

<p>THEOREM</p> <p><b>Alternative Formula for Expected Value</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Transformation of Random Variable</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Sample Mean and Variance</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Delta Method</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Multivariate Delta Method</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Geometric Series</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Gaussian Tail Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Markov's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Chebyshev's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Hoeffding's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>

<p>Let <math>X</math> be a random variable and define <math>Y = g(X)</math>, where <math>g</math> must be a monotonic function. Then <math>Y</math> has pdf</p> $p_Y(y) = p_X(g^{-1}(y)) \left  \frac{dg^{-1}(y)}{dy} \right $	<p>Let <math>X</math> be a non-negative continuous random variable. Then</p> $\mathbb{E}(X) = \int_0^\infty P(X > t) dt.$ <p>Analogously, if <math>X</math> is a discrete non-negative random variable, we have</p> $\mathbb{E}(X) = \sum_{k=1}^\infty P(X \geq k).$
<p>if <math>X \sim N(\mu, \sigma^2)</math> and <math>Y = g(X)</math> with <math>\sigma^2</math> small, we have</p> $Y \approx N(g(\mu), \sigma^2(g'(\mu))^2)$ <p>Proof: Start with Taylor Expansion of <math>g(X)</math> around <math>\mu</math>.</p>	<p>Given a random sample <math>X_1, \dots, X_n \sim N(\mu, \sigma^2)</math>, the following statements are true:</p> <ul style="list-style-type: none"> <li>• <math>\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)</math></li> <li>• <math>\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2</math></li> <li>• <math>\bar{X}</math> and <math>S^2</math> are independent</li> </ul>
<p>For <math>r \in (0, 1)</math>,</p> $a + ar + ar^2 + \dots = \frac{a}{1-r}.$ <p>A partial geometric series <math>a + ar + ar^2 + \dots + ar^{n-1}</math> sums up to <math>\frac{a(1-r^n)}{1-r}</math></p>	<p>Suppose that <math>Y_n = (Y_{n1}, \dots, Y_{nk})</math> is a sequence of random variables such that</p> $\sqrt{n}(Y_n - \mu) \rightsquigarrow N(0, \Sigma).$ <p>For a function <math>g: \mathbb{R}^k \rightarrow \mathbb{R}</math> and the gradient of <math>g</math> with respect to <math>y</math> be <math>\nabla g(y) = \left(\frac{\partial g}{\partial y_1}, \dots, \frac{\partial g}{\partial y_k}\right)^\top</math>. Denoting with <math>\nabla_\mu</math> the gradient evaluated at <math>y = \mu</math> where all elements are assumed to be non-zero, we have</p> $\sqrt{n}(g(Y_n) - g(\mu)) \rightsquigarrow N(0, \nabla_\mu^\top \Sigma \nabla_\mu).$
<p>If <math>X</math> is a non-negative random variable with existing expectation <math>\mathbb{E}(X)</math>, then</p> $P(X > \varepsilon) \leq \frac{\mathbb{E}(X)}{\varepsilon}$ <p>Proof: Trivially, the inequality <math>\varepsilon \mathbb{1}(X &gt; \varepsilon) \leq X</math> holds. Taking the expectation on both sides and rearranging yields Markov's inequality.</p>	<p>If <math>X \sim N(0, 1)</math>, then</p> $P( X  > \varepsilon) \leq \frac{2}{\varepsilon} e^{-\varepsilon^2/2}.$ <p>In general, for a sample <math>X_1, \dots, X_n</math> where <math>X_i \sim N(\mu, \sigma^2)</math>, we have</p> $P( \bar{X}_n - \mu  > \varepsilon) \leq \frac{2\sigma}{\varepsilon\sqrt{n}} \exp\left(-\frac{n\varepsilon^2}{2\sigma^2}\right)$
