

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

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

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

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

<p>Let \mathcal{A} be a class of sets and F be a finite set $\{x_1, \dots, x_k\}$. Let G be some subset of F. \mathcal{A} picks out G if $A \cap F = G$ for some $A \in \mathcal{A}$. The set F is shattered if $s(\mathcal{A}, F) = 2^k$, i.e. if all subsets can be picked out by \mathcal{A}.</p>	<p>Let $\hat{\theta}$ be the MLE of θ and let $\tau(\theta)$ be a continuous function of θ. Under regularity conditions, for every $\epsilon > 0$ and $\theta \in \Theta$,</p> $\lim_{n \rightarrow \infty} P_{\theta} \left(\left \tau(\hat{\theta}) - \tau(\theta) \right \geq \epsilon \right) = 0,$ <p>i.e. $\tau(\hat{\theta})$ is a consistent estimator of $\tau(\theta)$. 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 d 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 \mathcal{F}_k denotes all finite sets with k elements. Fact: $s_k(\mathcal{A}) \leq 2^k$.</p>
<p>X_n converges almost surely to X, written $X_n \xrightarrow{a.s.} X$, if, for every $\varepsilon > 0$,</p> $P \left(\lim_{n \rightarrow \infty} X_n - X < \varepsilon \right) = 1.$ <p>Almost sure convergence of X_n to X 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 T is any function of the data X_1, \dots, X_n, i.e. $T = g(X_1, \dots, X_n)$.</p>
<p>A sequence of random variables X_n converges to X in quadratic mean (L_2 convergence) if</p> $\mathbb{E}(X_n - X)^2 \rightarrow 0$ <p>as $n \rightarrow \infty$. We write $X_n \xrightarrow{q.m.} X$.</p>	<p>X_n converges to X in probability ($X_n \xrightarrow{P} X$), if</p> $\forall \varepsilon > 0 : P(X_n - X > \varepsilon) \rightarrow 0$ <p>as $n \rightarrow \infty$ (notice that we thus have $X_n - X = o_P(1)$).</p>
<p>Between the different convergence definitions, the following relationships hold:</p> <ul style="list-style-type: none"> • $X_n \xrightarrow{a.s.} X$ implies that $X_n \xrightarrow{P} X$. • $X_n \xrightarrow{q.m.} X$ implies that $X_n \xrightarrow{P} X$. • $X_n \xrightarrow{P} X$ implies that $X_n \rightsquigarrow X$. • If $X_n \rightsquigarrow X$ and if X has a point mass distribution, i.e. $P(X = c) = 1$ for some c, then $X_n \xrightarrow{P} X$. 	<p>X_n converges to X in distribution if</p> $\lim_{n \rightarrow \infty} F_n(t) = F(t)$ <p>at all t for which F is continuous. We write $X_n \rightsquigarrow X$.</p>

<p>THEOREM</p> <p>Continuous Mapping Theorem</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>Slutsky's Theorem</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>The Weak Law of Large Numbers (WLLN)</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>The Strong Law of Large Numbers (SLLN)</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>The Central Limit Theorem (CLT)</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>Estimate σ in CLT</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Multivariate Central Limit Theorem</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Loss Function</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Risk of an Estimator</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Minimax Risk</p> <p>INTERMEDIATE STATISTICS</p>

