

General conventions: functions as `function()`, packages as `{package}` and objects as `object`

Data
<p>General Workflow</p> <ol style="list-style-type: none"> 1. Initialize a Project. 2. Open RMarkdown File and save in project directory. 3. Load packages you'll need in first chunk. <ul style="list-style-type: none"> ◦ e.g., <code>{tidyverse}</code>, <code>{rio}</code>, <code>{ggeffects}</code>, <code>{car}</code> 4. Write prose and R code in RMarkdown file. <p>Add-on Packages</p> <ul style="list-style-type: none"> • <code>install.packages(<package name>)</code> downloads the package from the web to your computer - only need to do once per R version. • <code>library(<package name>)</code> makes the package functions available for your R session - need to do once per R session. <p>Importing Data</p> <p>Use the <code>{rio}</code> package to import data from Stata, SPSS, etc...:</p> <ul style="list-style-type: none"> • <code>dat <- import(<path>)</code> where <code><path></code> is the path to a file. • <code>export(dat, <filename>)</code> will export the data to another proprietary format (Stata, SPSS, etc...). <p>Data Types</p> <ul style="list-style-type: none"> • Numeric - numbers without labels • Factors - categorical variables - numbers with labels. <ul style="list-style-type: none"> ◦ <code>factor()</code> and <code>as.factor()</code> can turn numeric and character data into factors. ◦ <code>factorize()</code> from the <code>{rio}</code> package will take labels read in with the data from Stata or SPSS and apply them to the variables. • Strings - character strings with no attached numerical information. <p>Investigating Properties of Data For a dataset named <code>data</code> and a variable <code>x</code>.</p> <ul style="list-style-type: none"> • <code>str(data)</code> - shows the properties of all the variables in <code>data</code>. • <code>str(data\$x)</code> - shows the properties of the <code>x</code> variable in the <code>data</code>. • <code>names(data)</code> - shows the names of all the variables in <code>data</code>. • <code>head(data)</code> - shows the first six rows of all the variables in <code>data</code>. • <code>tail(data)</code> - shows the last six rows of all the variables in <code>data</code>. <p>{dplyr} Operations</p> <ul style="list-style-type: none"> • <code>%>%</code> pipe chains functions together. • <code>select()</code> chooses columns based on logical expression, name or number. • <code>filter()</code> chooses observations based on logical expression. • <code>group_by()</code> groups observations so operations are performed within groups. • <code>mutate()</code> adds or modifies variables to existing dataset. • <code>summarise()</code> collapses/aggregates across observations (within group). <ul style="list-style-type: none"> ◦ common operations: <code>first()</code>, <code>last()</code>, <code>mean()</code>, <code>sd()</code>, <code>median()</code>, <code>min()</code>, <code>max()</code>, <code>quantile()</code>, <code>n()</code>

Data, Viz and Linear Models
<p>Recoding Variables</p> <ul style="list-style-type: none"> • Using the <code>recode()</code> function from the <code>{car}</code> package: <ul style="list-style-type: none"> ◦ <code>recode(x, <recode statements>)</code> ◦ <code>lo</code> and <code>hi</code> stand in for the minimum and maximum of the variable, <code>else</code> is short for all other values not explicitly mentioned. ◦ Example: <code>college_ed = recode(educ, "lo:12='None'; 13:15='Some'; 16:hi='Degree'", as.factor=TRUE, levels=c("None", "Some", "Degree"))</code> • Using <code>case_when()</code> from the <code>{dplyr}</code> package. <ul style="list-style-type: none"> ◦ Example: <code>college_ed = case_when(educ <= 12 ~ 'None', educ > 12 & educ < 16 ~ 'Some', educ >= 16 ~ 'Degree', TRUE ~ NA_character_)</code> <p>Data Viz</p> <p>In the <code>{ggplot2}</code> framework, data can be mapped to the following aesthetics - <code>shape</code>, <code>colour</code>, <code>fill</code>, <code>alpha</code>, <code>linetype</code>, <code>size</code> and the properties controlled with the <code>scale_*_manual()</code> functions.</p> <ul style="list-style-type: none"> • <code>geom_bar()</code> to make bar plots, use <code>stat="identity"</code> if you are supplying the height of the bar with <code>y</code> and <code>position=position_dodge()</code> for side-by-side bar plots. • <code>geom_line()</code> for lines connecting points. • <code>geom_histogram()</code> for histograms, use <code>position="identity"</code> for overlapping (superposed) histograms. • <code>geom_tile()</code> for heatmap plots, you can add the text onto the tiles with <code>geom_text()</code>. • <code>geom_mosaic()</code> from the <code>{ggsomaic}</code> package for making mosaic plots. <p>Other useful functions:</p> <ul style="list-style-type: none"> • <code>coord_flip()</code> to swap <code>x</code> and <code>y</code> axes. • <code>facet_wrap(~group)</code> to make different panels for each value of <code>group</code>. • <code>theme(axis.text.x=element_text(rot=45, hjust=1))</code> to rotate <code>x</code>-axis labels 45 degrees. <p>Model Formulas</p> <ul style="list-style-type: none"> • Generally take the form of <code><dependent variable> ~ <independent variables></code> <ul style="list-style-type: none"> ◦ <code>x1 + x2</code> - additive relationship. ◦ <code>x1 * x2</code> - multiplicative including main effects. ◦ <code>x1:x2</code> - multiplicative excluding main effects. ◦ <code>log(x)</code> - natural logarithm of <code>x</code>. ◦ <code>I(x^p)</code> - raise <code>x</code> to the <code>p</code> power. ◦ <code>poly(x, p)</code> - include a <code>p</code> degree polynomial in <code>x</code>.

