

<div>Probability</div> <div><math display="block">P(A \cup B) = P(A) + P(B) - P(A \cap B)</math><math display="block">P(A \cap B) = P(A) + P(B) - P(A \cup B)</math><math display="block">P(A \cup B) \leq P(A) + P(B), P(A) = 1 - P(A^c) \geq 0</math><b>Δ Bayes Theorem</b><math display="block">P(A \mid B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A, B)}{P(B)} = \frac{P(A)P(B \mid A)}{P(B)}</math><math display="block">P(A \cap B) = P(B)P(A \mid B)</math><b>Δ Law of Total Probability</b><math display="block">P(B) = P(A_1 \cap B) + \dots + P(A_n \cap B)</math><math display="block">= P(A_1)P(B \mid A_1) + \dots + P(A_n)P(B \mid A_n)</math><b>If Independent:</b><math display="block">P(A \mid B) = P(A), P(A \cap B \mid C) = P(A \mid C)P(B \mid C)</math><b>Pairwise independent A,B,C:</b><math display="block">P(A \cap B) = P(A)P(B), P(B \cap C) = P(B)P(C), P(A \cap C) = P(A)P(C)</math><b>Mutually/Fully independent A,B,C:</b><math display="block">P(A \cap B \cap C) = P(A)P(B)P(C)</math></div>	<div>Variance</div> <div><math display="block">Var(X) = E\left[(X - E[X])^2\right] = E\left[X^2\right] - (E[X])^2 = \sigma^2, \sigma = \sqrt{\sigma^2}</math><math display="block">Var\left(X^2\right) = E\left[\left(X^2\right)^2\right] - \left(E\left[X^2\right]\right)^2</math><b>Law of Total Variance</b><math display="block">Var(X) = E[Var(X \mid Y)] + Var(E[X \mid Y])</math>decompose <math>X</math> into <math>X \mid Y</math>, <math>Y</math> should be something that influences <math>X</math>.<b>If independent:</b><math display="block">Var(X \pm Y) = Var(X) + Var(Y)</math><b>If dependent:</b><math display="block">Var(X_1 + \dots) = \sum_{i=1}^n Var(X_i) + \sum_{i,j} Cov(X_i, X_j)</math><b>Scaling of Var(X)</b>extract the constant and <math>a^2</math>:<math display="block">Var(aX + b) = a^2 Var(X)</math><b>Sample Variance</b><math display="block">S_n = \frac{1}{n} \sum \left(X_i - \overline{X_n}\right)^2, E[S_n] = \frac{n-1}{n} \sigma^2</math><b>Unbiased Sample Variance</b><math display="block">\widetilde{S}_n = \frac{n}{n-1} S_n = \frac{1}{n-1} \sum \left(X_i - \overline{X_n}\right)^2 = \frac{1}{n-1} \left(\sum X^2 - n(\overline{X})^2\right), E\left[\widetilde{S}_n\right] = \sigma^2</math><b>Covariance</b><math display="block">Cov(X, Y) = E[(X - E[X])(Y - E[Y])]</math><math display="block">= E[XY] - E[X]E[Y]</math><math display="block">Cov(X, Y) = 0</math> means <math>X, Y</math> indepdent, only for <b>Normal</b> distributions<b>Correlation Coefficient</b><math display="block">\rho_{XY} = \frac{Cov(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1]</math></div>	<div>POISSON PROCESS</div> <div><b>Poisson</b> Distribution: <math>X</math> = # of events in a fixed time interval. Idea: Random # of <math>k</math> events in <math>t</math> units of time with rate <math>\lambda</math>, exact timing of <math>k</math> is random. <math>F(x)</math> = prob of <math>x</math> events happending in a fixed <math>t</math> time interval? General PMF: <math>P(X = k) = \frac{(\lambda t)^k e^{-\lambda t}}{k!}</math><math display="block">E[X_t] = Var(X_t) = \lambda t, I = \frac{1}{\lambda}, \lambda = \frac{E[X_t]}{t}</math> For 1 unit of time <math>t = 1</math>: <math>P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}</math><math display="block">E[X] = Var(X) = \lambda, E\left[X^2\right] = \lambda(1 + \lambda)</math><math display="block">I = \frac{1}{\lambda}, L(\lambda) = \frac{\lambda \sum X_i e^{-n\lambda}}{(x_1! \dots x_n!)}</math>Probability of <b>0</b> event in <b>1</b> unit time = <math>e^{-\lambda}</math>; <b>1</b> event = <math>\lambda e^{-\lambda}</math>, <b>2</b> events = <math>\frac{\lambda^2 e^{-\lambda}}{2}</math> When dealing with <b>time</b> intervals between events, use the <math>\sim \exp(\lambda)</math>: <b>Exponential</b> Distribution: <math>X</math> = unit time until <b>next</b> event in Poisson Process, aka the inter-arrival time. Idea: Random time len between Fixed Events <math>f_X(x) = \lambda e^{-\lambda x}, F_X(x) = 1 - e^{-\lambda x}</math> The CDF is: <math>P(\text{1}^{st} \text{ event at exactly } x \text{ time})</math>.<math display="block">E[X] = \frac{1}{\lambda}, E\left[X^2\right] = \frac{2}{\lambda^2}, Var(X) = (E[X])^2 = I(\lambda) = \frac{1}{\lambda^2}</math><math display="block">P(X \geq a) = \int_a^\infty \lambda e^{-\lambda x} dx = e^{-\lambda a}, \exp(\lambda) = \sum_{k=0}^\infty \frac{\lambda^k}{k!}</math><math display="block">\widehat{\lambda}^{MLE} = \frac{1}{\overline{X_n}} = \frac{n}{\sum X}, Var\left(\widehat{\lambda}^{MLE}\right) = \frac{\lambda^2}{n}</math><math display="block">L(\lambda) = \lambda^n e^{-\lambda} \sum X_i, l(\lambda) = n \ln \lambda - \lambda \sum X_i</math>Typically when we use <math>T_1, T_2</math> etc. this means "inter-arrival time". <math>T_2</math> = time interval from <math>T_1</math> to <math>T_2 \sim \exp(\lambda)</math> <b>Properties</b> - In a Poisson Process, events must be independent, and <math>\lambda</math> must be the same. - Conversely, each interval can be broken up into independent subintervals. - This Memoryless property means, given a fixed time <math>t</math>, any sub-intervals <math>\sim \text{Unif}[0, t]</math> <b>Sum</b> of Poisson Process: (<math>\lambda_i</math> must be independent) <math>X_1 \sim \text{Poi}(\lambda_1), X_2 \sim \text{Poi}(\lambda_2)</math> <math>X = X_1 + X_2 : X \sim \text{Poi}(\lambda_1 + \lambda_2)</math>. for 1 unit of time, you'd expect to see <math>\lambda_A + \lambda_B</math> events, so at anytime, you expect to see <math>\frac{\lambda_A}{\lambda_A + \lambda_B}</math> type <math>\lambda_A</math> events. see <math>\frac{\lambda_A}{\lambda_A + \lambda_B}</math> type <math>\lambda_A</math> events. <math>P(\text{out of } n \text{ arrivals are total of } k \text{ type A events}) = \binom{n}{k} \left(\frac{\lambda_A}{\lambda_A + \lambda_B}\right)^k \left(\frac{\lambda_A}{\lambda_A + \lambda_B}\right)^{n-k}</math><math display="block">Erlang(k) = Erlang\left(\frac{k}{2}\right) + Erlang\left(\frac{k}{2}\right)</math> <b>Split</b> of Poisson Process: <math>Ber(p), \lambda_A = \lambda p, \lambda_B = \lambda(1 - p)</math> <b>Order Statistics:</b> expected of the <math>T^{th}</math> event = <math>\frac{1}{\lambda} \sum_{i=1}^k \frac{1}{n - i + 1}</math> e.g. 3 engines, <math>E[T_3] = \frac{1}{3\lambda} + \frac{1}{2\lambda} + \frac{1}{\lambda}</math> <math>f_{T_3}(t) = F_{T_3}(t)</math> means the Prob that the 3<sup>rd</sup> order stat happened before <math>t</math>, which is the prob that 1<sup>st</sup> and 2<sup>nd</sup> event all each happened</div>	<div>Bernoulli Process</div> <div>Trials must be independent and <math>p</math> constant. <math>X \in \{0, 1\}</math> <b>Single Bernoulli Trial::</b> <math>P(X = 1) = E[X] = p, Var(X) = p(1 - p)</math><math display="block">I(p) = \frac{1}{p(1 - p)}, \widehat{p}^{MLE} = \frac{\text{Sum}}{n} = \overline{X_n}</math><math display="block">P(X = x) = \begin{cases} p &amp; \text{if } x = 1 \\ 1 - p &amp; \text{if } x = 0 \end{cases}</math> or equivocally: <math display="block">P(X = x) = p^x (1 - p)^{1-x}, I(p) = \frac{1}{p(1 - p)}</math><math display="block">L(X) = \prod p^{X_i} (1 - p)^{1-X_i} = p^{\sum X_i} (1 - p)^{n - \sum X_i}</math><b>n Bernoulli Trials:</b> The <math>P()</math> of <math>k</math> successes in <math>n</math> trials: <math display="block">P(X = k) = L(p) = \binom{n}{k} p^k (1 - p)^{n-k}</math>Sum of <math>k</math> success in <math>n</math> trials: <math>X = k = X_1 + \dots + X_n</math> <math>E[X] = np, Var(X) = np(1 - p)</math><b>Time until 1<sup>st</sup> success:</b> 1<sup>st</sup> success at trial <math>k</math> means trials <math>k - 1</math> all failed, and trial <math>k</math> succeed (Geometric Distro): <math display="block">P(X_1 = k) = (1 - p)^{k-1} p</math><math display="block">F_X(k) = P(X \leq k) = 1 - (1 - p)^k</math><b>Time of <math>k^{th}</math> success:</b> this is an extended case of the above (Negative Geometric). <math display="block">P(Y_k = n) = \binom{n-1}{k-1} p^k (1 - p)^{n-k}</math>means <math>P(k^{th} \text{ success at } n \text{ time})</math> For Geometrics: <math>E[Y_k] = \frac{k}{p}</math>, this means "E[time] until <math>k^{th}</math> success", Assume memorylessness: <math>P(A \mid B) = P(A)</math>. <math