<p>Let X_n and Y_n be sequences of random variables and let X be a simple random variable and c a constant.</p> <p>We have</p> <ul style="list-style-type: none"> • If $X_n \rightsquigarrow X$ and $Y_n \rightsquigarrow c$, then $X_n + Y_n \rightsquigarrow X + c$. • If $X_n \rightsquigarrow X$ and $Y_n \rightsquigarrow c$, then $X_n Y_n \rightsquigarrow cX$. <p>In general, $X_n \rightsquigarrow X$ and $Y_n \rightsquigarrow Y$ does not imply that $X_n + Y_n \rightsquigarrow X + Y$.</p>	<p>Let X_n and Y_n be sequences of random variables. Also, let X and Y be simple random variables. For a continuous function g, we have</p> <ol style="list-style-type: none"> 1. If $X_n \xrightarrow{P} X$, then $g(X_n) \xrightarrow{P} g(X)$. 2. If $X_n \rightsquigarrow X$, then $g(X_n) \rightsquigarrow g(X)$.
<p>Let X_1, \dots, X_n be iid with mean μ. Then we have $\bar{X}_n \xrightarrow{a.s.} \mu$.</p>	<p>Given a random sample X_1, \dots, X_n iid, the sample mean \bar{X}_n converges in probability to μ. Therefore, $\bar{X}_n - \mu = o_p(1)$.</p>
<p>Let X_1, \dots, X_n be an iid sample where $\mathbb{E}(X_i) = \mu$ and $\mathbb{V}(X_i) = \sigma^2$. Let $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$. Denote the sample variance with $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. Then</p> $T_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \rightsquigarrow N(0, 1).$ <p>Proof: We have that $T_n = Z_n W_n$, where $Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow N(0, 1)$ and $W_n = \frac{\sigma}{S_n} \xrightarrow{P} 1$. The result then follows from Slutsky's Theorem.</p>	<p>Let X_1, \dots, X_n be an iid sample where $\mathbb{E}(X_i) = \mu$ and $\mathbb{V}(X_i) = \sigma^2$. Let $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$. 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 $Z \sim N(0, 1)$.</p>
<p>A loss function $L(\theta, \hat{\theta}) : \Theta^2 \rightarrow [0, \infty)$ measures the cost associated with the value of an estimator $\hat{\theta}$ not being equal to the true parameter θ. Common loss functions are</p> <ol style="list-style-type: none"> 1. Squared Loss 2. Absolute Loss 3. Zero-One Loss 	<p>Let X_1, \dots, X_n be a sample of iid random vectors where $X_i = (X_{1i}, \dots, X_{ki})^\top$ with mean $\mu = (\mu_1, \dots, \mu_k)^\top$ and covariance matrix Σ. Let $\bar{X} = (\bar{X}_1, \dots, \bar{X}_k)^\top$ where $\bar{X}_j = n^{-1} \sum_{i=1}^n X_{ji}$. Then,</p> $\sqrt{n}(\bar{X} - \mu) \rightsquigarrow N(0, \Sigma)$
<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 $\hat{\theta}$. It follows that an estimator $\hat{\theta}$ 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 $\hat{\theta}$ 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>DEFINITION</p> <p>Bayes Risk</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Posterior Risk</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Bayes Risk (in terms of posterior risk)</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>Common Bayes Estimators</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Minimax of Bayes Estimator</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>Bayes Estimator with Constant Risk</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>p-value</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Likelihood Function</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Equivariance Property of MLE</p> <p>INTERMEDIATE STATISTICS</p>	<p>THEOREM</p> <p>Mean Squared Error (MSE)</p> <p>INTERMEDIATE STATISTICS</p>

<p>The posterior risk of an estimator $\hat{\theta}(x^n)$ 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 $\hat{\theta}$ with prior distribution π 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 θ: 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"> • Under squared error loss, the Bayes estimator is the posterior mean $\mathbb{E}(\theta X = x^n)$. • Under absolute loss, the Bayes estimator is the posterior median $F_{\theta X}^{-1}(\frac{1}{2})$. • Under 0-1-loss, the Bayes estimator is the posterior mode of $\pi(\theta x^n)$. 	<p>The Bayes risk $B_{\pi}(\hat{\theta})$ can also be expressed as</p> $B_{\pi}(\hat{\theta}) = \int r(\hat{\theta} x^n)m(x^n)dx^n,$ <p>where $m(x^n)$ is the marginal distribution of the data (sometimes called the <i>evidence</i>). An estimator $\hat{\theta}$ which minimizes the posterior risk is therefore a Bayes estimator since the integrand in $B_{\pi}(\hat{\theta})$ will be minimal at all x.</p>
<p>Let $\hat{\theta}$ be the Bayes estimator under some prior distribution π. If the risk is constant (with respect to θ) then this estimator is minimax. Proof: We have that $R(\theta, \hat{\theta}) = c$, where c is some constant. It follows that $B_{\pi}(\hat{\theta}) = \int r(\hat{\theta} x^n)m(x^n)dx^n = c$ as well and hence $R(\theta, \hat{\theta}) \leq B_{\pi}(\hat{\theta})$ holds for all θ. By the “Minimax of Bayes Estimator” Theorem, this implies that the estimator is minimax.</p>	<p>Let $\hat{\theta}$ be the Bayes estimator under some prior π. 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 $\hat{\theta}$ is the minimax estimator and π is called a least favorable prior. Proof: By contradiction. Assume that $\hat{\theta}$ 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 $X^n = (X_1, \dots, X_n)$ have joint density $p(x^n; \theta)$ where $\theta \in \Theta$. The likelihood function $L: \Theta \rightarrow [0, \infty)$ is the joint density regarded as a function of parameter θ, 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 $W(X^n) > c$. Then the p-value when $X^n = x^n$ is</p> $p(x^n) = \sup_{\theta \in \Theta_0} P_{\theta}(W(X^n) \geq W(x^n))$
<p>The mean squared error (MSE) is</p> $MSE = \mathbb{E}_{\theta}[(\hat{\theta} - \theta)^2] = \int (\hat{\theta}(x^n) - \theta)^2 p(x^n; \theta) dx^n.$ <p>The MSE can be decomposed into variance and bias squared, i.e.</p> $MSE = \mathbb{V}_{\theta}(\hat{\theta}) + Bias^2,$ <p>where $Bias = \mathbb{E}_{\theta}(\hat{\theta}) - \theta$.</p>	<p>Let $\hat{\theta}$ be the MLE. If $\eta = g(\theta)$, then the MLE of η is $\hat{\eta} = g(\hat{\theta})$. Proof: Suppose g is invertible so $\eta = g(\theta)$ and $\theta = g^{-1}(\eta)$. Define $L^*(\eta) = L(\theta)$ where $\theta = g^{-1}(\eta)$. Hence,</p> $L^*(\hat{\eta}) = L(\hat{\theta}) \geq L(\theta) = L^*(\eta)$ <p>and thus $\hat{\eta}$ maximizes $L^*(\eta)$. For non-invertible functions, this still holds if we define</p> $L^*(\eta) = \sup_{\theta: \tau(\theta) = \eta} L(\theta).$