<p>Let <math>X_1, \dots, X_n</math> be independent observations such that <math>\mathbb{E}[X_i] = \mu</math> and <math>a \leq X_i \leq b \quad \forall i</math>. Then, for any <math>\varepsilon &gt; 0</math>, we have</p> $P( \bar{X}_n - \mu  > \varepsilon) \leq 2 \exp\left(\frac{-2n\varepsilon^2}{(b-a)^2}\right)$	<p>Let <math>X</math> be a random variable with mean <math>\mu</math> and variance <math>\sigma^2</math>. Then</p> $P( X - \mu  > \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}$ <p>Proof: Define <math>Z = (X - \mu)^2</math>. By Markov's Inequality, we have <math>P(Z &gt; \varepsilon^2) \leq \frac{\mathbb{E}(Z)}{\varepsilon^2}</math> which is equivalent to</p> $P( X - \mu  > \varepsilon) \leq \frac{\mathbb{E}[(X - \mu)^2]}{\varepsilon^2} = \frac{\sigma^2}{\varepsilon^2}$

<p>THEOREM</p> <p><b>Bernstein's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>McDiarmid's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Jensen's Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Cauchy-Schwartz Inequality</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Little <math>o</math></b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Big <math>O</math></b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Little <math>o_p</math></b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Big <math>o_p</math></b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Consistent Estimator</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Consistency Conditions</b></p> <p>INTERMEDIATE STATISTICS</p>

<p>Let <math>X_1, \dots, X_n</math> be independent random variables. Suppose that</p> $\sup_{x_1, \dots, x_n, x'_i}  g(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - g(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)  \leq c_i \quad (1)$ <p>for <math>i = 1, \dots, n</math>. Then</p> $P(g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)] \geq \varepsilon) \leq \exp - \frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}.$	<p>Let <math>X_1, \dots, X_n</math> be independent observations such that <math>\mathbb{E}[X_i] = 0</math>, <math> X_i  \leq M</math>, and <math>\mathbb{V}(X_i) \leq \sigma^2</math>. Then, for every <math>\varepsilon &gt; 0</math>, we have</p> $P( \bar{X}_n  \geq \varepsilon) \leq 2 \exp\left(-\frac{n\varepsilon^2}{2\sigma^2 + \frac{2}{3}M\varepsilon}\right).$
<p>Let <math>X</math> and <math>Y</math> be two random variables with finite variance. Then</p> $\mathbb{E} XY  \leq \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}.$	<p>Let <math>X</math> be a random variable and <math>g</math> a convex function. Then <math>\mathbb{E}g(X) \geq g(\mathbb{E}(X))</math>. On the other hand, if <math>g</math> is concave, we have <math>\mathbb{E}g(X) \leq g(\mathbb{E}(X))</math>. Example: From Jensen's inequality, it follows that <math>\mathbb{E}(X^2) \geq \mathbb{E}(X)^2</math>, since <math>g(x) = x^2</math> is convex.</p>
<p><math>a_n = O(b_n)</math> if for large <math>n &gt; n_0</math>, there exists some constant <math>C &gt; 0</math> such that <math>a_n &lt; Cb_n</math></p>	<p><math>a_n = o(b_n)</math> means that <math>\forall C</math> and <math>n &gt; n_0</math>,</p> $a_n < Cb_n$ <p>(<math>a_n</math> is bounded from above by <math>b_n</math>)</p>
<p><math>Y_n</math> is <math>O_p(1)</math> if, for any <math>\varepsilon &gt; 0</math>, there exists some finite constant <math>C &gt; 0</math> such that</p> $P( Y_n  > C) \leq \varepsilon$ <p>for all <math>n</math> (stochastically bounded from above).  <math>Y_n = O_p(a_n)</math> means that <math>\frac{Y_n}{a_n} = O_p(1)</math>.</p>	<p><math>Y_n</math> is <math>o_p(1)</math> if, for every <math>\varepsilon &gt; 0</math>, we have</p> $P( Y_n  > \varepsilon) \rightarrow 0$ <p>or equivalently</p> $P( Y_n  \leq \varepsilon) \rightarrow 1.$ <p><math>Y_n = o_p(a_n)</math> means that <math>\frac{Y_n}{a_n} = o_p(1)</math>.</p>
