

Special Matrix	<p>Symmetric matrix A is +ve/(-ve) definite if for any vector x, $x^T A x > (<) 0$. \geq (non -ve); \leq (non +ve) definite</p> <p>If A is a sq matrix and $A^T A = I$, then A is orthogonal matrix \rightarrow rows/cols of A form orthonormal basis for $R^n \rightarrow$ all rows/cols are pairwise orthogonal ($u \cdot v = 0$) & all rows/cols are unit vector</p> <p>If AB and BA are both compatible $\rightarrow \text{tr}(AB) = \text{tr}(BA)$ & AB and BA have the same non-zero eigenvalues</p> <p>If A is non-negative definite, $\max_x \frac{x^T A x}{x^T x} = \lambda_{\max}(A)$, where λ_{\max} denotes largest eigenvalue</p> <p>If A is non-negative definite and B is non-singular (i.e. has inverse, $\det(B) \neq 0$), $\max_x \frac{x^T A x}{x^T B x} = \lambda_{\max}(AB^{-1})$</p>
Descriptive quantities	<p>Var of X is $\sigma^2 = E(X - E(X))^2 = \int (x - E(X))^2 dF(x)$. Sample var is $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$</p> <p>Let Y be another r.v. w observed sample y_1, \dots, y_n. Covariance btw X and Y is $\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$</p> <p>If $\text{Cov}(X, Y) = 0$, then X and Y are un-correlated. Sample covariance is $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})$</p> <p>Suppose X has a cts dist. The $(1 - \alpha)$-quantile (upper α-quantile), q_α is defined s.t. $F(q_\alpha) = 1 - \alpha$ or $1 - F(q_\alpha) = \alpha$</p> <p>Mean: $E(\sum_{j=1}^m c_j X_j) = \sum_{j=1}^m c_j E(X_j)$. Var: $\text{Var}(X) = E(X^2) - [E(X)]^2$</p> <p>$\text{Var}(\sum_{j=1}^m a_j X_j) = \sum_{i=1}^m \sum_{j=1}^m a_i a_j \text{Cov}(X_i, X_j)$ in general; $\text{Var}(\sum_{j=1}^m a_j X_j) = \sum_{j=1}^m a_j^2 \text{Var}(X_j)$ if X_j's are un-correlated</p> <p>Covariance: $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$. $\text{Cov}(\sum_{i=1}^m a_i X_i, \sum_{j=1}^n b_j Y_j) = \sum_{i=1}^m \sum_{j=1}^n a_i b_j \text{Cov}(X_i, Y_j)$</p> <p>Sample var: $\frac{1}{n} \sum_{i=1}^n (x_i)^2 - (\bar{X})^2$. Sample cov: $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y}) = \frac{1}{n} \sum_{i=1}^n (x_i y_i) - \bar{X}\bar{Y}$</p> <p>Let A, B be constant matrices and b, c be vectors. Let X, Y be random vectors.</p> <p>$E(AX + \mathbf{b}) = AE(X) + \mathbf{b}$ where $E(X) = (E(X_1), \dots, E(X_p))^T$</p> <p>$\text{Var}(AX + \mathbf{b}) = A\text{Var}(X)A^T$ where $\text{Var}(X) = \text{Cov}(X_i, X_j) = E(X - E(X))(X - E(X))^T$. Note $\text{Cov}(X_i, X_i) = \text{Var}(X_i)$</p> <p>$\text{Cov}(AX + \mathbf{b}, BX + \mathbf{c}) = A\text{Cov}(X, Y)B^T$ where $\text{Cov}(X, Y) = \text{Cov}(X_i, Y_j) = E(X - E(X))(Y - E(Y))^T$</p>
Matrix op on Exp and Var	<p>r.v. X has uni-variate normal dist if pdf is $f(y \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$ where μ is the mean and σ^2 is the var of the dist</p> <p>Normal dist is denoted by $N(\mu, \sigma^2)$. If $\mu = 0, \sigma^2 = 1$, it is the standard normal dist. Normal dist is symmetric about its mean</p> <p>If $X \sim N(\mu, \sigma^2)$, then $Z = (X - \mu)/\sigma \sim N(0, 1) :=$ standardisation</p>
Uni-variate Normal dist	<p>Let $\mathbf{Y} = (Y_1, \dots, Y_m)^T$ be random vector whose components are iid standard normal variables. Let A be $q \times m$ constant matrix and $\boldsymbol{\mu}$ a constant vector. Dist of $\mathbf{X} = \mathbf{A}\mathbf{Y} + \boldsymbol{\mu}$ is a q-dimensional multivariable normal dist w mean $E(\mathbf{X}) = \boldsymbol{\mu}$ and var matrix $\Sigma = \mathbf{A}\mathbf{A}^T$, and denoted by $N(\boldsymbol{\mu}, \Sigma)$</p> <p>$\Sigma = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_q) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_q) \\ \dots & \dots & \dots & \dots \\ \text{Cov}(X_q, X_1) & \text{Cov}(X_q, X_2) & \dots & \text{Var}(X_q) \end{pmatrix}$. Normal dist is uniquely determined by its mean and var matrix Σ</p> <p>pdf of $N(\boldsymbol{\mu}, \Sigma)$ is $f(\mathbf{x} \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{q/2} \Sigma ^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}$</p> <p>If X is a multivariate normal vector, then for any constant matrix B, $\mathbf{B}\mathbf{X}$ has a multivariate normal dist $N(\mathbf{B}\boldsymbol{\mu}, \mathbf{B}\Sigma\mathbf{B}^T)$</p> <p>$\forall$ constant vector c, the LC $\mathbf{c}^T \mathbf{X}$ has a univariate normal dist $N(\mathbf{c}^T \boldsymbol{\mu}, \mathbf{c}^T \Sigma \mathbf{c})$ and any component of \mathbf{X} is an univariate normal var</p> <p>If $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, then $\mathbf{Z} = \Sigma^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim N(\mathbf{0}, I)$, i.e. components of \mathbf{Z} are iid $N(0, 1)$ variables</p>
Multi-variate normal dist	<p>Let $\mathbf{Z} = (Z_1, \dots, Z_m)^T$. Suppose Z_j's are iid $N(0, 1)$ variables. Dist of $\mathbf{Z}^T \mathbf{Z} = \sum_{j=1}^m Z_j^2$ is called the χ^2-dist w d.f. m and denoted by χ_m^2</p> <p>Suppose $Z \sim N(0, 1)$, $U \sim \chi_m^2$, Z and U are indep. The dist of $Z/(\sqrt{U/m})$ is called the t-dist w d.f. m and denoted by t_m</p> <p>Suppose $U \sim \chi_m^2$, $V \sim \chi_n^2$, U and V are indep. Dist of $\frac{U/m}{V/n}$ is called the F-dist w d.f. m and n and denoted by $F_{m,n}$</p> <p>If $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$, then $W = (\mathbf{X} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim \chi_m^2$</p> <p>If $Z \sim N(0, I)$, A is symmetric and idempotent, then $\mathbf{Z}^T \mathbf{A} \mathbf{Z} \sim \chi_r^2$, where $r = \text{rk}(A) = \text{tr}(A)$</p>
χ^2 , t and F - dist	<p>Let θ be param of interest. Suppose $\theta = g(m_1, m_2, \dots)$ (1^{st} moment, 2^{nd} moment, ...). Then MME of θ is obtained by replacing the theoretical moments in the fn w the corresponding sample moments, i.e. $\hat{\theta}_{MME} = g(\hat{m}_1, \hat{m}_2, \dots)$</p> <p>For any dist, the var $\sigma^2 = m_2 - m_1^2 = E(X^2) - [E(X)]^2$, its MME is given by $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{X}^2$. MME is not unique.</p>
Mtd of Moment estimation (MME)	<p>If X has pdf $f(x, \theta)$, given the observation x_1, \dots, x_n of a random sample, the log likelihood fn of θ is defined as $\ell(\theta) = \sum_{i=1}^n \ln f(x_i, \theta)$. The MLE of θ is value of θ that maximizes the log-likelihood fn.</p> <p>E.g. let (x_1, \dots, x_n) be observation of random sample from $N(\mu, \sigma^2)$.</p> <p>Likelihood fn: $\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}$ (μ is prediction here, i.e. $x_i - \mu = y_i - \hat{y}_i$)</p> <p>The log likelihood fn of (μ, σ^2) is then $\ell(\mu, \sigma^2) = \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$</p> <p>MLE of μ and σ^2 are obtained by maximizing $\ell(\mu, \sigma^2)$ and are given by $\hat{\mu} = \bar{X}$, $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$</p>
Maximum likelihood estimation (MLE)	<p>Null hypothesis H_0 and alternative hypothesis H_1. The 2 hypotheses are mutually exclusive</p> <p>A test statistic, $T(x)$ is used. If $T(x) \geq c$ for a predetermined constant $c \rightarrow$ reject H_0. If not \rightarrow don't reject H_0</p> <p>To control the type I error rate (reject H_0 when H_0 true) at a given level α, i.e. choose c s.t. $P(T(\mathbf{X}) \geq c H_0) \leq \alpha \rightarrow$ reject H_0</p> <p>H_0 should be s.t. type I error is more serious. If still not clear, H_0 should be a well-established theory OR opp of new guess</p>

Simple Linear Regression	<p>Correlation \neq causation. Correlation coefficient can only be btw -1 and 1, i.e. $-1 \leq \rho_{xy} \leq 1$</p> <p>Pearson's correlation, the theoretical correlation is $\rho_{xy} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$, the sample correlation is $\hat{\rho}_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} = \text{corr}(X, Y)$</p> <p>If $\rho_{xy} > 0$: move in same dir. If $\rho_{xy} < 0$: move in opp dir. Both DO NOT imply a causal r/s</p> <p>If $\rho_{xy} = 0$, X and Y have no LINEAR r/s. X and Y could have other kinds of r/s (e.g. quadratic)</p> <p>Y is the <i>response variable</i>. X = <i>covariate</i>/predictor</p> <p>A simple regression model (SRM) is $Y = \beta_0 + \beta_1 X + \epsilon$, $E(Y) = \beta_0 + \beta_1 X$ is called the regression function ($E(Y X)$ more technically correct way to write it)</p> <p>W observations (x_i, y_i), observed simple LRM is $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, $i = 1, \dots, n$. Error term is diff for ea observation</p> <p>Assumptions of simple linear regression model (LRM):</p> <table border="1"> <tr> <td>1. x_i and ϵ_i are indep. Note indep \rightarrow un-correlated, \nleftrightarrow not necessary. (Unless normal var, then indep \leftrightarrow un-correlated)</td><td></td></tr> <tr> <td>2. ϵ_i have mean 0</td><td>3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$</td></tr> <tr> <td>4. ϵ_i have common variance σ^2 (Homogeneity)</td><td>5. ϵ_i have a Normal dist (Normality)</td></tr> </table> <p>LRM is then called a Normal LRM. Assumptions can be simply stated, ϵ_i iid $\sim N(0, \sigma^2)$ in addition to (1)</p> <p>From $Y = \beta_0 + \beta_1 X + \epsilon$, 1) $E(Y) = \beta_0 + \beta_1 E(X)$ 2) $\text{Cov}(X, Y) = \beta_1 \text{Var}(X)$</p>	1. x_i and ϵ_i are indep. Note indep \rightarrow un-correlated, \nleftrightarrow not necessary. (Unless normal var, then indep \leftrightarrow un-correlated)		2. ϵ_i have mean 0	3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$	4. ϵ_i have common variance σ^2 (Homogeneity)	5. ϵ_i have a Normal dist (Normality)
1. x_i and ϵ_i are indep. Note indep \rightarrow un-correlated, \nleftrightarrow not necessary. (Unless normal var, then indep \leftrightarrow un-correlated)							
2. ϵ_i have mean 0	3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$						
4. ϵ_i have common variance σ^2 (Homogeneity)	5. ϵ_i have a Normal dist (Normality)						