display="block">Var(Y_k) = \frac{k(1 - p)}{p^2}</math> <b>Merging Bernoulli Processes</b> <math>Ber(p + q - pq) = P(p) \cup P(q) = P(\text{either or both})</math> occurs. <b>Splitting Bernoulli Processes</b> <math>A \sim Ber(pq), B \sim Ber(p(1 - q))</math> These streams are not independent <math>\text{Binom}(p) \approx N(np, np(1 - p))</math>: <math>P(X = 19) = P(18.5 \leq X \leq 19.5)</math></div>	<div>Normal Distribution</div> <div><b>Normal Distribution</b> <math>\sim N(\mu, \sigma^2)</math>: <math display="block">f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \sim N(\mu, \sigma^2)</math><math display="block">F(x) = P(X \leq x) = \int_{-\infty}^x \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(t - \mu)^2}{2\sigma^2}\right) dt</math><math display="block">E[X] = \mu, E\left[X^2\right] = \mu^2 + \sigma^2</math><math display="block">E\left[X^3\right] = \mu^3 + 3\mu\sigma^2, E\left[X^4\right] = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4</math><math display="block">Var(X) = E\left[(X - \mu)^2\right] = \sigma^2</math><math display="block">Var\left(X^2\right) = 4\mu^2\sigma^2 + 2\sigma^4, Var\left(\overline{X}\right) = \frac{\sigma^2}{n}</math><math display="block">L(\mu, \sigma^2 \mid X) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum (X_i - \mu)^2\right)</math><math display="block">l(L) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum (X_i - \mu)^2</math><math display="block">\widehat{\mu}^{MLE} = \frac{\sum X}{n} = \overline{X}, \widehat{\sigma^2}^{MLE} = \frac{\sum \left(X - \widehat{\mu}\right)^2}{n} = S_n</math><b>Standard Normal Distribution</b> <math>\sim N(0, 1)</math>: <math display="block">f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)</math><math display="block">F(x) = \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{t^2}{2}\right) dt</math><math display="block">E[X] = 0, E\left[X^2\right] = 1, E\left[X^3\right] = 0, E\left[X^4\right] = 3</math><math display="block">Var(X) = 1, Var\left(X^2\right) = 2, Var\left(\overline{X}\right) = \frac{1}{n}</math><math display="block">E\left[\overline{X_n}\right] = 0, Var\left(\overline{X_n}\right) = \frac{1}{n}</math><math display="block">E\left[\frac{\sum X_i^2}{n}\right] = 1, Var\left(\frac{\sum X_i^2}{n}\right) = \frac{2}{n}</math><math display="block">E\left[\left(\frac{\sum X_i}{n}\right)^2\right] = \frac{1}{n}, Var\left(\left(\frac{\sum X_i}{n}\right)^2\right) = \frac{2}{n^2}</math><math display="block">L(0, 1 \mid X) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2} \sum X_i^2\right)</math><math display="block">l(0, 1 \mid X) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum X_i^2</math><math display="block">\widehat{\mu}^{MLE} = \frac{\sum X}{n} = \overline{X}, \widehat{\sigma^2}^{MLE} = \frac{\sum \left(X - \overline{X}\right)^2}{n} = S_n</math><b>Δ wolfram's <math>\sigma</math> in <math>N(\mu, \sigma)</math> is <math>\sqrt{\sigma^2}</math></b> Linear Shift: <math>Y = aX + b, X \sim N(\mu, \sigma^2), Y = N(a\mu + b, a^2\sigma^2)</math> Independent Sum: <math>Z = N_1 + N_2</math>, <math>Z \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)</math> Independent Diff: <math>Z = N_1 - N_2</math>, <math>Z \sim N(\mu_X - \mu_Y, \sigma_X^2 + \sigma_Y^2)</math> Z-Table: <math>\Phi(-2) = P(Y \leq -2) = 1 - P(Y \leq 2) = 1 - \Phi(2)</math> Std to <math>N(0, 1)</math>: <math>\frac{X - \mu}{\sigma}, X = \mu + \sigma Z</math></div>
<div>Combinatorics:</div> <div><b># of k combos using n things</b> Using <b>{A,B,C}</b> (<b>3</b> choose <b>2</b>): <b>Combination(no order, no repeat):</b> <math display="block">\binom{n}{k} = \frac{n!}{k!(n-k)!} = 3: \{A, B\}, \{A, C\}, \{B, C\}</math>Using <b>{A,B,A,B,B,C,C}</b> (<b>6</b> choose <b>2</b>): <b>Combination(no order, repeats):</b> <math display="block">\binom{n+k-1}{k} = \frac{(n+k-1)!}{k!(n-1)!} = 6: \{A, A\}, \{A, B\} \dots</math>Using <b>{A,B,C}</b> (<b>3</b> choose <b>2</b>): <b>Permutation(unique, no repeat):</b> <math display="block">P(n, k) = \frac{n!}{(n-k)!} = 6: AB, BA, AC, CA, BC, CB</math><b>Permutation(unique, repeats):</b> <math display="block">P(n, k) = n^k = 9: AA, AB, AC, BA, BB, \dots</math><b>Subsets:</b> # subsets for <math>n</math> things: <math>2^n = 8: \{\}, \{A\}, \{A, B\}, \{A, B, C\} \dots</math> Idea: Think about what <math>\Theta</math> is</div>	<div>Uniform Distribution</div> <div><b>Discrete</b> <math display="block">p_X(x) = \frac{1}{b - a + 1}, F_X(k) = \frac{ k  - a + 1}{n}</math><math display="block">E[X] = \frac{a + b}{2}, Var(X) = \frac{(b - a + 1)^2 - 1}{12}</math>if <math>k</math> is an integer: <math display="block">F_X(x) = \frac{k - a + 1}{b - a + 1}</math><b>Continuous</b> <math display="block">f_X(x) = \frac{1}{b - a} x \in [a, b], F_X(x) = \frac{x - a}{b - a}</math><math display="block">E[X] = \frac{a + b}{2}, Var(X) = \frac{(b - a)^2}{12}</math><math display="block">L(b) = \frac{1}{b^n} 1\{\max X_i \leq b\} \text{ on } [0, b]</math></div>	<div>Meta</div> <div><b>step 1:</b> Understand exactly what the question is asking <b>step 2:</b> break the question into smaller parts. e.g. How to break up complex <math>P(A)</math> into smaller <math>P(A_1)P(A_2) \dots</math></div>	<div>MIXED RV</div> <div><math>P(X = Y) = p</math> <math>Y</math> ~discrete, <math>P(X = Z) = 1 - p</math> <math>Z</math> ~continuous CDF: <math>F(x) = p \cdot F_Y(x) + (1 - p) \cdot F_Z(x)</math> <math>E[X] = p \cdot E[Y] + (1 - p) \cdot E[Z]</math> Derive <math>f_X(x)</math> from <math>F_X(x)</math></div>	<div>Statistical Models</div> <div><math>(E, (P_\theta)_{\theta \in \Theta})</math> <math>E</math>: sample space <math>(X_1 \dots)</math> <math>P</math>: family of prob measures on <math>E</math></div>
<div>Expectations</div> <div><math display="block">E[X] = \sum x p_X(x), E[g(x)] = \sum_x g(x) p_X(x) = \mu</math><math display="block">E[X] = \int_{-\infty}^\infty x f_X(x) dx, E[g(x)] = \int_{-\infty}^\infty g(x) f_X(x) dx</math><math display="block">E\left[X^2\right] = Var(X) + (E[X])^2</math><b>Linearity of Expectation:</b> <math display="block">E[aX + bY + c] = aE[X] + bE[Y] + c</math><math display="block">E[X \pm Y] = E[X] \pm E[Y]</math><b>If Independent:</b> <math display="block">E[XY] = E[X]E[Y], E\left[(XY)^2\right] = E\left[X^2\right]E\left[Y^2\right]</math><b>Law of Total Expectation</b> <math display="block">E[X] = E[E[X \mid Y]] = \sum_y p_Y(y) E[X \mid Y = y]</math><math display="block">E[X] = \int_{-\infty}^\infty f_Y(y) E[X \mid Y = y] dy</math><b>Iterated Expectation</b> <math display="block">E[X] = E[E[X \mid Y]]</math><math display="block">E[XY] = E[E[XY \mid X]] = E[E[XY \mid Y]]</math><math display="block">E[XY \mid X] = X \cdot E[Y \mid X], E[XY \mid Y] = Y \cdot E[X \mid Y]</math><b>Conditional &amp; Joint</b> <math display="block">E[X \mid Y = y] = \sum_x x \cdot p_{X Y}(x \mid y) \text{ or for } g(x)</math><b>if X,Y are indepdent:</b> <math display="block">f_{X,Y}(x, y) = f(x) \cdot f(y), f(x) = \frac{f_{X,Y}(x, y)}{f(y)}</math><math display="block">f_{X Y}(x \mid y) = f(x)</math>deterministic r.v. are always indie</div>		<div>Conditional &amp; Joint</div> <div><math display="block">p(x \mid y) = \frac{p(x, y)}{p(y)} = \frac{p(x)p(y \mid x)}{p_Y(y)}</math>Fixed values <math>X = x</math> are different from <math>X</math>, <math>Y = X + N, N \sim N(0, 1) \rightarrow f(Y \mid X = x) \sim N(x, 1)</math> <b>Multiplication Rule</b> <math display="block">p(x, y) = p(y) \cdot p(x \mid y) = p(x) \cdot p(y \mid x)</math><math display="block">p(x, y, z) = p(x) \cdot p(y \mid x) \cdot p(z \mid x, y)</math><math display="block">p(x, y \mid z) = \frac{p(x, y, z)}{p(z)}</math>if <math>X, Y</math> cond. indie of <math>Z</math>: <math display="block">p(x, y \mid z) = p(x \mid z) \cdot p(y \mid z)</math>For Joint <b>Normal</b> <math>X, Y</math>: <math display="block">E[X \mid Y] = \mu_X + \rho \frac{\sigma_X}{\sigma_Y} (y - \mu_Y), Var(X \mid Y) = \sigma_X^2 \left(1 - \rho^2\right)</math><b>Marginal using Total Probability</b> <math display="block">p(x) = \sum_y p(x, y) = \sum_y p(x \mid y) p(y)</math><math display="block">f(x) = \int f(x, y) dy = \int f(x \mid y) f(y) dy</math>Prob of a joint in region <math>A</math>: <math display="block">P((X, Y) \in A) = \int \int_{(x, y) \in A} f(x, y) dx dy</math><b>Expected Value Rule</b></div>		