<p>Let <math>\theta_n</math> be a sequence of estimators of parameter <math>\theta</math> satisfying</p> $\lim_{n \rightarrow \infty} \mathbb{V}(\hat{\theta}_n) = 0$ <p>and Then <math>\hat{\theta}_n</math> is a consistent sequence of estimators of <math>\theta</math>.</p>	<p>A sequence of estimators <math>\theta_n</math> is consistent of the parameter <math>\theta</math> if, for every <math>\epsilon &gt; 0</math> and every <math>\theta \in \Theta</math> we have</p> $\lim_{n \rightarrow \infty} P_\theta( \theta_n - \theta  < \epsilon) = 1$ <p>or equivalently</p> $\lim_{n \rightarrow \infty} P_\theta( \theta_n - \theta  \geq \epsilon) = 0,$ <p>i.e. <math>\theta_n</math> converges in probability to <math>\theta</math>.</p>

<p>THEOREM</p> <p><b>Consistency of MLE</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Shattering</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Shatter Coefficient</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>VC Dimension</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Uniform Distribution</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Multinomial Distribution</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Statistic</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Almost Sure Convergence</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Convergence in Probability</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Convergence in Quadratic Mean</b></p> <p>INTERMEDIATE STATISTICS</p>

<p>Let <math>\mathcal{A}</math> be a class of sets and <math>F</math> be a finite set <math>\{x_1, \dots, x_k\}</math>. Let <math>G</math> be some subset of <math>F</math>. <math>\mathcal{A}</math> picks out <math>G</math> if <math>A \cap F = G</math> for some <math>A \in \mathcal{A}</math>. The set <math>F</math> is shattered if <math>s(\mathcal{A}, F) = 2^k</math>, i.e. if all subsets can be picked out by <math>\mathcal{A}</math>.</p>	<p>Let <math>\hat{\theta}</math> be the MLE of <math>\theta</math> and let <math>\tau(\theta)</math> be a continuous function of <math>\theta</math>. Under regularity conditions, for every <math>\epsilon &gt; 0</math> and <math>\theta \in \Theta</math>,</p> $\lim_{n \rightarrow \infty} P_{\theta} \left( \left  \tau(\hat{\theta}) - \tau(\theta) \right  \geq \epsilon \right) = 0,$ <p>i.e. <math>\tau(\hat{\theta})</math> is a consistent estimator of <math>\tau(\theta)</math>. The conditions are a) an iid random sample, b) identifiability of the parameter, c) common support and differentiability of the density and d) a parameter space which contains an open set of which the true parameter is an interior point.</p>
<p>The Vapnik-Chervonenkis (VC) Dimension is defined as</p> $d = d(\mathcal{A}) = \text{largest } k \text{ such that } s_k(\mathcal{A}) = 2^k.$ <p>This means that <math>d</math> is the size of the largest set that can be shattered.</p>	<p>The shatter coefficient is defined as</p> $s_k(\mathcal{A}) = \sup_{F \in \mathcal{F}_k} s(\mathcal{A}, F),$ <p>where <math>\mathcal{F}_k</math> denotes all finite sets with <math>k</math> elements. Fact: <math>s_k(\mathcal{A}) \leq 2^k</math>.</p>
<p>Multivariate version of Binomial. Draw all from urn with balls colored in <math>k</math> different colors. <math>p = (p_1, \dots, p_k)</math> where <math>\sum_j p_j = 1</math> and <math>p_j</math> is probability of drawing color <math>j</math>. Draw <math>n</math> balls from the urn with replacement and let <math>X = (X_1, \dots, X_n)</math> be the count of the number of balls of each color. Then <math>X</math> has a Multinomial distribution with pdf</p> $p(x) = \binom{n}{x_1 \dots x_k} p_1^{x_1} \dots p_k^{x_k}.$	<p>A continuous random variable <math>X</math> has a Uniform(<math>a, b</math>) distribution if its pdf is</p> $\begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$ <p>and CDF</p> $\begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 1 & \text{for } x \geq b \end{cases}.$ <p>The mean of <math>X</math> is <math>\frac{1}{2}(a+b)</math> and the variance <math>\frac{1}{12}(b-a)^2</math>.</p>
