# Plan

- Session 1
  - Workflow overview
  - Basic data manipulation

- Session 2
  - Join data tables
  - JMP graphing

- Session 3
  - **Modelling**
  - JMP Journal
  - JMP Scripting Language

# Modeling in JMP

- Linear regression (done)
  - predict a continuous variable

- Logistic regression
  - predict categorical variable (i.e. a classification problem)
    - binomial logistic regression: the categorical variable has binary outcomes (e.g., 0, 1)

- K-mean clustering
  - a method to partition observations (into clusters)

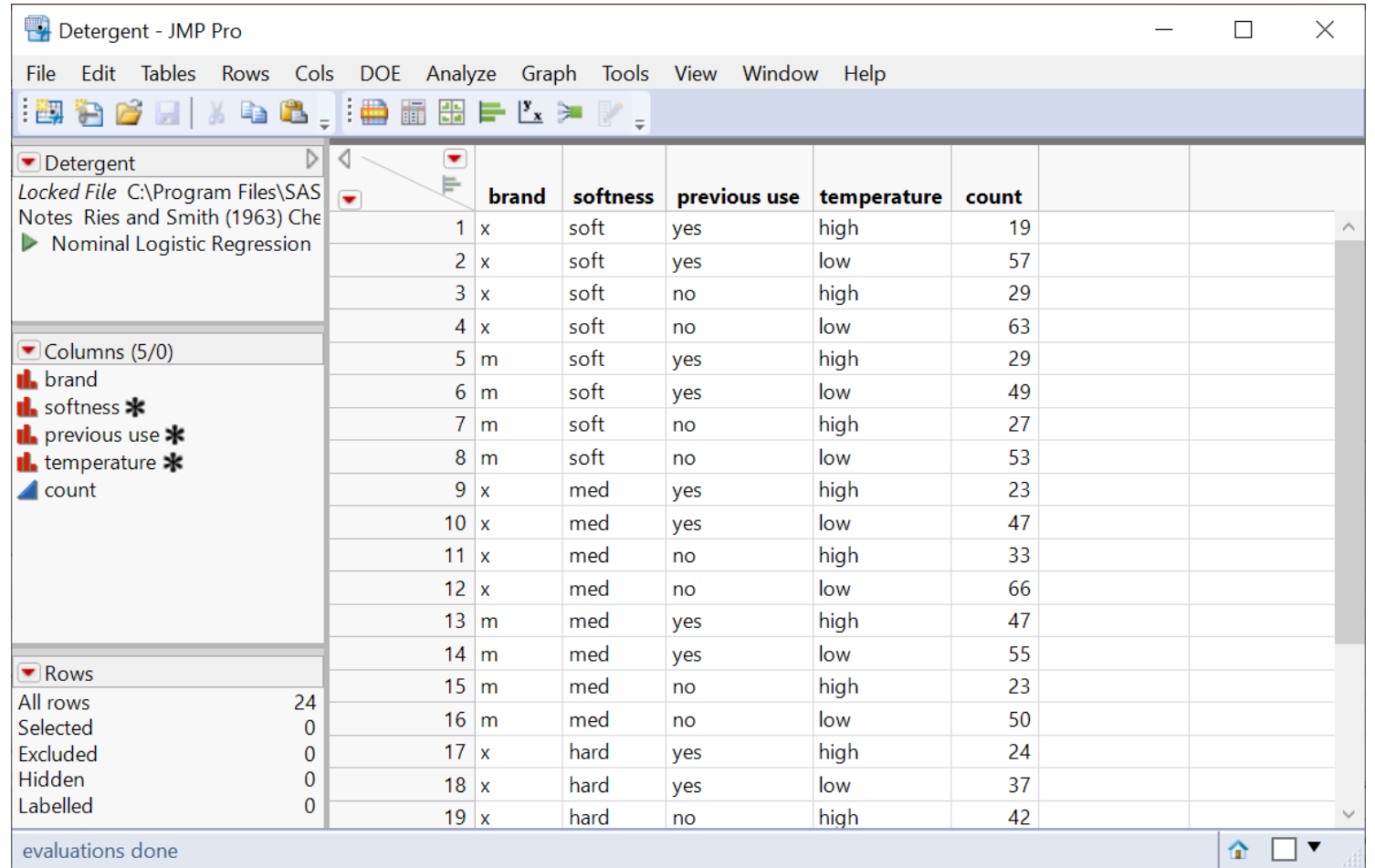- Outliers, missing values, and patterns

# Binomial Logistic Regression

- let $Y$ be the binary outcome variable
  - e.g. $\{0, 1\} = \{fail, \ success\}$

