
<ol style="list-style-type: none"> <li>1. Show that for any <math>a \in (0, 1), b \in \mathbb{R}, T</math> is Bayes for some prior. In particular, <math>\theta \sim \mathcal{N}(c, \tau^2)</math> works where <math>a = \frac{\tau^2}{\tau^2+1}, b = \frac{c}{\tau^2+1}</math>.</li> <li>2. Show that <math>T</math> is unique Bayes. So, for any other estimator <math>T'</math>, show that</li> </ol> $r_w(T') = r_w(T) \implies \mathbb{E} \left[ (T(X) - T'(X))^2 \right] = 0$ <p>where the expectation is with <math>\theta</math> integrated out, so with respect to some measure <math>P</math>.</p>	<p>Let <math>T_0 := X</math>.</p> <ol style="list-style-type: none"> <li>1. Suppose <math>a &gt; 1</math>. Then</li> </ol> $R(\theta, T) \geq \text{Var}_\theta(T) = a^2 > 1 = \text{Var}_\theta(T_0) = R(\theta, T_0)$ <p>so <math>T</math> is not admissible.</p> <ol style="list-style-type: none"> <li>2. Suppose alternatively that <math>a = 1, b \neq 0</math>. Then</li> </ol> $R(\theta, T) \geq \text{Bias}_\theta^2(T) + 1 = b^2 + 1 > \text{Var}_\theta(T_0) = R(\theta, T_0)$ <p>so once again <math>T</math> is inadmissible.</p>

<p>PROOF</p> <p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math> Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Sketch steps to prove that if <math>0 &lt; a &lt; 1</math> then <math>T</math> is admissible. (part-2)</p>	<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Knowing that <math>T</math> is Bayes for a prior <math>\theta \sim \mathcal{N}(c, \tau^2)</math> with <math>a = \frac{\tau^2}{\tau^2+1}, b = \frac{c}{\tau^2+1}</math>, show that it is unique Bayes. (part-1)</p>
<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Knowing that <math>T</math> is Bayes for a prior <math>\theta \sim \mathcal{N}(c, \tau^2)</math> with <math>a = \frac{\tau^2}{\tau^2+1}, b = \frac{c}{\tau^2+1}</math>, show that it is unique Bayes. (part-2)</p>	<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Knowing that <math>T</math> is Bayes for a prior <math>\theta \sim \mathcal{N}(c, \tau^2)</math> with <math>a = \frac{\tau^2}{\tau^2+1}, b = \frac{c}{\tau^2+1}</math>, show that it is unique Bayes. (part-3)</p>
<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math> Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Sketch steps to show that if <math>a = 1, b = 0</math> then <math>T</math> is admissible. (part-1)</p>	<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math> Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Sketch steps to show that if <math>a = 1, b = 0</math> then <math>T</math> is admissible. (part-2)</p>
<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math> Consider</p> $T := aX + b, a > 0, b \in \mathbb{R}.$ <p>Sketch steps to show that if <math>a = 1, b = 0</math> then <math>T</math> is admissible. (part-3)</p>	<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Show that <math>T = X</math> is extended Bayes. Keep in mind that for the prior <math>\theta \sim \mathcal{N}(0, m)</math>, the Bayes estimator is <math>T_m = \frac{m}{m+1}</math></p>
<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Knowing that <math>T = X</math> is extended Bayes for the prior <math>\theta \sim \mathcal{N}(0, m)</math>, show that it is also admissible. Keep in mind that the risk <math>r_{w_m}(T_m) = \frac{m}{m+1}</math>. (part-1)</p>	<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Knowing that <math>T = X</math> is extended Bayes for the prior <math>\theta \sim \mathcal{N}(0, m)</math>, show that it is also admissible. Keep in mind that the risk <math>r_{w_m}(T_m) = \frac{m}{m+1}</math>. (part-2)</p>

<p>Consider some estimator <math>T'</math> (4eq <math>\text{MSE} = \text{Var} + \text{Bias}^2</math>)</p> $ \begin{aligned} r_w(T') &= \mathbb{E}R(\theta, T') \\ &= \mathbb{E}[(\theta - T'(X))^2] \\ &= \mathbb{E}\mathbb{E}[(\theta - T'(X))^2   X] \\ &= \mathbb{E}\text{Var}(\theta   X) + \mathbb{E}[(\mathbb{E}(\theta   X) - T'(X))^2   X] \\ &= r_w(T) + \mathbb{E}[(T(X) - T'(X))^2], T = \mathbb{E}(\theta   X) \end{aligned} $	<p>3. Show that</p> $\mathbb{E}[(T(X) - T'(X))^2] = \int T(x) - T'(x) dP(x) = 0$ <p>implies <math>T = T'P_\theta</math>-almost surely for any <math>\theta</math>. To do this, show that <math>P</math> is a dominating measure (<math>\mathcal{N}(c, \tau^2 + 1)</math>) for <math>P_\theta</math> for any <math>\theta</math>.</p> <p>4. Conclude that <math>T</math> is unique Bayes and hence admissible.</p>
<p>Since all normal distributions dominate each other, <math>P</math> dominates <math>P_\theta</math> for all <math>\theta</math>, which means that <math>\mathbb{E}[(T(X) - T'(X))^2] = 0</math> integrated over <math>P \implies T = T'P_\theta</math>-almost surely for all <math>\theta</math>.</p> <p>So <math>T</math> is unique Bayes.</p>	<p>So, if <math>r_w(T') = r_w(T)</math> then <math>\mathbb{E}[(T(X) - T'(X))^2] = 0</math>. Now we want to show that this implies that indeed <math>T = T'</math> almost surely for all <math>\theta</math>. Note that the expectation <math>\mathbb{E}[(T(X) - T'(X))^2] = 0</math> is with respect to a measure <math>P</math> with <math>\theta</math> integrated out. So, we need to show that <math>P</math> dominates all <math>P_\theta</math>. <math>P</math> is the measure of <math>X</math>, which we can write as <math>X = \theta + \varepsilon</math> where <math>\theta \sim \mathcal{N}(c, \tau^2)</math> and <math>\varepsilon \sim \mathcal{N}(0, 1)</math>. So, <math>P</math> is the <math>\mathcal{N}(c, \tau^2 + 1)</math> distribution.</p>
<p>2. Confirm that the risk converges sufficiently quickly to fulfil the admissibility criterion, i.e. that</p> $\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0.$ <p>To do this, pick any open interval <math>U = (u, u + h)</math> and show by using a Taylor approx that for large <math>m</math>,</p> $\Pi_m(U) = \Phi\left(\frac{u+h}{\sqrt{m}}\right) - \Phi\left(\frac{u}{\sqrt{m}}\right) \approx \frac{1}{\sqrt{m}}\phi(0) \geq \frac{1}{\sqrt{m}C}$ <p>for some constant <math>C</math>.</p>	<p>1. Show that <math>T</math> is extended Bayes for prior <math>\theta \sim \mathcal{N}(0, m)</math>. The Bayesian estimator is then <math>T_m = \frac{m}{m+1}X</math>, while the risk of <math>T</math> is always 1, so</p> $r_{w_m}(T) - r_{w_m}(T_m) = 1 - \frac{m}{m+1} \xrightarrow{m \rightarrow \infty} 0.$
$ \begin{aligned} R(\theta, T_m) &= \left(\frac{m}{m+1}\theta - \theta\right)^2 + \frac{m^2}{(m+1)^2}, (\text{MSE} = \text{B}^2 + \text{Var}) \\ &= \frac{1}{(m+1)^2}\theta^2 + \frac{m^2}{(m+1)^2} \\ \mathbb{E}R(\theta, T_m) &= r_{w_m}(T_m) = \frac{m}{(m+1)^2} + \frac{m^2}{(m+1)^2} (\mathbb{E}\theta^2 = m) \\ &= \frac{m}{(m+1)} \\ r_w(T) &= 1. \end{aligned} $ <p>So we can already see that <math>T</math> is extended Bayes, because</p> $r_{w_m}(T) - r_{w_m}(T_m) = 1 - \frac{m}{m+1} \xrightarrow{m \rightarrow \infty} 0$	<p>3. At this point we get</p> $\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \leq \frac{1 - \frac{m}{m+1}}{\frac{1}{\sqrt{m}C}} = \frac{\frac{1}{m+1}}{\frac{1}{\sqrt{m}}} C \xrightarrow{m \rightarrow \infty} 0$ <p>and so, since the loss function is continuous, <math>T</math> is admissible.</p>
<p>Take an open interval <math>U = (u, u + h)</math>. Now</p> $ \begin{aligned} \Pi_m(U) &= \Phi\left(\frac{u+h}{\sqrt{m}}\right) - \Phi\left(\frac{u}{\sqrt{m}}\right) \\ &= \frac{1}{\sqrt{m}}\phi\left(\frac{u}{\sqrt{m}}\right)h + o(1/\sqrt{m}) \\ &\approx \frac{1}{\sqrt{m}}\phi(0) \quad (\text{large } m) \\ &\geq \frac{1}{\sqrt{m}C} \text{ for some constant } C \end{aligned} $	<p>Since we know <math>T</math> is extended Bayes, we know <math>r_{w_m}(T) - r_{w_m}(T_m) \xrightarrow{m \rightarrow \infty} 0</math>. We just need to check if it's sufficiently quick so that for any <math>U \subset \mathbb{R}</math> open it's true that</p> $\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \xrightarrow{m \rightarrow \infty} 0$ <p style="text-align: center;">Taylor approx.</p> $\left( F(x+h) = F(x) + F'(x)h + o(h) \right)$



<p>Suppose <math>X \sim \mathcal{N}(\theta, 1)</math> where <math>\theta \in \Theta = \mathbb{R}</math>. Let <math>R(\theta, T) := \mathbb{E}_\theta(T(X) - \theta)^2</math>. Knowing that <math>T = X</math> is extended Bayes for the prior <math>\theta \sim \mathcal{N}(0, m)</math>, show that it is also admissible. Keep in mind that the risk <math>r_{w_m}(T_m) = \frac{m}{m+1}</math>. (part-3)</p>	<p>DEFINITION</p> <p>What is the least squares estimator?</p>
<p>LEMMA</p> <p>State lemma about constructing least squares estimators using linear algebra operations.</p>	<p>Let <math>f = \mathbb{E}Y</math>. What is the best linear approximation of <math>f</math> ?</p>
<p>LEMMA</p> <p>State lemma about estimation and misspecification errors of least squares models.</p>	<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also that the noise is uncorrelated, so <math>\mathbb{E}\epsilon\epsilon^T = \sigma^2 I</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show that then</p> $\mathbb{E}\hat{\beta} = \beta^*$ <p>with <math>\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}</math></p>
<p>Let <math>\mathbb{E}(Y) = f</math>, let <math>\epsilon := Y - f</math> be the noise. Suppose also that the noise is uncorrelated, so <math>\mathbb{E}\epsilon\epsilon^T = \sigma^2 I</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show</p> $\mathbb{E} \left\  X \left( \hat{\beta} - \beta^* \right) \right\ _2^2 = \sigma^2 p$ <p>(part-1)</p>	<p>Let <math>\mathbb{E}(Y) = f</math>, let <math>\epsilon := Y - f</math> be the noise. Suppose also that the noise is uncorrelated, so <math>\mathbb{E}\epsilon\epsilon^T = \sigma^2 I</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show</p> $\mathbb{E} \left\  X \left( \hat{\beta} - \beta^* \right) \right\ _2^2 = \sigma^2 p$ <p>(part-2)</p>
<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also that the noise is uncorrelated, so <math>\mathbb{E}\epsilon\epsilon^T = \sigma^2 I</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show that then</p> $\mathbb{E} \ X\hat{\beta} - f\ _2^2 = \sigma^2 p + \ X\beta^* - f\ _2^2$	<p>LEMMA</p> <p>State the lemma about the distribution of linear approximation errors.</p>