<p><math>X_n</math> converges almost surely to <math>X</math>, written <math>X_n \xrightarrow{a.s.} X</math>, if, for every <math>\varepsilon &gt; 0</math>,</p> $P \left( \lim_{n \rightarrow \infty}  X_n - X  < \varepsilon \right) = 1.$ <p>Almost sure convergence of <math>X_n</math> to <math>X</math> is equivalent to</p> $\forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P \left( \sup_{m \geq n}  X_m - X  \leq \varepsilon \right) = 1.$	<p>A statistic <math>T</math> is any function of the data <math>X_1, \dots, X_n</math>, i.e. <math>T = g(X_1, \dots, X_n)</math>.</p>
<p>A sequence of random variables <math>X_n</math> converges to <math>X</math> in quadratic mean (<math>L_2</math> convergence) if</p> $\mathbb{E}(X_n - X)^2 \rightarrow 0$ <p>as <math>n \rightarrow \infty</math>. We write <math>X_n \xrightarrow{q.m.} X</math>.</p>	<p><math>X_n</math> converges to <math>X</math> in probability (<math>X_n \xrightarrow{P} X</math>), if</p> $\forall \varepsilon > 0 : P( X_n - X  > \varepsilon) \rightarrow 0$ <p>as <math>n \rightarrow \infty</math> (notice that we thus have <math>X_n - X = o_P(1)</math>).</p>

<p>DEFINITION</p> <p><b>Convergence in Distribution</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Convergence Relationships</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Continuous Mapping Theorem</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Slutsky's Theorem</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>The Weak Law of Large Numbers (WLLN)</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>The Strong Law of Large Numbers (SLLN)</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>The Central Limit Theorem (CLT)</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Estimate <math>\sigma</math> in CLT</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Multivariate Central Limit Theorem</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Loss Function</b></p> <p>INTERMEDIATE STATISTICS</p>

<p>Between the different convergence definitions, the following relationships hold:</p> <ul style="list-style-type: none"> <li>• <math>X_n \xrightarrow{a.s.} X</math> implies that <math>X_n \xrightarrow{P} X</math>.</li> <li>• <math>X_n \xrightarrow{q.m.} X</math> implies that <math>X_n \xrightarrow{P} X</math>.</li> <li>• <math>X_n \xrightarrow{P} X</math> implies that <math>X_n \rightsquigarrow x</math>.</li> <li>• If <math>X_n \rightsquigarrow X</math> and if <math>X</math> has a point mass distribution, i.e. <math>P(X = c) = 1</math> for some <math>c</math>, then <math>X_n \xrightarrow{P} X</math>.</li> </ul>	<p><math>X_n</math> converges to <math>X</math> in distribution if</p> $\lim_{n \rightarrow \infty} F_n(t) = F(t)$ <p>at all <math>t</math> for which <math>F</math> is continuous. We write <math>X_n \rightsquigarrow X</math>.</p>
<p>Let <math>X_n</math> and <math>Y_n</math> be sequences of random variables and let <math>X</math> be a simple random variable and <math>c</math> a constant. We have</p> <ul style="list-style-type: none"> <li>• If <math>X_n \rightsquigarrow X</math> and <math>Y_n \rightsquigarrow c</math>, then <math>X_n + Y_n \rightsquigarrow X + c</math>.</li> <li>• If <math>X_n \rightsquigarrow X</math> and <math>Y_n \rightsquigarrow c</math>, then <math>X_n Y_n \rightsquigarrow cX</math>.</li> </ul> <p>In general, <math>X_n \rightsquigarrow X</math> and <math>Y_n \rightsquigarrow Y</math> does not imply that <math>X_n + Y_n \rightsquigarrow X + Y</math>.</p>	<p>Let <math>X_n</math> and <math>Y_n</math> be sequences of random variables. Also, let <math>X</math> and <math>Y</math> be simple random variables. For a continuous function <math>g</math>, we have</p> <ol style="list-style-type: none"> <li>1. If <math>X_n \xrightarrow{P} X</math>, then <math>g(X_n) \xrightarrow{P} g(X)</math>.</li> <li>2. If <math>X_n \rightsquigarrow X</math>, then <math>g(X_n) \rightsquigarrow g(X)</math>.</li> </ol>