(Generalized) Linear Models
<p>Linear Models</p> <ul style="list-style-type: none"> • Estimate with <code>model <- lm(formula, data=data)</code> • Presentation <ul style="list-style-type: none"> ◦ <code>stargazer(model)</code> from the <code>{stargazer}</code> package to make publication-ready tables. ◦ <code>ggpredict(model)</code> from the <code>{ggeffects}</code> package for making predicted effects plots. ◦ <code>intQualQuant()</code> from the <code>{DAMisc}</code> package for categorical-continuous interactions. <ul style="list-style-type: none"> ▪ <code>type='slopes'</code> with <code>plot=FALSE</code> prints the simple slopes. ◦ <code>DAintfun2()</code> and <code>DAintfun()</code> from the <code>{DAMisc}</code> package for continuous-continuous interaction conditional effect graphs. <ul style="list-style-type: none"> ▪ <code>changeSig()</code> to identify changes in the significance of conditional effects. • Diagnostics <ul style="list-style-type: none"> ◦ Linearity/Functional Form <code>crPlots(model)</code> from the <code>{car}</code> package. ◦ Heteroskedasticity - <code>ncvTest(model)</code> from the <code>{car}</code> package or residual vs fitted plot. ◦ Normality of Residuals - <code>shapiro.test(model\$residuals)</code> or density plot of residuals. ◦ Outliers and Influential Observations - <code>InfluencePlot(model)</code> from the <code>{car}</code> package. <ul style="list-style-type: none"> ▪ <code>Leverage - hatvalues(model)</code> ▪ <code>Discrepancy - rstudent(model)</code> or <code>rstandard(model)</code> ▪ <code>Cook's D - cooks.distance(model)</code> ▪ <code>DFBETAS - dfbetas(model)</code> <p>Generalized Linear Models</p> <ul style="list-style-type: none"> • Estimate with <code>glm(formula, data, family)</code> • Model fit <ul style="list-style-type: none"> ◦ <code>binfit()</code> from the <code>{DAMisc}</code> package for pseudo-(R^2) measures. ◦ <code>pre()</code> from the <code>{DAMisc}</code> package to calculate the proportional reduction in error. ◦ <code>AIC()</code> and <code>BIC()</code> in base R to calculate information criterion measures. ◦ <code>logLik()</code> to display the log-likelihood. ◦ <code>lrtest()</code> from the <code>{lmtest}</code> or <code>anova()</code> from base R to test nested models. ◦ <code>clarke_test()</code> from the <code>{clarkeTest}</code> package to test non-nested models.

General conventions: functions as `function()`, packages as `{package}` and objects as `object`

GLMs and Categorical Models

Generalized Linear Models

- Calculating Effects.
 - Odds Ratios for Logit: `exp(code(model))`
 - First difference at reasonable values - `glmChange()` from the `{DAMisc}` package.
 - Find "central" values of all values except x_j
 - Calculate predicted probabilities for $x_j = x_0 + \delta$ and $x_j = x_0$
 - Subtract the latter from the former
 - Average first difference - `glmChage2()` from the `{DAMisc}` package.
 - Calculate the predicted probabilities for all observations.
 - Calculate the predicted probabilities for all observations setting $x_j = x_j + \delta$
 - Subtract the former from the latter to get the individual first differences and take the average.
- Marginal Effect - `margins()` from the `{margins}` package.
 - Calculate the partial first difference of the predicted probabilities w.r.t. x_j for each observation.
 - Average across all of the marginal effects.
- Plotting Effects
 - At Reasonable Values - `ggpredict()` from the `{ggeffects}` package
 - Hold all variables except x_j at central values.
 - Calculate predicted probabilities at $x_j = \{min(x_j), \dots, max(x_j)\}$
 - Plot predicted probabilities on y against $\{min(x_j), \dots, max(x_j)\}$ on x .
 - Average Effect - `aveEffPlot()` from the `{DAMisc}` package.
 - Calculate predicted probability for every observation at each value of $x_j = \{min(x_j), \dots, max(x_j)\}$
 - For each value of $x_j = \{min(x_j), \dots, max(x_j)\}$, calculate the average predicted probability.
 - Plot the average predicted probabilities against $\{min(x_j), \dots, max(x_j)\}$.