<p>The least squares estimator is</p> $\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \arg \min_{\begin{pmatrix} a \\ b \end{pmatrix} \in \mathbb{R}^{p+1}} \ Y - a - Xb\ _2^2$ <p>where <math>\ V\ _2^2 = V^T V = \sum_{i=1}^n V_i^2, V \in \mathbb{R}^n</math>. Sometimes <math>\hat{\alpha}</math> is replaced instead by a constant term appended to <math>X</math>. Then <math>\tilde{X} \rightarrow X, p+1 \rightarrow p, \begin{pmatrix} a \\ b \end{pmatrix} \rightarrow b</math>, and we write</p> $\hat{\beta} = \arg \min_{b \in \mathbb{R}^p} \ Y - Xb\ _2^2.$	<p>So</p> $\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \leq \frac{1 - \frac{m}{m+1}}{\frac{1}{\sqrt{mC}}} = \frac{\frac{1}{m+1}}{\frac{1}{\sqrt{m}}} C \xrightarrow{m \rightarrow \infty} 0.$ <p>Since the above criterion is fulfilled and the loss function is continuous, <math>T</math> is admissible.</p>
<p>We set <math>\beta^* := (X^T X)^{-1} X^T f</math>, and then the best linear approximation of <math>f</math> is <math>X\beta^*</math>.</p>	<p>Suppose <math>X \in \mathbb{R}^{n \times p}</math> has rank <math>p</math>. Then the least squares estimator is</p> $\hat{\beta} = (X^T X)^{-1} X^T Y$
$\begin{aligned} \mathbb{E}\hat{\beta} &= \mathbb{E} \left[ (X^T X)^{-1} X^T Y \right] \\ &= (X^T X)^{-1} X^T \mathbb{E}[\epsilon] + (X^T X)^{-1} X^T f = \beta^* \\ \text{Cov}(\hat{\beta}) &= \text{Cov} \left( (X^T X)^{-1} X^T \epsilon \right) \\ &= (X^T X)^{-1} X^T \underbrace{\text{Cov}(\epsilon)}_{\sigma^2 I} X (X^T X)^{-1} \\ &= (X^T X)^{-1} \sigma^2. \end{aligned}$	<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also that the noise is uncorrelated, so <math>\mathbb{E}\epsilon\epsilon^T = \sigma^2 I</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>.</p> <ol style="list-style-type: none"> <li>1. <math>\mathbb{E}\hat{\beta} = \beta^*, \text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}</math></li> <li>2. <math>\mathbb{E} \ X(\hat{\beta} - \beta^*)\ _2^2 = p\sigma^2</math></li> <li>3. <math>\mathbb{E} \ X\hat{\beta} - f\ _2^2 = \underbrace{p\sigma^2}_{\text{estimation error}} + \underbrace{\mathbb{E} \ X\beta^* - f\ _2^2}_{\text{misspecification error}}.</math></li> </ol>
<p>Now just continue the equations we had above</p> $\begin{aligned} \mathbb{E} \ X(\hat{\beta} - \beta^*)\ _2^2 &= \mathbb{E} \ PP^T \epsilon\ _2^2 \\ &= \mathbb{E} \ P^T \epsilon\ _2^2 \quad (\text{Eqn. above}) \\ &= \mathbb{E} \ V\ _2^2 \\ &= \mathbb{E} V^T V \\ &= \text{Cov}(V) \quad (\mathbb{E} V = 0) \\ &= p\sigma^2 \end{aligned}$	<p>We have</p> $\begin{aligned} X(X^T X)^{-1} X^T &= PP^T \\ \implies \ X(\hat{\beta} - \beta^*)\ _2^2 &= \ PP^T \epsilon\ _2^2 \\ &= \epsilon^T PP^T PP^T \epsilon \\ &= \epsilon^T PP^T \epsilon = \ P^T \epsilon\ _2^2. \end{aligned}$ <p>Let <math>V := P^T \epsilon</math>. Knowing that <math>\mathbb{E} P^T \epsilon = 0</math>, we get</p> $\text{Cov}(V) = P^T \text{Cov}(\epsilon) P = I \sigma^2$
<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also <math>\epsilon \sim \mathcal{N}(0, \sigma^2 I)</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Then</p> <ol style="list-style-type: none"> <li>1. <math>\hat{\beta} \sim \mathcal{N}(\beta^*, \sigma^2 (X^T X)^{-1})</math>,</li> <li>2. <math>\frac{\ x(\hat{\beta} - \beta^*)\ _2^2}{\sigma^2} \sim \chi_p^2.</math></li> </ol>	$X\hat{\beta} - f = \underbrace{X(\hat{\beta} - \beta^*)}_{\in \text{column space of } X} + \underbrace{(X\beta^* - f)}_{\text{orthogonal to } X}$ <p>So by Pythagoras for all <math>b</math> it's true that</p> $\ X\hat{\beta} - f\ _2^2 = \ X(b - \beta^*)\ _2^2 + \ X\beta^* - f\ _2^2$ <p>and, taking the expectation,</p> $\mathbb{E} \ X\hat{\beta} - f\ _2^2 = \underbrace{\mathbb{E} \ X(b - \beta^*)\ _2^2}_{= \sigma^2 p \text{ from lemma}} + \mathbb{E} \ X\beta^* - f\ _2^2$

<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also <math>\epsilon \sim \mathcal{N}(0, \sigma^2 I)</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show that then</p> $\hat{\beta} \sim \mathcal{N}\left(\beta^*, \sigma^2 (X^T X)^{-1}\right).$	<p>Let <math>\mathbb{E}(Y) = f</math>, and let <math>\epsilon := Y - f</math> be the noise. Suppose also <math>\epsilon \sim \mathcal{N}(0, \sigma^2 I)</math>. Also let <math>\hat{\beta}</math> be the least squares estimator, and <math>X\beta^*</math> with <math>X \in \mathbb{R}^{n \times p}</math> be the best linear approximation of <math>f</math>. Show that then</p> $\frac{\ X(\hat{\beta} - \beta^*)\ _2^2}{\sigma^2} \sim \chi_p^2$
<p>Suppose <math>Y = X\beta + \epsilon</math> for <math>\epsilon \sim \mathcal{N}(0, \sigma^2 I)</math> with a known <math>\sigma^2</math>. What is a pivot you could use to test for the hypothesis <math>H_0 : \beta = \beta_0</math> ?</p>	<p>LEMMA</p> <p>State lemma about testing a linear hypothesis with a restriction.</p>
<p>What does it mean for <math>Z_n</math> to converge to <math>Z</math> in probability?</p>	<p>What does it mean for <math>Z_n</math> to converge to <math>Z</math> almost surely?</p>
<p>What does it mean for <math>Z_n</math> to converge to <math>Z</math> in distribution?</p>	<p>How are convergence in probability and in distribution related?</p>
<p>THEOREM</p> <p>State the central limit theorem.</p>	<p>THEOREM</p> <p>State the Cramér-Wold device.</p>

$\begin{aligned}\ X(\hat{\beta} - \beta^*)\ _2^2 &= \ PP^T\epsilon\ _2^2 \\ &= \ PV\ _2^2 \quad \left( \text{Let } P^T\epsilon = V \right)\end{aligned}$ <p>Now <math>\mathbb{E}V = 0, \text{Cov}(V) = \sigma^2 I</math></p> $\ PV\ _2^2 = V^T P^T P V = \ V\ _2^2 = \sum_{j=1}^p V_j^2$ <p>where <math>V_1, \dots, V_p</math> are i.i.d. <math>\mathcal{N}(0, \sigma^2)</math>. Hence, <math>\frac{\sum_{j=1}^p V_j^2}{\sigma^2} \sim \chi_p^2</math></p>	$\hat{\beta} = \underbrace{(X^T X)^{-1} X^T f}_{\beta^*} + \underbrace{(X^T X)^{-1} X^T \epsilon}_{\mathcal{N}(0, (X^T X)^{-1} \sigma^2)}$
<p>Let the model be <math>Y = X\beta + \epsilon</math> with <math>\epsilon \sim \mathcal{N}(0, \sigma^2 I)</math>. Our hypothesis is <math>H_0 : B\beta = 0</math> where <math>B \in \mathbb{R}^{q \times p}</math> is a matrix of restrictions. Estimator <math>\hat{\beta}_0</math> defined as</p> $\hat{\beta}_0 = \arg \min_{b \in \mathbb{R}^p, Bb=0} \ Y - Xb\ _2^2.$ <p>Now under the null hypothesis</p> $\frac{\ Y - X\hat{\beta}_0\ _2^2 - \ Y - X\hat{\beta}\ _2^2}{\sigma^2} \sim \chi_q^2.$	<p>We could use the distribution function of the <math>\chi_p^2</math> distribution. This is because we know</p> $\frac{\ X(\hat{\beta} - \beta)\ _2^2}{\sigma^2} \sim \chi_p^2.$
<p><math>Z_n</math> converges to <math>Z</math> almost surely if</p> $\mathbb{P}\left(\lim_{n \rightarrow \infty} Z_n = z\right) = 1$ <p>We then write <math>Z_n \xrightarrow{\text{a.s.}} Z</math>.</p>	<p><math>Z_n</math> converges to <math>Z</math> in probability if for all <math>\varepsilon &gt; 0</math>,</p> $\lim_{n \rightarrow \infty} \mathbb{P}(\ Z_n - Z\  > \varepsilon) = 0$ <p>We then write <math>Z_n \xrightarrow{\mathbb{P}} Z</math></p>
<p>Convergence in probability implies convergence in distribution, but not the other way around.</p>	<p><math>Z_n</math> converges in distribution to <math>Z</math> if</p> $\lim_{n \rightarrow \infty} \mathbb{E}f(Z_n) = \mathbb{E}f(Z)$ <p>for all <math>f : \mathbb{R}^p \rightarrow \mathbb{R}</math> bounded and continuous. We then write <math>Z_n \xrightarrow{\mathcal{D}} Z</math>.</p>
<p>Let <math>(\{Z_n\}, Z)</math> be a collection of <math>\mathbb{R}^p</math>-valued random variables. Then</p> $Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} Z \iff a^T Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} a^T Z \forall a \in \mathbb{R}^p.$	<p>If <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \in \mathbb{R}</math> with <math>\mathbb{E}X = \mu, \text{Var}(X) = \sigma^2</math>, and <math>\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i</math>, then</p> $\frac{(\bar{X}_n - \mu)}{\sigma/\sqrt{n}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$