- Let $p = prob(Y = 1)$; $\frac{p}{1-p}$ is then the odds of being 1 (or success)

- Binary logistic regression models the logit-transformed probability as a linear relationship with the predictor variables
  - maximum likelihood estimation

$$logit(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k.$$

https://stats.idre.ucla.edu/other/mult-pkg/faq/general/faq-how-do-i-interpret-odds-ratios-in-logistic-regression/

# Binomial Logistic Regression (Demo): Data

- Preference for a brand of detergent (Ries and Smith, 1963)
  - Detergent.jmp

- Survey Questions
  - which brand do you prefer, x or m
  - water softness
  - previous user of m
  - water temperature

# Binomial Logistic Regression (Demo): Fit

- Analyze > Fit Model
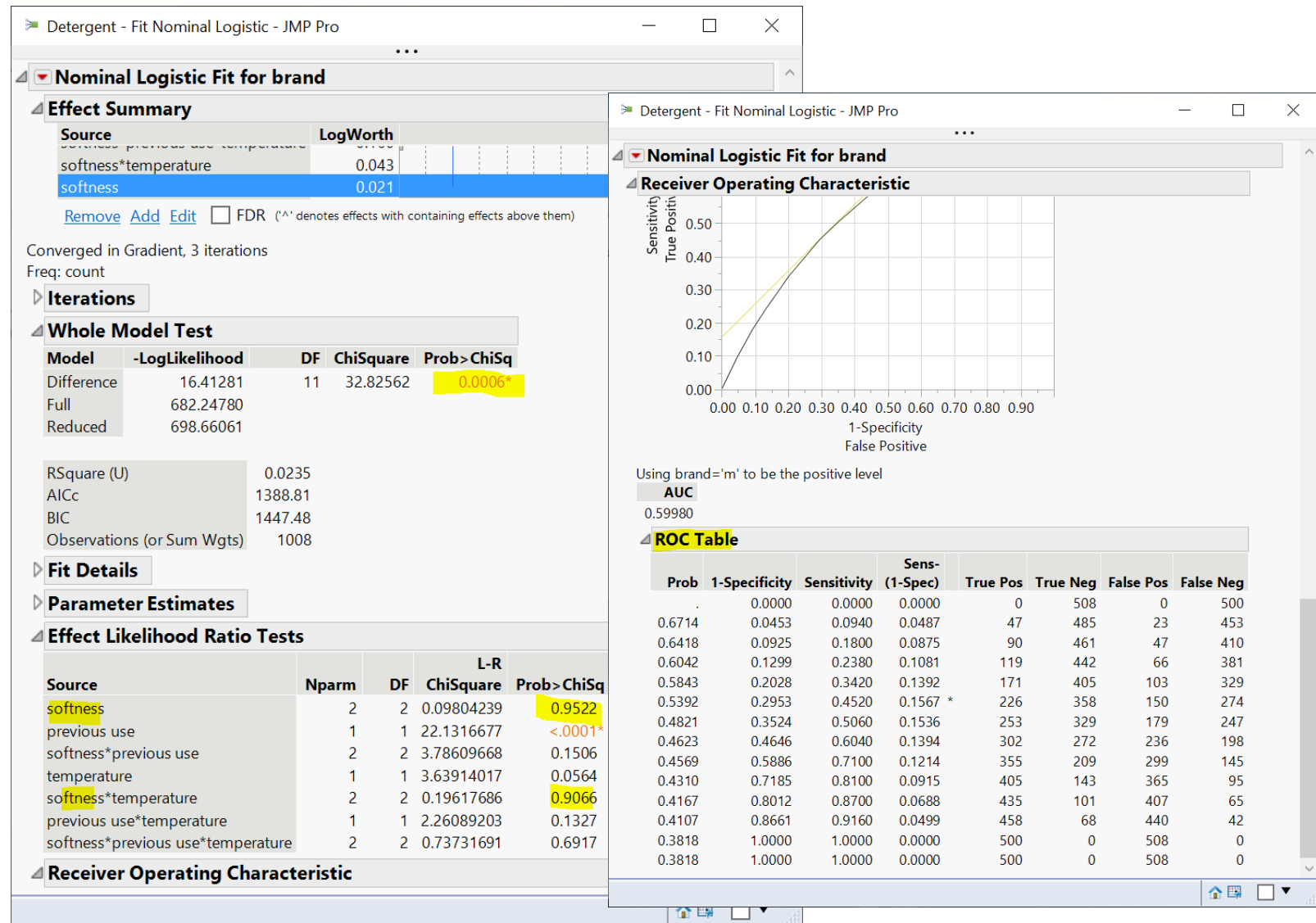
# Binomial Logistic Regression (Demo): Report

