

Matrix properties:

- $A = \begin{bmatrix} 1 & -1 \\ 2 & 3 \\ 3 & 4 \end{bmatrix}$
- $A^T = \begin{bmatrix} 1 & 2 & 3 \\ -1 & 3 & 4 \end{bmatrix}$
- $A^T A = \begin{bmatrix} 1 & -1 \\ 2 & 3 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ -1 & 3 & 4 \end{bmatrix} = \begin{bmatrix} 10 & 10 & 10 \\ 10 & 22 & 22 \\ 10 & 22 & 22 \end{bmatrix}$
- $A^T A^{-1} = I$
- $A^{-1} = \frac{1}{10} \begin{bmatrix} 4 & -3 \\ -3 & 1 \end{bmatrix} = \begin{bmatrix} 0.4 & -0.3 \\ -0.3 & 0.1 \end{bmatrix}$
- $AB = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31} & a_{11}b_{12} + a_{12}b_{22} + a_{13}b_{32} \\ a_{21}b_{11} + a_{22}b_{21} + a_{23}b_{31} & a_{21}b_{12} + a_{22}b_{22} + a_{23}b_{32} \end{bmatrix}$
- $(AB)^T = B^T A^T$
- $\text{Symmetry: } A = A^T$
- $\text{Diagonality: } \text{Off-diagonal elements are zero}$
- $\text{Identity: } \text{Diagonal elements 1, others zero}$
- $\text{Scalar matrix: } \text{Diagonal elements some scalar, others zero}$
- $\text{rank}(A) = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} = n$ Rules for matrix ranks:
 - $\text{rank}(A) = 0 \iff A = 0 \text{ matrix}$
 - $\text{if } A \neq 0, 1 \leq \text{rank}(A) \leq \text{rank}(c)$
 - $\text{rank}(Ac) = \text{rank}(A)$
 - $\text{rank}(A+B) = \text{rank}(AA) = \text{rank}(AA) \text{ Hat matrix}$
 - $\text{rank}(AB) \leq \text{rank}(A) + \text{rank}(B)$ $\leq \min(\text{rank}(A), \text{rank}(B))$
 - $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$ when two matrices multiplied
- $A^{-1} A = A A^{-1} = I$
- $\text{Idempotency: } AA = A$

Variance - stabilizing terms of $Y \rightarrow$ if var \uparrow with \hat{Y}_i , \downarrow with the latter and vice versa

Multicollinearity

- Two or more predictors highly correlated w/ each other \rightarrow If perfectly uncorrelated, marginal and partial relationships w/ outcome variable will be the same
- Effects of multicollinearity:

predictors	effect on $\hat{\beta}$	effect on $SE(\hat{\beta})$
uncorrelated	unchanged	unchanged or reduced
highly corr.	change, can be large	Increase, potentially large one
- What's going on w/ coefficient estimates?
 - In an MLR model, they are partial coefficients (holding other predictors constant / controlling)
 - In a SLR model, they are marginal coefficients (ignoring any other potential predictors)
 - In MLR, marginals and partials can even have opposite signs because you may be holding constant a predictor that naturally correlates w/ the one in question (i.e. BSA and weight). Model may not have much info about that rare situation.
- What's going on with SEs ?
 - Intuitively, if X_1 and X_2 are highly correlated and we hold one constant, there is little variation in the other. Thus, the data contain little information about the partial effect of one on the outcome variable holding the other constant. The SEs of the partial coeffs will be large, reflecting that there is little information about the situation in the data.
 - Mathematically, if X is not full rank, $X^T X$ isn't either, meaning $\det(X^T X) = 0$, so inverse does not exist. If predictors close to linear dependence, $\det(X^T X)$ close to 0 and $\sqrt{\det(X^T X)}$ large. Since $\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$, the SEs will be high.
 - Off estimates, SE and p-values can all be dramatically different if other predictors in model when multicollinearity is present.

* Note: $Y^T Y = \sum Y_i^2$

$X^T X = \begin{bmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{bmatrix}$

$X^T Y = \begin{bmatrix} \sum Y_i \\ \sum X_i Y_i \end{bmatrix}$

Inverse exists when:

- Full rank, $\det \neq 0$
- Rank $= p$

$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ $A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$

Basic Results

- $A+B = B+A$
- $(AB)+C = A+(B+C)$
- $(AB)C = A(BC)$
- $C(A+B) = CA+CB$
- $K(A+B) = KA+KB$
- $(A')' = A$
- $(A+B)' = A'+B'$
- $(AB)' = B'A'$
- $(ABC)' = C'B'A'$
- $(AB)^{-1} = B^{-1}A^{-1}$
- $(ABC)^{-1} = C^{-1}B^{-1}A^{-1}$
- $(A')^{-1} = A$
- $E\{Y\} = E\{Y - E\{Y\}\} + E\{Y\}$
- $\sigma^2 \{Y\} = E\{(Y - E\{Y\})^2\}$
- $\text{cov}(Y) = E\{YY^T\} - E\{Y\}E\{Y^T\}$
- $\text{var}(Y) = E\{Y^T Y\} - E\{Y\}E\{Y^T\}$
- $\text{cov}(XY) = E\{XY^T\} - E\{X\}E\{Y^T\}$
- $\text{var}(X) = E\{X^T X\} - E\{X\}E\{X^T\}$
- $\text{cov}(X^T Y) = E\{X^T Y\} - E\{X^T\}E\{Y\}$
- $\text{var}(Y|X) = E\{Y^T Y|X\} - E\{Y|X\}E\{Y^T|X\}$
- $\text{cov}(X^T Y|X) = E\{X^T Y|X\} - E\{X^T|X\}E\{Y|X\}$

SLR

$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$ $X = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$ $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$ $\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$

$Y = X\beta + \varepsilon$ *Note $E\{Y\} = X\beta$

Normal equations:

$SSTO = Y^T Y - (\frac{1}{n}) Y^T J Y$

$SSE = e^T e = (Y - \hat{Y})^T (Y - \hat{Y}) = Y^T Y - \hat{Y}^T Y$

$SSR = \hat{Y}^T X \hat{Y}$

$HH = H \text{ (idempotent)}$

$H^T = H \text{ (symmetric)}$

$SSTO = Y^T [I - (\frac{1}{n}) J] Y$

$SSE = Y^T [I - H] Y$

$SSR = Y^T [H - (\frac{1}{n}) J] Y$

$\text{To obtain the estimated regression coefficients: } (X^T X)^{-1} X^T \hat{Y} = (X^T X)^{-1} X^T Y$

$\text{Fitted: } \hat{Y} = X\hat{\beta} \quad \text{or} \quad \hat{\beta} = (X^T X)^{-1} X^T Y$

$\sigma^2 \{\varepsilon\} = \sigma^2 (I-H) \text{ and is estimated by } \hat{\sigma}^2 \{\varepsilon\} = MSE(I-H)$

$\text{Because } \sigma^2 \{\varepsilon\} = (I-H) \sigma^2 \{Y\} (I-H)^T$

$= (I-H) \sigma^2 I (I-H)^T \quad \sigma^2 \{\hat{\beta}\} = \sigma^2 (X^T X)^{-1}$

$= \sigma^2 (I-H)(I-H)^T \quad \sigma^2 \{\hat{\beta}\} = MSE(X^T X)^{-1}$

$= \sigma^2 (I-H) \quad \sigma^2 \{\hat{\beta}\} = MSE(I-H)$

More ANOVAs:

General MLR

$Y = X\beta + \varepsilon$; $\sigma^2 \{\varepsilon\} = \sigma^2 I$

$E\{Y\} = X\beta$; $\sigma^2 \{Y\} = \sigma^2 I$

$X^T \hat{Y} = X^T Y$; $\hat{Y} = X\hat{\beta}$; $e = Y - \hat{Y} = Y - X\hat{\beta}$ So $\hat{Y}_h = X_h^T \hat{\beta}$

$\hat{\beta} = (X^T X)^{-1} (X^T Y)$; $\hat{Y} = X\hat{\beta}$; $H = X(X^T X)^{-1} X^T$

$SSTO = Y^T Y - (\frac{1}{n}) Y^T J Y$; $e = (I-H)Y$; $\sigma^2 \{\hat{Y}_h\} = \sigma^2 X_h^T (X^T X)^{-1} X_h$

$= Y^T [I - (\frac{1}{n}) J] Y$; $\sigma^2 \{\varepsilon\} = \sigma^2 (I-H)$; $= X_h^T \sigma^2 \{\varepsilon\} X_h$

$SSE = e^T e = (Y - \hat{Y})^T (Y - \hat{Y})$; $\sigma^2 \{\varepsilon\} = MSE(I-H)$

$= Y^T Y - \hat{Y}^T Y = Y^T (I-H)Y$; $MSR = \frac{SSR}{P-1}$

$SSR = \hat{Y}^T X \hat{Y} = Y^T [H - (\frac{1}{n}) J] Y$; $MSE = \frac{SSE}{n-p}$

Diagnosing Multicollinearity:

- Pairwise correlation not sufficient
- Need to consider degree of linear relationship with the full set of other predictors in the model
- Use the R^2 obtained from regressing x_j on the other predictors (R_{jj}^2)
- Recall $\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{(1-R_{jj}^2)} S_{jj}^2$ when model only includes x_j .
- Now, with other predictors in the model,
- $\text{Var}(\hat{\beta}_j) = \frac{1}{1-R_{jj}^2} \frac{\sigma^2}{\sum_{i \neq j} R_{ij}^2}$
- $\frac{1}{1-R_{jj}^2} = \text{Variance Inflation Factor (VIF)}$ for predictor j .
- When R_{jj}^2 close to zero, VIF $_j$ close to 1 and SE not much affected, and vice versa
- $\sqrt{VIF_j}$ gives impact on $SE(\hat{\beta}_j)$. VIF $_j$ to a problem, possibly smaller if interacting w/ relevant
- Two considerations for effect on SE:
 - $\text{Var}(\hat{\beta}_j) = \frac{1}{1-R_{jj}^2} \frac{\sigma^2}{(n-1)S_{jj}^2}$
 - Centering variables before creating interaction terms and polynomial terms can help reduce correlation and improve interpretation.
 - Handling naturally correlated predictors
 - If not interfering w/ important inferences, do nothing
 - Consider changing predictors (often several variables represent same underlying construct)
 - Often, when trying to explain or control for confounding, we can just choose one and not the other.