	<p>So $\beta_1 = \frac{Cov(X,Y)}{Var(X)} = \rho_{xy} \frac{\sigma_y}{\sigma_x}$ and $\beta_0 = \mu_y - \beta_1 \mu_x$ (exact values of β_0, β_1)</p> <p>(β_0, β_1) will then be the soln to $\min_{b_0, b_1} E(Y - b_0 - b_1 X)^2$ (i.e. error term minimised)</p> <p>So regression fn $\beta_0 + \beta_1 X$ is best linear approximation of X to Y</p>																																	
Estimation of LRM	<p>LSE of Property (β_0, β_1) minimizes $E(Y - b_0 - b_1 X)^2$ gives rise to the least square estimation (LSE), which estimate the params by minimizing the sum of squares, $Q = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = SSE$</p> <p>To find LSE of β_0, β_1, $\frac{\partial Q}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$, and $\frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i)$ and equate both to 0</p> <p>Summary: $\hat{\beta}_0 = \bar{Y} - \beta_1 \bar{X}$, $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$, where $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{Y} = \frac{1}{n} \sum_{i=1}^n y_i$</p>																																	
	<p>LSE of $\sigma^2 = E(\epsilon^2)$. $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$</p> <p>SSE (sum of square of error) = $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$ and σ^2 estimated by $\hat{\sigma}^2 = s^2 = SSE/(n-2)$</p> <p>Estimator of σ is $\hat{\sigma}$ is called the residual standard error</p>																																	
	Fitted/estimated regression fn: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$		Residuals: $e_i = y_i - \hat{y}_i$, $i = 1, \dots, n$																															
	Fitted values: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$, $i = 1, \dots, n$		Value $\hat{Y}^* = \hat{\beta}_0 + \hat{\beta}_1 X^*$ is called the predicted value of a response at X^*																															
	<p>Variation of Y is estimated by total sum of squares: $SST = SSR + SSE$</p> <p>SSE is variation explained by X; regression sum of square</p> <p>SSE is variation caused by random errors; residual sum of squares. Note $y_i = \hat{y}_i + \epsilon_i$</p>																																	
	ANOVA table for SRM																																	
	<table><tr><th>source</th><th>df (deg of freedom)</th><th>SS (sum of sq)</th><th>MS (mean of sq)</th><th>F (f ratio)</th></tr><tr><td>Regression</td><td>1</td><td>$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2$</td><td>$MSR = SSR/1$</td><td>$MSR/MSE$</td></tr><tr><td>Error</td><td>n-2</td><td>$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$</td><td>$MSE = SSE/(n-2)$</td><td></td></tr><tr><td>Total</td><td>n-1</td><td>$SST = \sum_{i=1}^n (y_i - \bar{Y})^2$</td><td></td><td></td></tr></table>					source	df (deg of freedom)	SS (sum of sq)	MS (mean of sq)	F (f ratio)	Regression	1	$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2$	$MSR = SSR/1$	MSR/MSE	Error	n-2	$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	$MSE = SSE/(n-2)$		Total	n-1	$SST = \sum_{i=1}^n (y_i - \bar{Y})^2$											
	source	df (deg of freedom)	SS (sum of sq)	MS (mean of sq)	F (f ratio)																													
	Regression	1	$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2$	$MSR = SSR/1$	MSR/MSE																													
	Error	n-2	$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	$MSE = SSE/(n-2)$																														
Total	n-1	$SST = \sum_{i=1}^n (y_i - \bar{Y})^2$																																
<p>Coefficient of determination, R^2 = proportion of variation in Y explained by X. Measures strength of correlation btw Y and X</p>																																		
<p>Adjusted R^2: Less biased estimate is given by $R_a^2 = \frac{n-1}{n-p-1} R^2 - \frac{p}{n-p-1}$ where p is num of predictors and equals 1 for simple LRM</p> <p>R^2 is strictly increasing as p increases. But R_a^2 does not necessarily incr as p incr</p>																																		
<p>$\hat{\beta}_0 = \bar{Y} - \beta_1 \bar{X}$, $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$, $SSR = \beta_1^2 \sum_{i=1}^n (x_i - \bar{X})^2$, $R^2 = \text{corr}(Y, \hat{Y})^2 = \text{corr}(Y, X)^2 = \rho_{xy}^2 = \frac{\beta_1^2 \sigma_x^2}{\sigma_y^2} = \frac{\beta_1^2 \sigma_x^2}{\beta_1^2 \sigma_x^2 + \sigma^2} = \frac{SSR}{SST}$</p> <p>Summary results explained:</p> <table><tr><th></th><th>Estimate</th><th>Std. Error</th><th>t value</th><th>Pr(> t)</th></tr><tr><td>(Intercept)</td><td>$\hat{\beta}_0$</td><td>$s(\hat{\beta}_0)$</td><td>$\frac{\hat{\beta}_0}{s(\hat{\beta}_0)} = T_{\beta_0}$</td><td>p-value for 2-sided test on $\hat{\beta}_1$</td></tr><tr><td>X</td><td>$\hat{\beta}_1$</td><td>$s(\hat{\beta}_1)$</td><td>$\frac{\hat{\beta}_1}{s(\hat{\beta}_1)} = T_{\beta_1}$</td><td>p-value for 2-sided test on $\hat{\beta}_1$</td></tr><tr><td colspan="3">Residual standard error:</td><td colspan="2">$\hat{\sigma} = \sqrt{MSE}$</td></tr><tr><td>Multiple R-squared:</td><td colspan="2">$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = \frac{MSR}{MSE}$</td><td>Adjusted R-squared:</td><td>R_a^2</td></tr><tr><td>F-statistic:</td><td colspan="2">$MSR/MSE = (T_{\beta_1})^2$</td><td>p-value</td><td>p-value of the significant F test</td></tr></table>						Estimate	Std. Error	t value	Pr(> t)	(Intercept)	$\hat{\beta}_0$	$s(\hat{\beta}_0)$	$\frac{\hat{\beta}_0}{s(\hat{\beta}_0)} = T_{\beta_0}$	p-value for 2-sided test on $\hat{\beta}_1$	X	$\hat{\beta}_1$	$s(\hat{\beta}_1)$	$\frac{\hat{\beta}_1}{s(\hat{\beta}_1)} = T_{\beta_1}$	p-value for 2-sided test on $\hat{\beta}_1$	Residual standard error:			$\hat{\sigma} = \sqrt{MSE}$		Multiple R-squared:	$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = \frac{MSR}{MSE}$		Adjusted R-squared:	R_a^2	F-statistic:	$MSR/MSE = (T_{\beta_1})^2$		p-value	p-value of the significant F test
	Estimate	Std. Error	t value	Pr(> t)																														
(Intercept)	$\hat{\beta}_0$	$s(\hat{\beta}_0)$	$\frac{\hat{\beta}_0}{s(\hat{\beta}_0)} = T_{\beta_0}$	p-value for 2-sided test on $\hat{\beta}_1$																														
X	$\hat{\beta}_1$	$s(\hat{\beta}_1)$	$\frac{\hat{\beta}_1}{s(\hat{\beta}_1)} = T_{\beta_1}$	p-value for 2-sided test on $\hat{\beta}_1$																														
Residual standard error:			$\hat{\sigma} = \sqrt{MSE}$																															
Multiple R-squared:	$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = \frac{MSR}{MSE}$		Adjusted R-squared:	R_a^2																														
F-statistic:	$MSR/MSE = (T_{\beta_1})^2$		p-value	p-value of the significant F test																														
<p>$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$. Using R^2, can explain how much of variation of Y is due to X</p> <p>$\hat{\beta}_1$ is estimated amt of change in expectation of Y when X incr by an unit amt</p> <p>Can predict Y^* w new observation X^* using $\hat{\beta}_0 + \hat{\beta}_1 X^*$</p> <p>Fitted regression function: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$. We can say Y incr/decr as X incr. When X incr by 1 unit, Y incr/decr by $\hat{\beta}_1$</p> <p>X explains about $(R^2 * 100)\%$ of the variation in Y</p>																																		
Theoretical Properties of LSE	<p>Unbiasedness of $s^2 / \hat{\sigma}^2$ as an estimator of σ^2: Thus $E s^2 = \sigma^2$</p>																																	
	<p>Properties of estimated/fitted regression fn:</p> <p>The estimated regression fn $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ is an unbiased estimator of $EY = \beta_0 + \beta_1 X$, i.e. $E(\hat{Y}) = E(Y)$</p> <p>$\text{var}(\hat{Y}) = \left[\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right] \sigma^2$. OR using $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$, $\text{var}(\hat{Y}) = \text{var}(\hat{\beta}_0) + 2X \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) + X^2 \text{var}(\hat{\beta}_1)$</p> <p>For the prediction \hat{Y} at a new value X, the prediction mean square error is: $E(Y - \hat{Y})^2 = \text{var}(\epsilon) + \text{var}(\hat{Y}) = \left[1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right] \sigma^2$</p>																																	
	<p>Dist of LS estimators: Normality condition is assumed. $\hat{\beta}_0, \hat{\beta}_1$ are LC of Y_1, \dots, Y_n</p> <p>By property of normal dist, $\hat{\beta}_0, \hat{\beta}_1$ are normally distributed: $\hat{\beta}_0 \sim N\left(\beta_0, \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right] \sigma^2\right)$, $\hat{\beta}_1 \sim N\left(\beta_1, \left[\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right]\right)$</p> <p>$(n-2) \frac{s^2}{\sigma^2}$ (where $s^2 = \hat{\sigma}^2$) $\sim \chi^2$-dist w df n-2 & is indep from $\hat{\beta}_0, \hat{\beta}_1$ (verification given later)</p>																																	
	<p>Significant test: test whether or not there is a linear regression r/s btw the response var & the covariate</p> <p>Null hypothesis, H_0: no such r/s. Alternative hypothesis, H_1: r/s exist</p> <p>test statistic is F-statistic, where $F = \frac{MSR}{MSE}$ (intuitively, if variation caused by covariate > variation caused by error \Rightarrow r/s exist)</p> <p>Under H_0, $F \sim F$-dist with df 1 and n-2, since SSR and SSE indep and each follows a χ^2-dist</p> <p>If $F >$ upper α quantile $f_{1, n-2}(\alpha)$, H_0 is rejected at level α, otherwise not rejected OR $P(F > f_{1, n-2}(\alpha)) < \alpha = P(\text{type 1 error})$</p>																																	
	<p>t-statistic for the inference on β_1</p> <p>Dist of $\hat{\beta}_1$ cannot be used directly for inference, since its var involves σ^2 which is unknown</p> <p>Let $s^2(\hat{\beta}_1) = \frac{s^2}{\sum_{i=1}^n (x_i - \bar{X})^2}$. Define $T_{\beta_1} = \frac{\hat{\beta}_1 - \beta_1}{s(\hat{\beta}_1)} = t\text{-dist} \sim t_{n-2}$.</p>																																	
	Statistical Inference for simple LRM	<p>Two-sided test for β_1: $H_0: \beta_1 = 0$ and $H_1: \beta_1 \neq 0$</p> <p>Test-statistic: $T_{0\beta_1} = \frac{\hat{\beta}_1}{s(\hat{\beta}_1)} \sim t_{n-2}$ under H_0</p> <p>Significance level α (usually taken as 0.05 or 0.01). If $T_{0\beta_1} > t_{n-2}(\alpha/2)$, then reject H_0. Otherwise, don't reject H_0</p> <p>OR p-value = $p = 2P(t_{n-2} > T_{0\beta_1})$ (T is observed value here, above was r.v). If $p < \alpha$, reject H_0; otherwise, don't reject H_0</p> <p>For simple LRM, the two-sided test for $\beta_1 = p$-value of F-statistic. F-statistic here = (t-statistic)²</p>																																
		<p>One-sided test for β_1</p> <p>1) $H_0: \beta_1 \leq 0$, $H_1: \beta_1 > 0$. (Equal sign w H_0). If $T_{0\beta_1} > t_{n-2}(\alpha)$, OR $p = P(t_{n-2} > T_{0\beta_1})$: reject H_0; otherwise don't</p> <p>2) $H_0: \beta_1 \geq 0$, $H_1: \beta_1 < 0$. If $T_{0\beta_1} < -t_{n-2}(\alpha)$, OR $p = P(t_{n-2} < T_{0\beta_1})$: reject H_0; otherwise don't</p>																																
		<p>For both 2-sided, and 1-sided test, value 0 in hypotheses can be replaced by any constant c, then test statistic = $\frac{\hat{\beta}_1 - c}{s(\hat{\beta}_1)}$</p>																																

	<p>From dist of $T_{\hat{\beta}_1}$, $100(1 - \alpha)\%$ CI for β_1 is $[\hat{\beta}_1 - t_{n-2}(\alpha/2)s(\hat{\beta}_1), \hat{\beta}_1 + t_{n-2}(\alpha/2)s(\hat{\beta}_1)]$. CI are used for 2-sided tests</p> <p>$100(1 - \alpha)\%$ lower confidence bound: $\beta_1 \geq \hat{\beta}_1 - t_{n-2}(\alpha)s(\hat{\beta}_1)$. This corresponds to 1-sided test ($H_0: \beta_1 \leq 0, H_1: \beta_1 > 0$)</p> <p>$100(1 - \alpha)\%$ upper confidence bound: $\beta_1 \leq \hat{\beta}_1 + t_{n-2}(\alpha)s(\hat{\beta}_1)$. This corresponds to 1-sided test ($H_0: \beta_1 \geq 0, H_1: \beta_1 < 0$)</p> <p>The inference on β_0</p> <p>$T_{\hat{\beta}_0} = \frac{\hat{\beta}_0 - \beta_0}{s(\hat{\beta}_0)}$, where $s^2(\hat{\beta}_0) = \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right] \sigma^2$</p> <p>Dist of $T_{\hat{\beta}_0}$ also t_{n-2}. Inference is similar to β_1. E.g. 95% CI for β_0 is $\hat{\beta}_0 \pm t_{n-2}(\alpha/2)s(\hat{\beta}_0)$</p> <p>Prediction</p> <p>CI of $E(Y)$ when predictor value is x_h: $\hat{y}_h \pm t_{n-2}(\alpha/2)\sqrt{\text{var}(\hat{Y})}$, where $\text{var}(\hat{Y}) = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right]$ where \hat{y} is the predicted value for x_h</p> <p>Prediction interval of Y_{new} when predictor value is x_h: $\hat{y}_h \pm t_{n-2}(\alpha/2)\sqrt{\text{var}(Y_{\text{new}})}$, where $\text{var}(Y_{\text{new}}) = \hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{X})^2} \right]$</p> <p>$\text{var}(Y_{\text{new}}) = \text{var}(\hat{y}_h) + \hat{\sigma}^2$, where $\hat{y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h$</p> <p>Manual computation: $\text{var}(\hat{y}_h) = \text{var}(\hat{\beta}_0) + x_h^2 \text{var}(\hat{\beta}_1) + 2x_h \text{cov}(\hat{\beta}_0, \hat{\beta}_1)$. values and $\hat{\sigma}^2$ can be found from summary results</p>
	<p>LSE as Method of moments estimation (MME): just replace theoretical moments w sample moments</p> <p>From $\beta_1 = \rho_{xy} \frac{\sigma_y}{\sigma_x}$, $\beta_0 = \mu_y - \beta_1 \mu_x \Rightarrow$ LSE is then $\hat{\beta}_1 = \hat{\rho}_{xy} \frac{\hat{\sigma}_y}{\hat{\sigma}_x}$ and $\hat{\beta}_0 = \hat{\mu}_y - \hat{\beta}_1 \hat{\mu}_x$, where a qty w a hat = sample version of that qty</p>
	<p>LSE as Maximum Likelihood estimation (MLE): Under the Normality assumption, the log-likelihood is</p> <p>$-\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} Q \Rightarrow$ MLE of σ^2 is given by $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{n-2}{n} s^2$</p> <p>Maximizing the log-likelihood to obtain the MLE of the regression coefficients is equivalent to minimizing Q</p>