- Overall model fit is significant

- Softness doesn't seem to contribute too much

- ROC Table
  - [sensitivity, specificity, etc.](#)

# Your Turn (Hands-on)

- Do the same analysis without the softness variable

- Save the analysis script in the data table

- Challenge: How to construct a table of correct classification rate at each probability cutoff

$$correct\ classification\ rate = \frac{true\ positive + true\ negative}{total\ \#\ of\ predictions}$$
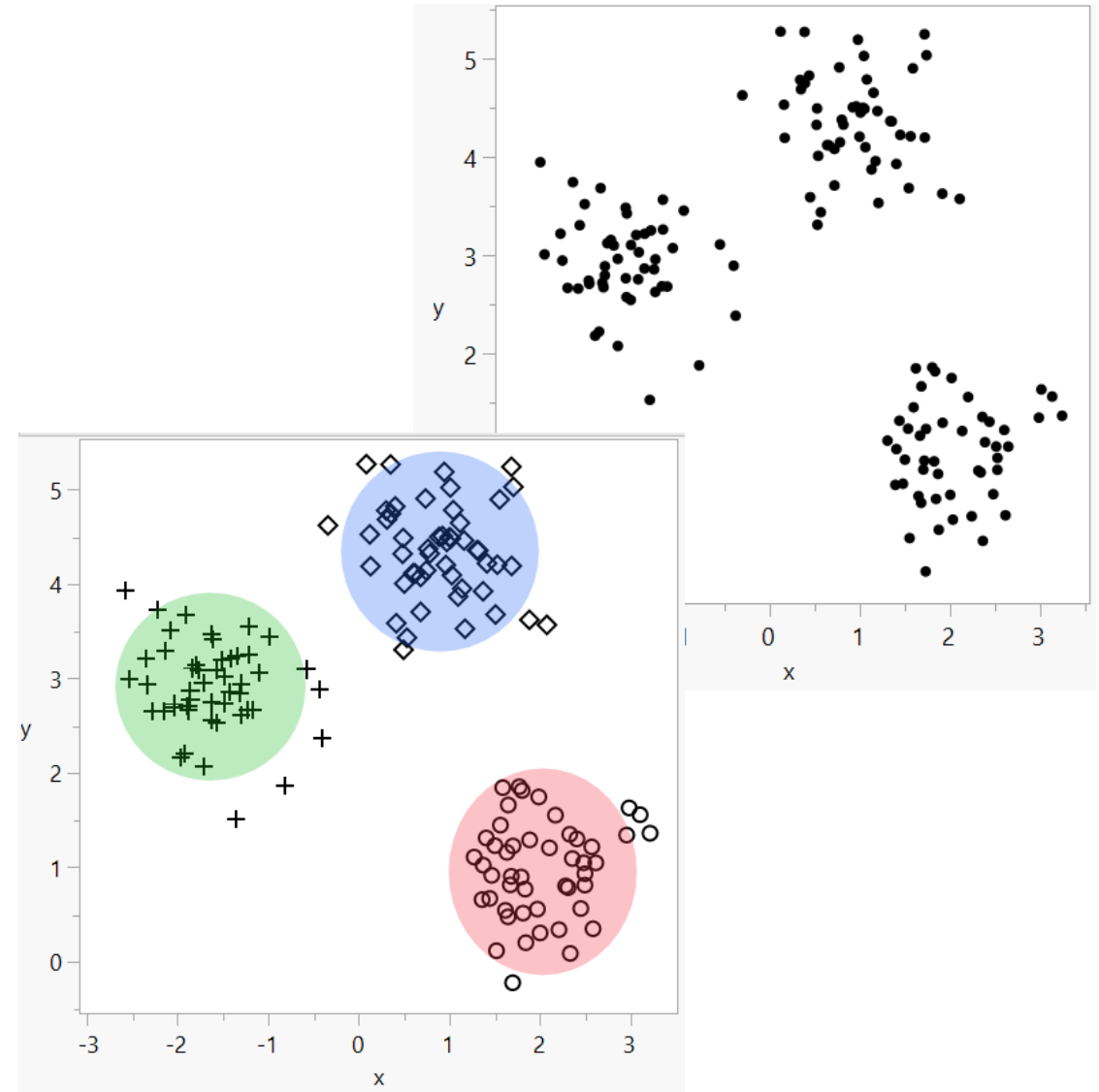
# K-Means Cluster Analysis