<p>THEOREM</p> <p><i>State the Portmanteau theorem.</i></p>	<p>Suppose <math>(z_n)_{n \geq 1}</math> is a sequence of <math>\mathbb{R}^p</math> vectors, <math>(r_n)_{n \geq 1}</math> is a sequence of positive constants. What do these order symbols mean?</p> <ol style="list-style-type: none"> <li>1. <math>z_n = \mathcal{O}(1)</math></li> <li>2. <math>z_n = o(1)</math></li> <li>3. <math>z_n = \mathcal{O}(r_n)</math></li> <li>4. <math>z_n = o(r_n)</math></li> </ol>
<p>What does it mean for a sequence of random variables <math>(Z_n)_{n \geq 1}</math> to be bounded in probability?</p>	<p>Suppose <math>(Z_n)_{n \geq 1}</math> is a sequence of <math>\mathbb{R}^p</math> random vectors, <math>(r_n)_{n \geq 1}</math> is a sequence of positive constants. What do these symbols mean?</p> <ol style="list-style-type: none"> <li>1. <math>Z_n = \mathcal{O}_{\mathbb{P}}(1)</math></li> <li>2. <math>Z_n = o_{\mathbb{P}}(1)</math></li> <li>3. <math>Z_n = \mathcal{O}_{\mathbb{P}}(r_n)</math></li> <li>4. <math>Z_n = o_{\mathbb{P}}(r_n)</math></li> </ol>
<p>What does <math>Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} Z</math> imply about the order of <math>Z_n</math>?</p>	<p>THEOREM</p> <p><i>State Slutsky's theorem.</i></p>
<p>PROOF</p> <p><i>Prove the Slutsky's theorem</i></p>	<p>Let <math>X_1, \dots, X_n</math> be i.i.d. <math>X \in \mathbb{R}</math> with mean <math>\mu</math> and variance <math>\sigma^2</math>. Does <math>\sqrt{n}(\bar{X}_n^2 - \mu^2)</math> converge to something? If so, to what, and with what type of convergence?</p>
<p>What does it mean for an estimator <math>T_n</math> to be consistent?</p>	<p>What does it mean for an estimator <math>T_n</math> to be called asymptotically normal?</p>

<table border="1"> <thead> <tr> <th>Order symbol</th><th>Meaning</th></tr> </thead> <tbody> <tr> <td>(1) <math>z_n = \mathcal{O}(1)</math></td><td><math>\lim_{n \rightarrow \infty} \ z_n\  &lt; \infty</math></td></tr> <tr> <td>(2) <math>z_n = o(1)</math></td><td><math>\lim_{n \rightarrow \infty} \ z_n\  = 0</math></td></tr> <tr> <td>(3) <math>z_n = \mathcal{O}(r_n)</math></td><td><math>\lim_{n \rightarrow \infty} \ z_n/r_n\  &lt; \infty</math></td></tr> <tr> <td>(4) <math>z_n = o(r_n)</math></td><td><math>\lim_{n \rightarrow \infty} \ z_n/r_n\  = 0</math></td></tr> </tbody> </table>	Order symbol	Meaning	(1) $z_n = \mathcal{O}(1)$	$\lim_{n \rightarrow \infty} \ z_n\  < \infty$	(2) $z_n = o(1)$	$\lim_{n \rightarrow \infty} \ z_n\  = 0$	(3) $z_n = \mathcal{O}(r_n)$	$\lim_{n \rightarrow \infty} \ z_n/r_n\  < \infty$	(4) $z_n = o(r_n)$	$\lim_{n \rightarrow \infty} \ z_n/r_n\  = 0$	<p>The following are equivalent:</p> <ul style="list-style-type: none"> <li>• <math>\mathbb{E}f(Z_n) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathbb{E}f(Z)</math> for all <math>f: \mathbb{R}^p \rightarrow \mathbb{R}</math> bounded and continuous,</li> <li>• <math>\mathbb{E}f(Z_n) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathbb{E}f(Z)</math> for all <math>f: \mathbb{R}^p \rightarrow \mathbb{R}</math> bounded and Lipschitz,</li> <li>• <math>\lim_{n \rightarrow \infty} \mathbb{P}(Z_n \leq z) = \mathbb{P}(Z \leq z) =: G(z)</math> for all <math>G</math>-continuity points <math>z</math>.</li> </ul>
Order symbol	Meaning										
(1) $z_n = \mathcal{O}(1)$	$\lim_{n \rightarrow \infty} \ z_n\  < \infty$										
(2) $z_n = o(1)$	$\lim_{n \rightarrow \infty} \ z_n\  = 0$										
(3) $z_n = \mathcal{O}(r_n)$	$\lim_{n \rightarrow \infty} \ z_n/r_n\  < \infty$										
(4) $z_n = o(r_n)$	$\lim_{n \rightarrow \infty} \ z_n/r_n\  = 0$										
<table border="1"> <thead> <tr> <th>Order symbol</th><th>Meaning</th></tr> </thead> <tbody> <tr> <td>(1) <math>Z_n = \mathcal{O}_{\mathbb{P}}(1)</math></td><td>(<math>\star</math>)</td></tr> <tr> <td>(2) <math>Z_n = o_{\mathbb{P}}(1)</math></td><td><math>Z_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0</math></td></tr> <tr> <td>(3) <math>Z_n = \mathcal{O}_{\mathbb{P}}(r_n)</math></td><td><math>Z_n/r_n = \mathcal{O}_{\mathbb{P}}(1)</math></td></tr> <tr> <td>(4) <math>Z_n = o_{\mathbb{P}}(r_n)</math></td><td><math>Z_n/r_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0</math></td></tr> </tbody> </table> <p>With (<math>\star</math>) = <math>\lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \mathbb{P}(\ Z_n\  &gt; M) = 0</math></p>	Order symbol	Meaning	(1) $Z_n = \mathcal{O}_{\mathbb{P}}(1)$	( $\star$ )	(2) $Z_n = o_{\mathbb{P}}(1)$	$Z_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$	(3) $Z_n = \mathcal{O}_{\mathbb{P}}(r_n)$	$Z_n/r_n = \mathcal{O}_{\mathbb{P}}(1)$	(4) $Z_n = o_{\mathbb{P}}(r_n)$	$Z_n/r_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$	<p>A sequence of random variables <math>(Z_n)_{n \geq 1}</math> is bounded in probability if</p> $\lim_{M \rightarrow \infty} \lim_{n \rightarrow \infty} \mathbb{P}(\ Z_n\  > M) = 0$
Order symbol	Meaning										
(1) $Z_n = \mathcal{O}_{\mathbb{P}}(1)$	( $\star$ )										
(2) $Z_n = o_{\mathbb{P}}(1)$	$Z_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$										
(3) $Z_n = \mathcal{O}_{\mathbb{P}}(r_n)$	$Z_n/r_n = \mathcal{O}_{\mathbb{P}}(1)$										
(4) $Z_n = o_{\mathbb{P}}(r_n)$	$Z_n/r_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0$										
<p>Let <math>(Z_n)_{n \geq 1}</math>, <math>Z</math>, and <math>(A_n)_{n \geq 1}</math> be (sequences of) random variables in <math>\mathbb{R}^p</math>. Furthermore, let <math>a \in \mathbb{R}^p</math> be a constant vector. Then</p> $\left\{ \begin{array}{l} Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} Z \\ \text{and} \\ A_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} a \end{array} \right\} \implies A_n^T Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} a^T Z.$	$Z_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} Z \text{ implies } Z_n = \mathcal{O}_{\mathbb{P}}(1)$										
<p>We know that</p> <ol style="list-style-type: none"> <li>1. by the CLT, <math>\sqrt{n}(\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \sigma^2)</math>,</li> <li>2. by the LLN, <math>(\bar{X}_n + \mu) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 2\mu</math>.</li> </ol> <p>So <math>\sqrt{n}(\bar{X}_n^2 - \mu^2) = \sqrt{n}(\bar{X}_n - \mu)(\bar{X}_n + \mu) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, 4\mu^2\sigma^2)</math>. (Slutsky)</p>	<p>See pg 132 skript.</p>										
<p><math>T_n</math> is called asymptotically normal if</p> $\sqrt{n}(T_n - g(\theta)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, V_\theta)$ <p>where <math>V_\theta</math> is the asymptotic covariance matrix.</p>	<p><math>T_n</math> is called consistent if</p> $T_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}_\theta} g(\theta)$										

<p>What does it mean for <math>T_n</math> to be called asymptotically linear?</p>	<p>How are asymptotic linearity and normality related?</p>
<p>Suppose <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \in \mathbb{R}</math> with <math>\mathbb{E}_\theta X = \mu</math>, <math>\text{Var}_\theta(X) = \sigma^2 &lt; \infty</math>. Let <math>g(\theta) = \mu</math>, and <math>T_n = \bar{X}_n</math>. What is the influence function in this case?</p>	<p>Suppose <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \in \mathbb{R}</math> with <math>\mathbb{E}_\theta X = \mu</math>, <math>\text{Var}_\theta(X) = \sigma^2 &lt; \infty</math>. Let <math>g(\theta) = \mu^2</math>, and <math>T_n = \bar{X}_n^2</math>. What is the influence function and asymptotic variance in this case?</p>
<p>THEOREM</p> <p>State the <math>\delta</math>-method theorem. (part-1)</p>	<p>THEOREM</p> <p>State the <math>\delta</math>-method theorem. (part-2)</p>
<p>Suppose <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \sim \text{Bernoulli}(\theta)</math>. What is the asymptotic distribution of</p> $\sqrt{n} \left( \log \frac{\bar{X}_n}{1 - \bar{X}_n} - \log \frac{\theta}{1 - \theta} \right)?$ <p>(part-1)</p>	<p>Suppose <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \sim \text{Bernoulli}(\theta)</math>. What is the asymptotic distribution of</p> $\sqrt{n} \left( \log \frac{\bar{X}_n}{1 - \bar{X}_n} - \log \frac{\theta}{1 - \theta} \right)?$ <p>(part-2)</p>
<p>DEFINITION</p> <p>What is the empirical risk?</p>	<p>DEFINITION</p> <p>What is an M-estimator?</p>