We can also get the slope of \hat{y}_i as mean of y 's change depending on value of x_i .
 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$
 Relative \hat{y}_i is change when you graph x_i is slope in \hat{y}_i when $x_i = 0$
 $\hat{\beta}_{x_i}$ is slope in \hat{y}_i when $x_i = 1$
 $\hat{\beta}_{x_i}$ is change in slope associated w/
 one-unit increase in the other predictor
 predictors i need to be (for interp purposes)

Leverage:
 x_i is i th row of design matrix $H = X(X'X)^{-1}X'$
 Elements of H are h_{ij} , column $X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$
 $h_{i,j} = x_i'(X'X)^{-1}x_j = x_j'(X'X)^{-1}x_i = h_{j,i}$
 $h_{ii} = x_i'(X'X)^{-1}x_i$
 Imagine SLR, where $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2}$
 \hookrightarrow Leverage is a type of relative distance measure
 Leverage for obs i is proportion of the total squared deviations of the x_i from their mean \bar{x} that is contributed by the observation (plus a $\frac{1}{n}$ term same for all obs.)
 • Far from mean, high leverage
 • Close to mean, low leverage
 Generally, $\bar{x} = [x_1, \dots, x_{p-1}]$ is mean vector, or centroid of predictor variables. $x^{**} = x_i - \bar{x}$ is mean-centered version of covariate raw vector for observation i
 $X^* = \begin{bmatrix} x_1^{**} \\ \vdots \\ x_n^{**} \end{bmatrix} = \begin{bmatrix} (x_{11} - \bar{x}_1) & \dots & (x_{1,p-1} - \bar{x}_{1,p-1}) \\ \vdots & \vdots & \vdots \\ (x_{n1} - \bar{x}_1) & \dots & (x_{n,p-1} - \bar{x}_{n,p-1}) \end{bmatrix}$
 $\frac{1}{n-1} X^* X^*$ is sample variance-covariance matrix of predictors $h_{ii} = \frac{1}{n-1} x_i^{**} (X^* X^*)^{-1} x_i^{**}$, so h_{ii} measures distance from centroid taking into account covariance structure of x 's

Partial regression plots (added-variable)
 are plots of residuals from partial regressions
Example plot of $e_{1|2,3,\dots,p-1}$ versus $x_{1|2,3,\dots,p-1}$, along with partial regression line (from SLR regressing residuals on residuals)
 • Helps with identifying influential obs., checking for non-linearity, suggesting monotonicity or not.

Step functions
 • Break X into bins K1 indicator variables
 • Create cutpoints C_1, \dots, C_K called knots
 \hookrightarrow piecewise constant function
 Piecewise linear reg:
 Range X
 $x_i \leq C_1 \quad Y_i = \beta_0 + \beta_1 x_i + e_i$
 $C_1 < x_i \leq C_2 \quad Y_i = \beta_0 + \beta_1 x_i + \beta_2 (x_i - C_1) + e_i$
 \vdots
 (continuous at knots)

Regression spline
 \hookrightarrow piecewise polynomial constrained to be smooth and cont. at knots.
 $\text{Plot } C_1(x_i) + \beta_1 C_1(x_i) x_i + \beta_2 C_1(x_i) x_i^2 + \dots$
 not continuous at knots
 Dropping β_0 terms (intercepts) gives continuity, but not smooth
 Dropping second and third linear terms gives equal first derivatives (duh)
 Dropping quadratic terms gives equal curvature

Total diagnostics
 • Conditional y -outliers vs. marginal y -outliers
 \hookrightarrow If grouped w/ rest of y data, but y_i is outlier in conditional dist. of y given x , we have conditional y -outlier.
 \hookrightarrow If not grouped w/ rest of y data, but y_i falls within conditional dist. of y given x , we have a marginal y -outlier

• Influence depends on x -outlyingness (leverage) and conditional y -outlyingness (residual). Also sample size.

• Leverage values are bounded $\frac{1}{n} \leq h_{ii} \leq 1$
 • Also, $\sum h_{ii} = p$, where $p = \text{rank}(X)$
 \hookrightarrow So $\bar{h} = \frac{p}{n}$, and 2 \bar{h} or 3 \bar{h} are used as rules of thumb for high leverage.

Residuals:
 Using raw residuals for model diagnostics has limitations
 • Raw res. do not all have same variance
 • True errors: $\text{Var}(\varepsilon) = \sigma^2 I$; Res: $\text{Var}(\hat{\varepsilon}) = \sigma^2 (I - H)$
 \hookrightarrow So $\text{Var}(\varepsilon_i) = \sigma^2$, but $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1-h_{ii})$
 \hookrightarrow Var of res. for obs. i depends on its leverage.
 • We standardize them: $e_i = \frac{\varepsilon_i}{\sqrt{1-h_{ii}}}$
 For constant variance 1,
 \hookrightarrow but recall $\hat{\sigma} = \text{root MSE} = \sqrt{\frac{1}{n-p} \sum e_i^2}$ which depends on the raw residual
 1) If obs i is a regression outlier, and has an unusual res., it will inflate est of β_i .
 2) Numerator and denominator not independent, so statistic doesn't follow a t dist.
 • Studentized deleted residual: $d_i = \frac{e_i}{\sqrt{1-h_{ii}}}$ where $\hat{\sigma}_{-i}$ is estimate of root MSE based on regression with observation i deleted.
 • With observation i deleted, error terms normally dist, then std.

Measures of influence:
 • Case deletion - how to calculate $\hat{\beta}$ with i th case deleted?
 \hookrightarrow Denote as $\hat{\beta}_{(-i)}$, a $p \times 1$ vector
 $\hookrightarrow X$ is $n \times p$ design matrix w/ linearly independent columns
 $\hookrightarrow X_{(-i)}$ is design matrix w/ i th obs deleted, a $(n-1) \times p$ matrix
 $\hookrightarrow Y_{(-i)}$ is vector of observed values w/ i th obs deleted, a $(n-1) \times 1$
 • We want to compute $\hat{\beta}_{(-i)} = [X_{(-i)}' X_{(-i)}]^{-1} X_{(-i)}' Y_{(-i)}$
 • Claim: we can compute $\hat{\beta}_{(-i)}$ for each i , to the original model output obtained by fitting the model to all n obs, without refitting the model.

$[X_{(-i)}' X_{(-i)}]^{-1} = (X'X)^{-1} + \frac{(X'X)^{-1} x_i x_i' (X'X)^{-1}}{1-h_{ii}}$
 where $h_{ii} = x_i' (X'X)^{-1} x_i$
 $X_{(-i)}' Y_{(-i)} = X'Y - x_i' Y_i$
 $\hat{\beta}_{(-i)} = [(X'X)^{-1} + (X'X)^{-1} x_i x_i' (X'X)^{-1}] \frac{1-h_{ii}}{1-h_{ii}} [X'Y - x_i' Y_i]$
 \hookrightarrow $\hat{\beta}_{(-i)}$ is i th residual for obs. i
 \hookrightarrow can obtain case-deleted MSE, fitted values, etc. using $\hat{\beta}_{(-i)}$. Don't refit the model every time.

DFFITS_i: measures influence of i th obs. on i th fitted value \hat{y}_i
 $\hookrightarrow \text{DFFITS}_i = \frac{\hat{y}_i - \hat{y}_{(-i)}}{\sqrt{\sigma^2(1-h_{ii})}}$ \hookrightarrow fitted val for obs. i when i th obs is omitted when fitting regression model
 $\text{Var}(\hat{y}_i) = \text{Var}(\hat{y}_i | \hat{\beta})$
 $= x_i' \text{Var}(\hat{\beta}) x_i$
 $= \sigma^2 x_i' (X'X)^{-1} x_i$
 $= \sigma^2 h_{ii}$
 \hookrightarrow SD of \hat{y}_i , but using estimate of MSE when i th obs. is omitted when fitting the model.

DFFITS_i vs. \hat{y}_i :
 • Gives influence on fitted value in SD units
 • Measure how well the model predicts each obs, when it is not used to fit the model. Think of them as measures of prediction accuracy.
 • Note we can calculate using entire dataset:
 $\text{DFFITS}_i = t_i (1-h_{ii})^{1/2}$ where t_i is std. deleted res. for obs. i and h_{ii} is leverage value.

Cook's Distance: summary index of influence of the i th obs. on all fitted values:
 $D_i = \frac{\sum (\hat{y}_i - \hat{y}_{(-i)})^2}{p \hat{\sigma}^2} = \frac{(\hat{y}_i - \hat{y}_{(-i)})^2}{p \hat{\sigma}^2}$
 $= \frac{(\hat{\beta}_{(-i)} - \hat{\beta})' X' X (\hat{\beta}_{(-i)} - \hat{\beta})}{p \hat{\sigma}^2}$
 $= \frac{e_i^2 h_{ii}}{p \hat{\sigma}^2 (1-h_{ii})^2}$

DFBETAS_{ij}: provide changes in j th estimated regression coeff when i th obs. omitted from data set
 $\text{DFBETAS}_{ij} = \frac{\hat{\beta}_j - \hat{\beta}_{(-i)}}{\sqrt{C_{jj}}}$
 where C_{jj} is j th diagonal element of $(X'X)^{-1}$. Units in standard errors.

Addressing influential obs:
 • Check for incorrect data
 • Perhaps observations don't belong to pop. of interest.
 • May need to change model
 • Consider impact of sample size.
 • Sensitivity analysis → estimate model with and without point

Model evaluation
 $AIC = -2 \log L(\hat{\theta}) + 2p$ for p parameters
 $\downarrow AIC$ is better $\Rightarrow n \log(\hat{\theta}^2) + 2p$
 $\hat{\theta}^2 = \frac{\text{SSE}}{n-p}$
 $BIC = -2 \log L(\hat{\theta}) + \log(n)p$
 \downarrow indicates good model $C_p = \frac{\text{SSE}(n+1) - (n-p)}{n-p}$
 $AIC \leq 3-7 > 10$

SLR

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \begin{array}{l} \text{note, centered model gives} \\ \hat{Y}_i = \bar{Y} + \hat{\beta}_1(X_i - \bar{X}) \\ \text{so intercept becomes } \bar{Y} \end{array}$$

- β_0 = intercept; when range of data includes $X=0$, β_0 is the mean of the dist. of Y at $X=0$. Otherwise, not meaningful
- β_1 = slope of regression line; indicates change in conditional mean of Y per one unit increase in X
- σ^2 = conditional variance of Y given X ; estimated by MSE

Assumptions:

- Linearity: $E(\varepsilon_i) \equiv E(\varepsilon_i | X_i) = 0$
 $E(Y_i) \equiv E(Y_i | X_i) = E(\beta_0 + \beta_1 X_i + \varepsilon_i)$
 $= \beta_0 + \beta_1 X_i + E(\varepsilon_i)$
 $= \beta_0 + \beta_1 X_i$

- Constant variance: $\text{Var}(\varepsilon | X_i) = \sigma^2$
 $\text{Var}(Y | X_i) = E[(Y_i - \mu_{Y|X_i})^2]$
 $= E[(Y_i - \beta_0 - \beta_1 X_i)^2]$
 $= E(\varepsilon_i^2) = \sigma^2$

- Normality: $\varepsilon_i \sim N(0, \sigma^2)$
 $Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$

- Independence: Observations sampled independently. Must be evaluated by data collection procedure evaluation.

- Fixed X , or X measured without error and independent of errors.

- X is not invariant \rightarrow must have variation in X

MER

For k predictors, model becomes:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \varepsilon_i$$

\hookrightarrow Note $k+1$ parameters

\hookrightarrow β_j is change in $E(Y)$ for a one-unit increase in X_j when all other predictors are held constant.