<b>Differences Squared Measures</b> (X-c)^2: Emphasize bigger differences, makes finding critical points easier. E[(X-c)^2]: average deviations. E[(X-θ)^2]: <b>MSE</b> of an estimator, aka Quadratic Risk: R(θ̂) = Var(θ̂) + (Bias(θ̂))^2 e.g. MSE of μ-X̄ is standard error^2: E[(μ-X̄)^2] = (σ/√n)^2 = σ^2/n <b>Variance &amp; Sample Variance:</b> E[(X-E[X])^2] = Var(X), E[(X-X̄)^2] → Sn <b>LMS:</b> Least Mean Squared Error, aka Bayes Estimator, minimizes posterior, is usually the μ=E[X], or the E[X Y] of a joint distribution: E[(X-X̂ <sup>LMS</sup> )^2] = E[(X-E[X Y])^2] = Var(X Y) <b>Test Statistics</b> Typically measures the differences between an estimator λ̂ and the hypothesis λ0: λ̂-λ0
---

<b>ESTIMATORS</b> <b>Asym. Normal if:</b> √n(θ̂-θ) → N(0,σ^2) <b>Consistent if:</b> θ̂ → θ as n → ∞ <b>Bias:</b> Bias(θ̂) = E[θ̂] - θ <b>Quadratic Risk:</b> R(θ) = E[ θ̂-θ ^2] <b>Unbiased Estimator</b> θ̂ such that bias(θ̂)=0, which means θ̂ → θ if a θ̂ is biased, use linearity of expectations to create a new estimator such that E[θ̂] = c · E[θ̂] = θ <b>Continuous Mapping Theorem</b> Apply continuous g(x) to a sequence of {Xi} → X, then the transformed {g(Xi)} will also → g(X). As in, applying g(x) preserves the convergence property. <b>Note:</b> plug in estimators in operations, e.g. plug-in λ̂ <sup>MLE</sup> = 1/X into l(λ)
---

<b>MLE Estimator</b> Since L(X θ) means the likelihood of observing Xi assuming θ is true, MLE aim to find a θ̂ that makes the data most probable: θ̂ <sup>MLE</sup> = arg max_θ L(X θ).  Does not incorp prior p(θ), unlike <b>MAP</b> . The MLE calculated for the prior p(θ) is called a constrained MLE <b>wolfram:</b> d/dθ ln(∏ f_θ(x)) = 0 Likelihood: all of L(X θ) happening together: L(X θ) = ∏_{i=1}^n f(Xi θ).  <b>Properties</b> If support of P_θ doesn't depend on θ, θ* is not at boundary, I(θ) is invertible: <b>Consistent:</b> As n ↑, MLE → θ* → to true θ <b>Asymptotic Normality:</b> For large n, √n(θ̂ <sup>MLE</sup> - θ*) → N(0, 1/I(θ*)), Var(θ̂) = 1/I(θ),
--

Θ: Param set. Well specified if the true θ* ∈ Θ Δ sample space must not depend on parameter Δ sample space must be the support for the distribution. i.e. ([0,∞), {N(μ,σ^2)}) is not valid because the sample space for a N is all R <b>Identifiability</b> θ identifiable only if mapping θ ∈ Θ → P_θ is injective (injective: θ ≠ θ' ⇒ P_θ ≠ P_θ') <b>LLN</b> if iid: lim_{n→∞} X̄_n = E[X] = μ <b>CLT</b> if iid, and n is large: √n/σ (X̄_n - μ) = (X̄_n - μ)/σ/√n → N(0,1) √n(X̄ - μ) → N(0,σ^2): X̄ ~ N(μ, σ^2/n) need to enlarge X̄ - μ by √n, otherwise as n ↑, X̄ - μ will → 0
--

<b>Critical Values</b> (z-score) a = 0.025, 2-tails: a/2 = 0.0125 in each tail q_{a/2} = 2.24, 1-tail: q_a = 1.96 a = 0.05, 2-tails: a/2 = 0.025 in each tail q_{a/2} = 1.96, 1-tail: q_a = 1.645. a = 0.1, 2-tails: a/2 = 0.05 in each tail q_{a/2} = 1.645, 1-tail: q_a = 1.28.
--

<b>Composition of P(A∩B∩C)</b> P(A^c) → P(A) → P(B A) ↑ P(A∩B) → P(C A∩B) → P(A∩B∩C)
---

<b>Parameters</b> θ: True but unknown population θ θ*: Asym of θ, is the lim of θ_n θ̂: point estimator of θ θ0: the θ that H0 hypothesizes θ̂: alt estimator, estimator of θ θn: sequence of parameters θ̄: mean of θ, or θ_n θ̂ <sup>MLE,MAP</sup> : specific types of θ̂
---

<b>Bayesian</b> Treats θ as random variables. Uses prior knowledge. Provide posterior distribution Makes probabilistic statements about θ. <b>Frequentist</b> Treats θ as fixed but unknown. Relies solely on data. Provides θ̂ and confidence intervals. Makes probabilistic statements about the data.
---

<b>QQ Plot</b> Visually access what distribution plotted dataset follows, by comparing its quantiles to the quantiles of a
---

prob that 1^st, 2^nd, 3^rd event all each happened before time t, which means F_{T_3}(t) = P(T_1 ≤ t, T_2 ≤ t, T_3 ≤ t) = (1 - e^{-λt})^3, differentiate this to get f_{T_3}(t)
---

<b>Gamma Distribution</b> f(x; α, β) = β^α / Γ(α) · x^{α-1} e^{-βx}, Var(X) = α/β^2 x̂ <sup>Bayes</sup> = E[X] = α/β, E[X^2] = α(α+1)/β^2 Γ(α) normalizes the distribution to 1. if 'α' integer: Γ(α) = (α-1)!, posterior π(θ X) ∝ Γ(θ; α, β), θ̂ <sup>MAP</sup> = Mode = α-1/β X1 ~ Γ(α1, β), X2 ~ Γ(α2, β), X1 + X2 ~ Γ(α1+α2, β) For use with time-related: until k <sup>th</sup> , not limited to just next event: f(x) = (λ^k x^{k-1} e^{-λx}) / Γ(k) · E[X] = k/λ, Var(X) = k/λ^2 <b>Erlang:</b> Gamma but k must be integer: f(x) = (λ^k x^{k-1} e^{-λx}) / (k-1)!, E[X] = k/λ, Var(X) = k/λ^2
---

<b>Beta Distribution</b> Continuous distribution defined on [0,1]. Models proportions, used as prior. f(x; α, β) = x^{α-1} (1-x)^{β-1} / B(α, β) B(α, β) = ∫_0^1 t^{α-1} (1-t)^{β-1} dt = Γ(α)Γ(β) / Γ(α+β) E[X] = α / (α+β), Var(X) = αβ / (α+β)^2 (α+β+1) Mode = (α-1) / (α+β-2) When α=β=1, Beta is a <b>Uniform</b> .
---

<b>Confidence Interval</b> An interval that will contain the true parameter with a likelihood of 1-α, cannot depend on unknown parameter. CI = [Estimate - Critical Value · Standard Error, Estimate + Critical Value · Standard Error] P(X̄_n - q_{α/2} · σ/√n ≤ μ ≤ X̄_n + q_{α/2} · σ/√n) = 1 - α Δ√σ^2 = σ <b>Standard Errors Forms</b> σ/√n, √(σ^2/n), √(p(1-p)/n), √(σ1^2/n1 + σ2^2/n2) for X̄1 - X̄2 Plug-in Method: use the θ̂ in variance formulas for standard error.
---

<b>Hypothesis Tests</b> Δ fail to reject H0 is not accepting H0 θ, Xi, θ̂: Given sample space, we use samples Xi to construct θ̂ for the θ we are interested in inferencing. H0, H1: We hypothesized a value for θ we call θ0, H1 could be θ0 ≠ θ, (2-tailed) or θ0 > or < θ (1-tail).  Tn: test statistic, usually (θ̂ - θ0) / SE(θ̂) Tn   H0 ~ a known distribution, the goal is to leverage the known properties of distributions to infer H0, e.g. Z-test, t-test, χ^2 test, Wald's, LRT. <b>1 tail:</b> P(X > or < Tn) is the tail to the right or left; <b>2-tails:</b> P( X  >  Tn ) for tails on both side.
--