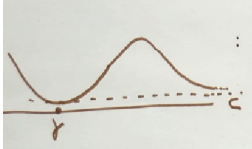
<p>Asymptotic linearity implies asymptotic normality</p>	<p><math>T_n</math> is called asymptotically linear if there exists a function <math>l_\theta : \mathcal{X} \rightarrow \mathbb{R}^p</math> with <math>\mathbb{E}_\theta l_\theta(X) = 0, \text{Cov}(l_\theta(X)) := V_\theta &lt; \infty</math> such that</p> $T_n - g(\theta) = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbb{P}_\theta} \left( \frac{1}{\sqrt{n}} \right)$
<p>We have</p> $T_n - g(\theta) = \bar{X}_n^2 - \mu^2 = 2\mu(\bar{X}_n - \mu) + \underbrace{(\bar{X}_n - \mu)^2}_{=O_{\mathbb{P}_\theta}(\frac{1}{n}) = o_{\mathbb{P}_\theta}(\frac{1}{\sqrt{n}})}$ <p>hence the influence function is</p> $l_\theta(x) = 2\mu(x - \mu)$ <p>and the asymptotic variance is <math>V_\theta = \text{Cov}(l_\theta(X)) = 4\mu^2\sigma^2</math></p>	<p>We have</p> $\bar{X}_n - \mu = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)$ <p>so the influence function is</p> $l_\theta(x) = x - \mu$
<p>If moreover <math>T_n</math> is an asymptotically linear estimator of <math>\gamma</math> with influence function <math>l_\theta</math>, then <math>h(T_n)</math> is an asymptotically linear estimator of <math>h(\gamma)</math> with influence function</p> $\dot{h}(\gamma)^T l_\theta,$ <p>so we can write</p> $h(T_n) - h(\gamma) = \frac{1}{n} \sum_{i=1}^n \dot{h}(\gamma)^T l_\theta(X_i) + o_{\mathbb{P}_\theta} \left( \frac{1}{\sqrt{n}} \right).$	<p>Suppose <math>T_n</math> is an asymptotically normal estimator of <math>\gamma = g(\theta) \in \mathbb{R}^p</math> with asymptotic covariance matrix <math>V_\theta</math>. Suppose that a function <math>h : \mathbb{R}^p \rightarrow \mathbb{R}</math> is differentiable at <math>\gamma</math>. Then <math>h(T_n)</math> is an asymptotically normal estimator of <math>h(\gamma)</math> with asymptotic variance</p> $\dot{h}(\gamma)^T V_\theta \dot{h}(\gamma),$ <p>so we can write</p> $\sqrt{n}(h(T_n) - h(\gamma)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, \dot{h}(\gamma)^T V_\theta \dot{h}(\gamma)\right).$
$\dot{h}(\theta) = \frac{1}{\theta} + \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)},$ $\begin{aligned} \dot{h}(\theta)^T V_\theta \dot{h}(\theta) &= (\dot{h}(\theta))^2 V_\theta \\ &= \frac{\theta(1-\theta)}{(\theta(1-\theta))^2} = \frac{1}{\theta(1-\theta)} \end{aligned}$ <p>and hence</p> $\sqrt{n} \left( \log \frac{\bar{X}_n}{1 - \bar{X}_n} - \log \frac{\theta}{1 - \theta} \right) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N} \left( 0, \frac{1}{\theta(1-\theta)} \right)$	<p>We know <math>T_n = X_n</math> is an asymptotically linear estimator of <math>\gamma = g(\theta) = \theta</math> with</p> $l_\theta(x) = x - \theta, V_\theta = \mathbb{E}l_\theta(X)^2 = \theta(1 - \theta).$ <p>Furthermore, we know the function <math>h(\theta) = \log \frac{\theta}{1-\theta}</math> is differentiable for <math>\theta \in (0, 1)</math>. So, we can use the <math>\delta</math>-technique to show that <math>h(T_n)</math> is an asymptotically linear estimator of <math>h(\gamma)</math> with asymptotic variance <math>\dot{h}(\gamma)^T V_\theta \dot{h}(\gamma)</math>. So</p> $h(\theta) = \log \frac{\theta}{1 - \theta},$
<p>An M-estimator is the estimator which minimizes the empirical risk, so</p> $\hat{\gamma}_n = \arg \min_{c \in \Gamma} \hat{R}_n(c).$ <p>It's also called an empirical risk minimizer (ERM).</p>	<p>Let <math>X_1, \dots, X_n, \dots</math> be i.i.d. copies of <math>X</math>. Then the empirical risk is</p> $\hat{R}_n(c) := \frac{1}{n} \sum_{i=1}^n \rho_c(x_i)$



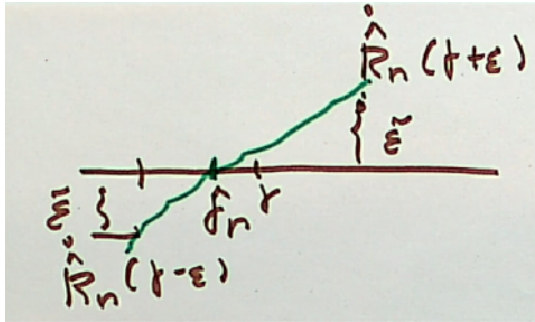
<p>Suppose <math>\gamma = \mathbb{E}X</math>. What is an M-estimator for <math>\gamma</math> ?</p>	<p>DEFINITION</p> <p>What is a Z-estimator?</p>
<p>Find a Z-estimator <math>\gamma</math> for the situation <math>X \in \mathbb{R}, c \in \mathbb{R}, \rho_c(x) = (x - c)^2</math>.</p>	<p>Find a Z-estimator <math>\gamma</math> for the situation <math>X \in \mathbb{R}, c \in \mathbb{R}, \rho_c(x) =  x - c </math>. Describe the ideal, i.e. non-empirical version. (part-1)</p>
<p>Find a Z-estimator <math>\gamma</math> for the situation <math>X \in \mathbb{R}, c \in \mathbb{R}, \rho_c(x) =  x - c </math>. Describe the ideal, i.e. non-empirical version. (part-2)</p>	<p>Find a Z-estimator <math>\gamma</math> for the situation <math>X \in \mathbb{R}, c \in \mathbb{R}, \rho_c(x) =  x - c </math>. Describe the empirical version. (part-1)</p>
<p>Find a Z-estimator <math>\gamma</math> for the situation <math>X \in \mathbb{R}, c \in \mathbb{R}, \rho_c(x) =  x - c </math>. Describe the empirical version. (part-2)</p>	<p>Consider the problem of finding an MLE estimator where the density of <math>X</math> is <math>p_\theta</math>. What is the loss function <math>\rho_\theta(x)</math> in this case?</p>
<p>What is the Kullback-Leibler information?</p>	<p>Show that if the loss function is <math>\rho_\theta(x) = -\log p_\theta(x)</math> and the true parameter is <math>\theta \in \Theta</math>, then for any other <math>\vartheta \in \Theta</math> the Kullback-Leibler information is non-negative</p> $K(\vartheta \mid \theta) = R(\vartheta) - R(\theta) \geq 0$ <p>holds.</p>

<p>Suppose the parameter space is an open set <math>\Gamma \subset \mathbb{R}^p</math>. Suppose the partial derivative</p> $\psi_c(x) = \frac{\partial}{\partial c} \rho_c(x) = \dot{\rho}_c(x)$ <p>exists for all <math>x</math>. Then the Z-estimator <math>\hat{\gamma}_n</math> is a solution for which the partial derivative of the empirical risk is 0, so <math>\dot{\hat{R}}_n(c) \Big _{c=\hat{\gamma}_n} = 0</math> where the derivative <math>\dot{\hat{R}}_n</math> is</p> $\dot{\hat{R}}_n(c) = \frac{\partial}{\partial c} \hat{R}_n(c) = \frac{1}{n} \sum_{i=1}^n \psi_c(x_i)$	<p>Note that</p> $\gamma = \mathbb{E}X = \arg \min_{c \in \Gamma} \mathbb{E} [(X - c)^2]$ <p>so the M-estimator is simply the sample average, i.e.</p> $\hat{\gamma}_n = \arg \min_{c \in \Gamma} \frac{1}{n} \sum_{i=1}^n (X_i - c)^2 = \bar{X}_n$
<p>Let <math>F</math> be the distribution function <math>F(\cdot) = P(X \leq \cdot)</math>. Now the risk is (Partial integration)</p> $\begin{aligned} \mathbb{E} X - c  &= \int_{x \leq c} (c - x) dF(x) + \int_{x > c} (x - c) dF(x) \\ &= \int_{x \leq c} (c - x) dF(x) - \int_{x > c} (x - c) d(1 - F(x)) \\ &= \overbrace{\left[ (c - x)F(x) \right]_{-\infty}^c}^{=0} - \overbrace{\left[ (x - c)(1 - F(x)) \right]_c^{\infty}}^{=0} \\ &\quad + \int_{x \leq c} F(x) dx + \int_{x > c} (1 - F(x)) dx. \end{aligned}$	<p>We have</p> $\begin{aligned} \psi_c(x) &= \frac{\partial}{\partial c} \rho_c(x) = -2(x - c), \\ \dot{\hat{R}}_n(c) \Big _{c=\gamma} &= -2(\bar{X}_n - c) \Big _{c=\gamma} \triangleq 0 \\ \gamma &= \bar{X}_n. \end{aligned}$
<p>The empirical distribution function is</p> $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n 1\{x_i \leq x\}$ <p>Then we have the empirical risk as</p> $\hat{R}_n(c) = \int_{x \leq c} \hat{F}_n(x) dx + \int_{x > c} (1 - \hat{F}_n(x)) dx$ <p>with the derivative <math>\dot{\hat{R}}_n(c) = 2\hat{F}_n(c) - 1</math></p>	<p>So, the derivative of risk is</p> $\dot{R}(c) = F(c) - (1 - F(c)) = 2F(c) - 1.$ <p>Putting that to 0, we get</p> $\dot{R}(\gamma) = 2F(\gamma) - 1 \triangleq 0 \implies \gamma = F^{-1}\left(\frac{1}{2}\right).$
<p>It's <math>\rho_\theta(x) = -\log p_\theta(x)</math></p>	<p>So, we need to set</p> $\hat{F}_n(\hat{\gamma}_n) \triangleq \frac{1}{2}$ <p>which makes <math>\hat{\gamma}_n</math> the median, so by convention</p> $\hat{\gamma}_n := \begin{cases} X_{(\frac{n+1}{2})} & \text{for } n \text{ odd} \\ \frac{X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)}}{2} & \text{for } n \text{ even} \end{cases}$
<p>Take any <math>\tilde{\vartheta} \in \Theta</math></p> $\begin{aligned} R(\tilde{\vartheta}) - R(\theta) &= -\mathbb{E}_\theta \log \frac{p_{\tilde{\vartheta}}(x)}{p_{\tilde{\theta}}(x)} \\ &\geq -\log \mathbb{E}_\theta \frac{p_{\tilde{\vartheta}}(X)}{p_\theta(X)} \quad (\text{Jensen}) \\ &= -\log \left[ \int \frac{p_{\tilde{\vartheta}}}{p_\theta} p_\theta \right] = -\log 1 = 0. \end{aligned}$	<p>Suppose <math>\Theta \subset \mathbb{R}^p</math> and that the densities <math>p_\theta = dP_\theta/d\nu</math> exist with respect to some <math>\sigma</math>-finite measure <math>\nu</math>. The Kullback-Leibler information is then</p> $K(\tilde{\theta} \mid \theta) = \mathbb{E}_\theta \log \left( \frac{p_\theta(X)}{p_{\tilde{\theta}}(X)} \right)$

<p>LEMMA</p> <p>State the lemma about uniform convergence of empirical M-estimators.</p>	<p>Show that if we have the uniform convergence</p> $\sup_{c \in \Gamma} \left  \hat{R}_n(c) - R(c) \right  \xrightarrow[n \rightarrow \infty]{} 0, \mathbb{P}_{\theta^-} \text{ a.s. } .$ <p>Then for the empirical risk minimizer <math>\hat{\gamma}_n</math></p> $R(\hat{\gamma}_n) \xrightarrow[n \rightarrow \infty]{} R(\gamma), \mathbb{P}_{\theta^-} \text{ a.s.}$
<p>What does it mean for <math>\gamma</math> to be well-separated?</p>	<p>If <math>\gamma</math> is well-separated, what does</p> $R(\hat{\gamma}_n) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} R(\gamma)$ <p>imply regarding <math>\hat{\gamma}_n</math> and <math>\gamma</math> ?</p>
<p>What does the notation below concerning the (empirical) measure mean? Let <math>f : \mathcal{X} \rightarrow \mathbb{R}^p</math>.</p> <ol style="list-style-type: none"> <li>1. <math>\hat{P}_n f</math>,</li> <li>2. <math>P f</math>.</li> </ol>	<p>Write <math>R(c)</math> and <math>\hat{R}_n(c)</math> using the <math>P</math> and <math>\hat{P}_n</math> notation.</p>
<p>THEOREM</p> <p>State the uniform law of large numbers</p>	<p>Show that if</p> <ol style="list-style-type: none"> <li>(i) <math>\Gamma</math> is compact,</li> <li>(ii) <math>c \mapsto \rho_c(x)</math> is continuous for all <math>x</math>,</li> <li>(iii) <math>\mathbb{E} \sup_{c \in \Gamma}  \rho_c(X)  &lt; \infty</math>.</li> </ol> <p>then for all <math>\varepsilon</math> there exists a finite ”” <math>\varepsilon</math>-bracketing set”” <math>\{[f_j^L, f_j^U]\}_{j=1}^N</math> such that (continued)</p>
<p>(continued) for <math>f_j^L : \mathcal{X} \rightarrow \mathbb{R}, f_j^U : \mathcal{X} \rightarrow \mathbb{R}</math></p> <ol style="list-style-type: none"> <li>1. <math>f_j^U \geq f_j^L</math></li> <li>2. <math>P(f_j^U - f_j^L) \leq \varepsilon</math>, and</li> <li>3. for all <math>c \in \Gamma</math> there exists a <math>j \in \{1, \dots, N\}</math> such that</li> </ol> $f_j^L \leq \rho_c \leq f_j^U$	<p>THEOREM</p> <p>State the uniform law of large numbers</p>