- A method to partition $n$ observations into $k$ clusters
  - such that total within-cluster sum of squares (between observations to cluster centroid) is minimized

- A **cluster** refers to a collection of data points aggregated together because of certain similarities

- Need to set $k$
  - There are methods to help you decide the value of $k$

# K-Means Cluster Analysis: An Example

- Observations: 150 2-d points

- Set $k = 3$
  - partition each observation to one of the 3 clusters $S = \{S_1, S_2, S_3\}$

- K-means clustering algorithm finds 3 clusters such that

$$\underset{S}{\mathrm{argmin}} \sum_{i=1}^{3} \sum_{x \in S_i} \|x - \mu_i\|^2$$

Within-cluster sum of squares

# K-mean Cluster Platform (Demo)

# Your Turn (Hands-on)

- Import the country_risk.xlsx data (data/basics/country_risk.xlsx)
  - note that it's an Excel file and column header starts at row 2

- Perform a pair-wise correlation analysis across the following 5 variables
  - Corruption, Peace, Legal, GDP Growth, Population
  - Note that Corruption and Legal variables are highly correlated
  - hint: use the Multivariate platform
    - Menu: ***Analyze -> Multivariate Methods - > Multivariate***

- Perform a K-means cluster analysis
  - As a start, use Peace, Legal and GDP Growth as factors; and set k=3
  - Produce a scatterplot matrix
  - Can you label each cluster (high-risk, medium-risk, etc.)?

# Outliers, Missing Values, and Patterns



Note: I should perhaps put this slide somewhere between data manipulation and modelling

# Plan

- Session 1
  - Workflow overview
  - Basic data manipulation

- Session 2
  - Join data tables
  - JMP graphing

- Session 3
  - Modelling
  - **JMP Journal**
  - **JMP Scripting Language**

# JMP Journal – Communicate Your Results

- Create a JMP journal when you want to present your results

- A JMP journal combine two kind of presentations
  - Static: embed output of JMP (graphs and reports), fixed at a moment in time
  - Dynamic: built from outlines containing text and buttons (links) that organize data tables and reports

- Getting-started resources
  - Dmitry's video about JMP Journal on Quercus (6 mins)
  - Creating, Using and Sharing JMP Journals (45 mins)

# JMP Script Language (JSL)

# Connect to SAS and Database Systems

- SAS
  - Learning resource: Using SAS from JMP

- Database