<p>Let <math>X_1, \dots, X_n</math> be iid with mean <math>\mu</math>. Then we have <math>\bar{X}_n \xrightarrow{a.s.} \mu</math>.</p>	<p>Given a random sample <math>X_1, \dots, X_n</math> iid, the sample mean <math>\bar{X}_n</math> converges in probability to <math>\mu</math>. Therefore, <math>\bar{X}_n - \mu = o_p(1)</math>.</p>
<p>Let <math>X_1, \dots, X_n</math> be an iid sample where <math>\mathbb{E}(X_i) = \mu</math> and <math>\mathbb{V}(X_i) = \sigma^2</math>. Let <math>\bar{X}_n = n^{-1} \sum_{i=1}^n X_i</math>. Denote the sample variance with <math>S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2</math>. Then</p> $T_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \rightsquigarrow N(0, 1).$ <p>Proof: We have that <math>T_n = Z_n W_n</math>, where <math>Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow N(0, 1)</math> and <math>W_n = \frac{\sigma}{S_n} \xrightarrow{P} 1</math>. The result then follows from Slutsky's Theorem.</p>	<p>Let <math>X_1, \dots, X_n</math> be an iid sample where <math>\mathbb{E}(X_i) = \mu</math> and <math>\mathbb{V}(X_i) = \sigma^2</math>. Let <math>\bar{X}_n = n^{-1} \sum_{i=1}^n X_i</math>. Then</p> $Z_n \equiv \frac{\bar{X}_n - \mu}{\sqrt{\mathbb{V}(\bar{X}_n)}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow Z,$ <p>where <math>Z \sim N(0, 1)</math>.</p>
<p>A loss function <math>L(\theta, \hat{\theta}) : \Theta^2 \rightarrow [0, \infty)</math> measures the cost associated with the value of an estimator <math>\hat{\theta}</math> not being equal to the true parameter <math>\theta</math>. Common loss functions are</p> <ol style="list-style-type: none"> <li>1. Squared Loss</li> <li>2. Absolute Loss</li> <li>3. Zero-One Loss</li> </ol>	<p>Let <math>X_1, \dots, X_n</math> be a sample of iid random vectors where <math>X_i = (X_{1i}, \dots, X_{ki})^\top</math> with mean <math>\mu = (\mu_1, \dots, \mu_k)^\top</math> and covariance matrix <math>\Sigma</math>. Let <math>\bar{X} = (\bar{X}_1, \dots, \bar{X}_k)^\top</math> where <math>\bar{X}_j = n^{-1} \sum_{i=1}^n X_{ji}</math>. Then,</p> $\sqrt{n}(\bar{X} - \mu) \rightsquigarrow N(0, \Sigma)$



<p>DEFINITION</p> <p><b>Risk of an Estimator</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Minimax Risk</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Bayes Risk</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Posterior Risk</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Bayes Risk (in terms of posterior risk)</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Common Bayes Estimators</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Minimax of Bayes Estimator</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Bayes Estimator with Constant Risk</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>p-value</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Likelihood Function</b></p> <p>INTERMEDIATE STATISTICS</p>

<p>The minimax risk is defined as</p> $R_n = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}).$ <p>It is the risk of the estimator whose maximal risk is lowest among all competing estimators <math>\hat{\theta}</math>. It follows that an estimator <math>\hat{\theta}</math> is minimax if</p> $\sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}).$	<p>The risk of an estimator <math>\hat{\theta}</math> is the expected value of the associated loss function, where the expectation is taken over all sample variables:</p> $R(\theta, \hat{\theta}) = \mathbb{E}(L(\theta, \hat{\theta})) = \int L(\theta, \hat{\theta}(x^n)) p(x^n; \theta) dx^n$ <p>Under squared error loss, the risk is equal to the mean squared error.</p>