"controlling for all other predictors"

1. Regress Y on X_2, \dots, X_k
2. Regress X_1 on X_2, \dots, X_k and get residuals $\varepsilon_{i1}, y_{1|X_2, \dots, X_k}$
3. SLR of $\varepsilon_{i1}, y_{1|X_2, \dots, X_k}$ on X_1 gives $\hat{\beta}_1$

* We think of each partial regression coefficient as providing an estimate of the linear relationship between \hat{Y} and that predictor, after linear relationship between Y and all other predictors has been accounted (controlled) for.



models w/ higher R², where R² just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

models w/ higher R², where R² really just an artifact of more predictors

$$\text{var}(\hat{Y}_n) = \hat{\sigma}^2 = \left(\frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right); \text{ CI: } \hat{Y}_n \pm t_{1-\alpha/2, n-2} \text{ SE}(\hat{Y})$$

prediction interval: $\hat{Y}_n \pm t_{1-\alpha/2, n-2} \sqrt{1 + \frac{1}{n} + \frac{(X_n - \bar{X})^2}{\sum(X_i - \bar{X})^2}}$

$$\text{var}(\hat{\beta}_0) = \hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum(X_i - \bar{X})^2} \right) \Rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\text{var}(\hat{\beta}_0)}} \sim t_{n-2}; \text{ CI: } \hat{\beta}_0 \pm t_{1-\alpha/2, n-2} \text{ SE}(\hat{\beta}_0)$$

$$\bullet \hat{\beta}_1 = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2} = r \frac{s_y}{s_x}$$

$$\text{var}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum(X_i - \bar{X})^2} \Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\text{var}(\hat{\beta}_1)}} \sim t_{n-2}; \text{ CI: } \hat{\beta}_1 \pm t_{1-\alpha/2, n-2} \text{ SE}(\hat{\beta}_1)$$

$$\bullet \hat{\sigma}^2 = \frac{1}{n-2} \sum(Y_i - \hat{Y}_i)^2 \stackrel{\text{(mean square error, error squares)}}{=} \frac{\text{RSS}}{\text{sum of squares}}$$

Mean square error, error mean square, residual mean sum of squares (RSS)

$\bullet \hat{\sigma}$ has same units as Y and is estimated std. dev. of observations around their mean.
 \hookrightarrow root MSE

\bullet Precision = $\frac{1}{\text{variance}}$

$\bullet R^2 = \frac{\text{Reg SS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$ (coefficient of determination)
 $\hookrightarrow 0 \leq R^2 \leq 1$

\hookrightarrow % variation explained by model w/ regression
 \hookrightarrow may not capture strong non-linear relationship

\hookrightarrow high R^2 not a guarantee for good fit -

$\bullet r = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}} = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum(X_i - \bar{X})^2} \sqrt{\sum(Y_i - \bar{Y})^2}}$
 \hookrightarrow for SLR, $r^2 \pm R^2$
 \hookrightarrow so for SLR, slope coefficient is a scaled version of r by std. devs. of x and y : $\hat{\beta}_1 = r \frac{s_y}{s_x}$

$\hookrightarrow r_{xy}$ values all given in Pearson correlation matrix

for MLR

tips us about individual corr. matrix

Degrees of freedom:
 $K = \# \text{ predictors}$
 $K+1 = \# \text{ parameters (Intercept)}$

Source	df	ss	ms	f
Model	1	Reg SS	MS Reg = $\frac{\text{Reg SS}}{1} = \frac{\text{Reg SS}}{n-2}$	$F_{1, n-2}$
Residual	n-2	RSS	$\text{MSE} = \frac{\text{RSS}}{n-2}$	
Total	n-1	TSS		

Note:
 $E(\text{MSE}) = \sigma^2$
 $E(\text{MS Reg}) = \sigma^2 + \beta_1^2 \frac{\sum(X_i - \bar{X})^2}{n-2}$
 $\text{so as } \beta_1 \rightarrow 0, \text{ MS Reg} \rightarrow \text{MSE}$

F-test gives same result as t-test for $\beta_1 = 0$
 \hookrightarrow large F therefore supports $\beta_1 \neq 0$.

$$\text{From practice: } \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0)} = \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0) + 2 \text{cov}(\hat{\beta}_j, \hat{\beta}_0)} = \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0) + 2 \text{Var}(\hat{\beta}_0) R^2} = \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0)(1 + 2R^2)}$$

$$\text{For } \hat{\beta}_j: \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0)} = \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0) + 2 \text{cov}(\hat{\beta}_j, \hat{\beta}_0)} = \frac{\text{Var}(\hat{\beta}_j)}{\text{Var}(\hat{\beta}_0)(1 + 2R^2)}$$

Test value of single coefficients:
 $H_0: \beta_j = \beta_j^*$
 $H_A: \beta_j \neq \beta_j^*$
 $t^* = \frac{\hat{\beta}_j - \beta_j^*}{\sqrt{\text{var}(\hat{\beta}_j)}} \sim t_{n-2}$

$$\text{Var}(\hat{\beta}_j) = \left(\frac{1}{1-R^2} \right) \left(\frac{\sigma^2}{\sum(X_{ij} - \bar{X})^2} \right)$$

Partial regression of X_j on all other X_i

$$\text{Eq: } \hat{\beta}_j^* = \hat{\beta}_j - \frac{\sum(X_{ij} - \bar{X})}{\sum(X_{ij} - \bar{X})^2} * (\text{SE}(\hat{\beta}_j))$$

Source	df	ss	ms	f
Model	K	Reg SS	MS Reg = $\frac{\text{Reg SS}}{K} = \frac{\text{Reg SS}}{n-K-1}$	$F_{K, n-K-1}$
Residual	n-K-1	RSS	$\text{MSE} = \frac{\text{RSS}}{n-K-1}$	
Total	n-1	TSS		

$$\hookrightarrow F = \frac{\text{Reg SS}/K}{\text{RSS}/(n-K-1)} = \frac{(n-K-1)}{K} \left(\frac{R^2}{1-R^2} \right)$$

If we want to test whether adding a set of predictors X_m, X_{m+1}, \dots, X_j to a model that already contains X_1, X_2, \dots, X_{m-1} has a sig. effect in terms of explaining variation in Y , we consider diff. in RSS or Reg SS between models with and without X_m, \dots, X_j :

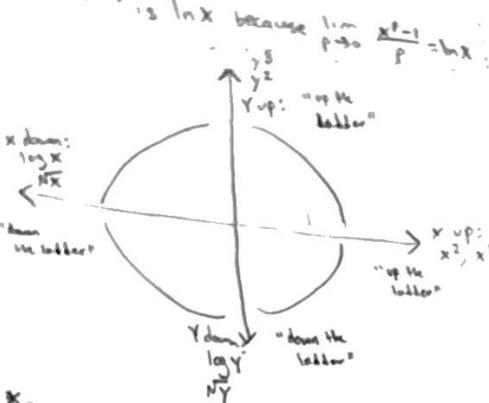
$$\text{SSE}(X_1, \dots, X_{m-1}) - \text{SSE}(X_1, \dots, X_m, \dots, X_j) = \text{Reg SS}(X_1, \dots, X_j) - \text{Reg SS}(X_1, \dots, X_{m-1})$$

* We are comparing a Full and a Reduced model! $H_0: \beta_m = \dots = \beta_j = 0$

link transformation:

$$x \rightarrow x^p = \frac{x^p - 1}{p}$$

x^{10} is $\ln x$ because $\lim_{p \rightarrow 0} \frac{x^p - 1}{p} = \ln x$



Transform X or Y?

- transforming Y generally has larger effect on dist. of error terms.
- think about interpretation

Model Diagnostics:

- Goal is to determine whether assumptions are met.
- Many assumptions are about the errors, which are not observable so we use residuals analysis instead.



Q-Q	Interp	Assumption	Methods for Checking	Consequence of violation	Remedies
	normal	Constant Variance	Residuals analysis, plot of e vs Y	increased standard errors	Variance stabilizing transformation of y , weighted least squares, nothing if not bad
	fat tails (heavy)				
	positive skew	$E(\epsilon_i) = 0$	Residuals analysis, plot of e vs. \hat{Y}	mean function is misspecified, regression coefficients are wrong	Transformations, deal w/ nonlinearities, (i.e., polynomial reg)
	negative skew	Normality (worst violation)	Univariate residuals plot; Q-Q, P-P	important for validity of tests and CI in small samples	Transformations of Y
Q-Q has range of data P-P has $[0,1] \times [0,1]$	Independent obs.		Understand data collection	SEs are wrong (\downarrow) if we model correlated values as though independent; parameter estimates biased	methods for correlated data

* Power transformations can help linearize when relationship is monotone and "simple" (i.e., no point of inflection). And all values are positive → add a start if need be

In general, for $\ln(Y)$,

effect of a one-unit increase in X_j is to multiply mean of Y by e^{β_j}

Interpreting natural log transformations: When we ln-transform a predictor, we can say that a 1% increase in X is associated with an increase in the conditional mean of y by about 0.01β . When we ln-transform the outcome variable, the effect of a one-unit increase in y is to multiply y by $\exp[\hat{\beta}_j]$.

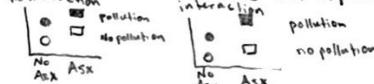
Suppose y and x are related by $y = \alpha x^\beta$.

If we increase x by 1%, the change in y is equivalent to multiplying it by 1.01^β , which is roughly equal to increasing y by β percent.

When y and x are both ln-transformed, β is here called an elasticity coefficient.

Suppose Y_i is measure of breathing difficulty, $X_{1i}=1$ if asthma, $X_{2i}=1$ if exposed to pollution

Interaction model:



$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{1i}X_{2i} + \epsilon_i$$

↳ For children w/ asthma, mean change due to air pollution exposure is $\beta_2 + \beta_3$

↳ For no asthma, mean change due to air pollution exposure is β_2

* We say "asthma modifies (moderates) the effect of air pollution on breathing difficulties" or vice versa.

OR "Effect of air pollution on breathing difficulties is conditional on ast status"

Interaction terms are symmetric in the two variables:

$$Y_i = \beta_0 + \beta_2 X_{2i} + (\beta_1 + \beta_3 X_{2i}) X_{1i} + \epsilon_i$$

• When $X_{2i} = 0$, effect of X_{1i} on $E(Y_i)$ is β_1 .

• When $X_{2i} = 1$, effect of X_{1i} on $E(Y_i)$ is $\beta_1 + \beta_3$

$$Y_i = \beta_0 + \beta_1 X_{1i} + (\beta_2 + \beta_3 X_{1i}) X_{2i} + \epsilon_i$$

• When $X_{1i} = 0$, effect of X_{2i} on $E(Y_i)$ is β_2 .