<b>Expected Value Rule</b> E[g(X, Y)] = ∑_x ∑_y g(x, y) p(x, y) E[g(X, Y)] = ∫ E[g(x, y)   Y = y] f(y) dy E[g(X, Y)   Y = y] = ∫ g(x, y) f(x   y) dy <b>CDF:</b> F(x, y) = P(X ≤ x, Y ≤ y) = ∫_{-∞}^y ∫_{-∞}^x f(s, t) ds dt <b>Convolution</b> Find the Probability of the Sum of variables: Z = X + Y p_Z(z) = ∑_x p_X(x) p_Y(z - x) f_Z(z) = ∫_{-∞}^∞ f_X(x) f_Y(z - x) dx <b>Random # * Random Amounts</b> N: # of events Xi: Amount of each event T: Total Amount: T = ∑_{i=1}^N Xi E[T] = E[N] E[X] = E[N] · μ by law of total variance: Var(T) = E[Var(T   N)] + Var(E[T   N]) Var(T) = σ^2 · E[N] + μ^2 · Var(N)
--

<b>Fisher Information</b> Δ use only <u>ONE</u> observation: 1/Var(X) I(θ) = E[(d/dθ ln f(X; θ)]^2 = -E[d^2/dθ^2 ln f(X; θ)] I(θ) = Var(l'(θ)) = -E[l''(θ)], ∑ I_i(θ) = n · I(θ) Cramer-Rao Bound: Var(θ̂) ≥ 1/I(θ) Invariance re-para: I(φ) = I(θ) (dθ/dφ)^2 Matrix: I_{ij}(θ) = E[d/dθ_i ln f(X; θ) · d/dθ_j ln f(X; θ)] not well-defined if distro support depends on unknown paramenter (e.g. shifted exp(λ)) d^2/dθ^2 of log-likelihood must exist: I(θ) = Var(l'(θ)) = -E[l''(θ)]
--

<b>Wald's Test</b> Accesses θ̂ <sup>MLE</sup> with H0, typically for H1: θ̂ ≠ θ0. Assumes θ̂ ~ N(). Z-form: T_n^Wald = (θ̂ - θ0) / SE(θ̂) = (θ̂ - θ0) / √Var(θ̂) ~ N(0,1),
--

reject H0 when  Tn  > z_{a/2} (2-tailed). χ1^2 form: T_n^Wald = ((θ̂ - θ0) / SE(θ̂))^2 = ((θ̂ - θ0) / √Var(θ̂))^2 ~ χ1^2 re-para g(θ0) =, wald will not be robust. Multi-variate: n(θ̂ <sup>MLE</sup> - θ0)^T · (θ̂ <sup>MLE</sup> - θ0) · I(θ̂ <sup>MLE</sup> ) → χ_d^2 Tn =    √n I(θ0)^{1/2} (θ̂ <sup>MLE</sup> - θ0)   ^2 → χ_d^2
---

<b>Likelihood Ratio Test</b> Purpose: test if data fits model better with θ0 or without. H0: selectively set hypothesized values to part or all parameters of the target
--

if any part of the ∫ "doesn't converge", means infinite and can't be normalized. - Evaluating at specific points gives the <b>relative density</b> of that point, not the actual probabilities, to get the actual probability in that range: integrate between the bounds: P(a ≤ X ≤ b) = ∫_a^b f(x) dx, this gives the CDF for that range.
---

<b>CDF</b> (gives actual Prob) → 0 as x → -∞, → 1 as x → ∞. Gives the Prob of the interval which the CDF was integrated up to F_X(x) = 1 when x ≥ upper bound non-decreasing, right-continuous. P(a < X ≤ b) = F_X(b) - F_X(a) Integrate up to x for the CDF of X F_X(x) = ∫_{-∞}^x f_X(t) dt = P(X ≤ x) Inversely, to get the PDF of X, differentiate f_X(x) = d/dx F_X(x) Evaluating specific value gives the area to the left: P(X ≤ x), for the area to right: 1 - P(X ≤ x), which gives P(X > x). Y = F(x), F_Y ~ Uni(0,1) <b>Empirical CDF:</b> step fn that estimates population F(t) on each 1/n step. F_n(t) = 1/n ∑ 1{Xi ≤ t}, eCDF is a step fn so discont, jumps with 1/n at each Xi, converge to true F(t). Asym Normal: √n(F_n(t) - F(t)) → N(0, F(t)(1 - F(t))), F(t)(1 - F(t)) is the var of binom, due to the indicator.
--

<b>PMF</b> Non-negative: P(X = x) ≥ 0, Must Sum to 1: ∑_x p(x) = 1, p(x) = P(X = x) Linear Transform Y = aX + b, p_Y(y) = p_X(y - b/a), f_Y(y) = 1/ a  f_X(y - b/a) g is monotonic: f_Y(y) = f_X(h(y))  dh(y)/dy , where h is inverse of g 1) find CDF: F_Y(y) = P(g(X) ≤ y) 2) derive CDF for y to find PDF
--

<b>Misc</b> (a - b)^2 = a^2 + b^2 - 2ab (a + b)^2 = a^2 + b^2 + 2ab <b>Log</b> ln(mn) = ln(m) + ln(n), ln(m/n) = ln(m) - ln(n) ln(m^r) = r ln(m), ln(1) = 0, ln(x ≤ 0) undefined remove or add e: e^{ln(x)} = x e^{a+b} = e^a · e^b, e^{ab} = (e^a)^b, e^{-x} = 1/e^x, e^0 = 1 add ln to both sides to remove e: e^x = 2 → x = ln(2) add e to both sides to remove ln: ln(x) = y is e^{ln(x)} = e^y then x = e^y Wolfram: complete the square x^2.. <b>Exponent</b> (ab)^x = a^x b^x, (a^x)^y = a^{xy}, a^x a^y = a^{x+y}
---

$\sigma^2 \cdot (g(\theta))^2$  if re-parameterized. MLE minimizes **KL diver**, can be Biased.  $f(x)$  must be continuously differentiable to find critical values. No sharp corner, vertical tangents, or discontinuous functions.

**MAP Estimator**  
**MAP:** Maximum Posterior Estimator. Incorporates prior  $p(\theta)$ , unlike **MLE**.  
 $\hat{\theta}^{MAP} = \arg \max \pi(\theta | X) = \arg \max L_n(X | \theta) \pi(\theta)$ .  
Find critical value(s) via  $\frac{d}{d\lambda} \pi(\theta | X) = 0$ , verify min or max by  $\frac{d^2}{d\lambda^2} \pi(\theta | X)$  and plug back into  $\lambda^{MAP}$ .  
MAP is the mode of  $\arg \max \pi(\theta | X)$ , which means  $\hat{\theta}^{MAP}$  is one of the values of  $\theta$  if  $\theta$  is a discrete set.

**Example:**  
given  $g(x) = \lambda e^{-\lambda} x, \lambda \sim \exp(\theta)$ , find  $\lambda^{MAP}$ .  
Solution 1.  
1) Posterior  $\propto$  Likelihood\*Prior. We ignore the marginal  $p(X)$  because it acts as a constant here so it will not change  $\arg \max$  of  $\lambda$ .  
2)  $\pi(\lambda | X) \propto L(X | \lambda) \pi(\lambda) = \lambda^n e^{-\lambda} \sum X \cdot (\theta e^{-\theta \lambda})$   
3) differentiate this posterior with respect to the parameter we're trying to estimate, set it to 0 and solve.  
4) wolfram:  $\frac{d}{d\lambda} \lambda^n e^{-\lambda} \sum X \cdot (\theta e^{-\theta \lambda}) = 0$ , we get  $\lambda^{MAP} = \frac{\sum X + \theta}{\sum X + \theta}$   
5) wolfram:  $\frac{d^2}{d\lambda^2} \ln(\lambda^n e^{-\lambda} \cdot (\theta e^{-\theta \lambda})) = -\frac{n}{\lambda^2}$   
6) plug this back into  $\lambda^{MAP}$  we get  $-\frac{(\sum X + \theta)^2}{n}$ , since  $(\sum X + \theta)^2$  must  $> 0$ , and  $n > 0$ , 2nd derivative test is neg, which means  $\lambda^{MAP}$  is the maximum.