<p>THEOREM</p> <p>Rao-Blackwell Theorem</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Sufficiency</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Factorization Theorem</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Minimal Sufficiency</p> <p>INTERMEDIATE STATISTICS</p>
<p>THEOREM</p> <p>Find Minimal Sufficient Statistic</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Empirical CDF</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Kernel Density Estimator</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Uniform Distribution</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Normal Distribution</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Multivariate Normal Distribution</p> <p>INTERMEDIATE STATISTICS</p>

<p>Suppose that we have a random sample $X_1, \dots, X_n \sim p(x; \theta)$. An estimator T is sufficient for θ if the conditional distribution of $X_1, \dots, X_n T$ does not depend on θ. Thus</p> $p(x_1, \dots, x_n t, \theta) = p(x_1, \dots, x_n t).$	<p>Let W be an unbiased estimator of $\tau(\theta)$ and let T be a sufficient statistic. Define $W' = \mathbb{E}(W T)$. Then W' is unbiased with variance $\mathbb{V}_\theta(W') \leq \mathbb{V}_\theta(W) \quad \forall \theta$.</p>
<p>T is a minimal sufficient statistic for θ if it is sufficient and if it is a function of any other sufficient statistic U, i.e. $T = g(U)$ for some function g.</p>	<p>An estimator $T(X^n)$ is sufficient for θ if the joint pdf of X^n 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 $\frac{1}{n}$ 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 $\hat{F}_n(x) \sim \text{Bernoulli}(F_X(x))$. 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 $\sup_x \hat{F}_n(x) - F(x) \xrightarrow{P} 0$.</p>	<p>An estimator T 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>A continuous random variable X has a Uniform(a, b) 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 X is $\frac{1}{2}(a+b)$ and the variance $\frac{1}{12}(b-a)^2$.</p>	<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 $h > 0$ is the bandwidth and K, the kernel, is a symmetric density with mean zero.</p>
<ul style="list-style-type: none"> Let $X \in \mathbb{R}^d$. Then $X \sim N(\mu, \Sigma)$ if $f_X(x) = \frac{1}{(2\pi)^{d/2} \Sigma ^{1/2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}$ <ul style="list-style-type: none"> $M_X(t) = \exp(\mu^\top t + \frac{t^\top \Sigma t}{2})$ $\mathbb{E}[X] = \mu, \text{cov}[X] = \Sigma$ 	<ul style="list-style-type: none"> $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/(2\sigma^2)}$ $-\infty < x < \infty, -\infty < \mu < \infty, \sigma^2 > 0$ $M_X(t) = e^{\mu t + \sigma^2 t^2/2}$ $\mathbb{E}[X] = \mu, \mathbb{V}[X] = \sigma^2$

<p>THEOREM</p> <p>Multivariate Normal Transformations</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Multinomial Distribution</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Binomial Distribution</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Bernoulli Distribution</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Geometric Distribution</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Poisson Distribution</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Exponential Distribution</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Beta Distribution</p> <p>INTERMEDIATE STATISTICS</p>
<p>DEFINITION</p> <p>Gamma Distribution</p> <p>INTERMEDIATE STATISTICS</p>	<p>DEFINITION</p> <p>Cauchy Distribution</p> <p>INTERMEDIATE STATISTICS</p>