$ \begin{aligned} 0 &\leq R(\hat{\gamma}_n) - R(\gamma) \\ &= - \left[ \hat{R}_n(\hat{\gamma}_n) - R(\hat{\gamma}_n) \right] + \left[ \hat{R}_n(\gamma) - R(\gamma) \right] \\ &\quad + \underbrace{\hat{R}_n(\hat{\gamma}_n) - \hat{R}_n(\gamma)}_{\leq 0(\hat{\gamma}_n \text{ is ERM})} \\ &\leq 2 \sup_c \left  \hat{R}_n(c) - R(c) \right  \xrightarrow{n \rightarrow \infty} 0, \mathbb{P}_{\theta-\text{a.s.}} \end{aligned} $	<p>Suppose the uniform convergence</p> $\sup_{c \in \Gamma} \left  \hat{R}_n(c) - R(c) \right  \xrightarrow{n \rightarrow \infty} 0, \mathbb{P}_{\theta-\text{a.s.}}$ <p>Then for the empirical risk minimizer <math>\hat{\gamma}_n</math></p> $R(\hat{\gamma}_n) \xrightarrow{n \rightarrow \infty} R(\gamma), \mathbb{P}_{\theta-\text{a.s.}}$
<p>If <math>\gamma</math> is well-separated, then</p> $R(\hat{\gamma}_n) \xrightarrow[n \rightarrow \infty]{\mathbb{P}} R(\gamma)$ <p>implies</p> $\ \hat{\gamma}_n - \gamma\  \xrightarrow[n \rightarrow \infty]{P} 0$	<p><math>\gamma</math> is well-separated if for all <math>\varepsilon &gt; 0</math> we have</p> $\inf\{R(c) : \ c - \gamma\  > \varepsilon\} > R(\gamma)$  <p>(Figure: <math>\gamma</math> not well separated)</p>
<ol style="list-style-type: none"> <li>1. <math>R(c) = P\rho_c</math></li> <li>2. <math>\hat{R}_n(c) = \hat{P}\rho_c</math></li> </ol>	<ol style="list-style-type: none"> <li>1. <math>\hat{P}_n f = \frac{1}{n} \sum_{i=1}^n f(X_i)</math></li> <li>2. <math>Pf = \mathbb{E}f(X)</math></li> </ol>
<p>Let <math>\delta &gt; 0</math> and define</p> $w(x, \delta, c) = \sup_{\tilde{c} \in \Gamma : \ \tilde{c} - c\  < \delta}  \rho_c(x) - \rho_{\tilde{c}}(x) .$ <p>By continuity (ii) <math>\lim_{\delta \downarrow 0} w(x, \delta, c) = 0</math> for all <math>x</math>.  By assumption (iii), the functions are dominated by their supremum which is finite, so we can use dominated convergence to show</p> $\lim_{\delta \downarrow 0} Pw(\cdot, \delta, c) = 0$ <p>for all <math>c</math>.</p>	<p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma</math> is compact,</li> <li>2. <math>c \mapsto \rho_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E} \sup_{c \in \Gamma}  \rho_c(X)  &lt; \infty</math>.</li> </ol> <p>Then we have the uniform convergence</p> $\sup_{c \in \Gamma} \left  \hat{R}_n(c) - R(c) \right  \xrightarrow{n \rightarrow \infty} 0, P_{\theta-\text{a.s.}}$
<p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma</math> is compact,</li> <li>2. <math>c \mapsto \rho_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E} \sup_{c \in \Gamma}  \rho_c(X)  &lt; \infty</math>.</li> </ol> <p>Then we have the uniform convergence</p> $\sup_{c \in \Gamma} \left  \hat{R}_n(c) - R(c) \right  \xrightarrow{n \rightarrow \infty} 0, P_{\theta-\text{a.s.}}$	<p>So, <math>\forall \varepsilon &gt; 0</math> there exists a <math>\delta_c</math> such that <math>Pw(\cdot, \delta_c, c) \leq \frac{\varepsilon}{2}</math>.  Define the open cover <math>B_c = \{\tilde{c} : \ \tilde{c} - c\  &lt; \delta_c\}</math> of <math>\Gamma</math>.  By compactness (i) of <math>\Gamma</math>, there exists a finite subcover <math>\{B_{c_j}\}_{j=1}^N</math>. Let <math>c \in \Gamma</math> be arbitrary, then there exists a <math>j : c \in B_{c_j}</math> such that</p> $ \begin{aligned} \rho_c &\leq \rho_{c_j} + w(\cdot, d_{c_j}, c_j) =: f_j^U \\ \rho_c &\geq \rho_{c_j} - w(\cdot, d_{c_j}, c_j) =: f_j^L \end{aligned} $ <p>Then <math>P(f_j^U - f_j^L) = 2Pw(\cdot, \delta_{c_j}, c_j) \leq \varepsilon</math> as we wished.</p>

<p>Suppose we wish to show that if</p> <ul style="list-style-type: none"> <li>(i) <math>\Gamma</math> is compact,</li> <li>(ii) <math>c \mapsto \rho_c(x)</math> is continuous for all <math>x</math>,</li> <li>(iii) <math>\mathbb{E} \sup_{c \in \Gamma}  \rho_c(X)  &lt; \infty</math>.</li> </ul> <p>then we have the uniform convergence</p> $\sup_{c \in \Gamma}  \hat{R}_n(c) - R(c)  \xrightarrow{n \rightarrow \infty} 0, P_{\theta^-} \text{ a.s.}$ <p>(continued)</p>	<p>(continued) Show that to show the above it is sufficient that there exists a finite "ε-bracketing set"</p> <p><math>\{[f_j^L, f_j^U]\}_{j=1}^N</math> such that for <math>f_j^L : \mathcal{X} \rightarrow \mathbb{R}, f_j^U : \mathcal{X} \rightarrow \mathbb{R}</math></p> <ol style="list-style-type: none"> <li>1. <math>f_j^U \geq f_j^L</math></li> <li>2. <math>P(f_j^U - f_j^L) \leq \varepsilon</math>, and</li> <li>3. for all <math>c \in \Gamma</math> there exists a <math>j \in \{1, \dots, N\}</math> such that</li> </ol> $f_j^L \leq \rho_c \leq f_j^U$
<p>Suppose we wish to show that if</p> <ul style="list-style-type: none"> <li>(i) <math>\Gamma</math> is compact,</li> <li>(ii) <math>c \mapsto \rho_c(x)</math> is continuous for all <math>x</math>,</li> <li>(iii) <math>\mathbb{E} \sup_{c \in \Gamma}  \rho_c(X)  &lt; \infty</math>.</li> </ul> <p>then we have the uniform convergence</p> $\sup_{c \in \Gamma}  \hat{R}_n(c) - R(c)  \xrightarrow{n \rightarrow \infty} 0, P_{\theta^-} \text{ a.s.}$ <p>What are the key steps of the proof? (part-1)</p>	<p>LEMMA</p> <p>State lemma about existence of Z-estimators.</p>
<p>PROOF</p> <p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma \subset \mathbb{R}</math>,</li> <li>2. <math>c \mapsto \psi_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E}_{\theta}  \psi_c(x)  &lt; \infty</math> for all <math>c</math></li> <li>4. there exists a <math>\delta &gt; 0</math> such that <math>\dot{R}(c) &gt; 0</math> for <math>c \in (\gamma, \gamma + \delta)</math> and <math>\dot{R}(c) &lt; 0</math> for <math>c \in (\gamma - \delta, \gamma)</math>.</li> </ol> <p>Then <math>\mathbb{P}_{\theta}</math>-a.s. there is a solution <math>\hat{\gamma}_n</math> of <math>\dot{R}_n(\hat{\gamma}_n) = 0</math> and this solution is consistent. (part-1)</p>	<p>PROOF</p> <p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma \subset \mathbb{R}</math>,</li> <li>2. <math>c \mapsto \psi_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E}_{\theta}  \psi_c(x)  &lt; \infty</math> for all <math>c</math></li> <li>4. there exists a <math>\delta &gt; 0</math> such that <math>\dot{R}(c) &gt; 0</math> for <math>c \in (\gamma, \gamma + \delta)</math> and <math>\dot{R}(c) &lt; 0</math> for <math>c \in (\gamma - \delta, \gamma)</math>.</li> </ol> <p>Then <math>\mathbb{P}_{\theta}</math>-a.s. there is a solution <math>\hat{\gamma}_n</math> of <math>\dot{R}_n(\hat{\gamma}_n) = 0</math> and this solution is consistent. (part-2)</p>
<p>PROOF</p> <p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma \subset \mathbb{R}</math>,</li> <li>2. <math>c \mapsto \psi_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E}_{\theta}  \psi_c(x)  &lt; \infty</math> for all <math>c</math></li> <li>4. there exists a <math>\delta &gt; 0</math> such that <math>\dot{R}(c) &gt; 0</math> for <math>c \in (\gamma, \gamma + \delta)</math> and <math>\dot{R}(c) &lt; 0</math> for <math>c \in (\gamma - \delta, \gamma)</math>.</li> </ol> <p>Then <math>\mathbb{P}_{\theta}</math>-a.s. there is a solution <math>\hat{\gamma}_n</math> of <math>\dot{R}_n(\hat{\gamma}_n) = 0</math> and this solution is consistent. (part-3)</p>	<p>Let <math>X \in \mathbb{R}, \theta \in \mathbb{R}</math>. Take the density</p> $p_{\theta}(x) = \frac{\exp(x - \theta)}{(1 + \exp(x - \theta))^2}, \text{ for } x \in \mathbb{R}.$ <p>Let the loss <math>\rho_{\theta}</math> be</p> $\rho_{\theta}(x) = -(x - \theta) + 2 \log(1 + \exp(x - \theta))$ <p>with</p> $\frac{\partial}{\partial c} \rho_{\theta}(x) = \frac{1 - \exp(x - \theta)}{1 + \exp(x - \theta)}.$ <p>Does there exist a consistent Z-estimator for this problem?</p>
<p>Consider a function <math>f : \mathcal{X} \rightarrow \mathbb{R}^p</math>. What does <math>P f f^T</math> mean?</p>	<p>State the multidimensional central limit theorem using (empirical) measure notation, i.e. using <math>\hat{P}_n f</math> and <math>P f</math>.</p>