GLM Interactions

- Two separable questions (that are inseparable in the linear model).
 - Is the product term needed?
 - Evaluate with statistical significance of interaction term(s).
 - Is there an interaction?
 - Calculate the second difference/derivative
 - `seconDiff()` from `{DAMisc}` calculates this.
 - For a cross-partial derivative (difference in first derivatives), you could do it by hand.

Categorical, Survey and Multilevel Models

Ordinal Models

- Estimate with `c1m()` from the `{ordinal}` package or `polr()` from the `{MASS}` package.
- `ordfit()` and `pre()` from the `{DAMisc}` package, `AIC()` and `BIC()` from base R as well as `lrtest()` from `{lmtest}` and `clarke_test()` from `{clarkeTest}` all work.
- `ordChange()` and `ordChange2()` from `{DAMisc}` calculate first differences using the 'reasonable values' and 'observed values' approaches, respectively.
- `ggpredict()` from `{ggeffects}` and `ordAveEffPlot()` from `{DAMisc}` plot effects for the 'reasonable values' and 'observed values' approaches, respectively.
- `nominal_test()` from `{ordinal}` for `c1m()` objects and `poTest()` from `{car}` for `polr()` objects do the Brant test.

Multinomial Models

- Estimate with `multinom()` from the `{nnet}` package or `mlogit()` from the `{mlogit}` package.
 - `mlogit()` also allows estimation of the conditional logit model.
- `mnlgit()` and `pre()` from the `{DAMisc}` package, `AIC()` and `BIC()` from base R as well as `lrtest()` from `{lmtest}` and `clarke_test()` from `{clarkeTest}` all work.
- `mnlChange()` and `mnlChange2()` from `{DAMisc}` calculate first differences using the 'reasonable values' and 'observed values' approaches, respectively.
- `ggpredict()` from `{ggeffects}` and `mnlAveEffPlot()` from `{DAMisc}` plot effects for the 'reasonable values' and 'observed values' approaches, respectively.

Complex Survey Data

- You can set the survey design properties with `svydesign()` from the `{survey}` package or `as_survey_design()` from the `{srvyr}` package.
 - Use `as_survey()` from `{srvyr}` to convert an object created with the `{survey}` package into a 'survey tibble' to be used with `{dplyr}` and other `{srvyr}` functions.
- `sumStats()` from the `{DAMisc}` package makes summary statistics (optionally by the values of other variables).
- `xt()` from `{DAMisc}` makes weighted cross-tabulations with measure of fit for survey data.
- `svyglm()` and `svyolr()` from `{survey}` will estimate GLM and ordinal models. `svymle()` will estimate any MLE model for complex survey data.
- `ggpredict()` from `{ggeffects}` will generate predicted probabilities as will `glmChange2()` and `probc1()` from the `{DAMisc}` package.
 - `glmChange2()` and `probc1()` use unweighted averages of individual first differences.

Multilevel and Measurement Models

Multilevel Models

- Estimated with `lmer()` or `glmer()` from `{lme4}` for (GL)Ms. For ordinal models, use `c1mm()` from `{ordinal}` and for multinomial models use `mlogit()` from `{mlogit}`.
 - For `lmer()` and `glmer()`, including in the formula `(1 | group)` will make random intercepts by group and `(1 + l2var | group)` will make random intercepts and random slopes for `l2var` by group.
- `coef()` will provide group-level intercepts and coefficients from the model.
- `ggpredict()` from `{ggeffects}` will graph either fixed overall effects or group-level effects - with confidence intervals that include the fixed-effect (and optionally random-effect) variances.
- `brm()` from `{brms}` will estimate a Bayesian version of the model using Stan.
 - NOTE: this could take a long time depending on the size of the model/data.

PCA/Exploratory FA

- `princomp()` from base R does principal components analysis.
- `ggbiplot()` from the `{ggbiplot}` package makes a biplot for the PCA.
- `fa()` from the `{psych}` package does EFA with many options.
 - `rotate=<rotation>` will rotate the solution see `?fa` for options.
- `scree()` from `{psych}` makes a scree plot for evaluating dimensionality.

CFA, SEM and Growth Models

- Estimated with `cfa()`, `sem()` or `growth()` from the `{lavaan}` package with primary operators:
 - `lv ~ ind1 + ind2` specifies that `ind1` and `ind2` are indicators of latent variable `lv`.
 - `dv ~ iv1 + iv2` specifies a regression of `dv` on `iv1` and `iv2`
 - `ac := a*c` specifies that the computed value `ac` is the product of `a` and `c`.
 - `ind1 ~ ind2` indicates that the covariance between `ind1` and `ind2` should be free and estimated by the model.
- `modificationIndices()` print the modification indices for `{lavaan}` models.
- `summary(model)` summarizes the model object with the following arguments:
 - `standardize=TRUE` will give the standardized solution.
 - `fit.measures=TRUE` give more scalar measures of fit for the model.
- `predict(model)` predicts the latent variables for the observed dataset used in the model.