• When $X_{1i} = 1$, effect of X_{2i} on $E(Y_i)$ is $\beta_2 + \beta_3$

We can have interactions among 3 or more variables. Here's a model with 3 dichotomous (0/1) variables and a

3-way interaction:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{1i}X_{2i} + \beta_5 X_{1i}X_{3i}$$

If one predictor was categorical w/ 3 levels? Let Z_1 and Z_2 be dummy variables:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 Z_{1i} + \beta_3 Z_{2i} + \beta_4 X_{1i}Z_{1i} + \beta_5 X_{1i}Z_{2i} + \epsilon_i$$

When we have more than one interaction term, we conduct a partial F-test as to whether or not an interaction exists: $H_0: \beta_4 = 0, \beta_5 = 0$

Interpretation of lower-order terms in model w/ interaction term

$$Y_i = \beta_0 + \beta_1 d_i + (\beta_2 + \beta_3 d_i) X_{1i} + \epsilon_i$$

β_0 is int. for men; β_1 is diff. in mean outcome for individuals who differ on d_i by 1 unit

β_2 is slope for women; β_3 is slope for men vs. women

β_3 is diff. in mean outcome for individuals who have $X_{1i} = 0$ → if this doesn't make sense, center X_{1i} at its mean!

Interaction:

$$Y_i = \beta_0 + \beta_1 d_i + \beta_2 X_{1i} + \beta_3 d_i X_{1i} + \epsilon_i$$

↳ men: $\beta_0 + \beta_2 X_{1i} + \epsilon_i$

↳ men: $(\beta_0 + \beta_1) + (\beta_2 + \beta_3) X_{1i} + \epsilon_i$

* Allows for difference in slopes. β_0 is int. for women; β_1 is int. for men vs. women

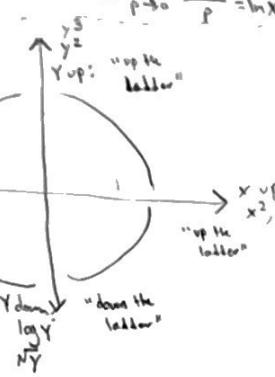
β_2 is slope for women

β_3 is slope for men vs. women

Box-Cox transformation:

$$x \rightarrow x^{(p)} = \frac{x^p - 1}{p}$$

$x^{(0)}$ is $\ln x$ because $\lim_{p \rightarrow 0} \frac{x^p - 1}{p} = \ln x$

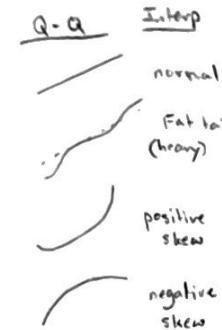
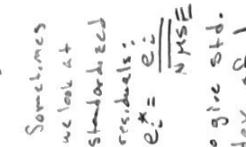


Transform X or Y?

- transforming Y generally has larger effect on dist. of error terms.
- think about interpretation

Model Diagnostics:

- Goal is to determine whether assumptions are met.
- Many assumptions are about the errors, which are not observable so we use residuals analysis instead.



Q-Q has range of data
P-P has $[0,1] \times [0,1]$

Assumption	Methods for Checking	Consequence of violating	Remedies
constant variance	Residuals analysis, plot of e vs y	increased standard errors	Variance stabilizing transformation of y , weighted least squares, nothing if not bad
$E(e_i) = 0$	Residuals analysis, plot of e vs. \hat{Y}	mean function is misspecified, regression coefficients are wrong	Transformations, deal w/ non-linears (i.e., polynomial reg.)
Normality (w/ or + violation)	Univariate residuals plot; Q-Q, P-P	important for validity of tests and CI in small samples	Transformation of y
Independent obs.	Understand data collection	S_E s are wrong (\downarrow) if we model correlated values as though independent; parameter estimates biased	methods for correlated data

* Power transformations can help linearize when relationship is monotone and "simple" (i.e., no point of inflection). And all values are positive \rightarrow add a start if need be

In general, for $\ln(y)$, effect of a one-unit increase in X_j is to multiply mean of Y by $e^{\hat{\beta}_j}$

Interpreting natural log transformations: When we ln-transform a predictor, we can say that a 1% increase in X is associated with an increase in the conditional mean of y by about $0.01\hat{\beta}$. When we ln-transform the outcome variable, the effect of a one-unit increase in X is to multiply y by $\exp[\hat{\beta}]$.

Suppose y and x are related by $y = \alpha x^\beta$. If we increase x by 1%, the change in y is equivalent to multiplying it by 1.01^β , which is roughly equal to increasing y by β percent. \rightarrow When y and x are both ln-transformed, β is here called an elasticity coefficient

Suppose Y_i is measure of breathing difficulty, $X_{1i}=1$ if asthma, $X_{2i}=1$ if exposed to pollution

Interaction model:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{1i}X_{2i} + \varepsilon_i$$

\rightarrow For children w/ asthma, mean change due to air pollution exposure is $\beta_2 + \beta_3$

\rightarrow For no asthma, mean change due to air pollution exposure is β_2

* We say "asthma modifies (moderates) the effect of air pollution on breathing difficulties" or vice versa.

OR "Effect of air pollution on breathing difficulties is conditional on ast status"

Interaction terms are symmetric in the two variables;

$$\rightarrow Y_i = \beta_0 + \beta_2 X_{2i} + (\beta_1 + \beta_3 X_{2i}) X_{1i} + \varepsilon_i$$

* When $X_{2i} = 0$, effect of X_{1i} on $E(Y_i)$ is β_1 .

* When $X_{2i} = 1$, effect of X_{1i} on $E(Y_i)$ is $\beta_1 + \beta_3$

$$\rightarrow Y_i = \beta_0 + \beta_1 X_{1i} + (\beta_2 + \beta_3 X_{1i}) X_{2i} + \varepsilon_i$$

* When $X_{1i} = 0$, effect of X_{2i} on $E(Y_i)$ is β_2 .

* When $X_{1i} = 1$, effect of X_{2i} on $E(Y_i)$ is $\beta_2 + \beta_3$

We can have interactions among 3 or more variables.
Here's a model with 3 dichotomous (0/1) variables and a

3-way interaction:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{1i}X_{2i} + \beta_5 X_{1i}X_{3i}$$

If one predictor was categorical w/ 3 levels? Let Z_1 and Z_2 be dummy variables:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 Z_{1i} + \beta_3 Z_{2i} + \beta_4 X_{1i}Z_{1i} + \beta_5 X_{1i}Z_{2i} + \varepsilon_i$$

When we have more than one interaction term, we conduct a partial F-test as to whether or not an interaction exists: $H_0: \beta_4 = 0, \beta_5 = 0$

Interpretation of lower-order terms in model w/ interaction term

$$Y_i = \beta_0 + \beta_1 d_i + (\beta_2 + \beta_3 d_i) X_{1i} + \varepsilon_i$$

* β_1 is diff. in mean outcome for individuals who differ on d_i by 1 unit and have $X_{1i} = 0$ \rightarrow if this doesn't make sense, center X_{1i} at its mean!

Interaction between categorical and continuous variable

No interaction: $d_i=1$ for men, $d_i=0$ for women; X_{1i} yes/no

$$Y_i = \beta_0 + \beta_1 d_i + \beta_2 X_{1i} + \varepsilon_i$$

\rightarrow women: $\beta_0 + \beta_2 X_{1i} + \varepsilon_i$

\rightarrow men: $(\beta_0 + \beta_1) + \beta_2 X_{1i} + \varepsilon_i$ * common slope model

β_1 is difference in intercepts, constant vertical dist. b/wn lines expected difference in mean income b/wn men and women when education held constant.

For

is int. for women

men vs. women

vs. women

Interaction:

$$Y_i = \beta_0 + \beta_1 d_i$$

Properties

- Cauchy-Schwarz**: $E[XY] \leq E[X]E[Y] \leq (E[X]^2)^{1/2}(E[Y]^2)^{1/2}$
- For any r.v.s X and Y , $E[XY] = E[X]E[Y]$
- $E[XY] \leq E[X]E[Y] \leq (E[X]^2)^{1/2}(E[Y]^2)^{1/2}$ So if X and Y have means μ_X and μ_Y and variances σ_X^2 and σ_Y^2 , we have $(\text{Cov}(X, Y))^2 \leq \sigma_X^2 \sigma_Y^2$. This implies that the conditionals equal the marginals when $X \perp Y$.
- Jensen's: If $f(x)$ is differentiable, convex iff $f'(x) \leq f''(y)$ for all $x < y$; if $f(x)$ is convex, then $E[g(x)] \geq g(E[X])$; if $f(x)$ is concave, then $E[g(x)] \leq g(E[X])$. Equality holds if and only if, for every line $at+bx$ tangent to $g(x)$ at $x = EX$, $P(g(x) = at+bx) = 1$. One immediate application tells us $EX^2 \geq (EX)^2$. Strict inequality holds if $g(x)$ is strictly convex since $g(x) = x^2$ is convex.
- Chebychev's**: Let X be a r.v. and let $g(x)$ be nonnegative. Then for any $r \geq 0$, $P(g(X) \geq r) \leq \frac{E[g(X)]}{r}$. Most common use says to let $g(x) = \frac{(x-\mu)^2}{\sigma^2}$. For convenience, rewrite $r = t^2$. Then $P(\frac{(X-\mu)^2}{\sigma^2} \geq t^2) \leq \frac{1}{t^2} E(\frac{(X-\mu)^2}{\sigma^2}) = \frac{1}{t^2}$.
- Markov's**: Let $u(x)$ be a nonnegative function of the r.v. X . If $E[u(X)]$ exists, then for every positive constant c , $P[u(X) \geq c] \leq E[u(X)]/c$.

Discrete Uniform: pmf: $p(x) = \frac{1}{N}; x=1, 2, \dots, N$; $\sigma^2 = \frac{N+1}{2}; \sigma^2 = \frac{1}{N(N+1)(N-1)}$

Binomial: pmf: $p(x) = \binom{n}{x} p^x (1-p)^{n-x}$; mgf: $M(t) = (1-p + pe^t)^n$; $E[e^{tX}] = \sum_{x=0}^n \binom{n}{x} (e^t)^x (1-p)^{n-x} e^{-n}$; $M(t) = E(e^{tX}) = \sum_{x=0}^n \binom{n}{x} (e^t)^x (1-p)^{n-x} e^{-n}$; $\sigma^2 = np; \sigma^2 = np(1-p)$. # successes in n independent Bernoulli trials.

Exponential (β): Frequently a model for waiting times. pmf: $f(x) = \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha) \beta^\alpha}, 0 \leq x < \infty$; $\alpha, \beta > 0$; mgf: $M(t) = (\frac{1}{1-\beta t})^\alpha, t < \frac{1}{\beta}$. $\alpha=1$ gives exponential $\lambda = \alpha/\beta$; $\alpha=2, \beta=2$ gives χ^2 $\sigma^2 = \alpha\beta^2$.