<p>Hence</p> $\sup_{c \in \Gamma} \left  \left( \hat{P}_n - P \right) \rho_c \right  \leq \underbrace{\max_{1 \leq j \leq N} \max \left\{ \left  \left( \hat{P}_n - P \right) f_j^L \right , \left  \left( \hat{P}_n - P \right) f_j^U \right  \right\}}_{\xrightarrow{n \rightarrow \infty} 0, P_\theta\text{-a.s.}} + \varepsilon$	<p>If we have this, then</p> $\begin{aligned} \left( \hat{P}_n - P \right) \rho_c &\leq \hat{P}_n f_j^U - P f_j^U + \underbrace{P \left( f_j^U - \rho_c \right)}_{\leq \varepsilon}, \\ \left( \hat{P}_n - P \right) \rho_c &\geq \hat{P}_n f_j^L - P f_j^L + \underbrace{P \left( f_j^L - \rho_c \right)}_{\geq -\varepsilon}. \end{aligned}$
<p>Suppose</p> <ol style="list-style-type: none"> <li>1. <math>\Gamma \subset \mathbb{R}</math>,</li> <li>2. <math>c \mapsto \psi_c(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math>\mathbb{E}_\theta  \psi_c(x)  &lt; \infty</math> for all <math>c</math></li> <li>4. there exists a <math>\delta &gt; 0</math> such that <math>\dot{R}(c) &gt; 0</math> for <math>c \in (\gamma, \gamma + \delta)</math> and <math>\dot{R}(c) &lt; 0</math> for <math>c \in (\gamma - \delta, \gamma)</math>.</li> </ol> <p>Then <math>\mathbb{P}_\theta</math>-a.s. there is a solution <math>\hat{\gamma}_n</math> of <math>\dot{R}_n(\hat{\gamma}_n) = 0</math> and this solution is consistent.</p>	<ol style="list-style-type: none"> <li>1. Let <math>\delta &gt; 0</math> and define <math display="block">w(x, \delta, c) = \sup_{\tilde{c} \in \Gamma: \ \tilde{c} - c\  &lt; \delta}  \rho_c(x) - \rho_{\tilde{c}}(x) </math></li> <li>2. Use continuity and dominated convergence to show that for all <math>\varepsilon &gt; 0</math> there exists a <math>\delta_c</math> such that <math>Pw(\cdot, \delta_c, c) \leq \frac{\varepsilon}{2}</math></li> <li>3. Use the above to construct an open cover <math>B_c = \{\tilde{c} : \ \tilde{c} - c\  &lt; \delta_c\}</math> of <math>\Gamma</math> such that <math>\rho_c</math> is bounded from below and above by <math>\pm w(\cdot, \delta_c, c)</math>.</li> <li>4. Use compactness to show that there is a finite number of such bounds.</li> <li>5. Use previous card</li> </ol>
<p>On the set <math>A</math> we have</p> $\begin{aligned} \dot{R}_n(\gamma + \varepsilon) &= \underbrace{\dot{R}_n(\gamma + \varepsilon) - \dot{R}(\gamma + \varepsilon)}_{\geq -\tilde{\varepsilon}} + \underbrace{\dot{R}(\gamma + \varepsilon)}_{\geq 2\tilde{\varepsilon}} \geq \tilde{\varepsilon}, \\ \dot{R}_n(\gamma - \varepsilon) &= \underbrace{\dot{R}_n(\gamma - \varepsilon) - \dot{R}(\gamma - \varepsilon)}_{\leq \tilde{\varepsilon}} + \underbrace{\dot{R}(\gamma - \varepsilon)}_{\leq -2\tilde{\varepsilon}} \leq -\tilde{\varepsilon}. \end{aligned}$ <p>Now by continuity (ii) on the set <math>A</math> there is a <math>\hat{\gamma}_n \in (\gamma - \varepsilon, \gamma + \varepsilon)</math> such that <math>\dot{R}_n(\hat{\gamma}_n) = 0</math> and <math>\ \hat{\gamma}_n - \gamma\  \xrightarrow[n \rightarrow \infty]{P} 0</math></p>	<p>Let <math>0 &lt; \varepsilon &lt; \delta</math>. By (iv) there is an <math>\tilde{\varepsilon}</math> such that</p> $\dot{R}(\gamma + \varepsilon) \geq 2\tilde{\varepsilon} \text{ and } \dot{R}(\gamma - \varepsilon) \leq -2\tilde{\varepsilon}.$ <p>Let</p> $A = \left\{ \dot{R}_n(\gamma + \varepsilon) - \dot{R}(\gamma + \varepsilon) > -\tilde{\varepsilon}, \dot{R}_n(\gamma - \varepsilon) - \dot{R}(\gamma - \varepsilon) < \tilde{\varepsilon} \right\}$ <p>Then by the uniform law of large numbers <math>\mathbb{P}(A^C) \xrightarrow[n \rightarrow \infty]{} 0</math>.</p>
<p>The conditions for existence of consistent Z-estimators hold, so yes:</p> <ol style="list-style-type: none"> <li>1. <math>\vartheta \in \Gamma \subset \mathbb{R}</math></li> <li>2. <math>\vartheta \mapsto \psi_\vartheta(x) = \frac{\partial}{\partial \vartheta} \rho_\vartheta(x)</math> is continuous for all <math>x</math>,</li> <li>3. <math> \psi_\vartheta(x)  \leq 1</math>, so <math>\mathbb{E}_\theta  \psi_\vartheta(X)  &lt; \infty</math> for all <math>\vartheta</math>,</li> <li>4. since <math>\psi_\vartheta(x)</math> passes 0 and is continuous, for <math>\dot{R}(\vartheta) = \mathbb{E}_\theta \psi_\vartheta(X)</math> there exists a <math>\delta &gt; 0</math> such that</li> </ol> $\dot{R}(c) > 0 \text{ for } c \in (\gamma, \gamma + \delta) \text{ and } \dot{R}(c) < 0, c \in (\gamma - \delta, \gamma)$	
<p>Suppose <math>X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X \in \mathcal{X}</math>, and <math>f : \mathcal{X} \rightarrow \mathbb{R}^p</math> is a function with</p> $Pf < \infty \text{ and } \Sigma_f = Pff^T - (Pf)(Pf^T) < \infty,$ <p>then</p> $\sqrt{n} \left( \hat{P}_n - P \right) f \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma_f).$	$Pff^T = \begin{pmatrix} \mathbb{E} f_1^2(X) & \cdots & \mathbb{E} f_1(X) f_p(X) \\ \vdots & \ddots & \vdots \\ \mathbb{E} f_1(X) f_p(X) & \cdots & \mathbb{E} f_p^2(X) \end{pmatrix}$

<p>What is an empirical process?</p>	<p>Suppose</p> $\nu_n(c) = \sqrt{n} \left( \hat{P}_n - P \right) \psi_c$ <p>is an empirical process. What and how does <math>\nu_n(c)</math> converge if <math>P\psi_c\psi_c^T = \Sigma_c &lt; \infty</math> ?</p>
<p>What does it mean for an empirical process <math>\nu_n(\cdot)</math> to be asymptotically continuous?</p>	<p>THEOREM</p> <p>State the theorem about asymptotic linearity of Z-estimators.</p>
<p>Sketch proof: suppose</p> <ol style="list-style-type: none"> <li>1. <math>\hat{\gamma}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \gamma</math></li> <li>2. <math>\dot{\hat{R}}_n(\hat{\gamma}_n) = 0, \dot{R}(\gamma) = 0</math></li> <li>3. <math>\nu_n</math> is asymptotically continuous at <math>\gamma</math>,</li> <li>4. <math>M_\theta = \frac{\partial}{\partial \gamma^T} \dot{R}(\gamma)</math> exists and is invertible,</li> <li>5. <math>J_\theta := P\psi_\gamma\psi_\gamma^T &lt; \infty</math></li> </ol> <p>Then <math>\hat{\gamma}_n</math> has influence function <math>l_\theta = -M_\theta^{-1}\psi_\gamma</math>. (part-1)</p>	<p>Sketch proof: suppose</p> <ol style="list-style-type: none"> <li>1. <math>\hat{\gamma}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \gamma</math></li> <li>2. <math>\dot{\hat{R}}_n(\hat{\gamma}_n) = 0, \dot{R}(\gamma) = 0</math></li> <li>3. <math>\nu_n</math> is asymptotically continuous at <math>\gamma</math>,</li> <li>4. <math>M_\theta = \frac{\partial}{\partial \gamma^T} \dot{R}(\gamma)</math> exists and is invertible,</li> <li>5. <math>J_\theta := P\psi_\gamma\psi_\gamma^T &lt; \infty</math></li> </ol> <p>Then <math>\hat{\gamma}_n</math> has influence function <math>l_\theta = -M_\theta^{-1}\psi_\gamma</math>. (part-2)</p>
<p>Sketch proof: suppose</p> <ol style="list-style-type: none"> <li>1. <math>\hat{\gamma}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \gamma</math></li> <li>2. <math>\dot{\hat{R}}_n(\hat{\gamma}_n) = 0, \dot{R}(\gamma) = 0</math></li> <li>3. <math>\nu_n</math> is asymptotically continuous at <math>\gamma</math>,</li> <li>4. <math>M_\theta = \frac{\partial}{\partial \gamma^T} \dot{R}(\gamma)</math> exists and is invertible,</li> <li>5. <math>J_\theta := P\psi_\gamma\psi_\gamma^T &lt; \infty</math></li> </ol> <p>Then <math>\hat{\gamma}_n</math> has influence function <math>l_\theta = -M_\theta^{-1}\psi_\gamma</math>. (part-3)</p>	<p>What condition must be true to be able to write</p> $\dot{\hat{R}}_n(\hat{\gamma}_n) - \dot{R}(\hat{\gamma}_n) = \dot{\hat{R}}_n(\gamma) - \dot{R}(\gamma) + o_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right)?$
<p>Suppose <math>M_\theta = \frac{\partial}{\partial c^T} \dot{R}(\gamma)</math>. Give an asymptotic Taylor approximation to <math>\dot{R}(\hat{\gamma}_n) - \dot{R}(\gamma)</math>?</p>	<p>Suppose we know <math>J_\theta = P\psi_\gamma\psi_\gamma^T &lt; \infty</math>. What does that tell us about the order symbol of <math>\frac{1}{\sqrt{n}}\nu_n(\gamma)</math> ?</p>