Multiple Linear Regression	Matrix form of linear fns: $\sum_{i=1}^n a_i x_i = a_1 x_1 + \dots + a_n x_n = \mathbf{a}^T \mathbf{x}$, where $\mathbf{a} = (a_1, \dots, a_n)^T$, $\mathbf{x} = (x_1, \dots, x_n)^T$									
	Matrix form of quadratic fns: $\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where $\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}$									
	Matrix form of differentiation: Let $f(\mathbf{x})$ be a multivariate fn. Define $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right)^T$, then $\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$, and $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{A} \mathbf{x}$ if \mathbf{A} is symmetric, $(\mathbf{A} + \mathbf{A}^T) \mathbf{x}$ otherwise									
	In multiple LRM, find r/s btw response var Y and p predictor variables/covariates $\mathbf{X} = (X_1, \dots, X_p)$ Then multiple LRM is $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$, where ϵ is random error w mean 0 Then $EY = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ is the multiple linear regression fn. (Or more strictly speaking EY is $E(Y \mathbf{X})$) For n observations, $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$, $i = 1, \dots, n$									
	Let $\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$, $\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}$, $\boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$. Then $\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$ OR $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$. X aka design matrix									
<table><tr><td>Assumptions for multiple LRM (similar to simple LRM)</td><td>1. X_1, \dots, X_p and ϵ are indep.</td></tr><tr><td>2. ϵ_i have mean 0</td><td>3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$</td></tr><tr><td>4. ϵ_i have common variance σ^2 (Homogeneity)</td><td>5. ϵ_i have a Normal dist (Normality)</td></tr></table>					Assumptions for multiple LRM (similar to simple LRM)	1. X_1, \dots, X_p and ϵ are indep.	2. ϵ_i have mean 0	3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$	4. ϵ_i have common variance σ^2 (Homogeneity)	5. ϵ_i have a Normal dist (Normality)
Assumptions for multiple LRM (similar to simple LRM)	1. X_1, \dots, X_p and ϵ are indep.									
2. ϵ_i have mean 0	3. ϵ_i are pairwise un-correlated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$									
4. ϵ_i have common variance σ^2 (Homogeneity)	5. ϵ_i have a Normal dist (Normality)									
Let μ_Y denote mean of Y, $\boldsymbol{\mu}_X$ denote mean vector of $\mathbf{X} = (X_1, \dots, X_p)$, Σ_{XY} : covariance vector btw Y and X, Σ_{XX} : covariance matrix of X and $\boldsymbol{\beta}_1 = (\beta_1, \dots, \beta_p)$. By assumption 1, $\boldsymbol{\beta}_1 = \Sigma_{XX}^{-1} \Sigma_{XY}$, $\beta_0 = \mu_Y - \boldsymbol{\beta}_1^T \boldsymbol{\mu}_X$										
Least sq estimate for multiple LRM	Estimate $\beta_0, \beta_1, \dots, \beta_p$ by minimizing $Q = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 = \ \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\ ^2$ (norm) LSE of $\boldsymbol{\beta}$ is $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ $\hat{\sigma}^2 = \frac{SSE}{df} = \frac{\ \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\ ^2}{n-p-1} = \frac{\ \mathbf{y} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\ ^2}{n-p-1}$									
Hat matrix & properties	$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the hat matrix of X. It is the projection matrix of the linear space spanned by the cols of X $\mathbf{H} = \mathbf{H}^T$, $\mathbf{H} \mathbf{X} = \mathbf{X}$, $\mathbf{H}^2 = \mathbf{H}$, $(\mathbf{I} - \mathbf{H})^2 = \mathbf{I} - \mathbf{H}$, $\mathbf{X}^T \mathbf{H} = \mathbf{X}^T$ Residual vector: $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y}$. Vector of fitted values: $\hat{\mathbf{y}} = \mathbf{H} \mathbf{y}$ Hence $\mathbf{e}^T \mathbf{1} = 0$ (by assumption 2), $\mathbf{e}^T \mathbf{x}_j = 0$ (by assumption 1), mean of $\hat{\mathbf{y}} = \bar{y}$									
Decomposition of Sum of Squares	Matrices \mathbf{H}_T , \mathbf{H}_R and \mathbf{H}_E are all symmetric and idempotent (i.e. $\mathbf{H}^T \mathbf{H} = \mathbf{H}$) Hence $SST = \mathbf{y}^T \mathbf{H}_T \mathbf{y} = \mathbf{y}^T (\mathbf{I} - \frac{11^T}{n}) \mathbf{y}$. $SSR = \mathbf{y}^T \mathbf{H}_R \mathbf{y} = \mathbf{y}^T [\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \frac{11^T}{n}] \mathbf{y}$. $SSE = \mathbf{y}^T \mathbf{H}_E \mathbf{y} = \mathbf{y}^T [\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{y}$. $SST = SSR + SSE$									
Dist of Sum of Squares	Under assumptions of multiple regression models, (Note $\text{var}(\mathbf{A} \mathbf{Y}) = \mathbf{A} \text{var}(\mathbf{Y}) \mathbf{A}^T$, cov also similar) Under hypothesis $\beta_1 = \dots = \beta_p = 0$, SSR and SSE are indep. $\frac{SSE}{\sigma^2} \sim \chi_{n-p-1}^2$. $\frac{SSR}{\sigma^2} \sim \chi_p^2$									
ANOVA Table	Source of variation	SS	df	MS	F-statistic					
	Regression	SSR	p	MSR=SSR/p	MSR/MSE = $F_{p, n-p-1}$					
	Error	SSE	n-p-1	MSE=SSE/(n-p-1)						
	Total	SST	n-1							
Coefficient of multiple determination	coefficient of multiple determination for multiple LRM is $R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{SSR}{SST} = \text{corr}(\mathbf{y}, \hat{\mathbf{y}})^2$ Adjusted R^2 : $R_a^2 = R^2 - \frac{p}{n-p-1} (1 - R^2)$									
Multiple correlation coefficient	Correlation of response var Y w scalar covariate is measured by the Pearson's correlation coefficient The multiple correlation coefficient btw response var Y and vector \mathbf{z} of covariates $(X_1, \dots, X_p)^T$ is $\text{MCORR}(Y, \mathbf{z}) = \max_{\mathbf{a}} \text{CORR}(Y, \mathbf{a}^T \mathbf{z})$ where \mathbf{a} is a vector of constants, (i.e. $\mathbf{a}^T \mathbf{z}$ is LC of \mathbf{z}) $[\text{MCORR}(Y, \mathbf{z})]^2 = \frac{\sum_{yz} \sum_{zz}^{-1} \sum_{zy}}{\sigma_y^2}$, where $\sum_{yz} = (\text{Cov}(Y, X_1), \dots, \text{Cov}(Y, X_p))$, \sum_{zz} = variance matrix of \mathbf{z} and $\sum_{zy} = \sum_{yz}^T$, $\sigma_y^2 = \text{var}(Y)$ Let \mathbf{y} be vector of n observations of y, Z be the matrix of observed \mathbf{z} w ith row given by \mathbf{z}_i . $[\text{MCORR}(\mathbf{y}, \mathbf{Z})]^2 = \frac{\mathbf{y}^T (\mathbf{I} - \frac{11^T}{n}) \mathbf{Z} [\mathbf{Z}^T (\mathbf{I} - \frac{11^T}{n}) \mathbf{Z}]^{-1} \mathbf{Z}^T (\mathbf{I} - \frac{11^T}{n}) \mathbf{y}}{\mathbf{y}^T (\mathbf{I} - \frac{11^T}{n}) \mathbf{y}} = \frac{\text{cov}(\mathbf{Z}, \mathbf{y})^T [\text{var}(\mathbf{Z})]^{-1} \text{cov}(\mathbf{Z}, \mathbf{y})}{\text{var}(\mathbf{y})}$ In multiple LRM w vector of p covariates $\mathbf{z} = (X_1, \dots, X_p)^T$, $R^2 = [\text{MCORR}(\mathbf{y}, \mathbf{Z})]^2$									