**Canonical Exponential Family**  
**General:**  $f(x | \theta) = h(x) \exp(\eta(\theta) \cdot T(x) - A(\theta))$   
**1 parameter:**  $f_{\theta}(x) = \exp\left(\frac{x\theta - b(\theta)}{\phi} + c(x, \phi)\right)$   
 $= \exp\left(x \cdot \frac{\theta}{\phi} - \frac{b(\theta)}{\phi}\right) \cdot e^{c(x, \phi)}$  try to move all  $\theta$  related terms inside the  $\exp()$   
 $\theta = \eta(\theta)$ : main canonical parameter  
 $x = T(x)$ : sufficient stat, carry main random info on  $x$   
 $e^{c(x, \phi)} = h(x)$ : base measure, indie of  $\theta$   
 $\frac{b(\theta)}{\phi} = A(\theta)$ : log-partition function, normalizes the pdf to 1  
 $\phi$ : dispersion, effects variance (GLM)  
**Properties:**  
 $E[X] = b'(\theta)$ ,  $Var(X) = b''(\theta)\phi$ ,  $I(\theta) = \frac{nb''(\theta)}{\phi}$   
 $\Delta$  these formulas belongs to  $\exp(\eta(\theta) \cdot T(x) - A(\theta))$ , when using transformed para, must replace the transformed para with the equal  $\theta$ .  
**linear transformations** also canon.  
**Canon Link** relates  $\mu = E[X]$  to  $\theta$ :  
 $g(\mu(x)) = \eta = \theta = (b')^{-1}(\mu(x))$ , if  $\phi > 0$ , canon link is strictly increasing.  
**log-L:**  $l(\theta; x) = \sum \left[ \ln h(X_i) + \frac{T(X_i)\theta - b(\theta)}{\phi} \right]$

quantiles to the quantiles of a known distribution.  
Theoretical quantile-axis: the reference distribution to compare to, the points on the axis are its support.  
Sample quantile-axis: the distribution from the sample which we are interested in, values on axis are its support.

lighter tails  
lesser ext values  
such as uniform

fatter tails  
more ext values  
such as t-distro

**Canonical Exp Family Example:**  
Form:  $f_{\theta}(x) = h(x) \exp(\eta(\theta)T(x) - A(\theta))$   
 $f_{\lambda}(x) = \frac{x^4}{24\lambda^5} \exp\left(-\frac{x}{\lambda}\right) = \exp\left(-\frac{1}{\lambda} \cdot x\right) \cdot \frac{1}{\lambda^5} \cdot \frac{x^4}{24} = \exp\left(-\frac{1}{\lambda} \cdot x\right) \cdot \exp(-5 \ln(\lambda))$   
 $\exp(4 \ln(x) - \ln(24)) = \exp\left(-\frac{1}{\lambda} \cdot x - 5 \ln(\lambda) + 4 \ln(x) - \ln(24)\right)$

convert all to  $\theta$ :  $\theta = -\frac{1}{\lambda} \rightarrow \lambda = -\frac{1}{\theta}$ ,  
 $T(x) = x$ ,  $h(x) = 4 \ln(x) - \ln(24)$ ,  $\phi = 1$ ,  
 $A(\theta) = -5 \ln(\lambda) = -5 \ln\left(-\frac{1}{\theta}\right) = 5 \ln(-\theta)$   
since  $\ln(x \leq 0)$  is undefined, we must make sure  $-\theta > 0$ , since  $\theta = -\frac{1}{\lambda}$ , and  $\lambda > 0$ ,  $\theta < 0 \rightarrow -\theta > 0$ .  
 $E[X] = A'(\theta)$  is only true when  $\theta$  is  $\uparrow$  in the same direction as  $T(X)$ , and  $E[X]$  must  $\geq 0$ ,  $A'(\theta)$  cannot  $< 0$   
 $A'(\theta) = \frac{d}{d\theta} 5 \ln(-\theta) = \frac{5}{\theta}$ ,  $\theta < 0$ ,  $\left(\frac{5}{\theta}\right) < 0$  so this means

$E[X] = -A'(\theta) = -\frac{5}{\theta} = -\frac{5}{-\frac{1}{\lambda}} = 5\lambda$ ,  
 $Var(X) = A''(\theta)\phi = \frac{-5}{\theta^2} \cdot 1$ ,  $Var$  must be  $\geq 0$  so  $Var(X) = \frac{5}{\theta^2} \cdot 1 = \frac{5}{\left(-\frac{1}{\lambda}\right)^2} = 5\lambda^2$

$\hat{\lambda}^{MLE} = \frac{d}{d\lambda} \ln\left(\frac{x^4}{24\lambda^5} \exp\left(-\frac{x}{\lambda}\right)\right) = 0$   
for MLE of a single  $x$ :  
 $\frac{d}{d\lambda} \ln\left(\frac{x^4}{24\lambda^5} \exp\left(-\frac{x}{\lambda}\right)\right) = 0$ ,  $\hat{\lambda} = \frac{x}{5}$ , for the MLE of all the  $x$ ,  
 $\hat{\lambda}^{MLE} = \frac{1}{n} \cdot \frac{\sum X_n}{5} = \frac{\bar{X}_n}{5}$   
We already derived that  $E[X] = 5\lambda$ , since  $E[X] = \mu = \bar{X}_n = 1^{st}$  moment,  
 $\hat{\lambda}^{MM} = E\frac{X}{5} = \frac{\bar{X}_n}{5}$ .  $|\bar{X}_n - 5| > C_{0.05, n}$ :  
 $\bar{X}_n$  by CLT  $\rightarrow N\left(5, Var(\bar{X}_n)\right)$ ,  
 $Var(\bar{X}_n) = \frac{Var(X_i)}{n} = \frac{5\lambda_0^2}{n} = \frac{5 \cdot 1^2}{n} = \frac{5}{n}$ ,  
 $C_{\alpha, n} = q_{\frac{\alpha}{2}} \cdot \sqrt{Var(\bar{X}_n)} = 1.96 \cdot \sqrt{\frac{5}{n}}$

$95^{th}$  quantile means the area that makes up 95% of the area under curve, can also be expressed as:  $q(1-\alpha)$   
critical value:  $T_n$  that marks a certain area% under curve,  $q_{\alpha}$  or  $q_{\frac{\alpha}{2}}$   
we use  $1 - CDF$ , or  $1 - \phi(T_n)$  for the right tail probabilities.  
**p-value:** how likely the data "belongs" to  $H_0$ , the lower the p-value, the more unlikely  $H_0$  is true. Set a  $\alpha$  level, if  $p \leq \alpha$ , then  $H_0$  might not be true, reject  $H_0$ .  
Significance:  $\alpha$  = type I error, error rate for rejecting  $H_0$  when  $H_0$  is actually true. higher  $\alpha$  means more likely to reject  $H_0$ .  
 $\beta$ : type II error. Fail to reject  $H_0$  when  $H_1$  is actually true. Higher  $\beta \rightarrow$  more likely to fail to reject  $H_0$  when  $H_1$  is true.  
 $1 - \beta$ : power of test: prob of detecting true effect.

**Workflow**  
1) State  $H_0, H_1, \alpha$   
2) Compute an  $\hat{\theta}$  for  $H_0 : \theta_0$   
3) Compute an  $T_n$ , state its distribution  
4) Compute  $T_n$ 's p-value  
5) Find the critical value  $q_{\alpha}$  (1-tail) or  $q_{\frac{\alpha}{2}}$  (2-tails) of the distribution via wolfram  
6) Reject  $H_0$  if:  $T_n > q_{1-\alpha}$ ,  $|T_n| > q_{\frac{\alpha}{2}}$  or for z-test:  $p \leq \alpha$  (p is already std z-score)

**Bayesian Inference**  
Bayes Theorem: posterior =  $\frac{\text{prior} \cdot \text{likelihood}}{\text{evidence}}$   
**Prior:** previous belief about  $\theta$   
**Likelihood:** Prob of observing data assuming  $\theta$  is true  
**Evidence:** Marginal that ensure the posterior integrate to 1  
**Posterior:** Updated belief about  $\theta$  given new evidence  
 $\pi(\theta | X) = \frac{\pi(\theta) \cdot L_n(X | \theta)}{p(X)} \propto \pi(\theta) L_n(X | \theta)$ , remember to multiply the normalizing factor.  
 $p(X) = \int_{\theta} \pi(\theta) L_n(X | \theta)$   
**Conjugate Prior:** priors that, combined with  $L_n(X | \theta)$ , results in a posterior that ~ the same family of distributions as the prior. Prior:  $\theta \sim B\eta(\alpha, \beta)$ , Likelihood:  $P(X | \theta) = \theta^k (1 - \theta)^{n-k}$ , Posterior:  $\theta | X \sim B\eta(\alpha + k, \beta + n - k)$ .  $\Gamma \cdot \text{Exp} \sim \Gamma$

For proper priors:  $\sum \pi(\theta) = 1$ , or  $\int \pi(\theta) d\theta = 1$   
**Improper Prior:** Non-informative priors that don't integrate to 1, and are not valid distributions such as **Jeffreys**  
Prior:  $\pi_J(\theta) \propto \sqrt{\det I(\theta)}$ , invariant under reparameterization: if  $\eta = \Phi(\theta)$ , we need:  
 $\pi_J(\eta) \propto \pi_J(\theta) \cdot \left| \frac{d\theta}{d\eta} \right|$ , substituting  $\theta = \Phi^{-1}(\eta)$ , we have  $\pi_J(\eta) \propto \sqrt{\det I(\Phi^{-1}(\eta))} \cdot \left( \Phi'(\Phi^{-1}(\eta)) \right)^{-1}$   
Confidence region:  $P(a \leq \theta \leq b | X) = 1 - \alpha$ , Credible Interval:  $[a, b]$   
Under Bayesian frameworks, the convention is to use  $\pi$  to denote pmf or pdf

**Bayes Estimator**  
 $\hat{\theta}$  via its  $\pi(\theta | X)$  using a chosen loss function  $L(\theta, \hat{\theta})$ :  $\hat{\theta}^{\pi} = \arg \min E[L(\theta, \hat{\theta}) | X]$

use para of all parameters of the target model, and find  $L(\theta_0)$   
**H1:** use estimators such as MLE for all the paras of the model, and find  $L(\hat{\theta})$   
LRT is defined as  $\Lambda = \frac{L(\theta_0)}{L(\hat{\theta})}$ ,  
 $T_n = -2 \ln \Lambda = 2 \left( \ln L(\hat{\theta}) - \ln L(\theta_0) \right) \sim \chi_k^2$ , degrees of freedom = numbers of  $\theta_0 - \hat{\theta}$ , so if  $H_0$  had 1 restricted  $\lambda_0$ , for  $H_1$ 's 1 unrestricted  $\hat{\lambda}$ , then  $df = 1 - 0 = 1$ .  
 $\psi_{\alpha} = \mathbf{1}\{T_n > q_{\alpha}\}$