<p>If <math>P\psi_c\psi_c^T = \Sigma_c &lt; \infty</math>, then <math>\nu_n(c) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma_c)</math></p>	<p>We call</p> $\nu_n(c) = \sqrt{n} \left( \hat{P}_n - P \right) \psi_c$ <p>an empirical process indexed by <math>c \in \Gamma</math>.</p>
<ol style="list-style-type: none"> <li>1. <math>\hat{\gamma}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \gamma</math></li> <li>2. <math>\dot{\hat{R}}_n(\hat{\gamma}_n) = 0, \dot{R}(\gamma) = 0</math></li> <li>3. <math>\nu_n</math> is asymptotically continuous at <math>\gamma</math>,</li> <li>4. <math>M_\theta = \frac{\partial}{\partial \gamma^T} \dot{R}(\gamma)</math> exists and is invertible,</li> <li>5. <math>J_\theta := P\psi_\gamma\psi_\gamma^T &lt; \infty</math></li> </ol> <p>Then <math>\hat{\gamma}_n</math> has influence function <math>l_\theta = -M_\theta^{-1}\psi_\gamma</math>. Furthermore</p> $\sqrt{n}(\hat{\gamma}_n - \gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \underbrace{M_\theta^{-1} J_\theta M_\theta^{-1}}_{\text{" sandwich formula" }})$	<p>An empirical process <math>\nu_n(\cdot)</math> is called asymptotically continuous at <math>\gamma</math> if</p> $ \nu_n(\gamma_n) - \nu_n(\gamma)  \xrightarrow[n \rightarrow \infty]{\mathbb{P}} 0 \text{ as } \gamma_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \gamma$ <p>More formally, for all <math>\varepsilon &gt; 0</math> there exists a <math>\delta &gt; 0</math> such that</p> $\mathbb{P} \left( \sup_{\ c - \gamma\  \leq \delta}  \nu_n(c) - \nu_n(\gamma)  > \varepsilon \right) \xrightarrow[n \rightarrow \infty]{} 0$
$(*) = \frac{1}{\sqrt{n}} \nu_n(\gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) + M_\theta ([\hat{\gamma}_n - \gamma] [1 + o_{\mathbb{P}}(1)])$ <p>(Assum. (v) <math>\implies \nu_n(\gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, J_\theta) \implies \nu_n(\gamma) = \mathcal{O}_{\mathbb{P}}(1)</math>)</p> $= \mathcal{O}_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) + M_\theta ([\hat{\gamma}_n - \gamma] [1 + o_{\mathbb{P}}(1)])$ $\implies M_\theta ([\hat{\gamma}_n - \gamma] [1 + o_{\mathbb{P}}(1)]) = \mathcal{O}_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) \text{ by (iv) } M_\theta \text{ is invertible} \implies \ \hat{\gamma}_n - \gamma\  = \mathcal{O}_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$	$0 = \dot{\hat{R}}_n(\hat{\gamma}_n) \quad (2.)$ $= \underbrace{\dot{\hat{R}}_n(\hat{\gamma}_n) - \dot{R}(\hat{\gamma}_n)}_{= \frac{1}{\sqrt{n}} \nu_n(\hat{\gamma}_n)} + \dot{R}(\hat{\gamma}_n)$ $= \underbrace{\dot{\hat{R}}_n(\gamma) - \dot{R}(\gamma)}_{= \frac{1}{\sqrt{n}} \nu_n(\gamma)} + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) + \dot{R}(\hat{\gamma}_n) \quad (3. \text{asympt cont})$ $= \frac{1}{\sqrt{n}} \nu_n(\gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) + \dot{R}(\hat{\gamma}_n) - \underbrace{\dot{R}(\gamma)}_{=0}$
$\nu_n(\gamma) = \left( \hat{P}_n - P \right) \psi_\gamma = \dot{\hat{R}}_n(\gamma) - \dot{R}(\gamma)$ <p>must be asymptotically continuous at <math>\gamma</math>.</p>	$(*) = \frac{1}{\sqrt{n}} \nu_n(\gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right) + M_\theta (\hat{\gamma}_n - \gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$ $\implies M_\theta (\hat{\gamma}_n - \gamma) = \frac{1}{\sqrt{n}} \nu_n(\gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$ $(\hat{\gamma}_n - \gamma) = \frac{1}{\sqrt{n}} M_\theta^{-1} \nu_n(\gamma) + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$ $= \hat{P}_n M_\theta^{-1} \psi_\gamma + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$ $= \hat{P}_n l_\theta + o_{\mathbb{P}} \left( \frac{1}{\sqrt{n}} \right)$
<p><math>J_\theta = P\psi_\gamma\psi_\gamma^T &lt; \infty</math> implies <math>\nu_n(\gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, J_\theta)</math>, so the order of <math>\nu_n(\gamma) = \mathcal{O}(1)</math>, and so finally</p> $\frac{1}{\sqrt{n}} \nu_n(\gamma) = \mathcal{O} \left( \frac{1}{\sqrt{n}} \right)$	$\dot{\hat{R}}(\hat{\gamma}_n) - \dot{R}(\gamma) = M_\theta ([\hat{\gamma}_n - \gamma] [1 + o_{\mathbb{P}}(1)])$



<p>Consider the maximum likelihood estimator <math>\hat{\gamma}_n</math>, for <math>X \sim P_\theta</math> with density <math>p_\theta</math>, <math>\theta \in \mathbb{R}^p</math>, the loss function</p> $\rho_\theta(x) = -\log p_\theta(x)$ <p>and <math>\gamma = \theta</math>. Assuming regularity conditions hold, what can you say about the convergence of</p> $\sqrt{n}(\hat{\gamma}_n - \gamma)?$	<p>What are the steps to find the influence function and asymptotic variance of a Z-estimator?</p>
<p>Suppose</p> $\rho(x) = \begin{cases} x^2 & \text{if }  x  \leq k \\ 2( x  - k) & \text{if }  x  > k \end{cases}$ <p>and <math>\rho_c(x) = \rho(x - c)</math> for <math>c \in \mathbb{R}</math>, so <math>\rho_c</math> is the Huber loss. What is the influence function and asymptotic variance of an estimator <math>\hat{\gamma}_n</math> which minimizes the Huber loss? (part-1)</p>	<p>Suppose</p> $\rho(x) = \begin{cases} x^2 & \text{if }  x  \leq k \\ 2( x  - k) & \text{if }  x  > k \end{cases}$ <p>and <math>\rho_c(x) = \rho(x - c)</math> for <math>c \in \mathbb{R}</math>, so <math>\rho_c</math> is the Huber loss. What is the influence function and asymptotic variance of an estimator <math>\hat{\gamma}_n</math> which minimizes the Huber loss? (part-2)</p>
<p>Suppose</p> $\rho(x) = \begin{cases} x^2 & \text{if }  x  \leq k \\ 2( x  - k) & \text{if }  x  > k \end{cases}$ <p>and <math>\rho_c(x) = \rho(x - c)</math> for <math>c \in \mathbb{R}</math>, so <math>\rho_c</math> is the Huber loss. What is the influence function and asymptotic variance of an estimator <math>\hat{\gamma}_n</math> which minimizes the Huber loss? (part-3)</p>	<p>DEFINITION</p> <p>What is asymptotic relative efficiency?</p>
<p>Consider two estimators <math>T_{n,1}, T_{n,2}</math>. What does it mean if the asymptotic relative efficiency</p> $e_{2,1} > 1 \text{ ?}$	<p>What is the <math>(1 - \alpha)</math> asymptotic confidence interval based on an asymptotically normal estimator <math>T_n</math> of <math>\gamma</math> with asymptotic variance <math>V_\theta</math> ?</p>
<p>Consider two asymptotically normal estimators <math>T_{\theta,1}, T_{\theta,2}</math> with asymptotic relative efficiency <math>e_{2,1}</math>. What must be the ratio of sample size <math>n_1/n_2</math> for them to produce asymptotic confidence intervals of equal length?</p>	<p>Suppose <math>X = \mu + \varepsilon</math>, where <math>\varepsilon \sim F_0</math> symmetric around zero with density <math>f_0</math>, and <math>\text{Var}(\varepsilon) := \sigma^2 &lt; \infty</math>. Let <math>X_1, \dots, X_n, \dots</math> i. <math>\stackrel{i.i.d.}{\sim} X</math>. Consider the estimators</p> $T_{n,1} := \bar{X}_n \quad \text{w/ asymp. var. } V_{\theta,1} = \sigma^2$ $T_{n,2} := \text{sample median} \quad \text{w/ asymp. var. } V_{\theta,2} = \frac{1}{4f_0^2(0)}$ <p>If <math>F_0 = \Phi</math> is the standard normal distribution, what is the asymptotic relative efficiency and which estimator is more efficient? Recall that <math>\Phi'(x) = \phi(x) = \frac{1}{\sqrt{2\pi}} \exp[-\frac{1}{2}x^2]</math></p>

<ol style="list-style-type: none"> <li>Find <math>\psi_c = \frac{\partial}{\partial c} \rho_c(x)</math>.</li> <li>Compute <math>\dot{R}(c) = \mathbb{E} \psi_c(x)</math>, and find <math>\gamma</math> for which <math>\dot{R}(\gamma) = 0</math>.</li> <li>Determine <math>M_\theta = \left. \frac{\partial}{\partial c^T T} \dot{R}(c) \right _{c=\gamma}</math>.</li> <li>The influence function is <math>l_\theta = -M_\theta^{-1} \psi_\gamma</math>.</li> <li>Compute <math>J_\theta = P \psi_\gamma \psi_\gamma^T</math>.</li> <li>The asymptotic variance is <math>V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1}</math>.</li> <li>If the asymptotic variance depends on the unknown distribution and <math>\gamma</math>, by Slutsky's theorem we can plug in empirical estimates and retain the asymptotic properties.</li> </ol>	<p>Under regularity conditions</p> $\sqrt{n}(\hat{\gamma}_n - \gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, I^{-1}(\gamma))$
<ol style="list-style-type: none"> <li>The matrix <math>M_\theta</math> is then <math display="block">\left. \frac{\partial}{\partial c} \mathbb{E} \psi_c(x) \right _{c=\gamma} = 2F(c+k) - 2F(c-k) = M_\theta</math> </li> <li>The influence function is then <math display="block">l_\theta = -M_\theta^{-1} \psi_\gamma</math> <math display="block">= \frac{-1}{2F(\gamma+k) - 2F(\gamma-k)} \begin{cases} -2(x-\gamma) &amp;  x-\gamma  \leq k \\ 2k &amp; x-\gamma &lt; -k \\ -2k &amp; x-\gamma &gt; k \end{cases}</math> </li> </ol>	<ol style="list-style-type: none"> <li>First we need to find <math>\psi_c(x) = \frac{\partial}{\partial c} \rho_c(x)</math> as <math display="block">\psi(x) = \begin{cases} 2x &amp;  x  \leq k \\ -2k &amp; x &lt; -k \\ 2k &amp; x &gt; k, \end{cases}, \psi_c(x) = \begin{cases} -2(x-c) &amp;  x-c  \leq k \\ 2k &amp; x-c &lt; -k \\ -2k &amp; x-c &gt; k \end{cases}</math> </li> <li>Then we need to find <math>M_\theta = \left. \frac{\partial}{\partial c^T} \dot{R}(c) \right _{c=\gamma}</math>, so compute <math>\dot{R}(c)</math> as <math display="block">\dot{R}(c) = \mathbb{E} \psi_c(x) = (\text{long computations}) = 2 \left\{ \int_{c-k}^{c+k} F(x) dx - k \right\}</math> so we get <math>\gamma</math> by setting the above to 0. </li> </ol>
<p>Suppose <math>\gamma \in \Gamma \subset \mathbb{R}</math>, and for two estimators <math>T_{n,1}, T_{n,2}</math> we have</p> $\sqrt{n}(T_{n,1} - \gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}_\theta} \mathcal{N}(0, V_{\theta,1}),$ $\sqrt{n}(T_{n,2} - \gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}_\theta} \mathcal{N}(0, V_{\theta,2})$ <p>Then the asymptotic relative efficiency (a.r.e) is</p> $e_{2,1} := \frac{V_{\theta,1}}{V_{\theta,2}}$	<ol style="list-style-type: none"> <li>The asymptotic variance is <math>V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1}</math>, where <math>J_\theta = P \psi_\gamma \psi_\gamma^T</math>. So <math display="block">V_\theta = \frac{1}{(F(\gamma+k) - F(\gamma-k))^2} \times \left\{ \int_{ x-\gamma  \leq k} (x-\gamma)^2 dF(x) + k^2 F(\gamma-k) + k^2 (1 - F(\gamma+k)) \right\}.</math> </li> <li>The asymptotic variance depends on the unknown distribution <math>F</math>. However, by Slutsky's theorem the asymptotic variance holds even if we plug in the empirical estimates for <math>F</math> and <math>\gamma</math>.</li> </ol>
<p>The <math>(1 - \alpha)</math> asymptotic confidence interval is then</p> $\hat{\gamma}_n \pm \sqrt{\frac{V_\theta}{n}} \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right)$	<p>If <math>e_{2,1} &gt; 1</math>, then <math>T_{n,2}</math> is asymptotically more efficient than <math>T_{n,1}</math>.</p>
<p>The a.r.e is <math>e_{2,1} = 4\sigma^2 f_0^2(0)</math>. If <math>F_0 = \Phi</math>, then <math>f_0 = \phi</math>. Hence, <math>\sigma^2 = 1</math>, <math>f_0(0) = \frac{1}{\sqrt{2\pi}}</math>, so a.r.e. is <math>e_{2,1} = \frac{4}{2\pi} = \frac{2}{\pi} &lt; 1</math> so the estimator <math>T_{n,1} = \bar{X}_n</math> is more efficient.</p>	$\frac{\text{length}_1}{\text{length}_2} = \frac{2\sqrt{\frac{V_{\theta,1}}{n_1}} \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right)}{2\sqrt{\frac{V_{\theta,2}}{n_2}} \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right)} = \sqrt{\frac{n_2}{n_1}} \sqrt{\frac{V_{\theta,1}}{V_{\theta,2}}}$ $= \sqrt{\frac{n_2}{n_1}} \sqrt{e_{2,1}} \triangleq 1,$ $\frac{n_1}{n_2} = e_{2,1}.$