Partial correlation coefficient	<p>Let $X(-j)$ denote the sub-matrix of the design matrix X obtained by deleting col \mathbf{x}_j of X</p> <p>Let $\tilde{\mathbf{y}}$ be the residual of \mathbf{y} regressed on $X(-j)$, $\tilde{\mathbf{x}}_j$ is residual of \mathbf{x}_j regressed on $X(-j)$</p> <p>Correlation btw $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{x}}_j$ is called partial correlation btw \mathbf{y} and \mathbf{x}_j adjusting for effects of $X(-j)$, given by</p> $\text{CORR}(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}_j) = \frac{\sum_{i=1}^n (\tilde{y}_i - \bar{\tilde{y}})(\tilde{x}_{ij} - \bar{\tilde{x}}_j)}{\sqrt{\sum_{i=1}^n (\tilde{y}_i - \bar{\tilde{y}})^2 \sum_{i=1}^n (\tilde{x}_{ij} - \bar{\tilde{x}}_j)^2}} = \frac{\mathbf{y}^T (I - H_{-j}) \mathbf{y}}{\sqrt{\mathbf{y}^T (I - H_{-j}) \mathbf{y} \mathbf{x}_j^T (I - H_{-j}) \mathbf{x}_j}}$ <p>where H_{-j} is the hat matrix of $X(-j)$</p> <p>Since residual of \mathbf{y} regressed on $X(-j)$ is part of \mathbf{y} unexplained by $X(-j)$, the squared partial correlation btw \mathbf{y} and \mathbf{x}_j adjusting for effects of $X(-j)$ is the proportion of unexplained variation of \mathbf{y} which is explained by \mathbf{x}_j</p>
Explicit expression of $\hat{\beta}_j$	Explicit expression of $\hat{\beta}_j$ can be obtained by considering the minimization of $\ \mathbf{y} - X\boldsymbol{\beta}\ ^2$: 1 st) minimize wrt $\boldsymbol{\beta}_{-j}$ with β_j fixed, then minimize wrt β_j , where $\boldsymbol{\beta}_{-j}$ is the sub-vector of $\boldsymbol{\beta}$ eliminating β_j
	Minimizing w fixed β_j : $\min_{\boldsymbol{\beta}_{-j}} \ \mathbf{y} - X\boldsymbol{\beta}\ ^2 = \min_{\boldsymbol{\beta}_{-j}} \ (I - H_{-j})\mathbf{y} - \beta_j(I - H_{-j})\mathbf{x}_j\ ^2$
	Minimizing wrt β_j : $\hat{\beta}_j = \frac{\mathbf{x}_j^T (I - H_{-j}) \mathbf{y}}{\mathbf{x}_j^T (I - H_{-j}) \mathbf{x}_j}$ $\text{var}(\hat{\beta}_j) = \frac{\sigma^2}{\mathbf{x}_j^T (I - H_{-j}) \mathbf{x}_j}$
	$\frac{\hat{\beta}_j}{sd(\hat{\beta}_j)}$ (Standardization) = $\frac{\mathbf{x}_j^T (I - H_{-j}) \mathbf{y}}{\sigma \sqrt{\mathbf{x}_j^T (I - H_{-j}) \mathbf{x}_j}}$, (which tells us which covariate contribute more to variation in y) <p>Note $R_j^2 = \text{CORR}(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}_j) = \frac{\mathbf{y}^T (I - H_{-j}) \mathbf{y}}{\sqrt{\mathbf{y}^T (I - H_{-j}) \mathbf{y}} \sqrt{\mathbf{x}_j^T (I - H_{-j}) \mathbf{x}_j}} = \text{some constant} * \frac{\hat{\beta}_j}{sd(\hat{\beta}_j)}$</p>
Example	In ANOVA table, the sum of squares (SS) generated are sequential SS, ie. the SS associated w each covariate is the one after adjusting for the effects of the covariate preceding it. SSR of $y \sim x_1 + x_2$ is sum of all sequential SS
Properties of LSE	Dist of $\hat{\boldsymbol{\beta}}$ is normal w mean = $\boldsymbol{\beta}$, var = $\sigma^2 (X^T X)^{-1}$
	$(n-p-1) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p-1}^2$ $\hat{\sigma}^2$ and $\hat{\boldsymbol{\beta}}$ are indep
	Let $\hat{\Sigma} = \hat{\sigma}^2 (X^T X)^{-1}$. Let $\hat{\sigma}_j^2$ denote the jth diagonal elems of $\hat{\Sigma}$ (estimated var of $\hat{\beta}_j$) Note $\hat{\sigma}_j^2 = c_{jj} \hat{\sigma}^2$, where c_{jj} is the jth diagonal elem of $(X^T X)^{-1}$ Then $\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_j} \sim t_{n-p-1}$; & For any constant vector \mathbf{c} , $\frac{\mathbf{c}^T (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\sqrt{\mathbf{c}^T \hat{\Sigma} \mathbf{c}}} \sim t_{n-p-1}$.
	MSR and MSE are indep, since SSR and SSE are indep
Significance F-test	<p>Hypothesis: $H_0: \beta_1 = \dots = \beta_p = 0$ vs $H_1: \beta_j \neq 0$ for at least one of $j = 1, \dots, p$.</p> <p>Test statistic: $F = \text{MSR}/\text{MSE}$. Under H_0, $F \sim F_{p, n-p-1}$</p> <p>For a significance level α, reject H_0 if $F \geq f_{p, n-p-1}(\alpha)$ or the p-value $P(F_{p, n-p-1} \geq F) < \alpha$; otherwise, do not reject H_0</p>
Wald test statistic	<p>Let $\boldsymbol{\theta}$ be a vector of parameters, $\hat{\boldsymbol{\theta}}$ be its estimator, and $\hat{\Sigma}_{\boldsymbol{\theta}}$ the estimated var matrix of $\hat{\boldsymbol{\theta}}$</p> <p>Wald statistic for testing $H_0: \boldsymbol{\theta} = 0$ is given by $\hat{\boldsymbol{\theta}}^T \hat{\Sigma}_{\boldsymbol{\theta}}^{-1} \hat{\boldsymbol{\theta}}$</p> <p>Thus Wald test statistic for significance test $H_0: \boldsymbol{\beta}_1 = 0$, where $\boldsymbol{\beta}_1 = (\beta_1, \dots, \beta_p)^T$ is given by $W = \hat{\boldsymbol{\beta}}_1^T \hat{\Sigma}_{\boldsymbol{\beta}}^{-1} \hat{\boldsymbol{\beta}}_1$, where $\hat{\boldsymbol{\beta}}_1$ is estimator of $\boldsymbol{\beta}_1$ and $\hat{\Sigma}_{\boldsymbol{\beta}}^{-1}$ is estimated variance matrix of $\hat{\boldsymbol{\beta}}_1$</p> <p>In context of multiple LRM, $F = W/p$, where p is dimension of $\boldsymbol{\beta}_1$</p>
Individual t-test	<p>Answers qn: given other vars in model, does a particular predictor have a significant effect?</p> <p>Hypotheses: $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$.</p> <p>Test statistic: $T = \frac{\hat{\beta}_j}{\hat{\sigma}_j}$ where $\hat{\sigma}_j$ is the estimated SD of $\hat{\beta}_j$. Under H_0, $T \sim t_{n-p-1}$</p> <p>For a significance level α, reject H_0 if $T \geq t_{n-p-1}(\alpha/2)$ or p-value $2P(t_{n-p-1} \geq T) < \alpha$; otherwise do not reject H_0</p> <p>- p-value better than test statistic as not only can reject H_0, if p-value very small -> evidence supporting H_1 is strong</p> <p>- 1-sided test also same way as in simple LRM</p>
Testing general linear hypothesis	<p>General linear hypothesis : $H_0: \sum_{j=0}^p c_j \beta_j = 0$</p> <p>For testing linear hypothesis, test statistic is $T = \frac{\mathbf{c}^T \hat{\boldsymbol{\beta}}}{\sqrt{\mathbf{c}^T \hat{\Sigma} \mathbf{c}}}$. Under H_0, $t \sim t_{n-p-1}$</p> <p>If only a few components of \mathbf{c} are non-zero, T can be simplified.</p> <p>E.g. $\mathbf{c} = (c_1, c_2, 0, \dots, 0)^T$, the var $\mathbf{c}^T \hat{\Sigma} \mathbf{c}$ becomes $c_1^2 \text{var}(\hat{\beta}_1) + c_2^2 \text{var}(\hat{\beta}_2) + 2c_1 c_2 \text{cov}(\hat{\beta}_1, \hat{\beta}_2)$ and $\mathbf{c}^T \hat{\boldsymbol{\beta}}$ becomes $c_1 \hat{\beta}_1 + c_2 \hat{\beta}_2$</p> <p>For a significance level α, reject H_0 if $T \geq t_{n-p-1}(\alpha/2)$ or p-value $2P(t_{n-p-1} \geq T) < \alpha$; otherwise do not reject H_0</p>
CI and confidence bound	<p>A $100(1 - \alpha)\%$ CI for β_j is $[\hat{\beta}_j - \hat{\sigma}_j t_{n-p-1}(\alpha/2), \hat{\beta}_j + \hat{\sigma}_j t_{n-p-1}(\alpha/2)]$</p> <p>A $100(1 - \alpha)\%$ CI for $\mathbf{c}^T \boldsymbol{\beta}$ is $\left[\mathbf{c}^T \hat{\boldsymbol{\beta}} - \sqrt{\mathbf{c}^T \hat{\Sigma} \mathbf{c}} * t_{n-p-1}(\alpha/2), \mathbf{c}^T \hat{\boldsymbol{\beta}} + \sqrt{\mathbf{c}^T \hat{\Sigma} \mathbf{c}} * t_{n-p-1}(\alpha/2) \right]$</p> <p>The $100(1 - \alpha)\%$ confidence bounds for β_j are $\begin{cases} \text{Upper bound: } \beta_j \leq \hat{\beta}_j + \hat{\sigma}_j t_{n-p-1}(\alpha) \\ \text{Lower bound: } \beta_j \geq \hat{\beta}_j - \hat{\sigma}_j t_{n-p-1}(\alpha) \end{cases}$</p>
Prediction	<p>Given a new observation $\mathbf{x}_0 = (1, x_{01}, \dots, x_{0p})^T$, predicted value for both $E y_0$ and y_0 is $\hat{y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}}$</p> <p>Estimated variance of fitted value is $\hat{\sigma}_F^2(\hat{y}_0) = \mathbf{x}_0^T \hat{\sigma}^2 (X^T X)^{-1} \mathbf{x}_0$</p> <p>Estimated prediction error variance is $\hat{\sigma}_P^2(\hat{y}_0) = \hat{\sigma}^2 + \hat{\sigma}_F^2(\hat{y}_0)$</p> <p>Note $y_0 = \mathbf{x}_0^T \boldsymbol{\beta} + \epsilon$. But $E y_0 = \mathbf{x}_0^T \boldsymbol{\beta}$. Thats why prediction have extra error term</p> <p>CI for $E y_0$ is $\hat{y}_0 \pm \hat{\sigma}_F^2(\hat{y}_0) t_{n-p-1}(\alpha/2)$</p> <p>Prediction interval for y_0 is $\hat{y}_0 \pm \hat{\sigma}_P^2(\hat{y}_0) t_{n-p-1}(\alpha/2)$</p>

One-Way ANOVA	<p>One-Way ANOVA model deals w only one factor, A having a levels.</p> <p>$y_{ij} = \mu + \alpha_i + \epsilon_{ij}$, $i = 1, \dots, a$, $j = 1, \dots, n_i$, where y_{ij} is the value of Y for the j^{th} member of the i^{th} grp, μ and α_i are unknown params, ϵ_{ij} are iid random errors, and n_i is num of observations in i^{th} grp</p> <p>Let μ_i be mean response to i^{th} grp. $\mu_i = \mu + \alpha_i$. For identifiability, we restrict $\sum_{i=1}^a \alpha_i = 0$.</p> <p>Then μ represents the overall mean and α_i represents effect of i^{th} treatment</p> <p>Answers qn whether factor A has any effect (i.e. any diff btw its levels) by analyzing the components of the variation in Y</p> <p>SST = SSA + SSE. Total sum of squares = sum of squares attributable to factor A and sum of squares attributable to errors</p> <p>$\text{SST} = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2$, $\text{SSA} = \sum_{i=1}^a n_i (\bar{y}_{i.} - \bar{y}_{..})^2$, $\text{SSE} = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2$</p>
---------------	--

$$\bar{y}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}, \bar{y}_{..} = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^{n_i} y_{ij}, n = \sum_{i=1}^a n_i$$

One-way ANOVA table is:

Source	df	SS	MS	F-test
A	a-1	SSA	MSA	MSA/MSE
Error	n-a	SSE	MSE	
Total	n-1	SST		

Effect of the factor is tested by the F-statistic, $F = \text{MSA}/\text{MSE}$

Under hypothesis $\alpha_1 = \dots = \alpha_a = 0$, $F \sim F_{a-1, n-a}$. p-value for F-test use $\text{pf}(F\text{-statistic}, a-1, n-a, \text{lower.tail}=\text{FALSE})$

Two-way ANOVA

Allows analysis of 2 factors at the same time & analysis of interaction of 2 factors

Analyse effect of 2 factors A and B on response variable Y. Suppose A has a levels, B has b levels

$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$, $i = 1, \dots, a$, $j = 1, \dots, b$, $k = 1, \dots, n_{ij}$, (γ : interaction btw the 2 factors) where

$$\sum_{i=1}^a \alpha_i = 0, \sum_{j=1}^b \beta_j = 0, \sum_{i=1}^a \gamma_{ij} = \sum_{j=1}^b \gamma_{ij} = 0$$

Two-way ANOVA investigates whether there is interaction btw the 2 factors, whether the main effect (ave effects of 1 factor over all the levels of the other factor) is significant

$SST = SSA + SSB + SSE + SSAB$ (sum of squares due to interaction of factor A and B)

$$SST = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} (y_{ijk} - \bar{y}_{..})^2, \quad SSA = \sum_{i=1}^a n_i (\bar{y}_{i.} - \bar{y}_{..})^2, \quad SSB = \sum_{j=1}^b n_j (\bar{y}_{.j} - \bar{y}_{..})^2, \quad SSE = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} (y_{ijk} - \bar{y}_{ij.})^2$$

$$SSAB = \sum_{i=1}^a \sum_{j=1}^b n_{ij} (\bar{y}_{ij.} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2, \quad \bar{y}_{i.} = \frac{1}{n_i} \sum_{j=1}^b \sum_{k=1}^{n_{ij}} y_{ijk}, \quad n_i = \sum_{j=1}^b n_{ij}, \quad \bar{y}_{.j} = \frac{1}{n_j} \sum_{i=1}^a \sum_{k=1}^{n_{ij}} y_{ijk}$$

$$\bar{y}_{.j} = \frac{1}{n_j} \sum_{i=1}^a \sum_{k=1}^{n_{ij}} y_{ijk}, \quad n_j = \sum_{i=1}^a n_{ij}, \quad \bar{y}_{..} = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^{n_{ij}} y_{ijk}, \quad n = \sum_{i=1}^a n_i, \quad \bar{y}_{ij.} = \frac{1}{n_{ij}} \sum_{k=1}^{n_{ij}} y_{ijk}$$

Source	df	SS	MS	F-test
A	a-1	SSA	MSA = SSA/(a-1)	MSA/MSE
B	b-1	SSB	MSB = SSB/(b-1)	MSB/MSE
AB	(a-1)(b-1)	SSAB	MSAB = SSAB/((a-1)(b-1))	MSAB/MSE
Error	n-ab	SSE	MSE = SSE/(n-ab)	
Total	n-1	SST		

Test for interaction: $F = \text{MSAB}/\text{MSE}$. Under hypothesis of no interaction $F \sim F_{(a-1)(b-1), n-ab}$

Test for main effect: $F_1 = \text{MSA}/\text{MSE}$, $F_2 = \text{MSB}/\text{MSE}$. Under hypothesis of zero main effects, $F_1 \sim F_{a-1, n-ab}$, $F_2 \sim F_{b-1, n-ab}$

Factor w zero main effect \nRightarrow factor has no effect, unless the interaction effect DNE

p-value for $F_1 = \text{pf}(F_1, a-1, n-ab, \text{lower.tail}=\text{FALSE})$. p-value for $F_2 = \text{pf}(F_2, b-1, n-ab, \text{lower.tail}=\text{FALSE})$

p-value for $F = \text{pf}(F, (a-1)(b-1), n-ab, \text{lower.tail}=\text{FALSE})$

A and B have a significant interaction effect on Y; the main effect of A is significant, but main effect of B is not

Main effect & Contrasts

If change in 1 var elicits a change in another var, then the var has an effect on the other var

Effect of a factor is the differences of the expected response it causes among its levels

Effect can be measure by contrasts. Let μ_k denote expected response at level k.