Normal (μ, σ^2): pmf: $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$; $-\infty < x < \infty, -\infty < \mu < \infty, \sigma^2 > 0$; mgf: $M(t) = e^{\mu t + \frac{\sigma^2 t^2}{2}}$; $\mu = \mu, \sigma^2 = \sigma^2$. Let X_1, \dots, X_n be \perp r.v.s. Suppose $X_i \sim \chi^2(r_i)$. Let $Y = \sum_{i=1}^n X_i$. Then $Y \sim \chi^2(\sum_{i=1}^n r_i)$. We get a t -distribution when $W = N(0, 1), V \sim \chi^2(r)$. Let X_1, \dots, X_n be \perp r.v.s. We get an F-distribution when $U \sim \chi^2(r_1)$ and $V \sim \chi^2(r_2)$. So for $i=1, 2, \dots, n$, $X_i \sim N(\mu_i, \sigma_i^2)$ and we take $W = (U/r_1)/(V/r_2)$. Let $Y = \sum_{i=1}^n a_i X_i$, for constants a_1, \dots, a_n . Then $\bar{X} \sim N(\mu, \sigma^2/n)$. Let X_1, \dots, X_n be iid so $X_i \sim N(\mu, \sigma^2)$. Then $\bar{X} \sim N(\mu, \sigma^2/n)$. Let $X \sim \chi^2(r)$. If $k > r/2$, then $E(X^k)$ exists and is given by $E(X^k) = \frac{k!}{2} \Gamma(\frac{k+r}{2})$. Let X_1, \dots, X_n be iid. Suppose $X_i \sim \Gamma(\alpha_i, \beta)$. Let $Y = \sum_{i=1}^n X_i$. Then $Y \sim \Gamma(\sum_{i=1}^n \alpha_i, \beta)$.

Conditional: $f(y|x) = \frac{f(x,y)}{f(x)}$

General Mathematical Facts:

- $E[g(x)] = \int_{-\infty}^{\infty} g(x)p(x)dx$, if x is continuous
- $E[g(x)] = \sum_{x \in \mathbb{Z}} g(x)p(x)$, if x is discrete
- $E[g(x)+bg_2(x)+c] = aEg_1(x) + bEg_2(x) + c$
- $E[g(y)|x] = \int_{-\infty}^{\infty} g(y)f(y|x)dy$
- If $X \perp Y$, $E[g(x)h(y)] = [Eg(x)][Eh(y)]$
- If X and Y are any two r.v.s, $EX = E[E(X|Y)]$
- For any two r.v.s X and Y , $\text{Var}X = E[\text{Var}(X|Y)] + \text{Var}[E(X|Y)]$
- $\text{Cov}(X, Y) = E[(X-\mu_X)(Y-\mu_Y)] = E(XY) - \mu_X \mu_Y$
- $\rho_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$. If $X \perp Y$, then $\text{Cov}(X, Y) = \rho = 0$
- $\text{Var}(ax+by) = a^2 \text{Var}X + b^2 \text{Var}Y + 2ab \text{Cov}(X, Y)$ * So when $X \perp Y$, $\text{Var}(ax+by) = a^2 \text{Var}X + b^2 \text{Var}Y$

General Mathematical Facts:

- $\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n, |x| < 1$
- $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}, \frac{d}{dx} e^x = \sum_{n=0}^{\infty} (-1)^n x^n$
- $(a+b)^n = \sum_{x=0}^{\infty} \binom{n}{x} b^x a^{n-x}$
- $\ln(1+x) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1} x^n}{n}$
- $\frac{a}{1-t} = \sum_{n=1}^{\infty} ar^{n-1}$ for $|t| < 1$

MacLaurin Series for mgf: $M(t) = M(0) + \frac{M'(0)}{1!}t + \frac{M''(0)}{2!}t^2 + \dots + \frac{M^{(m)}(0)}{m!}t^m$

Geometric: A special case of the negative binomial distribution with $r=1$. geometric(p): $P(X=x) = \frac{1}{p} \frac{(1-p)^{x-1}}{x!}, x=0, 1, \dots$; $\sigma^2 = \frac{1-p}{p^2}$. mgf: $M(t) = \frac{1-(1-p)e^t}{1-(1-p)e^t - t}$, $t < -\log(1-p)$. * This distribution is memoryless; $P(X>s|X>t) = P(X>s-t)$.

Negative Binomial (r, p): Experiment results in one of k ways, C_k . Probability of result being an element of C_i . Experiment repeated n times. X_i is number of outcomes in C_i . If X_1, X_2, \dots, X_{n-r} are outcomes in corresponding set C_i , then the multinomial pmf is given by: $P(X) = \frac{n!}{x_1! \dots x_{n-r}!} p^{x_1} p^{x_2} \dots p^{x_{n-r}}$. Specifically, when $k=3$, we get the trinomial $P(X, Y) = \frac{n!}{x! y! (n-x-y)!} p^x p^y p^{n-x-y}$. $M(t_1, t_2) = (pe^{t_1} + p^2 e^{t_2} + p^3)^n$; $M(t_1, 0) = M_X(t_1) = [(1-p) + pe^{t_1}]^n$; $M(0, t_2) = M_Y(t_2) = [(1-p) + p^2 e^{t_2}]^n$.

Chi-squared (ν): Frequently a model for waiting time until first "success". pmf: $f(x) = \frac{1}{\Gamma(\nu/2)} x^{\nu/2-1} e^{-x/2}, 0 \leq x < \infty; \nu = 1, 2, 3, \dots$; mgf: $M(t) = (\frac{1}{1-2t})^{\nu/2}, t < \frac{1}{2}$. $\nu=1$ gives gamma $M=\beta$; $\nu=2$ gives geometric $\sigma^2 = \beta^2$.

Beta (α, β): pmf: $f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, 0 \leq x \leq 1, \alpha > 0, \beta > 0$; mgf: $M(t) = 1 + \sum_{k=1}^{\infty} \frac{(\frac{\alpha+k-1}{\alpha+k})^t k}{k!} \frac{x^k}{(1-x)^{k+1}}, \mu = \frac{\alpha}{\alpha+\beta}, \sigma^2 = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$. * The constant in beta pdf can be defined in terms of gamma functions, $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$.

F Related to χ^2 : $F_{v_1, v_2} = \frac{\chi_{v_1}^2}{\chi_{v_2}^2} = \frac{(x_{v_1})/(x_{v_2})}{(x_{v_2})/(x_{v_1})}$ where the χ^2 are independent. Also related $v_1 \rightarrow v_1 - 2$ to $v_2 \rightarrow v_2 - 2$ and $F_{v_1-2, v_2-2} = \frac{v_2}{v_1} F_{v_1, v_2}$.

Cauchy (θ, σ): pmf: $f(x) = \frac{1}{\pi \sigma} \frac{1}{1+(\frac{x-\theta}{\sigma})^2}, -\infty < x < \infty, \theta \in \mathbb{R}, \sigma > 0$. * Special case of t, when df=1. Also if $X, Y \sim i.i.d N(0, 1)$, then $\frac{X}{Y}$ is cauchy.

Uniform (a, b): pmf: $f(x) = \frac{1}{b-a}$ for $a \leq x \leq b$; mgf: $M(t) = e^{\frac{t(b-a)}{2}} - e^{\frac{t(a-b)}{2}} / (b-a)$. $\sigma^2 = (b-a)^2/12$. * If $a=0, b=1$, this is special case.

Bivariate normal: pmf: $f(x_1, y_1) = \frac{1}{2\pi\sigma_1\sigma_2} e^{-\frac{1}{2}(\frac{(x_1-\mu_1)^2}{\sigma_1^2} + \frac{(y_1-\mu_2)^2}{\sigma_2^2})}$ where $g = (\frac{x_1-\mu_1}{\sigma_1}, \frac{y_1-\mu_2}{\sigma_2})^T$. Bivariate trans. We start w/ $(x_1, y_1) \sim \mathcal{N}(0, I)$ and want (Y_1, Y_2) in terms of g . We make sure $x_1 = g_1$ and $y_1 = g_2$. Get covariance, $\text{cov}(x_1, y_1), \text{cov}(y_1, y_2), \text{cov}(x_1, y_2)$. old knowns: $\text{cov}(x_1, y_1) = \frac{1}{2}\sigma_1\sigma_2 \rho$; $\text{cov}(y_1, y_2) = \frac{1}{2}\sigma_2\sigma_2 \rho$; $\text{cov}(x_1, y_2) = \frac{1}{2}\sigma_1\sigma_2 \rho$.

Exponentials: $Ex^n = 0$ if $n < 0$ and odd; $Ex^n = \frac{1}{\Gamma(n)} \frac{1}{(1-\mu)^{n+1}}$ if $n < 0$ and even.

Commutativity: $P(\text{bankrupt}) = P(\text{bankrupt}) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$

$$A \cup B = B \cup A$$

$$A \cap B = B \cap A$$

Associativity:

$$A \vee (B \vee C) = (A \vee B) \vee C$$

$$A \wedge (B \wedge C) = (A \wedge B) \wedge C$$

Distributive laws:

$$A \wedge (B \vee C) = (A \wedge B) \vee (A \wedge C)$$

$$A \vee (B \wedge C) = (A \vee B) \wedge (A \vee C)$$

DeMorgan's laws:

$$(A \vee B)^c = A^c \wedge B^c$$

$$(A \wedge B)^c = A^c \vee B^c$$

Basic properties of $P(\cdot)$:

$$\begin{aligned} P(\emptyset) &= 0 \\ 0 &\leq P(A) \end{aligned}$$

$$\textcircled{3} P(A^c) = 1 - P(A)$$

$$\textcircled{4} \text{ If } A_1, A_2, \text{ then } P(A) \leq P(A_1 \cup A_2)$$

$$\textcircled{5} P(B \wedge A^c) = P(B) - P(A \wedge B)$$

$$\textcircled{6} P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

$$\hookrightarrow P(B \wedge A) \leq P(A) + P(B)$$

$$\hookrightarrow P(A^c \wedge B^c) \geq 1 - (P(A) + P(B))$$

Conditional probability and independence:

$$\text{Def'n: } P(A|B) = P(A)P(B)$$

↳ statistical independence

↳ from Bayes' w/ assumption that

given has no effect

↳ If A and B are independent, so are

• A and B^c

• A^c and B

• A^c and B^c

Def'n conditional prob: $P(A|B) = \frac{P(A \wedge B)}{P(B)}$

$$\hookrightarrow P(A|B) = P(B|A)P(A)$$

$$\hookrightarrow P(A|B) = P(B|A) \left(\frac{P(A)}{P(B)} \right) \text{ "Bayes"}$$

Series calculations: n terms in series (a, to an)

S_n = first n terms summed

d = diff between succ. terms

$$\text{Arithmetic: } a_n = a_1 + (n-1)d; a_i = \frac{a_{i-1} + a_{i+1}}{2}$$

$$S_n = \frac{a_1 + a_n}{2} \cdot n; S_n = \frac{2a_1 + (n-1)d}{2} \cdot n$$

Geometric: $a_n = a_1 \cdot q^{n-1}; a_i = \sqrt[n]{a_{i-1} \cdot a_{i+1}}$

$$S_n = \frac{a_1 - a_n}{q-1}; S_n = \frac{a_1(q^n - 1)}{q-1}$$

Powers of natural #s:

$$\sum_{k \geq 1} k = \frac{1}{2} n(n+1)$$

$$\sum_{k \geq 1} k^2 = \frac{1}{6} n(n+1)(2n+1)$$

$$\sum_{k \geq 1} k^3 = \frac{1}{4} n^2(n+1)^2$$

Let X have pmf $P_X(x) = \begin{cases} \frac{3!}{x!(3-x)!} \binom{3}{x}^2, & x=1, 2, 3 \\ 0, & \text{elsewhere} \end{cases}$

1st head appears odd, pay \$1; if even, win \$2. Game ends at 1st heads, not 1-to-1, so we split up and treat each

unique case differently.