<p>The posterior risk of an estimator <math>\hat{\theta}(x^n)</math> is</p> $r(\hat{\theta} x^n) = \int L(\theta, \hat{\theta}(x^n)) \pi(\hat{\theta} x^n) d\theta$	<p>The Bayes risk of an estimator <math>\hat{\theta}</math> with prior distribution <math>\pi</math> is</p> $B_{\pi}(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta.$ <p>Notice that the remaining uncertainty of the risk lies in different values for <math>\theta</math>: The risk already has dealt with uncertainty in the data, as</p> $R(\theta, \hat{\theta}) = \mathbb{E}_{\theta}(L(\theta, \hat{\theta}))$ <p>Estimators which minimize the Bayes risk are called <i>Bayes estimators</i>.</p>
<ul style="list-style-type: none"> <li>• Under squared error loss, the Bayes estimator is the posterior mean <math>\mathbb{E}(\theta X = x^n)</math>.</li> <li>• Under absolute loss, the Bayes estimator is the posterior median <math>F_{\theta X}^{-1}(\frac{1}{2})</math>.</li> <li>• Under 0-1-loss, the Bayes estimator is the posterior mode of <math>\pi(\theta x^n)</math>.</li> </ul>	<p>The Bayes risk <math>B_{\pi}(\hat{\theta})</math> can also be expressed as</p> $B_{\pi}(\hat{\theta}) = \int r(\hat{\theta} x^n) m(x^n) dx^n,$ <p>where <math>m(x^n)</math> is the marginal distribution of the data (sometimes called the <i>evidence</i>). An estimator <math>\hat{\theta}</math> which minimizes the posterior risk is therefore a Bayes estimator since the integrand in <math>B_{\pi}(\hat{\theta})</math> will be minimal at all <math>x</math>.</p>
<p>Let <math>\hat{\theta}</math> be the Bayes estimator under some prior distribution <math>\pi</math>. If the risk is constant (with respect to <math>\theta</math>) then this estimator is minimax. Proof: We have that <math>R(\theta, \hat{\theta}) = c</math>, where <math>c</math> is some constant. It follows that <math>B_{\pi}(\hat{\theta}) = \int r(\hat{\theta} x^n) m(x^n) dx^n = c</math> as well and hence <math>R(\theta, \hat{\theta}) \leq B_{\pi}(\hat{\theta})</math> holds for all <math>\theta</math>. By the “Minimax of Bayes Estimator” Theorem, this implies that the estimator is minimax.</p>	<p>Let <math>\hat{\theta}</math> be the Bayes estimator under some prior <math>\pi</math>. If its risk is always smaller than the Bayes risk, i.e. if</p> $R(\theta, \hat{\theta}) \leq B_{\pi}(\hat{\theta}) \quad \forall \theta,$ <p>then <math>\hat{\theta}</math> is the minimax estimator and <math>\pi</math> is called a least favorable prior. Proof: By contradiction. Assume that <math>\hat{\theta}</math> was not minimax. Then show that this would imply that the estimator did not minimize the Bayes risk in the first place (Hint: The average of a function is always less than or equal to its maximum).</p>
<p>Let <math>X^n = (X_1, \dots, X_n)</math> have joint density <math>p(x^n; \theta)</math> where <math>\theta \in \Theta</math>. The likelihood function <math>L: \Theta \rightarrow [0, \infty)</math> is the joint density regarded as a function of parameter <math>\theta</math>, i.e.</p> $L(\theta) = p(x^n; \theta).$ <p>The likelihood is not a pdf and defined only up to a constant of proportionality.</p>	<p>Suppose we have a test of the form: reject when <math>W(X^n) &gt; c</math>. Then the p-value when <math>X^n = x^n</math> is</p> $p(x^n) = \sup_{\theta \in \Theta_0} P_{\theta}(W(X^n) \geq W(x^n))$