Contrast is defined as $\sum_{k=1}^a c_k \mu_k$, w $\sum_{k=1}^a c_k = 0$, i.e. $\mathbf{c}^T \mathbf{1} = 0$, where $\mathbf{c} = (c_1, \dots, c_a)$ is called a contrast vector

If factor has no effect, then all contrasts are 0, i.e. $\mu_1 = \mu_2 = \dots = \mu_a \Leftrightarrow \sum_{k=1}^a c_k \mu_k = 0$, for all \mathbf{c}

Interaction Effect

If effect of factor A is diff when factor B is fixed at diff levels, or equivalently, effect of B is diff when A is fixed at diff levels, then it is said that the 2 factors have an interaction effect

E.g. If A and B only have 2 levels, the effects of A at the 2 levels of B are: $\mu_{21} - \mu_{11}, \mu_{22} - \mu_{12}$.

If both not same \rightarrow there is interaction btw A and B.

Interaction effect is measured by $(\mu_{22} - \mu_{12}) - (\mu_{21} - \mu_{11}) = \mu_{22} - \mu_{12} - \mu_{21} + \mu_{11}$

In general, effect of factor A at a fixed level, i, of B is measured by any contrasts $\sum_{k=1}^a c_k \mu_{ki}$. If there is at least one pair (i,j) and at least one contrast \mathbf{c} , s.t. $\sum_{k=1}^a c_k \mu_{ki} \neq \sum_{k=1}^a c_k \mu_{kj}$, then A and B have interaction effect

Interaction effect is measured by interaction contrasts: (i.e. contrast of the contrast)

$$\sum_{j=1}^b d_j [\sum_{k=1}^a c_k \mu_{kj}] = \sum_{j=1}^b \sum_{k=1}^a c_k d_j \mu_{kj}, \text{ where } \sum_{k=1}^a c_k = 0, \sum_{j=1}^b d_j = 0$$

Interaction contrast vector is of the form $(c_1 d_{11}, \dots, c_1 d_{1b}, \dots, c_a d_{a1}, \dots, c_a d_{ab})$

Components of the vector have restrictions: $\sum_{k=1}^a c_k d_j = 0$, $j=1, \dots, b$; $\sum_{j=1}^b d_j = 0$, $k=1, \dots, a$; $\sum_{k=1}^a \sum_{j=1}^b c_k d_j = 0$

Among the 1st b restrictions, only b-1 are indep, among 2nd a restrictions, only a-1 are indep, and altogether there are a+b-1 indep restrictions. Thus num of indep interaction contrasts is $ab-a-b+1 = (a-1)(b-1)$

Illustration of main and interaction effects



Limitation of ANOVA

Not convenient for detailed analysis of effects

In two-way ANOVA model, if n_{ij} (group sizes) are not the same, the SSAB does not measure the interaction effect, i.e.

in SSAB, the term $\bar{y}_{ij.} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..}$ is not the estimate of an interaction contrast

ANOVA by LRM

Factor predictor can be represented by dummy vars. If factor has a levels, it can be represented by a-1 dummy vars

Dummy vars: $u_k = \begin{cases} 1, & \text{if level } k \\ 0, & \text{otherwise} \end{cases}, k = 2, \dots, a$

LRM for one-way ANOVA

Model: $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$, expressed in dummy variables u_k : $y_{ij} = \beta_0 + \beta_2 u_2 + \dots + \beta_a u_a + \epsilon_{ij}$

In matrix form: $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where $\mathbf{y} = (y_{11}, \dots, y_{1n_1}, \dots, y_{a1}, \dots, y_{an_a})^T$, $\epsilon = (\epsilon_{11}, \dots, \epsilon_{1n_1}, \dots, \epsilon_{a1}, \dots, \epsilon_{an_a})^T$, $\mathbf{X} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{1}_{n_2} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{1}_{n_a} & \mathbf{0} & \dots & \mathbf{1}_{n_a} \end{pmatrix}$.

Alternatively, model can be expressed as $y_{ij} = \begin{cases} \beta_0 + \epsilon_{1j}, & \text{if level 1 } (i = 1) \\ \beta_0 + \beta_k + \epsilon_{kj}, & \text{if level } k \text{ } (i = k), k = 2, \dots, a \end{cases}$

Regression params and mean response at ea level have relation: $\mu_1 = \beta_0, \mu_k = \beta_0 + \beta_k, k=2, \dots, a$

	<p>Hence param $\beta_k = \mu_k - \mu_1$ (which is a contrast) is the diff btw expected values at level k and level 1</p> <p>LRM for two-way ANOVA</p> <p>2 cases: (i) no repeated observations at ea level combination, i.e. $n_{ij} = 1$; (ii) $n_{ij} > 1$</p> <p>(i): only main effects can be analyzed; (ii) both main effects and interaction effects can be analyzed</p> <p>Main effect models:</p> <p>Dummy var for 2 factors: $u_i = \begin{cases} 1, & \text{if level } i \\ 0, & \text{otherwise} \end{cases}, i = 2, \dots, a, v_j = \begin{cases} 1, & \text{if level } j \\ 0, & \text{otherwise} \end{cases}, j = 2, \dots, b$</p> <p>Main effect model: $y_{ijk} = \beta_0 + \sum_{i=2}^a \alpha_i u_i + \sum_{j=2}^b \beta_j v_j + \epsilon_{ijk}$</p> <p>Expectation of the main effect model at diff levels: $E y_{ijk} = \begin{cases} \mu_{11} = \beta_0, & \text{both A and B at level 1} \\ \mu_{i1} = \beta_0 + \alpha_i, & \text{A at level } i, \text{B at level 1} \\ \mu_{1j} = \beta_0 + \beta_j, & \text{A at level 1, B at level } j \\ \mu_{ij} = \beta_0 + \alpha_i + \beta_j, & \text{A at level } i, \text{B at level } j \end{cases}$</p> <p>$\alpha_i = u_{i1} - u_{11}$: diff of effect of A at level i and 1 when B is fixed at level 1, which is the same as $u_{ij} - u_{1j}$, the diff of effects of A at level i and 1 when B is fixed at level j. i.e. diff does not depend on j</p> <p>$\beta_j = u_{1j} - u_{11}$, diff of effects of B at level j and 1 when A is fixed at level 1, which is the same as $u_{ij} - u_{i1}$, the diff of effects of B at level j and 1 when A is fixed at level i. i.e. diff does not depend on i</p> <p>Interaction model: $y_{ijk} = \beta_0 + \sum_{i=2}^a \alpha_i u_i + \sum_{j=2}^b \beta_j v_j + \sum_{i=2}^a \sum_{j=2}^b \gamma_{ij} u_i v_j + \epsilon_{ijk}$</p> <p>Expectation of the main effect model at diff levels: $E y_{ijk} = \begin{cases} \mu_{11} = \beta_0, & \text{both A and B at level 1} \\ \mu_{i1} = \beta_0 + \alpha_i, & \text{A at level } i, \text{B at level 1} \\ \mu_{1j} = \beta_0 + \beta_j, & \text{A at level 1, B at level } j \\ \mu_{ij} = \beta_0 + \alpha_i + \beta_j + \xi_{ij}, & \text{A at level } i, \text{B at level } j \end{cases}$</p> <p>$\alpha_i = \mu_{i1} - \mu_{11}$: diff of effect of A at level i and 1 when B is fixed at level 1,</p> <p>$\beta_j = \mu_{1j} - \mu_{11}$, diff of effects of B at level j and 1 when A is fixed at level 1,</p> <p>$\alpha_i + \xi_{ij} = \mu_{ij} - \mu_{11}$: diff of effect of A at level i and 1 when B is fixed at level j,</p> <p>$\beta_j + \xi_{ij} = \mu_{ij} - \mu_{i1}$, diff of effects of B at level j and 1 when A is fixed at level i,</p> <p>$\xi_{ij} = (\mu_{ij} - \mu_{i1}) - (\mu_{1j} - \mu_{11}) = \gamma_{ij}$, interaction contrast (diff of effect of A at level i and 1 when B is at level j and 1)</p> <p>So main effect caused by level i and level 1 of A is the average over all the levels of B, i.e. $\alpha_i + \frac{1}{b} \sum_{j=2}^b \xi_{ij}$</p>
Inference on interaction effect	<p>For testing whether there is a significant overall interaction effect, under the regression model,</p> <p>$H_0: \xi_{ij} = 0, i=2, \dots, a; j=2, \dots, b$; vs H_1: at least one of $\xi_{ij} \neq 0$, where ξ_{ij} are the basis interaction contrasts</p> <p>Use F-test statistic to test hypothesis</p> <p>1) F-test statistic from table produced by anova</p> <p>2) F-test statistic from computation of Wald statistic: $W = \hat{\xi}^T \hat{\Sigma}_{\xi}^{-1} \hat{\xi}$, where $\hat{\xi}$ is vector of estimated ξ_{ij}'s and $\hat{\Sigma}_{\xi}$ is the estimated covariance matrix of $\hat{\xi}$. F-statistic is $F = \frac{W}{(a-1)(b-1)}$</p> <p>Under H_0 $F \sim F_{(a-1)(b-1), n-ab}$. Reject H_0 at level α, if $F > F_{(a-1)(b-1), n-ab}(\alpha)$ or p-value $\text{pf}(F, (a-1)(b-1), n-ab, \text{lower.tail}=\text{FALSE}) < \alpha$</p> <p>For testing of particular interaction effect: perform individual test on ξ_{ij}'s and on linear combination of ξ_{ij}'s</p> <p>E.g. 1) Test whether diff btw level j and 1 of Factor B is same at level i and level 1 of Factor A, i.e. $H_0: (u_{ij} - u_{i1}) - (u_{1j} - u_{11}) = \xi_{ij} = 0$</p> <p>2) Test whether diff btw level j and l of Factor B is same at level i and level 1 of Factor A, i.e. $H_0: (u_{ij} - u_{il}) - (u_{1j} - u_{1l}) = \xi_{ij} - \xi_{il} = 0$</p> <p>3) Test whether diff btw level j and k of B is same at level i and l of Factor A, i.e. $H_0: (u_{ij} - u_{ik}) - (u_{lj} - u_{lk}) = \xi_{ij} - \xi_{ik} - \xi_{lj} + \xi_{lk} = 0$</p> <p>General rule of thumb: if subscript of u contains 1: ignore; else convert to ξ</p> <p>Since any particular interaction contrast is a LC of ξ_{ij}'s, i.e. $c^T \xi$ where ξ is the vector of ξ_{ij}. Let $\hat{\xi}$ be the vector of estimated ξ_{ij}'s and $\hat{\Sigma}_{\xi}$ is the estimated covariance matrix of $\hat{\xi}$</p> <p>Test statistic for $c^T \xi = 0$ is $T_c = \frac{c^T \hat{\xi}}{\sqrt{c^T \hat{\Sigma}_{\xi} c}}$. Under H_0: $c^T \xi = 0$, $T_c \sim t_{n-ab}$</p> <p>E.g. for a particular interaction contrast, above formula can be simplified. $H_0: \xi_{ij} - \xi_{ik} - \xi_{lj} + \xi_{lk} = 0$. Only vector $\hat{\theta} = (\hat{\xi}_{ij}, \hat{\xi}_{ik}, \hat{\xi}_{lj}, \hat{\xi}_{lk})^T$ and its covariance matrix $\hat{\Sigma}_{\theta}$ are needed. And corresponding c reduces to $b = (1, -1, -1, 1)^T$</p>
Example main effect contrast	<p>Main effect contrasts cannot be conveniently calculated using interaction model</p> <p>Main effect model has correct estimates for main-effect contrasts. But estimate $\hat{\sigma}_M^2$ from main-effect model not correct estimate of error variance. $\text{SSE}_M = \text{SSE}_I + \text{SSAB}$ (error caused by interaction + interaction effect). Instead, $\hat{\sigma}_I^2$ from interaction model has correct variance</p>
	<p>Let X_m denote design matrix of main-effect model. The estimated var matrix from main-effect model is $V.m = \hat{\sigma}_M^2 (X_m^T X_m)^{-1}$</p> <p>It shld be adjusted to $V = \hat{\sigma}_I^2 (X_m^T X_m)^{-1} = (\hat{\sigma}_I^2 / \hat{\sigma}_M^2) V.m$</p> <p>Test statistic for main-effect contrast can be computed using coefficient from main-effect model and the adjusted estimated var matrix</p>
Remarks	<p>Insignificance of the main effect does not imply it has no effect if its interaction w another factor is significant</p> <p>When interaction is significant, levels of a factor shld be compared at ea level of the other factor</p> <p>In general, inference on main effect when interaction is significant is not very relevant</p>