**t Test**  
Test if  $\mu_1$  is significantly different from  $\mu_2$ , or between  $\mu$  and  $\mu_0$   
Assumes: iid, **X~N()**,  $\sigma_1^2 \approx \sigma_2^2$ ,  $n < 30$ , uses t-table, is non-asm  
**t-Distribution:** composed of  $N(0,1)$  rvs for small  $n$ ,  $n < 30$ .  
**Concept:** t-test is the Z-test but instead of  $Var(\hat{\theta})$  that is asym correct, use the unbiased  $Var(\text{Sample})$ , and introduce the extra variabilities.  
**1-sample 1-sided:**  $H_0: \mu \leq \mu_0$  vs  $H_1: \mu > \mu_0$   
 $T_n = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\sqrt{\hat{S}_n}} \sim t_{n-1}$ ,  $\psi_{\alpha} = \mathbf{1}\{T_n > q_{\alpha}\}$

**1-sample 2-sided:**  $H_0: \mu = \mu_0$  vs  $H_1: \mu \neq \mu_0$   
 $T_n = \frac{\sqrt{n}\bar{X}_n}{\sqrt{\hat{S}_n}} = \frac{\sqrt{n}\frac{\bar{X}_n - \mu_0}{\sigma}}{\sqrt{\frac{\hat{S}_n}{\sigma^2}}} \sim t_{n-1}$ ,  
 $\psi_{\alpha} = \mathbf{1}\{|T_n| > q_{\frac{\alpha}{2}}\}$   
**2-samples** (Welch's Test)  
special case of 2-sampled t-test where,  $\sigma_1^2 \neq \sigma_2^2$  (heteroscedasticity), and  $n_1 \neq n_2$   
 $\bar{X}_n \sim N\left(\mu_X, \frac{\sigma_X^2}{n}\right)$  and  $\bar{Y}_m \sim N\left(\mu_Y, \frac{\sigma_Y^2}{m}\right)$   
 $\frac{\bar{X}_n - \bar{Y}_m - (\mu_X - \mu_Y)}{\sqrt{\frac{\hat{\sigma}_X^2}{n} + \frac{\hat{\sigma}_Y^2}{m}}} \sim t_N$ ,  
 $N = \frac{\frac{\hat{\sigma}_X^2}{n}}{n^2(n-1)} + \frac{\frac{\hat{\sigma}_Y^2}{m}}{m^2(m-1)} \geq \min(n, m)$

**Bayes Estimator Example**  
1)  $\pi(p = 0.4) = 0.2, \pi(p = 0.7) = 0.8$   
2) Compute all likelihoods  $L(X | p)$ :  
 $L(X = 3 | p = 0.4) = \binom{6}{3} \cdot 0.4^3 \cdot (1 - 0.4)^{6-3} = 0.27648$ ,  $L(X = 3 | p = 0.7) = 0.18522$   
3) Compute all posteriors  $\pi(p | X)$ :  
 $\pi(p = 0.4 | X = 3) = \pi(p = 0.4) \cdot L(X = 3 | p = 0.4) = 0.2 \cdot 0.27648 = 0.055296$ ,  $\pi(p = 0.7 | X = 3) = 0.148176$   
4) Normalize  $\pi(p | X)$ :  
 $\pi(p = 0.4 | X = 3) = \frac{0.055296}{0.055296 + 0.148176} = 0.27176$ ,  $\pi(p = 0.7 | X = 3) = 0.72823$   
normalized  $\pi(p | X)$  should sum to 1.  
5)  $\hat{p}^{Bayes} = E[p | X = 3] = \sum p_i \cdot \pi(p | X) = (0.4 \cdot 0.27176) + (0.7 \cdot 0.72823) = 0.618$   
Even though  $p \in \{0.4, 0.7\}$ ,  $\hat{p}^{Bayes}$  thinks that  $E[p | X] = 0.618$  means the outcome is closer

**e limits**  
 $\lim_{n \rightarrow \infty} \left(1 - \frac{t}{n}\right)^n = e^{-t}$ ,  $\lim_{n \rightarrow \infty} \left(1 + \frac{t}{n}\right)^n = e^t$   
**Indicators**  
 $I_A = 1$  if  $P(A)$ ,  $I_A = 0$  if  $P(A^c)$ ,  $E[I_A] = P(A)$

**Wolfram Syntax**  
Distros as objects: NormalDistribution[ $\mu, \sigma$ ]  
**Δbeware**  $\sigma$  is std dev, not var!  
Commands: PDF[NormalDistribution[ $\mu, \sigma$ ], x], plots can be added on top  
Plot[PDF[NormalDistribution[2, 3], x], {x, -5, 10}], Mean[], Moment[], **Expected Value:** EV[f(x)]  
**Quantile**[NormalDistribution[0, 1], 0.975] for the Z-score of  $\alpha = 0.025$  of a standard normal, aka **inverse CDF**  
**Probability**  $x \leq 1.96$ , NormalDistribution[0, 1] for the area under curve.  
**Quantile**[StudentTDistribution[4], 1-a] for the critical value of the area to the left of  $\alpha$   
**CDF**[NormalDistribution[0, 1], 1.96] for area under curve upto x=1.96 of a standard normal  
**MLE:** all components in a single step:  $\frac{d}{d\theta} \ln\left(\prod f_{\theta}(x)\right) = 0$ , look for "real solutions" of  $\theta$  or  $x$ . Note: use  $x$  for  $\sum X_i$ . **Δ** Watch out for similar chars.

**Delta Method**  
Use CLT to  $\approx Var(\hat{\theta})$  of re-parameterized  $\hat{\theta}$ . Asym  
 $V(\hat{\theta}) = Var(\hat{\theta}) \cdot (g'(\theta))^2$ , this accounts for the change in variability induced by the transformation.  
 $g'(\theta)$  must be continous differentiable  
 $\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0, \sigma^2)$   
 $\sqrt{n}(g(\hat{\theta}) - g(\theta)) \rightarrow N(0, (g'(\theta))^2 \cdot \sigma^2)$   
**Example:** We want to estimate  $\sqrt{\mu}$  instead of the mean  $\mu$ . We define  $g(x) = \sqrt{x}$ , and use  $g(\bar{X}) = \sqrt{\bar{X}}$  as an estimator for  $\sqrt{\mu}$ . By delta method, we have:  $\sqrt{n}(g(\bar{X}) - g(\mu)) \rightarrow N(0, (g'(\mu))^2 \cdot \sigma^2)$ .  
Next, given  $g(\mu) = \sqrt{\mu}$ ,  $g'(\mu) = \frac{1}{2\sqrt{\mu}}$ , then,  
 $Var(\sqrt{\bar{X}}) \approx \frac{(g'(\mu))^2 \cdot \sigma^2}{n} = \frac{\left(\frac{1}{2\sqrt{\mu}}\right)^2 \cdot \sigma^2}{n} = \frac{\sigma^2}{4n\mu}$

**Delta Method Example 2**  
1) Given  $Y_i \sim Ber(\theta^t)$ ,  $\bar{Y}_n$  is an estimator for  $\theta^t$ , and  $\hat{\theta} = \bar{Y}_n^{\frac{1}{t}}$   
2) We need  $Var(\hat{\theta})$  as  $n \rightarrow \infty$   
3)  $Var(Y_i) = \theta^t \cdot (1 - \theta^t)$ ,  $Var(\bar{Y}_n) = \frac{\theta^t(1 - \theta^t)}{n}$   
4) Delta Method: Let  $g(x) = x^{\frac{1}{t}}$   
 $g'(x) = \left(\frac{1}{t}\right) x^{\frac{1}{t}-1}$ ,  $(g'(\theta^t))^2 = \left(\frac{1}{t^2}\right) \theta^{2-2t}$   
5)  $Var(\hat{\theta}) \approx (g'(\theta^t))^2 \cdot Var(\bar{Y}_n)$ :  
6)  $Var(\hat{\theta}) = \left(\frac{\theta^t(1 - \theta^t)}{n}\right) \cdot \left(\frac{1}{t^2}\right) \cdot \theta^{2-2t}$