<p>THEOREM</p> <p><b>Equivariance Property of MLE</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p><b>Mean Squared Error (MSE)</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Rao-Blackwell Theorem</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Sufficiency</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Factorization Theorem</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Minimal Sufficiency</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p><b>Find Minimal Sufficient Statistic</b></p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p><b>Empirical CDF</b></p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p><b>Kernel Density Estimator</b></p> <p>INTERMEDIATE STATISTICS</p>	

<p>The mean squared error (MSE) is</p> $MSE = \mathbb{E}_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] = \int (\hat{\theta}(x^n) - \theta)^2 p(x^n; \theta) dx^n.$ <p>The MSE can be decomposed into variance and bias squared, i.e.</p> $MSE = \mathbb{V}_{\theta}(\hat{\theta}) + Bias^2,$ <p>where <math>Bias = \mathbb{E}_{\theta}(\hat{\theta}) - \theta</math>.</p>	<p>Let <math>\hat{\theta}</math> be the MLE. If <math>\eta = g(\theta)</math>, then the MLE of <math>\eta</math> is <math>\hat{\eta} = g(\hat{\theta})</math>.  Proof: Suppose <math>g</math> is invertible so <math>\eta = g(\theta)</math> and <math>\theta = g^{-1}(\eta)</math>. Define <math>L^*(\eta) = L(\theta)</math> where <math>\theta = g^{-1}(\eta)</math>. Hence,</p> $L^*(\hat{\eta}) = L(\hat{\theta}) \geq L(\theta) = L^*(\eta)$ <p>and thus <math>\hat{\eta}</math> maximizes <math>L^*(\eta)</math>. For non-invertible functions, this still holds if we define</p> $L^*(\eta) = \sup_{\theta: \tau(\theta) = \eta} L(\theta).$
<p>Suppose that we have a random sample <math>X_1, \dots, X_n \sim p(x; \theta)</math>. An estimator <math>T</math> is sufficient for <math>\theta</math> if the conditional distribution of <math>X_1, \dots, X_n   T</math> does not depend on <math>\theta</math>. Thus</p> $p(x_1, \dots, x_n   t, \theta) = p(x_1, \dots, x_n   t).$	<p>Let <math>W</math> be an unbiased estimator of <math>\tau(\theta)</math> and let <math>T</math> be a sufficient statistic. Define <math>W' = \mathbb{E}(W T)</math>. Then <math>W'</math> is unbiased with variance <math>\mathbb{V}_{\theta}(W') \leq \mathbb{V}_{\theta}(W) \quad \forall \theta</math>.</p>
<p><math>T</math> is a minimal sufficient statistic for <math>\theta</math> if it is sufficient and if it is a function of any other sufficient statistic <math>U</math>, i.e. <math>T = g(U)</math> for some function <math>g</math>.</p>	<p>An estimator <math>T(X^n)</math> is sufficient for <math>\theta</math> if the joint pdf of <math>X^n</math> can be factored as</p> $p(x^n   \theta) = h(x^n) g(T(x^n); \theta).$
<p>The empirical cumulative distribution function (ECDF) puts mass <math>\frac{1}{n}</math> at each data point. It is defined as</p> $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x).$ <p>Notice that <math>\hat{F}_n(x) \sim \text{Bernoulli}(F_X(x))</math>. We also have that</p> $P \left( \sup_x  \hat{F}_n(x) - F(x)  > \varepsilon \right) \leq 2e^{-2n\varepsilon^2},$ <p>that is <math>\sup_x  \hat{F}_n(x) - F(x)  \xrightarrow{P} 0</math>.</p>	<p>An estimator <math>T</math> is minimal sufficient if and only if it has the following property:</p> $T(y^n) = T(x^n) \Leftrightarrow \frac{p(y^n; \theta)}{p(x^n; \theta)} \text{ does not depend on } \theta$
	<p>The kernel density estimator is a non-parametric estimator of the density function. It is defined as</p> $\hat{p}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K \left( \frac{x - X_i}{h} \right),$ <p>where <math>h &gt; 0</math> is the bandwidth and <math>K</math>, the kernel, is a symmetric density with mean zero.</p>