Family-wise type I error rate	<p>In multiple comparison problem, we investigate many contrasts which form a family.</p> <p>If a contrast is "actually" 0 but turn out to be "significant", these apparently significant contrasts are called artifacts</p> <p>To avoid claiming artifacts as significant contrasts, the family-wise type I error rate must be controlled</p> <p>Denote family of contrasts by \mathcal{C}. For any $j \in \mathcal{C}$, let T_j be test statistic & decision rule to reject H_0 be $T_j \geq c_j$ for a critical value c_j</p> <p>Family-wise type I error rate is $P(\cup_{j \in \mathcal{C}} \{T_j \geq c_j\})$</p>
General exploration	<p>Find if there is any contrasts of grp means which are statistically significant.</p> <p>For main effects: grp means are those at the level of a factor</p> <p>For interaction effect: grp means are those at the level combination of 2 factor</p> <p>So let v be the vector of the grp means. Let c denote a contrast vector of v and \mathcal{C} denote set of all possible contrasts.</p> <p>Then $H_0: c^T v = 0$ for all $c \in \mathcal{C}$. H_0 can be expressed in terms of the basis contrasts:</p> <p>Let θ be the vector of basis contrasts. H_0 is then equivalent to $a^T \theta = 0$ for any a</p> <p>In main effects: $\theta = (\alpha_2, \dots, \alpha_a)$ or $(\beta_2, \dots, \beta_b)$. In interaction effects: $\theta = (\xi_{ij})_{i=2, \dots, a; j=2, \dots, b}$</p> <p>Scheffe's solⁿ</p> <p>For an individual contrast $a^T \theta$, test statistic $T_a = \frac{a^T \hat{\theta}}{\sqrt{a^T \hat{\Sigma}_{\theta} a}}$, where $\hat{\theta}$ is estimate of θ, $\hat{\Sigma}_{\theta}$ is estimated var matrix of $\hat{\theta}$</p>

	<p>Need to find critical value c_α s.t. $P(\cup_{all a} \{ T_a \geq c_\alpha\} H_0) \leq \alpha$</p> <p>To control family-wise type I error rate at level α, Scheffe's soln is to take for all j, $c_j = c_\alpha = \sqrt{m_1 F_{m_1, m_2}(\alpha)}$, where m_1 is num of components of θ and m_2 is df of $\hat{\sigma}^2$, $F_{m_1, m_2}(\alpha)$ is the upper α quantile of the F_{m_1, m_2} dist. (c_α aka Scheffe's criterion)</p>
One-way ANOVA	<p>If multiple comparison is done directly w contrasts of level means (and not in terms of indep contrasts, i.e. Scheffe's), the test statistic for contrast $\sum_{j=1}^a c_j \mu_j$ is $T_c = \frac{\sum_{k=1}^a c_k \bar{Y}_k}{\sqrt{\hat{\sigma}^2 \sum_{k=1}^a \frac{c_k^2}{n_k}}}$ and $\max_{c \in \mathcal{R}^a} T_c^2 = \frac{1}{\hat{\sigma}^2} \max_{c \in \mathcal{R}^a} \frac{ \sum_{k=1}^a c_k \bar{Y}_k ^2}{\sum_{k=1}^a \frac{c_k^2}{n_k}}$. Note $\sum_{k=1}^a c_k \bar{Y}_k = \sum_{k=1}^a c_k (\bar{Y}_k - \bar{Y})$</p> <p>Let $\mathbf{b} = (\bar{Y}_1 - \bar{Y}, \dots, \bar{Y}_a - \bar{Y})^T$, $\mathbf{c} = (c_1, \dots, c_a)^T$, $\mathbf{D} = \text{diag}(n_1^{-1}, \dots, n_a^{-1})$. Then $\max_{c \in \mathcal{R}^a} \frac{ \sum_{k=1}^a c_k \bar{Y}_k ^2}{\sum_{k=1}^a \frac{c_k^2}{n_k}} = \max_{c \in \mathcal{R}^a} \frac{\mathbf{c}^T \mathbf{b} \mathbf{b}^T \mathbf{c}}{\mathbf{c}^T \mathbf{D} \mathbf{c}} = \mathbf{b}^T \mathbf{D}^{-1} \mathbf{b} = \text{SSA} = \sum_{i=1}^a n_i (\bar{y}_i - \bar{y})^2$</p> <p>And $\max_{c \in \mathcal{R}^a} T_c^2 = (a-1) \text{MSA} / \text{MSE} = \text{SSA} / \text{MSE}$ (since $\text{SSA} = (a-1) \text{MSA}$, $\hat{\sigma}^2 = \text{MSE}$) (where MSA / MSE is $F_{a-1, n_{ERR}}$)</p> <p>p-value = $P(F_{a-1, n-a} \geq F\text{-ratio})$</p> <p>Then the simultaneous p-value of 2-sided test is $P(\cup_{all a} \{ T_a \geq c_\alpha\}) = P(\max_{c \in \mathcal{R}^a} T_c^2 \geq c_\alpha) = P(F_{a-1, n-a} \geq T_c^2 / (a-1))$</p> <p>For 1-sided test: p-value is $(1/2)P(F_{a-1, n-a} \geq T_c^2 / (a-1))$ and critical value is $c_{2\alpha} = \sqrt{(a-1)F_{a-1, n-a}(2\alpha)}$,</p>
Individual & Simultaneous CI	<p>Simultaneous CI: CI which covers all params</p> <p>Simultaneous CI w confidence coefficient $1-\alpha$ for all j require that $P(\cap_{j=1}^m \{L_j \leq \theta_j \leq U_j\}) \geq 1-\alpha$</p> <p>Simultaneous CI for the contrasts $\sum_{k=1}^a c_k \mu_k$ is $\sum_{k=1}^a c_k \bar{Y}_k \pm \sqrt{\hat{\sigma}^2 \sum_{k=1}^a \frac{c_k^2}{n_k}} \sqrt{(a-1)F_{a-1, n_{ERR}}(\alpha)}$</p>
Approach for general exploring	<p>1) Conduct overall significance test to see if there is any effect on Y (if no -> stop)</p> <p>2) If significance test is significant, find particular significant effects (impossible to investigate all possible contrasts, so just look at rather diff \bar{Y}_k in summary data or look at estimated regression coefficients)</p>
Multiple comparison through LRM	<p>LRM for one-way ANOVA is $y_i = \beta_0 + \beta_2 u_{i2} + \dots + \beta_a u_{ia} + \epsilon_i$, $i = 1, \dots, n$ where $\beta_k = \mu_k - \mu_1$, $k \geq 2$</p> <p>For any contrast, we have $\sum_{k=1}^a c_k \mu_k = \sum_{k=2}^a c_k \mu_k - \sum_{k=2}^a c_k \mu_1 = \sum_{k=2}^a c_k (\mu_k - \mu_1) = \sum_{k=2}^a c_k \beta_k$</p> <p>Test statistic: $T_c = \frac{\sum_{k=1}^a c_k \hat{\beta}_k}{\sqrt{\tilde{\epsilon}^T \hat{\Sigma}_\beta \tilde{\epsilon}}}$, where $\tilde{\mathbf{c}} = (c_2, \dots, c_a)^T$, $\hat{\Sigma}_\beta$ is estimated covariance matrix of $(\hat{\beta}_2, \dots, \hat{\beta}_a)$</p> <p>$T_c$ is same as T_c when using level means</p>
Pairwise contrasts	<p>Only interested in certain pairwise contrasts, $\mu_k - \mu_j$, $1 \leq k < j \leq a$. Overall significance F-test not necessary</p> <p>Only need to find c_α s.t. $P\left(\max_{k,j} \frac{ \bar{Y}_k - \bar{Y}_j }{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n_k} + \frac{1}{n_j}\right)}} \geq c_\alpha\right) = \alpha$</p>
Studentized range dist	<p>Let \bar{Y}_i, $i=1, \dots, a$ be sample means of a samples w equal sizes ($n_i = n$)</p> <p>Studentized range statistic is $q_{a, n_{ERR}} = \sqrt{n}(\max \bar{Y}_i - \min \bar{Y}_i) / \sigma^2$</p> <p>Let $q_{a, n_{ERR}}$ denote upper α-quantile of the Studentized range dist, aka Tukey's criterion for pairwise comparison at level α</p>
Pairwise comparison procedure	<p>For studentized range dist, Q-statistic for $\mu_i - \mu_j$ is $Q_{ij} = \begin{cases} \frac{\sqrt{n} \bar{Y}_i - \bar{Y}_j }{\hat{\sigma}}, & \text{if } n_1 = \dots = n_a = n \\ \frac{\sqrt{\tilde{n}_{ij}} \bar{Y}_i - \bar{Y}_j }{\hat{\sigma}}, & \text{otherwise} \end{cases}$, where $\tilde{n}_{ij} = \frac{2n_i n_j}{n_i + n_j}$</p> <p>Contrast $\mu_i - \mu_j$ is significant at level α if $Q_{ij} > q_{a, n_{ERR}, \alpha}$</p>
Diff btw Q-statistic & t-statistic	<p>So $T_{ij} = Q_{ij} / \sqrt{2}$</p> <p>t-statistics can be used for pairwise comparison using Tukey's criterion, but T_{ij} must be compared w $q_{a, n_{ERR}, \alpha} / \sqrt{2}$</p> <p>i.e. contrast $\mu_i - \mu_j$ is significant at level α if $Q_{ij} > q_{a, n_{ERR}, \alpha}$ OR $T_{ij} > q_{a, n_{ERR}, \alpha} / \sqrt{2}$</p>
Equivalent form in terms of regression coefficient	<p>W conventional defn of dummy vars, the regression coefficients $\beta_i = \mu_i - \mu_1$, $i = 2, \dots, a$ OR $\mu_i - \mu_j = \beta_i - \beta_j$, $i, j > 1$</p> <p>t-statistics: $T_{ij} = \begin{cases} \frac{\hat{\beta}_i}{s.d(\hat{\beta}_i)}, & \text{if } j = 1 \\ \frac{\hat{\beta}_i - \hat{\beta}_j}{\sqrt{\text{var}(\hat{\beta}_i) + \text{var}(\hat{\beta}_j) - 2\text{cov}(\hat{\beta}_i, \hat{\beta}_j)}}, & \text{otherwise} \end{cases}$</p>
Tukey's simultaneous CI	<p>For $\beta_i = \mu_i - \mu_1$: $\hat{\beta}_i \pm \text{sd}(\hat{\beta}_i) q_{a, n_{ERR}, \alpha} / \sqrt{2}$</p> <p>For $\beta_i - \beta_j = \mu_i - \mu_j$, $i, j > 1$: $(\hat{\beta}_i - \hat{\beta}_j) \pm \sqrt{\text{var}(\hat{\beta}_i) + \text{var}(\hat{\beta}_j) - 2\text{cov}(\hat{\beta}_i, \hat{\beta}_j)} q_{a, n_{ERR}, \alpha} / \sqrt{2}$</p>
Bonferroni's Mtd	<p>If there are only k contrasts we are interest in, the overall type I error rate α for the k contrasts can be controlled by Bonferroni's Mtd: $\sum_{j=1}^k \alpha_j = \alpha$, where α_j is type I error rate for contrast j</p> <p>Each α_j can be specified, but in general just use $\alpha_j = \alpha/k$. Critical value $c_\alpha = T_{n-a}(\alpha/k*2)$. Check $T_j > c_\alpha$ for 2-sided test</p> <p>Rationale: $P(\cup_{j=1}^k \{ T_a \geq c_\alpha\}) \leq \alpha$, where $P(\cup_{j=1}^k \{ T_a \geq c_\alpha\}) \leq \sum_{j=1}^k P(T_a \geq c_{\alpha_j}) \leq \sum_{j=1}^k \alpha_j$. Thus Bonferroni mtd is a conservative mtd</p> <p>p-values: p-value for jth test is $p_j = kP(T_j \geq T_j^0)$ where T_j^0 is observed value of T_j and prob computed under dist of T_j</p> <p>If critical value for all individual test is set at T_j^0, then $\alpha_j = P(T_j \geq T_j^0)$ for $j = 1, \dots, k$ and overall type I error rate is $\alpha = k\alpha_j$</p> <p>i.e. p-value is the overall type I error rate when critical value is observed value of the statistic</p>
Summary	<p>General exploration: Scheffe's criterion</p> <p>Pairwise comparison: all 3 mtds can be used, but Tukey's Mtd more efficient than Scheffe's and Bonferroni's</p> <p>Pre-specified contrasts: Scheffe's and Bonferroni's. Bonferroni more efficient when num of pre-specified contrasts is small</p> <p>For pre-specified contrasts: can just compute both Scheffe's and Bonferroni's, smaller criteria value -> more efficient</p>

LRM w both factor & quant-itative predictor	<p>ANOCOV (analysis of covariance) model is a LRM w both factor & quantitative predictors, $y_i = \beta_0 + \sum_{j=2}^k \alpha_j u_{ij} + \gamma x_i + \epsilon_i$</p> <p>where u_j, $j = 2, \dots, k$ are dummy variables representing a factor predictor and x is a quantitative predictor, α_j is diff of Y btw type j and type 1 when X is same for both types, γ is effect of X on Y (i.e. when X change by a unit, Y has an expected change of γ)</p> <p>In general ANOCOV model could have more than 1 factor and quantitative predictor. There could be both main and interaction effects</p> <p>Traditionally, 1) ANOCOV = comparison of treatment effects of factor predictors, adjusting for effect of certain concomitant variables</p> <p>2) Comparison of regression fns: r/s btw response var and quantitative predictors is studied in diff categories.</p> <p>ANOCOV is based on adjusted SS's (which adjust for effect of concomitant vars). But adjusted SS's have complicated formulae and have similar limitations as traditional ANOVA</p> <p>W regression approach, when estimating factor, effect of concomitant vars is auto adjusted, & avoids drawbacks of traditional ANOVA</p>
---	--