$\log\text{-L: } \ell(\theta; x) = \sum_{i=1}^n \left[ \ln \pi(x_i) + \frac{-\ell(\theta; x_i)}{\phi} \right]$ Canonical Families: $N(\mu, \sigma^2), Ber(p), Poi(\lambda), Exp(\lambda), \Gamma(\alpha, \beta), \text{Binom}(k, p)$	
<b>Linear Regression</b> $Y = a + bX + \varepsilon$ . Y: dependent/response, X: Independent/Predictor, a: intercept, b: slope, $\varepsilon$ : error. Goal: Fit a line $y = a + bx$ , and find $a, b$ values that minimize some loss. OLS (vertical): $\sum (Y - a - bX)^2: \hat{a} = \bar{Y} - \hat{b}\bar{X}$ $\hat{b} = \frac{Cov(X, Y)}{Var(X)} = \frac{\sigma_{XY}}{\sigma^2_X} = \frac{\overline{XY} - \bar{X} \cdot \bar{Y}}{\overline{X^2} - (\bar{X})^2}$ wolfram: OLS: fit linear(x,y),(x,y)..., Cov(X,Y): covariance(x,x..)(y,y..), this gives unbiased, must $\cdot \frac{n-1}{n}$ to get the biased $\sigma_{XY}$ , same goes for variance(x,x..) horizontal distance: $\sum \left( X - \frac{Y-a}{b} \right)^2$ $\frac{1}{b} = \frac{\sigma_{XY}}{\sigma_Y^2} \rightarrow b = \frac{\sigma_Y^2}{\sigma_{XY}}$	
<b>Min(X) &amp; Max(X)</b> $P(\max > x) = 1 - P(\max < x) = 1 - [P(X_i < x)]^n$ $P(\max < x) = \prod_{i=1}^n P(X_i < x) = [P(X_i < x)]^n$ $P(\min > x) = [P(X_i > x)]^n = [1 - P(X_i < x)]^n$ $P(\min < x) = 1 - \prod_{i=1}^n P(X_i > x) = 1 - [P(X_i > x)]^n$ $P(\min < x) = 1 - P(\min > x)$	
<b>Total Variation Distance</b> Measures max distance between two distributions $P_\theta, P_{\theta'}$ $TV(P_\theta, P_{\theta'}) = \frac{1}{2} \sum_{x \in \mathcal{E}}  p_\theta(x) - p_{\theta'}(x) $ $TV(P_\theta, P_{\theta'}) = \frac{1}{2} \int_{-\infty}^{\infty}  f_\theta(x) - f_{\theta'}(x)  dx$ E is the joint support of $P_\theta, P_{\theta'}$ . <b>Properties</b> symmetric: $TV(P_\theta, P_{\theta'}) = TV(P_{\theta'}, P_\theta)$ positive: $0 \leq TV \leq 1$ definite: if $TV(P_\theta, P_{\theta'}) = 0$ then $P_\theta = P_{\theta'}$ triangle inequality: $TV(P_\theta, P_{\theta'}) \leq TV(P_\theta, P_{\theta''}) + TV(P_{\theta''}, P_{\theta'})$ if $TV = 1$ : then $P_\theta, P_{\theta'}$ disjoint, if $TV = 0$ : then $P_\theta, P_{\theta'}$ same.	
<b>Multiple Hypothesis Testing</b> Testing Multi- $\theta$ $\uparrow$ , we must control family-wise err=P(at least 1 type_I) $\leq \alpha$ <b>Bonferroni's test:</b> Divide $\alpha$ by # of tests $m$ : Reject $H_i$ if $p_i \leq \frac{\alpha}{m}$ Example: $m = 100$ at $\alpha = 0.05$ : $\text{FWER} \approx 1 - P(\text{type}_I = 0) = 1 - (1 - 0.05)^{100} \approx 0.994$ Very restrictive for large $m \rightarrow$ lower $\beta$ $\text{FDR} = \text{fraction of type}_I$ in all sig $\leq \alpha$ <b>Bonferroni corr:</b> reject if $m \cdot p\text{-value} \leq \alpha$ , controls FWER, but very conservative. <b>Holm-Bonferroni corr:</b> sort p-values, reject smallest $p_k: p_k > \frac{\alpha}{m-k+1}$ Bonferroni-Hochberg corr: sort p-values, threshold: $p_k \leq \frac{k \cdot \alpha}{m}, k = 1, \dots, m$	

<b>Log-Likelihood</b> $\ell(x   \theta) = \ln \left( \prod_{i=1}^n f_\theta(x_i) \right) = \sum_{i=1}^n \ln(f_\theta(x_i))$	
<b>Method of Moments Estimator</b> Analyze which aspect the parameter $\theta$ represents, then choose a moment that also represent the aspect as the values that minimize some loss. $\hat{\theta}_{MM} = \bar{X}_n^k = \frac{\sum X_i^k}{n}$ by LLN: $\hat{\theta}_k \rightarrow E_0[X^k]$ by CLT: $\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0, \Gamma(\theta))$ $\Gamma(\theta)$ : asym covariance matrix. $\hat{\theta}_1 = \mu = E[X] = \frac{\sum X_i}{n} = \bar{X}_n$ , $\hat{\theta}_2 = \mu^2 + \sigma^2 = E[X^2] = \frac{\sum X_i^2}{n}$ , $\sigma^2 = E[(X - \mu)^2] = \frac{\sum (x_i - \bar{X})^2}{n}$ -biased, $\mu'_k = E[X^k], \mu_k = E[(X - \mu)^k]$	
<b>KL Divergence</b> Measures the difference between 2 distributions: $KL(P_\theta, P_{\theta'}) = \sum p_\theta(x) \log \left( \frac{p_\theta(x)}{p_{\theta'}(x)} \right)$ $KL(P_\theta, P_{\theta'}) = \int f_\theta(x) \log \left( \frac{f_\theta(x)}{f_{\theta'}(x)} \right) dx$ <b>Properties</b> not symmetric: $KL(P_\theta, P_{\theta'}) \neq KL(P_{\theta'}, P_\theta)$ not negative: $KL(P_\theta, P_{\theta'}) \geq 0$ definite: $KL(P_\theta, P_{\theta'}) = 0$ if $P_\theta = P_{\theta'}$	
<b>KS Test:</b> Goodness of fit by comparing empirical CDF w/ CDF. Do not require bins. Provide numerical $T_n = \sqrt{n} \sup  F_n(t) - F_0(t) $ , p-value: $P(Z > T_n)$ assume $H_0$ is true, reject $H_0: T_n > q_\alpha$ <b>KL Test:</b> Similar to KS test but $\theta$ are estimated. More likely to reject than KS. If $\theta$ is given, not suitable. Test if Gaussian: $\sup  F_n(t) - \Phi_{\hat{\mu}, \hat{\sigma}^2}(t) $ Geometric Distribution $Geom(p)$ = similar to $\exp()$ but discrete, # of trials until 1 <sup>st</sup> $P(X = k) = (1 - p)^{k-1} p, k = 1, 2, \dots$	
<b>Significance Tests</b> is $j^{th} X_j$ significant to Y? $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$ $\gamma_j: j^{th}$ diag coeff $(\mathbb{X}^T \mathbb{X})^{-1} (\gamma_j > 0)$ $T_n^{(j)} = \frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\sigma}^2 \gamma_j}} \sim t_{n-p}, \hat{\sigma}^2 = \frac{\ Y - \mathbb{X}\hat{\beta}\ ^2}{n-p}$ Rejection Region: $R_{j,\alpha} = \left\{ \left  T_n^{(j)} \right  > q_{\frac{\alpha}{2}}(t_{n-p}) \right\}$	

$\text{function } L(\theta; y) : y = \text{sym}(\mathbb{X}^T L[\theta; y]) \mid \mathbb{X}$ Bayes is the posterior <b>mean</b> : $\hat{\theta}^\pi = E[\theta   X] = \int \theta \cdot \pi(\theta   X) d\theta$ , or $\sum \theta \cdot \pi(\theta   X)$ , if using $\alpha$ , divide each $\pi(\theta   X)$ by $\sum \pi(\theta   X)$ aka LMS, or conditional expectation. MAP is posterior <b>mode</b> : $\hat{\theta}^{MAP} = \arg \max_{\theta} \pi(\theta   X)$ Must use normalized $\pi(\theta   X)$ , not the $\alpha \pi(\theta) L(X   \theta)$ version without $p(X)$ . Or use the mean formula of known distribution types. asym $Var(\hat{\theta}) = \frac{1}{I(\theta)}$ . LMS Error: Error=Bayes-Asym_Value, $E[Err \text{ or }   X] = 0, Cov(Err \text{ or }, \hat{\theta}) = 0$ , $Var(\theta) = Var(\hat{\theta}) + Var(Err \text{ or })$ . Conditional MSE of Bayes is the posterior Var under quad loss: $E[(\hat{\theta} - \theta)^2   X] = Var(\theta   X)$	
<b>Multivariate Linear Regression</b> $\tilde{Y} = \mathbb{X} \tilde{\beta}^* + \tilde{\varepsilon}, \tilde{\beta} \in \mathbb{R}^p, \tilde{Y} \in \mathbb{R}^n, \mathbb{X} \in \mathbb{R}^{n \times p}$ . $X^T \beta = \mu(x) = E[Y   X = x] = \int y \cdot h(y   x) dy$ <b>LSE</b> $\hat{\tilde{\beta}}^{Bayes}: \hat{\tilde{\beta}} = \arg \min \ \tilde{Y} - \mathbb{X} \tilde{\beta}\ ^2$ , solution: $\hat{\tilde{\beta}} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \tilde{Y}$ , must be full rank: Rank( $X$ ) = $p$ , and $n \geq p$ . <b>Assumptions:</b> Design Matrix $\mathbb{X}$ deterministic, full rank, $\tilde{\varepsilon} \sim N(0, \sigma^2 I_n)$ iid: $\Rightarrow Y \sim N_n(\mathbb{X} \beta^*, \sigma^2 I_n) \Rightarrow I(\beta) = \frac{1}{\sigma^2} \mathbb{X}^T \mathbb{X}$ <b>Properties</b> of LSE LSE=MLE in homoscedastic Gaussian Distro of $\hat{\tilde{\beta}} \sim N_p(\beta^*, \sigma^2 (\mathbb{X}^T \mathbb{X})^{-1})$ , is asym. Quad Risk: $E[\ \hat{\tilde{\beta}} - \beta\ ^2] = \sigma^2 \text{trace}((\mathbb{X}^T \mathbb{X})^{-1})$ Prediction Error: $E[\ Y - \mathbb{X} \hat{\tilde{\beta}}\ ^2] = \sigma^2(n - p)$ Unbiased $\hat{\sigma}^2 = \frac{\ Y - \mathbb{X} \hat{\tilde{\beta}}\ ^2}{n - p} = \frac{1}{n - p} \sum \hat{\varepsilon}^2$ $(n - p) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2, \hat{B}, \hat{\sigma}^2$ orthogonal and indie	
<b>Generalized Linear Models "GLM"</b> Link $\mu$ of response $Y$ linearly to predictor $X$ : $\mu_i = b(\theta_i), \mu = E[Y   X]$ , via Link fn: $g(\mu_i) = \theta = a + bX = X_i^T \beta$ , inverse: $\mu = g^{-1}(\theta)$ . $\hat{\beta}$ : asym normal, find $\beta$ via MLE. $N: g(\mu) = \mu = X^T \beta$ , Binom: $g(\mu) = \ln \left( \frac{\mu}{n - \mu} \right)$ $Exp: g(\mu) = -\mu^{-1}, \mu = - (X^T \beta)^{-1}$ , $\Gamma: g(\mu) = \mu^{-2}, \mu = (X^T \beta)^{-\frac{1}{2}}$ , $Poi: g(\mu) = \ln(\mu), \mu = \exp(X^T \beta)$ , $Ber: g(\mu) = \ln \left( \frac{\mu}{1 - \mu} \right), \mu = \frac{1}{1 + \exp(-X^T \beta)}$	
<b>Bonferroni's Test</b> test group $\beta$ is sig at $\text{FWER} \leq \alpha$ , non-asym. $H_0: \beta_j = 0 \forall j \in S$ where $S \subseteq \{1, \dots, p\}$ $H_1: \exists j \in S$ where $\beta_j \neq 0$ $\psi = 1 \left\{ \frac{\max \left( \left  \hat{\beta}_1 \right , \left  \hat{\beta}_2 \right , \dots \right)}{\sqrt{Var(\hat{\beta}_j)}} > q_{\frac{\alpha}{2k}} \right\}$	