X = # of heads until heads appears. 0, 1, 2, 3

Random point in unit circle:

X = distance to origin

Probability of point lying in a smaller circle C T S P(C) = area/C

For 0 < x < 1, center $\sqrt{x^2 + y^2}$ is

point lying in circle of radius x.

So $P(X \leq x) = \frac{\pi x^2}{\pi r^2} = x^2$, so cdf is

$$F_X(x) = \begin{cases} 0, & x < 0 \\ x^2, & 0 \leq x < 1 \\ 1, & x \geq 1 \end{cases}$$

pdf is $f_X(x) = \begin{cases} 2x, & 0 \leq x < 1 \\ 0, & \text{elsewhere} \end{cases}$

$$P\left(\frac{1}{4} \leq X \leq \frac{1}{2}\right) = \int_{\frac{1}{4}}^{\frac{1}{2}} 2x dx = \frac{3}{16}$$

Call to switchboard

$$f(x) = \begin{cases} \frac{e^{-x}}{x}, & x > 0 \\ 0, & \text{elsewhere} \end{cases}$$

$$P(x > 4) = \int_4^\infty e^{-x} dx$$

$$f(x) = \frac{1}{x} F(x) = \frac{1}{x} - e^{-x}$$

is monotonically decreasing

for $x > 0$

cdf of X defined by $F_X(x) = P(X \leq x)$ for all x

$F_X(x)$ is a cdf iff

$$\lim_{x \rightarrow -\infty} F_X(x) = 0 \text{ and } \lim_{x \rightarrow \infty} F_X(x) = 1$$

b) $F_X(x)$ is a nondecreasing function of x

c) $F_X(x)$ is right-continuous: for every x_0 ,

$$\lim_{x \rightarrow x_0^+} F_X(x) = F(x_0)$$

The following two statements are equivalent:

1) r.v. X and Y are identically distributed

2) $F_X(x) = F_Y(x)$ for every x

pmf of a discrete r.v. X given by $f_X(x) = P(X=x)$ for all x

pdf of a continuous r.v. X is the function that

satisfies $f_X(x) = \int_x^\infty f_X(t) dt$ for all x

$P(X \leq x) = F_X(x) = \int_{-\infty}^x f_X(t) dt$ for all x

Thus, $\frac{d}{dx} F_X(x) = f_X(x)$

$f_X(x)$ is a pdf or pmf iff:

a) $f_X(x) \geq 0$ for all x

b) $\int_{-\infty}^\infty f_X(x) dx = 1$ (pmf) OR

$\int_{-\infty}^\infty f_X(x) dx = 1$ (pdf)

To find marginal cdfs, $F_{X_1}(x_1) = P(X_1 \leq x_1, -\infty < X_2 < \infty)$

To find marginal pmfs, hold marginal variable

and sum over the other. joint w/ respect to

other r.v.

Marginals: To find marginal pmfs, hold

marginal variable

and sum over the other.

To find marginal pdfs, $f_{X_1}(x_1) = \int_{-\infty}^\infty f_{(X_1, X_2)}(x_1, x_2) dx_2$

To find marginal pmfs, $f_{X_1}(x_1) = \lim_{x_2 \rightarrow \infty} F_{(X_1, X_2)}(x_1, x_2)$

Expectation:

$$E(x) = \int_{-\infty}^\infty x f(x) dx \text{ OR } E(x) = \sum x_i p(x_i)$$

$$E(k_1 g_1(x) + k_2 g_2(x) + c) = k_1 E(g_1(x)) + k_2 E(g_2(x)) + c$$

Joint US: Suppose X is a continuous r.v. defined by pdf $f(x)$.

Let Y be a r.v. $Y = g(x)$ for some function.

$$E(Y) = \int_{-\infty}^\infty g(x) f(x) dx$$

Suppose (X_1, X_2) is a continuous random vector.

Let Y be a r.v. defined by $Y = g(X_1, X_2)$

for some real-valued function g (i.e. $R^2 \rightarrow R$)

Then $E(Y) = \int_{-\infty}^\infty g(x_1, x_2) f_{(X_1, X_2)}(x_1, x_2) dx_1 dx_2$

$E(k_1 Y_1 + k_2 Y_2) = k_1 E(Y_1) + k_2 E(Y_2)$

Variance: Let X be a r.v. with finite mean

$M = E(x)$ and such that $E((x-M)^2) = \sigma^2$ is finite.

Then the variance of X is $\sigma^2 = E(x^2) - M^2$

$$= E(x^2) - 2Mx + M^2 = E(x^2) - 2M^2 + M^2 + M^2 = E(x^2) - M^2$$

Special power series: $\frac{1}{1-x} = 1 + x + x^2 + \dots$

$$\frac{1}{1+x} = 1 - x + x^2 - x^3 + \dots$$

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$$

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

$$\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$$

$$\tan x = x + \frac{x^3}{3} + \frac{x^5}{15} + \dots$$

From earlier results:

$$\text{Var}(X_2 | X_1) = E(X_2^2 | X_1) - (E(X_2 | X_1))^2$$

$$E(u(X_2) | X_1) = \int_{-\infty}^\infty u(x_2) f_{(X_1, X_2)}(x_1, x_2) dx_2$$

Finally: Let (X_1, X_2) be a random vector such that

$\text{Var}(X_2)$ is finite.

$$\hookrightarrow E(E(X_2 | X_1)) = E(X_2)$$

$$\hookrightarrow \text{Var}(E(X_2 | X_1)) \leq \text{Var}(X_2)$$

Derivatives: $\frac{d}{dx} (cf(x)) = cf'(x)$

$$(f(x) + g(x))' = f'(x) + g'(x)$$

$$\frac{d}{dx} (x^n) = nx^{n-1}; \frac{d}{dx} (c) = 0$$

$$\frac{d}{dx} (\ln(x)) = \frac{1}{x}$$

$$\frac{d}{dx} (\sin x) = \cos x; \frac{d}{dx} \cos x = -\sin x; \frac{d}{dx} \tan x = \sec^2 x$$

$$\frac{d}{dx} \sec x = \sec x \tan x; \frac{d}{dx} \csc x = -\csc x \cot x; \frac{d}{dx} \cot x = -\operatorname{cosec}^2 x$$

$$\frac{d}{dx} \sin^{-1} x = \frac{1}{\sqrt{1-x^2}}; \frac{d}{dx} \cos^{-1} x = -\frac{1}{\sqrt{1-x^2}}; \frac{d}{dx} \tan^{-1} x = \frac{1}{1+x^2}$$

$$\frac{d}{dx} (\ln(x)) = \frac{1}{x}; x > 0; \frac{d}{dx} (\ln|x|) = \frac{1}{x}, x \neq 0$$

$$\frac{d}{dx} (\log_a(x)) = \frac{1}{x \ln a}, x > 0$$

Integrals:

$$\int x^{-n} dx = \left(-\frac{1}{n-1}\right)x^{-n+1}, n \neq 1$$

$$\int \frac{1}{x+b} dx = \frac{1}{a} \ln|x+b|; \int x^p dx = \frac{x^{p+1}}{p+1}$$

$$\int a^x du = a^x; \int a^x du = \frac{a^x}{\ln a}; \int \ln u du = u \ln u - u$$

$$\int u^x du = (u-1)u^x + \int u^{x-1} du = u^x \ln u$$

Whenever you divide or multiply by a negative number, you must flip the inequality sign!

AUB = BUA
 ANB = BNA
 Associativity:
 AU(BVC) = (AUB)UC
 AN(BNC) = (ANB)NC
 Distributive laws:
 AN(BUC) = (ANB)U(ANC)
 AU(BNC) = (AUB)N(ANC)
 DeMorgan's laws:
 (AUB)' = A' ∩ B'
 (ANB)' = A' ∪ B'
 Basic properties of P(C):
 ① P(C) = 0
 ② 0 ≤ P(A) ≤ 1
 ③ P(A') = 1 - P(A)
 ④ If A, C, A₂, then P(A) ≤ P(A₂)
 ⑤ P(B ∩ A^c) = P(B) - P(A ∩ B)
 ⑥ P(A ∪ B) = P(A) + P(B) - P(A ∩ B)
 ↳ P(B ∪ A) ≤ P(A) + P(B) "Boole's"
 ↳ P(A' ∩ B') ≥ 1 - (P(A) + P(B))
 Conditional probability and independence:
 Def'n: P(A ∩ B) = P(A)P(B) means statistical independence
 ↳ from Bayes' w/ assumption that given has no effect
 ↳ If A and B are independent, so are A' and B'
 ↳ A and B
 ↳ A' and B'
 Defn conditional prob: P(A|B) = $\frac{P(A \cap B)}{P(B)}$
 ↳ P(A ∩ B) = P(A|B)P(B)
 ↳ P(A ∩ B) = P(B|A)P(A)
 ↳ P(A|B) = P(B|A) $\left(\frac{P(A)}{P(B)}\right)$ "Bayes"