Comparison of regression lines	Find whether regression lines have same intercept and whether have same slope -> Use ANOCOV model w interaction btw factor and quantitative predictor: $Y = \beta_0 + \sum_{i=2}^a \alpha_i u_i + \beta X + \sum_{i=2}^a \xi_i u_i X + \epsilon$, i.e. Category 1: $Y = \beta_0 + \beta X + \epsilon$. Category $i \geq 2$: $(\beta_0 + \alpha_i) + (\beta + \xi_i)X + \epsilon$. So just making inference on α_i , and ξ_i
LRM w non-linear predictor terms	In LRM: $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$, X need not be diff predictor variables, could be non-linear fns of predictor variables E.g. $Y = \beta_0 + \beta_1 Z + \beta_2 Z^2 + \dots + \beta_p Z^p + \epsilon$ $Y = \beta_0 + \beta_1 (1/X) + \epsilon$ (inverse model) $\log(Y) = \beta_0 + \beta_1 \log(x) + \beta_2 \log(v) + \epsilon$ (log model)
Polynomial regression models	$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i}^2 + \beta_p x_{2i}^2 + \epsilon_i$, $i = 1, \dots, n$ (model can have more than 1 predictor variables) $y_i = \beta_0 + \beta_1 (x_{1i} - \bar{x}_1) + \beta_2 (x_{2i} - \bar{x}_2) + \beta_3 (x_{1i} - \bar{x}_1)^2 + \beta_p (x_{2i} - \bar{x}_2)^2 + \epsilon_i$ (centralization to \downarrow multicollinearity of model)
Piece-wise linear models	R/s btw Y and X might be diff over diff ranges of X Hence use Auxiliary X truncated at point $X = c$: $\tilde{X} = (X - c)^+ = \begin{cases} X - c, & \text{if } X - c \geq 0 \\ 0, & \text{otherwise} \end{cases}$ Model from $y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ becomes $y_i = \beta_0 + \beta_1 X_i + \alpha_1 \tilde{X}_i + \epsilon_i$ (piece-wise linear in X) Slope of model changes from β_1 to $\beta_1 + \alpha_1$ at point $X = c$. Can continue adding more auxiliary terms if r/s diff at diff points

Cross Validation (CV)	Model M fitted using data (\mathbf{y}, \mathbf{X}) . To assess goodness of model, ideally use another data set, $(\tilde{\mathbf{y}}, \tilde{\mathbf{X}})$ and validate by prediction error $\ \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}\ ^2$. Practically, same dataset is split into training and testing data. So wastes information and is against principle of sufficiency Leave-out-one CV: Training data is n-1 observations, test data is 1 observation k-fold CV: Whole data divided into k parts, each time, one part is test data, remaining k-1 parts for training data Leave-out-one CV score: $CV = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(-i)})^2$, where $\hat{\boldsymbol{\beta}}^{(-i)}$ is the estimate of $\boldsymbol{\beta}$ by leaving out the i^{th} data point (y_i, \mathbf{x}_i) k-fold CV score: $CV_k = \frac{1}{k} \sum_{i=1}^k (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(-i)})^2$, where $(\mathbf{y}_i, \mathbf{x}_i)$ is j^{th} part of data (testing data) & $\hat{\boldsymbol{\beta}}^{(-i)}$ = estimate obtained by training data
Model selection strategy	Naive method: all subset selection. Not practical as for p predictors, there are 2^p possible models 1) Remove redundant predictors (when p is not very large) i. Fit full model, remove all predictors w p-value bigger than a certain level α ii. Fit full model, remove predictor w largest p-value which is $> \alpha$, then re-fit model w remaining predictors, repeat this until no predictor has p-value $> \alpha$ If covariates are not highly correlated, the 2 options produce the same selected model. If high correlation among covariates exists, ii. preferred 2) Forward selection (sequential procedure) Starts w null model M_0 w no predictors, then add predictors 1 at a time, choosing predictor having largest contribution to reduce the residual sum of squares Compare new model w old model by certain criterion. If new model better than old model -> continue; otherwise stop. Criterion can be any except R^2 and R_a^2 3) Backward Selection (inverted sequential method) Start w full model M_F w all predictors, then reduce model by removing predictors one at a time, choosing predictor w smallest contribution to reduce residual sum of squares to be removed. Compare reduced model w previous model by AIC. If AIC of new $<$ AIC of old -> continue; otherwise stop 4) Stepwise selection (mixture of forward & backward selection. Can be done upwards OR downwards) Upward stepwise selection; Start w null model. Add predictor to model, perform backward procedure until no predictor can be removed. Proceed to next forward step. Repeat Downward stepwise selection: Start w full model. Remove predictor from model, perform forward procedure until no predictors can be added. Proceed to next backward step. Repeat
Penalized likelihood approach	For LRM $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \epsilon$, the penalized likelihood approach select variables by minimizing $\frac{1}{n} \ \mathbf{y} - \mathbf{X}\boldsymbol{\beta}\ _2^2 + p_\lambda(\boldsymbol{\beta})$, where p_λ is the penalty function, λ is the penalty parameter whose value is to be chosen Procedure: specify sequence of λ values, at each value, carry out the penalized minimization, which yields a model w certain selected variables. Selection criteria is used to select the model If purpose to obtain model for prediction -> CV. If purpose to identify important variables -> EBIC
Common penalty functions	1) LASSO penalty: $p_\lambda(\boldsymbol{\beta}) = \lambda \sum_{j=1}^p \beta_j $ 2) Adaptive LASSO penalty: $\lambda \sum_{j=1}^p w_j \beta_j $, where w_j is taken as $ \hat{\beta}_j ^{-1}$. If $p \ll n$, $\hat{\beta}_j$ being the OLS (ordinary least square) estimator in multiple LRM. If p is close to or $> n$, $\hat{\beta}_j$ is the OLS estimate is the marginal LRM SCAD penalty: $p_\lambda(\beta) = \beta $ for $ \beta $ near 0, and equals a constant C for large $ \beta $, the two parts are connected by a smooth function MCP penalty: for large $ \beta $, it is a constant. Smoothly decreases to 0 w $p_\lambda(\beta) = \beta $ as its asymptote when $ \beta $ approaches 0
Rationale of penalized likelihood approach	E.g. LASSO (least absolute shrinkage and selection operator) estimates $\boldsymbol{\beta}$ by minimizing $\frac{1}{n} \ \mathbf{y} - \mathbf{X}\boldsymbol{\beta}\ _2^2 + \lambda \sum_{j=1}^p \beta_j $ If $\lambda = 0$, LASSO estimator is same as LSE. If $\lambda = \infty$, all components of $\boldsymbol{\beta}$ are estimated as 0 For a certain nonzero λ , some of the components will be estimated as nonzero, and others 0. The nonzero ones are shrunk version of LSE Variables w nonzero estimated coefficients are the selected variables

Model diagnostics	Assumptions made for model might not be true, leading to discrepancies. There are 2 types of discrepancies – systematic and local	
	Fitted values are $\hat{y}_i = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$, $i = 1, \dots, n$. Hat matrix is $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ Hat values (i^{th} diagonal elem of \mathbf{H}), $h_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i$ = hat value of i^{th} observation, where \mathbf{x}_i^T is the i^{th} row of \mathbf{X}	
	Partitioning \mathbf{X} as $\mathbf{X} = (\mathbf{1} \ \mathbf{Z})$, $(\mathbf{X}^T \mathbf{X})^{-1} = \begin{pmatrix} n & \mathbf{1}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{1} & \mathbf{Z}^T \mathbf{Z} \end{pmatrix}^{-1} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ $A_{22} = \left[\mathbf{Z}^T \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^T}{n} \right) \mathbf{Z} \right]^{-1}$ = sample covariance matrix, $A_{11} = \frac{1}{n} + \frac{\mathbf{1}^T \mathbf{Z}}{n} A_{22} \frac{\mathbf{Z}^T \mathbf{1}}{n}$, $A_{12} = -\frac{1}{n} \mathbf{1}^T \mathbf{Z} A_{22}$, $A_{21} = A_{12}^T$ Let i^{th} row vector \mathbf{x}_i of \mathbf{X} be $(1, \mathbf{z}_i)^T$ and $\bar{\mathbf{z}} = \frac{\mathbf{Z}^T \mathbf{1}}{n}$. Then $A_{11} = \frac{1}{n} + \bar{\mathbf{z}}^T A_{22} \bar{\mathbf{z}}$ Then $h_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i = \frac{1}{n} + (\mathbf{z}_i - \bar{\mathbf{z}})^T A_{22} (\mathbf{z}_i - \bar{\mathbf{z}}) \bar{\mathbf{z}} = \frac{1}{n} + \text{Mahalanobis dist of } i^{\text{th}} \text{ observation in } \mathbf{x}\text{-space to centroid of sample}$	
	Pearson's residuals are $e_i = y_i - \hat{y}_i$, $i = 1, \dots, n$. Let $\mathbf{e} = (e_1, \dots, e_n)^T$. Then $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H})\mathbf{y}$. $\mathbf{E}\mathbf{e} = \mathbf{0}$, $\text{var}(\mathbf{e}) = \sigma^2(\mathbf{I} - \mathbf{H})$. $\text{var}(e_i) = \sigma^2(1 - h_{ii})$ Under normality assumption, $\mathbf{e} \sim \mathbf{N}(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H}))$	
	Studentized residuals, $r_i' = \frac{e_i}{sd(e_i)} = \frac{y_i - \hat{y}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$ Studentized deleted residuals, $r_i^* = \frac{y_i - \hat{y}_{(i)}}{sd(y_i - \hat{y}_{(i)})} = \frac{y_i - \hat{y}_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}}$, where $\hat{y}_{(i)}$ is the predicted value by fitted model w i^{th} observation deleted, $\hat{\sigma}_{(i)}$ is the counter part of $\hat{\sigma}$ when i^{th} observation is deleted Cook's distance, $d_i = (\hat{\boldsymbol{\beta}}_{(i)} - \hat{\boldsymbol{\beta}})^T (\mathbf{X}^T \mathbf{X}) (\hat{\boldsymbol{\beta}}_{(i)} - \hat{\boldsymbol{\beta}}) / (p \hat{\sigma}^2)$, where $\hat{\boldsymbol{\beta}}_{(i)}$ is estimate of $\boldsymbol{\beta}$ when i^{th} observation is removed from data Variance Inflation Factor (VIF). Let R_k^2 be the coefficient of determination of model $\mathbf{X}_k = \beta_0 + \sum_{j \neq k} \beta_j \mathbf{X}_j + \epsilon$. VIF of \mathbf{X}_k is $\text{VIF}_k = \frac{1}{1 - R_k^2}$	
Systematic discrepancies	Caused by regression fn not linear, error terms don't have constant variance, error terms not indep, error terms don't have normal dist, impt predictors are omitted from model. Use residual plots to check for systematic error	
	If no discrepancies, residuals would appear like iid random errors w mean 0. In any residual plot, points are scattered evenly within a horizontal band around 0	
Check non-linearity	Plot Pearson's residual against fitted values Plot Pearson's residual against predictor variables Scatter plot of response against predictor variables	If any of the plots show a non-linear trend -> regression fn is not linear For residual vs predictors: no trend -> no obvious discrepancy in regression fn
Check homogeneity	Plot Pearson's residual against fitted values Plot Pearson's residual against predictor variables	If vertical range of residuals have obvious change along x-axis -> variances are not constant / not homogeneous

Check independence		Plot residual against time/space	If indep -> should be constant horizontal trend																
Check normality	Plot of studentized residuals through dist plot (box plot, histogram) OR normal probability plot of residuals OR QQ plot		<div>Skewed to the right</div> <div>Symetric but heavy tailed</div>																
	If normality holds, points in QQ plot shld fall on straight line $y = x$																		
Heavy tailed pattern	If dist of r.v. Y is skewed to the right (positive skew) relative to normal dist, then $P(Y \leq c) \leq P(Z \leq c)$ for all c Let y_q and z_q denote q-quantile of Y and Z, then $P(Y \leq y_q) = P(Z \leq z_q) \geq P(Y \leq z_q)$. Then $P(Y \leq y_q) \geq P(Y \leq z_q)$ Hence $y_q \geq z_q$ for all q. In Q-Q plot where y_q is plotted against z_q , the point (z_q, y_q) is above the point (z_q, z_q)		<div>skewed left</div> <div>skewed right</div>																
Check missing predictors		Plot residual against other predictors not included in the model	If 1 of the plot show a trend -> that predictor is missing																
Outliers	Leverage: whether point is far away from major cluster in x-space Since $\sum_i h_{ii} = \text{Tr}(H) = p$. Point is high leverage if $h_{ii} > \frac{2p}{n}$																		
	Consistency: whether point is consistent in terms of fitting in the (x,y)-space Studentized deletion residuals are the standardized prediction errors, Find points w highest $ r_i^* $ values -> possible outliers																		
	Influence: whether point highly affects fitting of model Find points w highest Cook's distance -> possible outliers																		
			<table><tr><td></td><td>a</td><td>b</td><td>c</td></tr><tr><td>Leverage</td><td>low</td><td>high</td><td>high</td></tr><tr><td>Consistency</td><td>No</td><td>Yes</td><td>No</td></tr><tr><td>Influence</td><td>low</td><td>low</td><td>high</td></tr></table>			a	b	c	Leverage	low	high	high	Consistency	No	Yes	No	Influence	low	low
	a	b	c																
Leverage	low	high	high																
Consistency	No	Yes	No																
Influence	low	low	high																
Assessment of outliers	Informal test of outliers: normal probability plots of studentized deletion residual, leverage h_i and Cook's distance d_i Formal test: To assess (y_i, \mathbf{x}_i) , introduce the dummy variable $u = \begin{cases} 1, & \text{for } i^{\text{th}} \text{ unit} \\ 0, & \text{otherwise} \end{cases}$. Significance of coefficient of u in the linear predictor indicates i^{th} point is an outlier.																		