<p>Multivariate version of Binomial. Draw all from urn with balls colored in k different colors.</p> <p>$p = (p_1, \dots, p_k)$ where $\sum_j p_j = 1$ and p_j is probability of drawing color j. Draw n balls from the urn with replacement and let $X = (X_1, \dots, X_n)$ be the count of the number of balls of each color. Then X 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>Assume that $X \in \mathbb{R}^d$ and that $X \sim N(\mu, \Sigma)$. Then the following statements are true:</p> <ul style="list-style-type: none"> • If X is multiplied with a scalar c, we have $cX \sim N(c\mu, c^2\Sigma)$ • If A is a $p \times n$ matrix and b is a $p \times 1$ column vector, then $AX + b \sim N(A\mu + b, A\Sigma A^T)$. • $(X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi_d^2$
<ul style="list-style-type: none"> • $P(X = x) = p^x (1 - p)^{1-x}$ • $x \in \{0, 1\}, 0 \leq p \leq 1$ • $M_X(t) = (1 - p) + pe^t$ • $\mathbb{E}[X] = p, \mathbb{V}[X] = p(1 - p)$ 	<p>Sequence of Bernoulli Trials</p> <ul style="list-style-type: none"> • $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$ • $k \in \{0, \dots, n\}, 0 \leq p \leq 1$ • $M_X(t) = (1 - p + pe^t)^n$ • $\mathbb{E}[X] = np, \mathbb{V}[X] = np(1 - p)$
<ul style="list-style-type: none"> • $P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$ • $x \in \{0, 1, \dots\}, 0 \leq \lambda < \infty$ • $M_X(t) = e^{\lambda(e^t - 1)}$ • $\mathbb{E}[X] = \lambda, \mathbb{V}[X] = \lambda$ 	<ul style="list-style-type: none"> • $P(X = x) = p(1 - p)^{x-1}$ • $x \in \{0, 1, \dots\}, 0 \leq p \leq 1$ • $M_X(t) = \frac{pe^t}{1 - (1-p)e^t}, t < -\log(1 - p)$ • $\mathbb{E}[X] = \frac{1}{p}, \mathbb{V}[X] = \frac{1-p}{p^2}$ • Only existing discrete distribution with memoryless property: $P(X > s X > t) = P(X > s - t)$
<ul style="list-style-type: none"> • $f_X(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1 - x)^{\beta-1}$ • $0 \leq x \leq 1, \alpha > 0, \beta > 0$ • $M_X(t) = 1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r} \right) \frac{t^k}{k!}$ • $\mathbb{E}[X] = \frac{\alpha}{\alpha+\beta}, \mathbb{V}[X] = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ • Recall that the beta function can be defined in terms of the Gamma function: $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ 	<ul style="list-style-type: none"> • $f_X(x) = \frac{1}{\beta} e^{(-x/\beta)}$ • $0 \leq x < \infty, \beta > 0$ • $M_X(t) = \frac{1}{1-\beta t}, t < \frac{1}{\beta}$ • $\mathbb{E}[X] = \beta, \mathbb{V}[X] = \beta^2$ • Only continuous distribution with memoryless property: $P(X > s X > t) = P(X > s - t)$. • Special case of Gamma distribution with $\alpha = 1$.
<ul style="list-style-type: none"> • $f_X(x) = \frac{1}{\pi\sigma} \frac{1}{1 + \left(\frac{x-\mu}{\sigma}\right)^2}$ • $-\infty < x < \infty, -\infty < \mu < \infty, \sigma > 0$ • The moments and mgf of the Cauchy do not exist (major consequence: CLT does not apply) • If $X, Y \sim N(0, 1)$, the ratio X/Y has the Cauchy distribution 	<ul style="list-style-type: none"> • $f_X(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{(-x/\beta)}$ • $0 \leq x < \infty, \alpha, \beta > 0$ • $M_X(t) = \left(\frac{1}{1-\beta t}\right)^\alpha, t < \frac{1}{\beta}$ • $\mathbb{E}[X] = \alpha\beta, \mathbb{V}[X] = \alpha\beta^2$ • When $\alpha = 1$, the Gamma becomes the Exponential distribution. With $\alpha = \frac{p}{2}$ and $\beta = 2$, the chi-squared distribution is recovered.