Series calculations: n terms in series (a₁ to a_n)
 S_n = sum of n terms summed
 d = diff between succ. terms
 Arithmetic: a_n = a₁ + (n-1)d; a_i = $\frac{a_{i-1} + a_{i+1}}{2}$
 S_n = $\frac{a_1 + a_n}{2} \cdot n$; S_n = $\frac{2a_1 + (n-1)d}{2} \cdot n$
 Geometric: a_n = a₁ * qⁿ⁻¹; a_i = $\sqrt[n]{a_1 \cdots a_{i-1} a_{i+1}}$
 S_n = $\frac{a_1 - a_n}{q-1}$; S_n = $\frac{a_1(q^n - 1)}{q-1}$
 Powers of natural #s:
 $\sum_{k=1}^{\infty} k = \frac{1}{2} n(n+1)$.
 $\sum_{k=1}^{\infty} k^2 = \frac{1}{6} n(n+1)(2n+1)$
 $\sum_{k=1}^{\infty} k^3 = \frac{1}{4} n^2(n+1)^2$
 Let X have pmf p_x(x) = $\begin{cases} \frac{3!}{x!(3-x)!} \binom{3}{x}^2, & x=1, 2, 3 \\ 0, & \text{elsewhere} \end{cases}$
 1st head appears odd, pay 1/2, if even, wins \$1, -\$3. Dist. Given by p_y(y) = $\begin{cases} \frac{1}{4} \text{ for } y=1, 3, 7, 11, \\ \frac{1}{2} \text{ for } y=2, 4, 6, 8, 10, \\ 0 \text{ elsewhere} \end{cases}$
 Now if Z is (X-2)² from game above, not 1-to-1, so we split up and treat each case differently.
 X = # of heads until 2nd appears. Odd is $\frac{1}{2}$. P_x(Z) = $\begin{cases} P_x(z=2) = \frac{1}{4} \text{ for } z=0, 4, 12, \dots \\ P_x(z=1) + P_x(z=3) = \frac{1}{2} \text{ for } z=1, 5, 9, \dots \\ P_x(z=2) = \frac{1}{4} \text{ for } z=2, 6, 10, \dots \end{cases}$
 random point in unit circle: X = distance to origin
 Probability of point lying in a smaller circle C is P(C) = $\frac{\pi r^2}{\pi R^2}$. For 0 < x < 1, point (x, x²) is point lying in circle of radius x. So P(x < x) = $\frac{\pi x^2}{\pi R^2} = x^2$, so cdf is F_x(x) = $\begin{cases} 0, & x \leq 0 \\ x^2, & 0 < x \leq 1 \\ 1, & x > 1 \end{cases}$
 pdf is f_x(x) = $\begin{cases} 2x, & 0 < x \leq 1 \\ 0, & \text{elsewhere} \end{cases}$
 P($\frac{1}{4} < x \leq \frac{1}{2}$) = $\int_{1/4}^{1/2} 2x dx = \frac{3}{16}$
 Call to switchboard: $\begin{cases} \frac{e^{-4y}}{4}, & 0 < y \leq 0 \\ 0, & \text{elsewhere} \end{cases}$
 P(x > y) = $\int_y^{\infty} f(x) dx$

CDF of X defined by F_x(x) = P(X ≤ x) for all x
 F_x(x) is a cdf iff
 a) $\lim_{x \rightarrow -\infty} F_x(x) = 0$ and $\lim_{x \rightarrow \infty} F_x(x) = 1$
 b) F_x(x) is a nondecreasing function of x
 c) F_x(x) is right-continuous: for every x₀, $\lim_{x \downarrow x_0} F_x(x) = F(x_0)$
 The following two statements are equivalent:
 1) r.v.s X and Y are identically distributed
 2) F_x(x) = F_y(x) for every x
 PMF of a discrete r.v. X given by f_x(x) = P(X=x) for all x
 PDF of a continuous r.v. X is the function that satisfies F_x(x) = $\int_{-\infty}^x f_x(t) dt$ for all x
 P(X ≤ x) = F_x(x) = $\int_x^{\infty} f_x(t) dt$ for all x
 Thus, $\frac{d}{dx} F_x(x) = f_x(x)$
 f_x(x) is a pdf or pmf iff:
 a) f_x(x) ≥ 0 for all x
 b) $\int_{-\infty}^{\infty} f_x(x) dx = 1$ (pmf) OR
 $\int_{-\infty}^{\infty} f_x(x) dx = 1$ (pdf)

Multivariate:
 Joint cdf: P(a₁ ≤ X₁ ≤ b₁, a₂ ≤ X₂ ≤ b₂) = F_{x₁x₂}(b₁, b₂) - F_{x₁x₂}(a₁, b₂) - F_{x₁x₂}(b₁, a₂) + F_{x₁x₂}(a₁, a₂)
 Joint pmf: p_{x₁x₂}(x₁, x₂) = P(X₁ = x₁, X₂ = x₂)
 Joint pdf: F_{x₁x₂}(x₁, x₂) = $\int_{-\infty}^{x_1} \int_{-\infty}^{x_2} f_{x_1x_2}(w_1, w_2) dw_1 dw_2$
 $\int_{-\infty}^{x_1} \int_{-\infty}^{x_2} f_{x_1x_2}(w_1, w_2) dw_1 dw_2 = f_{x_1x_2}(x_1, x_2)$
 Note that P(X₁, X₂) ∈ A = $\iint_A f_{x_1x_2}(x_1, x_2) dx_1 dx_2$, which is just the volume under the surface z = f_{x₁x₂}(x₁, x₂)
 Marginals: To find marginal pmfs, hold marginal variable and sum over the other. Joint wrt to other r.v.
 To find marginal pdfs, f_{x_i}(x_i) = $\int_{-\infty}^{\infty} f_{x_1x_2}(x_1, x_2) dx_2$
 Conditional distributions and expect.: $= \lim_{x_2 \rightarrow \infty} F(x_1, x_2)$
 Discrete: P(X₂ = x₂ | X₁ = x₁) = $P(X_1 = x_1, X_2 = x_2) / P(X_1 = x_1)$
 $= \frac{P(X_1 = x_1, X_2 = x_2)}{P(X_1 = x_1)} \rightarrow \text{conditional pmf}$
 $P_{x_1}(x_1) = P_{x_2|x_1}(x_2|x_1) \rightarrow$
 Continuous: $f_{x_2|x_1}(x_2|x_1) = \frac{f_{x_1x_2}(x_1, x_2)}{f_{x_1}(x_1)}$
 Note: $\int_{-\infty}^{\infty} f_{x_2|x_1}(x_2|x_1) dx_2 = \int_{-\infty}^{\infty} \frac{f_{x_1x_2}(x_1, x_2)}{f_{x_1}(x_1)} dx_2$
 $= \frac{1}{f_{x_1}(x_1)} \int_{-\infty}^{\infty} f_{x_1x_2}(x_1, x_2) dx_2 = \frac{f_{x_1}(x_1)}{f_{x_1}(x_1)} = 1$
 P(a \leq X₂ \leq b | X₁ = x₁) = $\int_a^b f_{x_2|x_1}(x_2|x_1) dx_2$
 P(c \leq X₁ \leq d | X₂ = x₂) = $\int_c^d f_{x_1|x_2}(x_1|x_2) dx_1$
 If u(X₂) is a function of X₂, the conditional expectation of u(X₂), given that X₁ = x₁, is E(u(X₂) | x₁) = $\int_{-\infty}^{\infty} u(x_2) f_{x_2|x_1}(x_2|x_1) dx_2$
 Conditional mean: E(X₂ | x₁)
 From earlier results:
 • Var(X₂ | x₁) = E(X₂² | x₁) - (E(X₂ | x₁))²
 • E(u(x₁) | x₂) = $\int_{-\infty}^{\infty} u(x_1) f_{x_1|x_2}(x_1|x_2) dx_1$
 Finally: Let (X₁, X₂) be a random vector such that Var(X₂) is finite.
 ↳ E(E(X₂ | x₁)) = E(X₂)
 ↳ Var(E(X₂ | x₁)) ≤ Var(X₂)

Special power series: $\frac{1}{1-x} = 1 + xt + xt^2 + \dots$
 $\frac{1}{1+x} = 1 - xt + xt^2 - xt^3 + \dots$
 $\ln(1+x) = -\frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$
 $\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$
 $\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$
 $\tan x = x + \frac{2x^3}{3} + \frac{2x^5}{15} + \dots$
 R_n = $\frac{f^{(n)}(\xi)(x-a)^n}{n!}$
 Taylor and MacLaurin:
 f(x) = f(a) + f'(a)(x-a) + $\frac{f''(a)(x-a)^2}{2!} + \dots + \frac{f^{(n-1)}(a)(x-a)^{n-1}}{(n-1)!} + R_n$
 Binomial series:
 (a+x)ⁿ = $a^n + \binom{n}{1} a^{n-1} x + \binom{n}{2} a^{n-2} x^2 + \binom{n}{3} a^{n-3} x^3 + \dots$
 We seek pmf p_y(y) of r.v. Y = X². Transformation y = g(x) = x² maps onto y = y^2 for y = 0, 1, 4, 9. In general, not a one-to-one, but here yes. x = ±y, so p_y(y) = p_x(±y) = $\frac{1}{4}$ for y = 0, 1, 4, 9, 16, ...
 A lot has 100 fuses, five to be tested. If all 5 are good, let accept. P_y(y) for y = 0, 1, 4, 9, 16, ...
 P_y(y) = $\begin{cases} \frac{1}{400}, & y=0 \\ \frac{1}{200}, & y=1 \\ \frac{1}{100}, & y=4 \\ \frac{1}{50}, & y=9 \\ \frac{1}{25}, & y=16 \end{cases}$
 Derivatives: $\frac{d}{dx} (cf(x)) = cf'(x)$
 $\frac{d}{dx} (f(x) + g(x)) = f'(x) + g'(x)$
 $\frac{d}{dx} (x^n) = nx^{n-1}$; $\frac{d}{dx} (c) = 0$
 $\frac{d}{dx} (\sin x) = \cos x$; $\frac{d}{dx} (\cos x) = -\sin x$; $\frac{d}{dx} (\tan x) = \sec^2 x$
 $\frac{d}{dx} (\sec x) = \sec x \tan x$; $\frac{d}{dx} (\csc x) = -\csc x \cot x$
 $\frac{d}{dx} (\sec^{-1} x) = \frac{1}{x\sqrt{1-x^2}}$; $\frac{d}{dx} (\csc^{-1} x) = -\frac{1}{x\sqrt{1-x^2}}$
 $\frac{d}{dx} (\ln x) = \frac{1}{x}$, x > 0; $\frac{d}{dx} (\ln|x|) = \frac{1}{x}$, x ≠ 0
 $\frac{d}{dx} (\log a x) = \frac{1}{x \ln a}$, x > 0
 Integrals:
 $\int x^{-n} dx = (-\frac{1}{n+1}) x^{-n-1}$, n ≠ 1
 $\int \frac{1}{x^{a+b}} dx = \frac{1}{a} \ln|x|^{a+b}$; $\int x^{a/b} dx = \frac{b}{a+1} x^{\frac{a+1}{b}}$
 $\int a^u du = e^u$; $\int a^{u/a} du = \frac{a}{\ln a}$; $\int a^{u-a} du = a^{u-a}$
 $\int e^{u-a} du = (u-a)e^u$; $\int \frac{1}{a^{u-a}} du = \ln|a|^{u-a}$
 Whenever you divide or multiply by a negative number, you must flip the inequality sign!