1. LRM w unequal variances	When ϵ_i 's don't have common variance, an unequal variance model is considered: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ has the variance matrix $\Sigma = \sigma^2 \begin{pmatrix} w_1^{-1} & 0 & \dots & 0 \\ 0 & w_2^{-1} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & w_n^{-1} \end{pmatrix}$, where w_i 's are unequal weights (if common variance: $\sigma^2 I$) If w_i 's are known, the unequal variance model can be transformed into an equal variance model. Let $\mathbf{W} = \begin{pmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & w_n \end{pmatrix}$. Multiply $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ by $\mathbf{W}^{1/2}$, then $\mathbf{W}^{1/2}\mathbf{y} = \mathbf{W}^{1/2}\mathbf{X}\boldsymbol{\beta} + \mathbf{W}^{1/2}\boldsymbol{\epsilon}$. Let $\tilde{\mathbf{y}} = \mathbf{W}^{1/2}\mathbf{y}$, $\tilde{\mathbf{X}} = \mathbf{W}^{1/2}\mathbf{X}$, $\tilde{\boldsymbol{\epsilon}} = \mathbf{W}^{1/2}\boldsymbol{\epsilon}$ Then $\text{var}(\tilde{\boldsymbol{\epsilon}}) = \mathbf{W}^{1/2}\Sigma\mathbf{W}^{1/2} = \sigma^2 I$. And model $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \tilde{\boldsymbol{\epsilon}}$ is a constant variance model Minimizing $\ \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\ _2^2$, we obtain the estimate of $\boldsymbol{\beta}$ as $\hat{\boldsymbol{\beta}}_W = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} = \text{weighted LSE (WLSE)}$ $\ \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\ _2^2$ can be expressed explicitly as $\sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$. The weight w_i reflect the relative importance of $(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$ in the estimation. The larger the variance of i^{th} term, the smaller the corresponding weight, since weight is inversely proportional to the variance For WLSE $\hat{\boldsymbol{\beta}}_W$, $E\hat{\boldsymbol{\beta}}_W = \boldsymbol{\beta}$, $\text{var}(\hat{\boldsymbol{\beta}}_W) = \sigma^2 (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$. And $\hat{\boldsymbol{\beta}}_W \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1})$ σ^2 is estimated as $\hat{\sigma}^2 = \frac{1}{n-p-1} \ \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}_W\ _2^2 = \frac{1}{n-p-1} \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_W)^2$. Inference on $\boldsymbol{\beta}$ is made in same way as in normal LRM
Estimation of unknown weights	If no replicates of predictor values, Weights can be estimated by the following procedure: 1. Fit regression model by unweighted least squares and obtain residuals \mathbf{r} and fitted values $\hat{\mathbf{y}}$ 2. Regress log of absolute residual on log of fitted values, e.g. $\ln r_i = \alpha_0 + \alpha_1 \ln \hat{y}_i + e_i$. If smallest $\hat{y}_i < 0$, replace \hat{y}_i by $\hat{y}_i - \min \hat{y}_i + c$, for some positive constant c 3. Weights are estimated as $w_i = \hat{y}_i^{-2\hat{\alpha}_1}$ OR $e^{-2\hat{\alpha}_0 \hat{y}_i^{-2\hat{\alpha}_1}}$. Constant $e^{-2\hat{\alpha}_0}$ don't really matter If there are replicates of predictor values, sample variance can be used in estimation of weights 0. Naive method: weight = $1/s^2$, s = sample SD 1. Estimate variance as a function of the mean by using the regression model: $\ln s_i = \alpha_0 + \alpha_1 \ln \bar{y}_i + e_i$, where s_i is sample sd, and \bar{y}_i is sample mean 2. The weight for the i^{th} predictor value is $w_i = \bar{y}_i ^{-2\hat{\alpha}_1}$
2. Multi-collinearity & its effects	Correlation among predictor variables such as pairwise correlation, linear dependence of one predictor on another predictor,... If 1 predictor is perfectly linearly dependent on the other predictors, $\mathbf{X}^T \mathbf{X}$ would be singular, i.e. non-invertible Practically, although perfect linear dependence will not occur, high multicollinearity can cause $\mathbf{X}^T \mathbf{X}$ to be nearly singular (i.e. having large condition number $\lambda_{\max}/\lambda_{\min}$), which renders the LSE extremely unstable Serious multicollinearity greatly increases variance of LSE and make LSE inaccurate and useless
Informal Diagnostics for multi-collinearity	Large changes in estimated regression coefficients when a predictor is added/deleted, or an observation is altered or deleted Nonsignificant results in individual test on the regression coefficients for known important predictor variables Estimated regression coefficients w an algebraic sign that is opp of that expected from theoretical consideration or prior experience Large coefficients of sample correlation btw pairs of predictors in the correlation matrix (linear trend, correlation > 0.6)
1) Formal diagnostic - VIF	\mathbf{X}_{-j} being the submatrix of \mathbf{X} w/o its j^{th} column and R_j^2 is the coefficient of multiple determination of \mathbf{x}_j is regressed on \mathbf{X}_{-j} If \mathbf{x}_j is uncorrelated w \mathbf{X}_{-j} -> $R_j^2 = 0$ and the variance of $\hat{\beta}_j = \sigma^2 / \text{SST}_j$ If \mathbf{x}_j is correlated w \mathbf{X}_{-j} -> the variance is inflated by a factor $\text{VIF}_j = (1 - R_j^2)^{-1}$ (Variance inflation factor)
2) Ridge regression	Since consequence of multicollinearity is $\mathbf{X}^T \mathbf{X}$ nearly singular, to remedy: add diagonal matrix to $\mathbf{X}^T \mathbf{X}$, i.e. $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ (which is invertible) Ridge regression estimator $\hat{\boldsymbol{\beta}}_{\text{RIG}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$, where $\lambda > 0$ is a parameter to be chosen $\hat{\boldsymbol{\beta}}_{\text{RIG}}$ is also the minimizer of the penalized sum of squares: $\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \ \mathbf{y} - \mathbf{X}\boldsymbol{\beta}\ ^2 + \lambda \ \boldsymbol{\beta}\ ^2$ $E\hat{\boldsymbol{\beta}}_{\text{RIG}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta} - \lambda (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\beta}$ -> estimator is biased. Bias is $\boldsymbol{\beta} - E\hat{\boldsymbol{\beta}}_{\text{RIG}} = \lambda (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \boldsymbol{\beta}$

Let Q be the orthogonal matrix s.t. $X^T X = Q \Delta Q^T$, where $\Delta = \text{Diag}(\tau_0, \tau_1, \dots, \tau_p)$. Thus $\beta - E\hat{\beta}_{RIG} = Q \begin{pmatrix} \frac{\lambda}{\tau_0 + \lambda} & 0 & \dots & 0 \\ 0 & \frac{\lambda}{\tau_1 + \lambda} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \frac{\lambda}{\tau_p + \lambda} \end{pmatrix} Q^T \beta$

Let $\tilde{\beta}_j$ denote the j^{th} component of $Q^T \beta$. We have $\|\beta - E\hat{\beta}_{RIG}\|^2 = \sum_{j=0}^p \left(\frac{\lambda}{\tau_j + \lambda}\right)^2 \tilde{\beta}_j^2$. Thus bias increase as λ increase

Thus $\text{tr}(\text{var}(\hat{\beta}_{RIG})) = \sigma^2 \sum_{j=0}^p \frac{\tau_j}{(\tau_j + \lambda)^2}$. The variance decrease as λ increase. All other cov = 0

Sum of mean squares errors of β is given by $MSE = \sum_{j=0}^p MSE(\hat{\beta}_j) = \|\beta - E\hat{\beta}_{RIG}\|^2 + \text{TR}(\text{Var}(\hat{\beta}_{RIG}))$.

Need to get balance btw bias and variance by minimizing MSE

MSE cannot be readily used as a criterion as it involves unknowns β and σ^2 . Can select λ by Cross validation (CV)

Let $\hat{\beta}_{-i}(\lambda)$ be the ridge regression estimate of β w parameter λ by deleting the i^{th} observations. The CV score is given by $CV(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - x_{(i)}^T \hat{\beta}_{-i}(\lambda))^2$, where $x_{(i)}^T$ is the i^{th} row vector of the design matrix X

Best λ is the minimizer of $CV(\lambda)$

Note ridge regression mainly used for building model for prediction. Cannot be used to assess importance or effects of predictor variables. The estimates from model cannot be used to construct CI or conduct hypo testing

If need to make inference on effects of predictors -> use strategy of removing predictors w large VIF

3. Non-normality (remedy w variable transformation)

If normal dist -> variance don't depend on mean (i.e. constant variance in regression models)

The violation of normality usually goes tgt w violation of constancy of variance

A variance stabilization transformation can help to rectify both discrepancy in normality and constancy of variance

If r/s btw variance and mean is known -> the desired variance stabilization transformation can be derived

To find the transformation -> use Box-Cox transformation

If variance σ^2 depend on mean μ , i.e. $\sigma^2 = V(\mu)$, a transformation can be found s.t. variance of transformed variable is \approx indep of mean

Let $h(Y)$ be a transformation. Use Taylor series to expand $h(Y)$ at μ , as $h(Y) \approx h(\mu) + h'(\mu)(Y - \mu)$, where $h(\mu)$ is a constant

Treating $h(\mu)$ as mean of $h(Y)$, $\text{var}(h(Y)) = E[h(Y) - h(\mu)]^2 \approx [h'(\mu)]^2 E[(Y - \mu)^2] = [h'(\mu)]^2 V(\mu)$

WLOG, setting $[h'(\mu)]^2 V(\mu) = 1$, $h'(\mu) = \frac{1}{\sqrt{V(\mu)}} \rightarrow h(\mu) = \int \frac{1}{\sqrt{V(\mu)}} d\mu$, aka Variance Stabilization transformation

Box-Cox transformation

If Variable transformation don't work?

For certain non-normal cts r.v., transformation cannot be determined solely by data type

Box-Cox transformation: $h(Y) = \frac{Y^\lambda - 1}{\lambda}$. λ can be determined by data

When $\lambda = 0$, the Box-Cox transformation is given by $h(Y) = \ln(Y)$, since $\lim_{\lambda \rightarrow 0} \frac{Y^\lambda - 1}{\lambda} = \ln(Y)$

Let $\sigma_i \propto \mu_i^\alpha$, where σ_i and μ_i are the SD and mean of i^{th} treatment effect. The λ in Box-Cox transformation is determined by variance stabilization transformation

$\sigma_i \propto \mu_i^\alpha$	α (just use closest half int, e.g. $\alpha = 0.04 \rightarrow$ use 0.05)	$\lambda = 1 - \alpha$	Transformation
$\sigma_i \propto \mu_i^3$	3	-2	reciprocal squared
$\sigma_i \propto \mu_i^2$	2	-1	reciprocal
$\sigma_i \propto \mu_i^{3/2}$	3/2	-1/2	reciprocal sqrt
$\sigma_i \propto \mu_i$	1	0	log
$\sigma_i \propto \mu_i^{1/2}$	1/2	1/2	sqrt
$\sigma_i \propto \text{constant}$	0	1	no transformation
$\sigma_i \propto \mu_i^{-1/2}$	-1/2	3/2	3/2 power
$\sigma_i \propto \mu_i^{-1}$	-1	2	square

Determination of α

Method 1: If observations are grouped, for each group, compute s_i and \bar{y}_i . Fit regression model $\ln s_i = \beta_0 + \beta_1 \ln \bar{y}_i + e_i$

If observations not group, fit $\ln |r_i| = \beta_0 + \beta_1 \ln \hat{y}_i + e_i$

For both, $\alpha = \text{estimate of } \beta_1$

Method 2: Only for grouped observations. Select a few α values, say α_k , $k = 1, \dots, K$. For each k , compute $R_k = \max_i \frac{s_i}{\bar{y}_i^{\alpha_k}} / \min_i \frac{s_i}{\bar{y}_i^{\alpha_k}}$

Select α_k with smallest R_k

Direct determination of λ

For grouped observations: select a few λ values, say λ_k , $k = 1, \dots, K$

For each k , make the transformation $y_{ij} \rightarrow y_{ij}^{\lambda_k}$. With the transformed data, compute $s_{\lambda_k i}^2$, $i = 1, \dots, g$

Select λ_k s.t. $\max_i s_{\lambda_k i}^2 / \min_i s_{\lambda_k i}^2$ is closest to 1

For non-grouped observations: select a few λ values, say λ_k , $k = 1, \dots, K$

For each k , make the transformation $y_{ij} \rightarrow y_{ij}^{\lambda_k}$. Analyze the regression models w $y_i^{\lambda_k}$ as response variable.

Select λ_k s.t. $MSE(\lambda_k)$ is smallest