<p>Suppose <math>X = \mu + \varepsilon</math>, where <math>\varepsilon \sim F_0</math> symmetric around zero with density <math>f_0</math>, and <math>\text{Var}(\varepsilon) := \sigma^2 &lt; \infty</math>. Let <math>X_1, \dots, X_n, \dots</math> i. i. d. <math>\sim X</math>. Consider the estimators</p> $T_{n,1} := \bar{X}_n \quad \text{w/ asymp. var. } V_{\theta,1} = \sigma^2$ $T_{n,2} := \text{sample median} \quad \text{w/ asymp. var. } V_{\theta,2} = \frac{1}{4f_0^2(0)}$ <p>If <math>F_0</math> is the Laplace distribution, what is the asymptotic relative efficiency and which estimator is more efficient?</p> $f_0(x) = \frac{1}{\sqrt{2}} \exp[-\sqrt{2} x ]$	<p>Take a random variable <math>Z \in \mathbb{R}^p</math> with <math>Z \sim \mathcal{N}(0, \Sigma), \Sigma &gt; 0</math>. What is the distribution of <math>Z^T \Sigma^{-1} Z</math> ?</p>
<p>DEFINITION</p> <p>What is an asymptotic pivot?</p>	<p>What are the two main ways to construct asymptotic pivots based on asymptotic covariance?</p>
<p>Suppose <math>\hat{\theta}_n</math> is a consistent estimator of <math>\theta</math>, and <math>\theta \mapsto V_\theta</math> is continuous. Give a consistent estimator of the asymptotic covariance <math>V_\theta</math>.</p>	<p>Consider an M-estimator with asymptotic variance <math>V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1}</math>. Give a consistent estimator of the asymptotic covariance <math>V_\theta</math>.</p>
<p>Suppose <math>\theta_n</math> is the maximum likelihood estimator for <math>\theta</math>. What's an asymptotic pivot for <math>\theta</math> that does not rely on the covariance <math>V_\theta</math> ?</p>	<p>Suppose <math>\hat{\theta}_n</math> is the maximum likelihood estimator for <math>\theta</math>, with density <math>p_\theta(x)</math> for <math>\theta \in \Theta \subset \mathbb{R}^p</math>. Show three possible asymptotic pivots for <math>H_0 : \theta = \theta_0</math>, and what you need to compute to derive them. What are their distributions? (part-1)</p>
<p>Suppose <math>\hat{\theta}_n</math> is the maximum likelihood estimator for <math>\theta</math>, with density <math>p_\theta(x)</math> for <math>\theta \in \Theta \subset \mathbb{R}^p</math>. Show three possible asymptotic pivots for <math>H_0 : \theta = \theta_0</math>, and what you need to compute to derive them. What are their distributions? (part-2)</p>	<p>LEMMA</p> <p>State the lemma about asymptotic convergence of restricted MLE. (part-1)</p>

<p>See that <math>\Sigma^{-\frac{1}{2}}Z \sim \mathcal{N}(0, I)</math>. Now we get <math>Z^T \Sigma^{-1} Z = \left\  \Sigma^{-\frac{1}{2}} Z \right\ ^2 \sim \chi_p^2</math></p>	<p>The a.r.e is <math>e_{2,1} = 4\sigma^2 f_0^2(0)</math>. Variance is</p> $\sigma^2 = \mathbb{E}\varepsilon^2 = \sqrt{2} \int_0^\infty x^2 \exp[-\sqrt{2}x] dx = 1.$ <p><math>f_0(0) = \frac{1}{\sqrt{2}}</math>, so a.r.e is <math>e_{2,1} = \frac{4}{2} = 2 &gt; 1</math>, so the median is more efficient than the mean.</p>
<p>1. Suppose the asymptotic variance <math>V_\theta = V(\gamma)</math> depends only on <math>\gamma</math>. Then an asymptotic pivot is</p> $Z_{n,1}(\gamma) := n (T_n - \gamma) V(\gamma)^{-1} (T_n - \gamma)$ <p>and we have <math>Z_{n,1}(\gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \chi_p^2</math>.</p> <p>2. If instead we have a consistent estimator of the variance <math>\hat{V}_n</math>, i.e. <math>\forall \theta</math> it's true that <math>\hat{V}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}_\theta} V_\theta</math>, then an asymptotic pivot is</p> $Z_{n,2}(\gamma) = n (T_n - \gamma) \hat{V}_n^{-1} (T_n - \gamma)$ <p>and using Slutsky's theorem we have again <math>Z_{n,2}(\gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}_\theta} \chi_p^2</math>.</p>	<p>Suppose <math>X \sim P_\theta, X_1, \dots, X_n, \dots</math> are i.i.d. copies of <math>X</math>, and we are testing. a hypothesis <math>H_0 : \gamma = \gamma_0</math>. We call a function <math>Z_n(\gamma) = Z_n(\gamma, X_1, \dots, X_n)</math> an asymptotic pivot if</p> $Z_n(\gamma) \xrightarrow[n \rightarrow \infty]{\mathcal{D}_\theta} Z$ <p>for all <math>\theta</math>, that is the distribution of <math>Z</math> does not depend on parameters <math>\theta</math>.</p>
<p>Just replace the true measure with the empirical measure. So instead of</p> $M_\theta = \frac{\partial}{\partial c^T} \mathbb{E} \psi_c(x) \Big _{c=\gamma} = P \dot{\psi}_\gamma,$ $J_\theta = \mathbb{E} \psi_\gamma(x) \psi_\gamma(x)^T = P \psi_\gamma \psi_\gamma^T,$ <p>use</p> $\hat{M}_n := \hat{P}_n \dot{\psi}_{\hat{\gamma}_n} = \frac{\partial^2}{\partial c \partial c^T} \frac{1}{n} \sum_{i=1}^n \rho_c(x_i) \Big _{c=\hat{\gamma}_n},$ $\hat{J} := \hat{P}_n \psi_{\hat{\gamma}_n} \psi_{\hat{\gamma}_n}^T$ $\hat{V}_n := \hat{M}_n^{-1} \hat{J}_n \hat{M}_n^{-1}$	$\hat{V}_n = V_{\hat{\theta}_n}$
<p>We need</p> $\rho_\vartheta(x) = -\log p_\vartheta(x),$ $\psi_\vartheta(x) = -s_\vartheta(x) = -\frac{\dot{\rho}_\vartheta(x)}{p_\vartheta(x)},$ $I(\theta) = -\mathbb{E} \dot{s}_\theta(x),$ $M_\theta = -P \dot{s}_\theta = I(\theta),$ $J_\theta = P s_\theta s_\theta^T = I(\theta)$ $V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1} = I(\theta)^{-1},$ $l_\theta = I(\theta)^{-1} s_\theta.$	<p>We can construct an asymptotic pivot using the twice log-likelihood ratio</p> $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) := 2 \sum_{i=1}^n \left[ \log p_{\hat{\theta}_n}(X_i) - \log p_\theta(X_i) \right].$ <p>Under regularity conditions, we have the convergence</p> $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \xrightarrow[n \rightarrow \infty]{\mathcal{D}_\theta} \chi_p^2$ <p>for all <math>\theta</math>.</p>
<p>Suppose we are testing the hypothesis <math>H_0 : R(\theta) = 0</math> where <math>R</math> is a vector of restrictions</p> $R(\theta) = (R_1(\theta) \cdots R_q(\theta))^T.$ <p>Let <math>\hat{\theta}_n</math> and <math>\hat{\theta}_n^0</math> be the unrestricted and restricted MLE, i.e.</p> $\hat{\theta}_n = \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_\vartheta(X_i)$ $\hat{\theta}_n^0 = \arg \max_{\vartheta \in \Theta: R(\vartheta)=0} \sum_{i=1}^n \log p_\vartheta(X_i).$	<p>Then we have the three asymptotic pivots</p> <ol style="list-style-type: none"> <li>1. <math>Z_{n,1}(\theta) = n (\hat{\theta}_n - \theta)^T I(\theta) (\hat{\theta}_n - \theta),</math></li> <li>2. <math>Z_{n,2}(\theta) = n (\hat{\theta}_n - \theta)^T I(\hat{\theta}_n) (\hat{\theta}_n - \theta),</math></li> <li>3. <math>Z_{n,3}(\theta) = 2 \left[ \mathcal{L}_n(\hat{\theta}_n) - \mathcal{L}_n(\theta) \right]</math> where <math>\mathcal{L}_n(\vartheta) = \sum_{i=1}^n \log p_\vartheta(x_i).</math></li> </ol> <p>They all asymptotically have the <math>\chi_p^2</math> distribution.</p>

<p>LEMMA</p> <p><i>State the lemma about asymptotic convergence of restricted MLE.</i> <i>(part-2)</i></p>	

	<p>Furthermore, let</p> $\mathcal{L}_n\left(\hat{\theta}_n\right)-\mathcal{L}_n\left(\hat{\theta}_n^0\right)=\sum_{i=1}^n\left[\log p_{\hat{\theta}_n}\left(X_i\right)-\log p_{\hat{\theta}_n^0}\right]$ <p>be the log-likelihood ratio for testing <math>H_0: R(\theta)=0</math>. Then under regularity conditions we have</p> $2\left[\mathcal{L}_n\left(\hat{\theta}_n\right)-\mathcal{L}_n\left(\hat{\theta}_n^0\right)\right] \stackrel{\mathcal{D}_\theta}{\underset{n \rightarrow \infty}{\rightarrow}} \chi_q^2$