Commutativity: $A \cup B = B \cup A$; $P(\text{turnhouse}) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2}$
 Associativity: $A \cup (B \cup C) = (A \cup B) \cup C$
 $A \cap (B \cap C) = (A \cap B) \cap C$
 Distributive laws:
 $A \cup (B \cap C) = (A \cup B) \cup (A \cup C)$
 $A \cap (B \cup C) = (A \cap B) \cap (A \cap C)$
 DeMorgan's laws:
 $(A \cup B)^c = A^c \cap B^c$
 $(A \cap B)^c = A^c \cup B^c$
 Basic properties of $P(\cdot)$:
 ① $P(\emptyset) = 0$
 ② $0 \leq P(A) \leq 1$
 ③ $P(A^c) = 1 - P(A)$
 ④ If A_1, A_2, \dots , then $P(A) \leq P(A_1 \cup A_2 \cup \dots)$
 ⑤ $P(B \cap A^c) = P(B) - P(A \cap B)$
 ⑥ $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 $\hookrightarrow P(B \cup A) \leq P(A) + P(B)$ "Boole's"
 $\hookrightarrow P(A^c \cap B^c) \geq 1 - (P(A) + P(B))$
 Conditional probability and independence:
 Def'n: $P(A|B) = P(A)P(B)$ means statistical independence
 Log from Bayes' w/ assumption that given has no effect
 \hookrightarrow If A and B are independent, so are A and BC, A^c and B, A^c and BC.
 Def'n conditional prob: $P(A|B) = \frac{P(A \cap B)}{P(B)}$
 $\hookrightarrow P(A \cap B) = P(A|B)P(B)$
 $\hookrightarrow P(A \cap B) = P(B|A)P(A)$
 $\hookrightarrow P(A|B) = P(B|A) \left(\frac{P(A)}{P(B)} \right)$ "Bayes"

Series calculations: n terms in series (a₁ to a_n)
 $S_n = \text{first } n \text{ terms summed}$
 $d = \text{diff between succ. terms}$
 Arithmetic: $a_n = a_1 + (n-1)d$; $a_i = \frac{a_1 + a_i + d}{2}$
 $S_n = \frac{a_1 + a_n}{2} \cdot n; S_n = \frac{2a_1 + (n-1)d}{2} \cdot n$
 Geometric: $a_n = a_1 \cdot q^{n-1}; a_i = \sqrt[n]{a_{i+1} \cdot a_{i+1}}$
 $S_n = \frac{a_1 - a_n}{q-1}; S_n = \frac{a_1(q^n - 1)}{q-1}$
 Powers of natural #s:
 $\sum_{k=1}^n k = \frac{1}{2} n(n+1)$.
 $\sum_{k=1}^n k^2 = \frac{1}{6} n(n+1)(2n+1)$
 $\sum_{k=1}^n k^3 = \frac{1}{4} n^2 (n+1)^2$
 Let X have pmf $p_X(x) = \begin{cases} \frac{3!}{x!(3-x)!} \binom{1/3}{x}^{3-x}, & x=1, 2, 3 \\ 0, & \text{elsewhere} \end{cases}$
 1st head appears odd, pay H1, if even with H1. Same game about H2, H3. Now if Z is $(X-2)^2$ from above, not $1-10-1$, so we split up and treat each case differently. X is # H1s until heads appears. $O \leq z \leq 3$. $P_Z(z) = \begin{cases} P_X(y)=y, & \text{for } z=0 \\ P_X(y)+P_X(y)=\frac{1}{2}, & \text{for } z=1 \\ P_X(y)+P_X(y)=\frac{1}{2}, & \text{for } z=2 \\ P_X(y)+P_X(y)=\frac{1}{2}, & \text{for } z=3 \end{cases}$

Random point in unit circle:
 $X = \text{distance to origin}$
 Probability of point lying in a smaller circle C is $P(C) = \frac{\text{area}}{\pi}$. For $0 \leq x \leq 1$, event $X \leq x$ is point lying in circle of radius x. So $P(X \leq x) = \frac{\pi x^2}{\pi} = x^2$, so cdf is $F_X(x) = \begin{cases} 0, & x < 0 \\ x^2, & 0 \leq x \leq 1 \\ 1, & x > 1 \end{cases}$
 pdf is $f_X(x) = \begin{cases} 2x, & 0 \leq x \leq 1 \\ 0, & \text{elsewhere} \end{cases}$
 $P(\frac{1}{4} \leq X \leq \frac{1}{2}) = \int_{1/4}^{1/2} 2x dx = \frac{3}{16}$

Call it to switchboard: $f(x) = \begin{cases} e^{-xy}, & 0 \leq x \leq 200 \\ 0, & \text{elsewhere} \end{cases}$

cdf of X defined by $F_X(x) = P(X \leq x)$ for all x
 $F_X(x)$ is a cdf iff
 a) $\lim_{x \rightarrow -\infty} F_X(x) = 0$ and $\lim_{x \rightarrow \infty} F_X(x) = 1$
 b) $F_X(x)$ is a nondecreasing function of x
 c) $F_X(x)$ is right-continuous: for every x_0 , $\lim_{x \rightarrow x_0^+} F_X(x) = F_X(x_0)$
 The following two statements are equivalent:
 1) r.v. X and Y are identically distributed
 2) $F_X(x) = F_Y(x)$ for every x
 PMF of a discrete r.v. X given by $f_X(x) = P(X=x)$
 pdf of a continuous r.v. X is the function that satisfies $F_X(x) = \int_{-\infty}^x f_X(t) dt$ for all x
 $P(X \leq x) = F_X(x) = \int_{-\infty}^x f_X(t) dt$ for all x
 Thus, $\frac{d}{dx} F_X(x) = f_X(x)$
 $f_X(x)$ is a pdf or pmf iff:
 a) $f_X(x) \geq 0$ for all x
 b) $\sum_x f_X(x) = 1$ (pmf) OR $\int_{-\infty}^{\infty} f_X(x) dx = 1$ (pdf)

Multivariate:
 Joint cdf: $P(a_1 \leq X_1 \leq b_1, a_2 \leq X_2 \leq b_2) = F_{X_1, X_2}(b_1, b_2) - F_{X_1, X_2}(a_1, b_2) - F_{X_1, X_2}(a_1, a_2) + F_{X_1, X_2}(a_1, a_2)$
 Joint pmf: $p_{X_1, X_2}(x_1, x_2) = P(X_1 = x_1, X_2 = x_2)$
 Joint pdf: $f_{X_1, X_2}(x_1, x_2) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} f_{X_1, X_2}(w_1, w_2) dw_1 dw_2$
 $\frac{\partial^2}{\partial x_1 \partial x_2} f_{X_1, X_2}(x_1, x_2) = f_{X_1, X_2}(x_1, x_2)$
 Note that $P((X_1, X_2) \in A) = \iint_A f_{X_1, X_2}(x_1, x_2) dx_1 dx_2$, which is just the volume under the surface $z = f_{X_1, X_2}(x_1, x_2)$
 Marginals: To find marginal pmfs, hold marginal variable and sum over the other. Marginal w/ respect to other r.v.
 To find marginal pdfs, $f_{X_1}(x_1) = \int_{-\infty}^{\infty} f_{X_1, X_2}(x_1, x_2) dx_2$
 To find marginal cdfs, $F_{X_1}(x_1) = P(X_1 \leq x_1, -\infty < X_2 < \infty)$
 Discrete: $P(X_2 = x_2 | X_1 = x_1) = \lim_{x_2 \rightarrow \infty} F_{X_1, X_2}(x_1, x_2)$
 $= P(X_1 = x_1, X_2 = x_2)$
 $= p_{X_1, X_2}(x_1, x_2)$ conditional pmf
 $P_{X_1}(x_1) = p_{X_2|x_1}(x_2 | x_1)$
 Continuous: $f_{X_2|x_1}(x_2 | x_1) = \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_1}(x_1)}$
 Note: $\int_{-\infty}^{\infty} f_{X_2|x_1}(x_2 | x_1) dx_2 = \int_{-\infty}^{\infty} \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_1}(x_1)} dx_2 = \frac{1}{f_{X_1}(x_1)} \int_{-\infty}^{\infty} f_{X_1, X_2}(x_1, x_2) dx_2 = \frac{f_{X_1}(x_1)}{f_{X_1}(x_1)} = 1$
 $P(a \leq X_2 \leq b | X_1 = x_1) = \int_a^b f_{X_2|x_1}(x_2 | x_1) dx_2$
 $P(c < X_1 < d | X_2 = x_2) = \int_c^d f_{X_1|x_2}(x_1 | x_2) dx_1$
 If $u(X_2)$ is a function of X_2 , the conditional expectation of $u(X_2)$, given that $X_1 = x_1$, is
 $E(u(X_2) | x_1) = \int_{-\infty}^{\infty} u(x_2) f_{X_2|x_1}(x_2 | x_1) dx_2$
 Conditional mean: $E(X_2 | x_1)$
 From earlier results:
 • $\text{Var}(X_2 | x_1) = E(X_2^2 | x_1) - (E(X_2 | x_1))^2$
 • $E(u(x_1) | x_2) = \int_{-\infty}^{\infty} u(x_1) f_{X_1|x_2}(x_1 | x_2) dx_1$
 Finally: Let (X_1, X_2) be a random vector such that $\text{Var}(X_2)$ is finite.
 $\hookrightarrow E(E(X_2 | X_1)) = E(X_2)$
 $\hookrightarrow \text{Var}(E(X_2 | X_1)) \leq \text{Var}(X_2)$
 Derivatives: $\frac{d}{dx} (cf(x)) = c f'(x)$
 $(f(x) \pm g(x))' = f'(x) \pm g'(x)$
 $\frac{d}{dx} (x^n) = nx^{n-1}; \frac{d}{dx} (c) = 0$
 $\frac{d}{dx} (f(g(x))) = f'(g(x))g'(x)$
 $\frac{d}{dx} \sin x = \cos x; \frac{d}{dx} \cos x = -\sin x; \frac{d}{dx} \tan x = \sec^2 x$
 $\frac{d}{dx} \sec x = \sec x \tan x; \frac{d}{dx} \csc x = -\csc x \cot x; \frac{d}{dx} \cot x = \frac{1}{\sin^2 x}$
 $\frac{d}{dx} \sin^{-1} x = \frac{1}{\sqrt{1-x^2}}; \frac{d}{dx} \cos^{-1} x = -\frac{1}{\sqrt{1-x^2}}; \frac{d}{dx} \tan^{-1} x = \frac{1}{1+x^2}$
 $\frac{d}{dx} \sec^{-1} x = \frac{1}{|x|\sqrt{x^2-1}}; \frac{d}{dx} \csc^{-1} x = \frac{1}{|x|\sqrt{1-x^2}}; \frac{d}{dx} \cot^{-1} x = \frac{1}{1+x^2}$
 $\frac{d}{dx} (a^x) = a^x \ln(a); \frac{d}{dx} e^x = e^x$
 $\frac{d}{dx} (\ln(x)) = \frac{1}{x}, x > 0; \frac{d}{dx} (\ln|x|) = \frac{1}{x}, x \neq 0$
 $\frac{d}{dx} (\log_a(x)) = \frac{1}{x \ln(a)}, x > 0$
 Integrals:
 $\int x^{-n} dx = \left(\frac{1}{n-1} \right) x^{-n+1}, n \neq 1$
 $\int \frac{1}{x^2} dx = \frac{1}{x} \ln|x|; \int x^p dx = \frac{1}{p+1} x^{p+1}$
 $\int a^x dx = a^x \frac{1}{\ln(a)}; \int \frac{1}{a^x} dx = \frac{1}{\ln(a)} a^x$
 $\int \ln x dx = x \ln x - x; \int \ln u du = u \ln u - u$
 $\int \ln u du = (u-1)e^u; \int \frac{1}{u} du = \ln|u|$
 Whenever you divide or multiply by a negative number, you must flip the inequality sign!