$p_j(x_i)$ means the outcome is closer to 0.7 than 0.4, $\hat{p}^{MAP} = 0.7$ because MAP is the mode and $P(p = 0.7) > P(p = 0.4)$ .	
<b>Categorical Likelihood</b> Are Zodiac signs $\sim \text{Unif}(a, b)$ ? $p_0 = \left( \frac{1}{12}, \frac{1}{12}, \dots \right)$ , prob of a $X_i \in j: p_j = j(P(X = a_j))$ , # of $X_i \in j: N_j = n \{X_i = a_j\}$ , likelihood: $L_n = p_1^{N_1} \cdot p_2^{N_2} \dots = \prod_{j=1}^k p_j^{N_j}$ with $\vec{p} = (p_1, p_2, \dots, p_k)$ , $\vec{p}^{MLE}: \hat{p}_j = \frac{N_j}{n}$ where $\sum p_j = 1$ , for a single $X_i$ : $P(X = a_j) = \prod_{j=1}^k p_i^{1(X_i=a_j)}$	
<b>Chi Squared Distribution</b> $\chi_k^2$ : $\sum$ of $k$ squared $N(0, 1)$ variables, with degrees of freedom = $k$ $E[X] = k, Var(X) = 2k, f(x) = \frac{1}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} e^{-\frac{x}{2}}$ $\chi^2 \sim \Gamma\left(\frac{k}{2}, 2\right), \chi^2 \sim N(k, \sqrt{2k})$ for large $k$ . As $n$ grows, $\chi_1^2$ (composed from a single normal) 95 <sup>th</sup> quantile value decreases, and the distribution is scaled by asym $\cdot q_\alpha + \frac{\text{any } q_\alpha - \text{asym } q_\alpha}{\sqrt{n}}$	
<b>Chi Squared Test</b> $\chi^2$ test is a goodness-of-fit test (match distribution, variable independence) that groups categorical data into bins, then compare observed data freq w/ expected data freq, requires counts in each bin $> 5$ . $H_0: \vec{p} = \vec{p}^0$ vs. $H_1: \vec{p} \neq \vec{p}^0$ $T_n = n \sum \left[ \frac{\left( \hat{p}_j - p_j^0 \right)^2}{p_j^0} \right] \rightarrow \chi_{k-1}^2$ , count of $j$ in $n: p_j = \frac{n_j}{n}, k: \#$ of categories. $H_0: p \in \{\text{Bin}(k, \theta)\}$ vs $H_1: p \notin \{\text{Bin}(k, \theta)\}$ $T_n = n \sum_{j=0}^k \frac{\left( \frac{N_j}{n} - f_\theta(j) \right)^2}{f_\theta(j)} \rightarrow \chi_{(k+1)-d-1}^2, \hat{\theta}$ : MLE, $f_\theta(j) = \binom{k}{j} \hat{\theta}^j (1 - \hat{\theta})^{k-j}, d: \#$ of para, $k - d - 1$ when start at $j = 1$	
<b>Optimization</b> global extremes could be on the supports (end points), if distribution defined on closed range. solve $h'(x) = 0$ for critical points. min/max critical points $h''(x) \leq 0$ : concave, max, $h'(x) \downarrow$ $h''(x) < 0$ : str concave, global max, $h'(x) \downarrow$ $h''(x) \geq 0$ : convex, min, $h'(x) \uparrow$ $h''(x) \geq 0$ : str convex, global min, $h'(x) \uparrow$ Multivariate min/max Gradient is the vector of partial deri: $\nabla h(\theta) = \left( \frac{dh}{d\theta_1}, \dots \right)$ , Hessian $Hh(\theta)$ is the matrix of 2 <sup>nd</sup> partial deri: $X^T Hh(\theta) X \leq 0$ : concave, max, vice versa If diagonal of Hessian are positive and Hessian symmetric: convex, minimum	

$\frac{\partial^2 -t(1 - \theta t)}{t^2} \setminus \begin{matrix} n \\ t^2 \end{matrix} \setminus \begin{matrix} t \\ t^2 \end{matrix}$ $= \frac{\partial^2 -t(1 - \theta t)}{t^2}$ We can remove $n$ to show that asym variance has large $n$ already.	
<b>Normal-like Functions</b> $f_X(x) = c \exp \left( - \left( ax^2 + bx + \gamma \right) \right)$ , complete the square: $-a \left( x + \frac{\beta}{2a} \right)^2 + c$ . therefore: $\mu = -\frac{\beta}{2a}, \sigma^2 = \frac{1}{2a}$ . peak of $f_X(x) = \min$ of $-(ax^2 + bx + \gamma)$ : $\frac{d}{dx} (-ax^2 - bx - \gamma) = -2ax - \beta, -2ax - \beta = 0, x = \mu = -\frac{\beta}{2a}$ . When posterior unimodal & symmetric, such as $\sim N()$ : $\hat{\theta}^{MAP} = \hat{\theta}^{LMS} = \mu$	
<b>Covariance Matrix</b> $\Sigma = \begin{pmatrix} Cov(X, X) & Cov(X, Y) \\ Cov(Y, X) & Cov(Y, Y) \end{pmatrix}$ $= E[(X - E[X])(Y - E[Y])]$ For single vector $X$ : $Var(X) = Cov(X)$ Linear: $Cov(\mathbf{A}X + \mathbf{B}) = Cov(\mathbf{A}X) = \mathbf{A}Cov(X)\mathbf{A}^T = \mathbf{A}\Sigma\mathbf{A}^T$ <b>Multivariate Gaussian</b> vector: defined by mean vector $\mu$ and $\Sigma$ $f_X(x) = \frac{1}{\sqrt{(2\pi)^d \det \Sigma}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$	
<b>Multivariate CLT</b> $X_i \sim R^d, E[X_i] = \mu, Cov(X_i) = \Sigma$ <b>Multivariate Delta Method</b> $\sqrt{n}(g(T_n) - g(\theta)) \rightarrow N(0, \nabla g(\theta)^T \Sigma \nabla g(\theta))$	
<b>Miscellaneous Theorems</b> <b>Markov Inequality:</b> $X \geq 0$ and $a \geq 0$ , provides upper bound on $P(X \geq a)$ . $P(X \geq a) \leq \frac{E[X]}{a}$ <b>Chebyshev Inequality:</b> upper bound on $P(X$ deviates from $\mu$ by $> c$ std dev) $P( X - \mu  > c) \leq \frac{\sigma^2}{c^2}$ <b>Slutsky's Theorem:</b> Relationships between a rv that $\rightarrow$ rv and a rv that $\rightarrow$ constant. Given $T_n \rightarrow T$ and $U_n \rightarrow c$ , then, $T_n + U_n \rightarrow T + c, T_n \cdot U_n \rightarrow T \cdot c, \frac{T_n}{U_n} \rightarrow \frac{T}{c}$ Convergence: $\rightarrow a$ In distribution, $T_n$ 's distro $\rightarrow T$ 's distro. $\rightarrow p$ : In probability, $U_n \rightarrow c$ : for any error $> 0, P( U_n - c  > \epsilon) \rightarrow 0$ <b>Cochran's Theorem:</b> Relationship between $S_n, \chi^2$ , typically in ANOVA $\frac{nS_n}{\sigma^2} \sim \chi_{n-1}^2$ or $S_n \sim \frac{\sigma^2}{n} \chi_{n-1}^2$ <b>Donsker's Theorem:</b> connects empirical CDF to Brownian bridge (A Gaussian Process). if F is cont CDF: $\sqrt{n} \sup  F_n(t) - F(t)  \rightarrow \sup_{t \in [0, 1]}  B(t) $ <b>Movire Laplace Correction:</b> using discrete rv to $\approx$ conti rv via CLT. $S_n \sim \text{Bin}(n, p): P(S_n \leq k) \approx P \left( Z \leq \frac{k - np + 0.5}{\sqrt{np(1 - p)}} \right)$	