

Today: What Is Data?
Graphics Principles

Sam Ventura
36-315

Department of Statistics
Carnegie Mellon University

January 23, 2017

1/8

structured (vs. unstructured)

What Is Data?

Data: Information organized in some fixed / standard /
easy-to-understand way (humans or computers)

Ex) temperature measurements, housing prices,
heights / weights / BMI / heart rate / cholesterol / etc
Image classification - MNIST → images, video
Tweets, stock prices

2/8

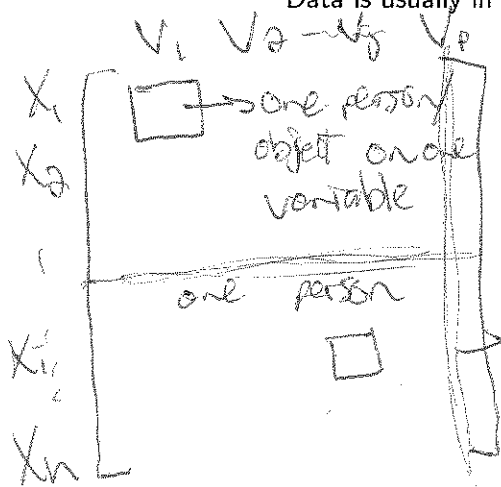
How Do We Describe Data?

rows
↑
etc

Two measurements used to describe datasets:

n : # of observations (people, objects, survey responses, days in the month)
 p or d : # of variables / covariates / features / questions → columns

Data is usually in matrix form:



* Rows = observations

↳ single row contains all information on all variables for one person/object

* Columns = variables

↳ one column has all answers to a single ? for all rows

Types of Data

Categorical:

Qualitative (usually), describes categories / classes / qualities of observation

CS: boolean, string/text, integer, factor

sub groups { ordinal / ordered: eg. "strongly agree", "agree", "neutral", ...
 nominal / unordered: e.g. race, gender, nationality

Continuous:

Quantitative (usually) real-valued, numerical data

CS: floats, integer, doubles, long

Graphics and Their Goals (from Tufte)

Graphics: visually display measured quantities by combining points, lines, coordinate system, numbers, symbols, words, shading, color

Goals: show data!

- ▶ induce viewer to think about substance, not graphical methodology
- ▶ avoid **distorting** the data
- ▶ present numbers in small space
- ▶ make large, complicated datasets more coherent
- ▶ encourage comparison of different pieces of data
- ▶ reveal data at several levels of detail
- ▶ **describe, explore, tabulate, identify relationships**
- ▶ be closely integrated with statistical/verbal descriptions of dataset

Graphs that do not meet these goals are not successful

Graphs leading viewers to make misleading conclusions should be avoided

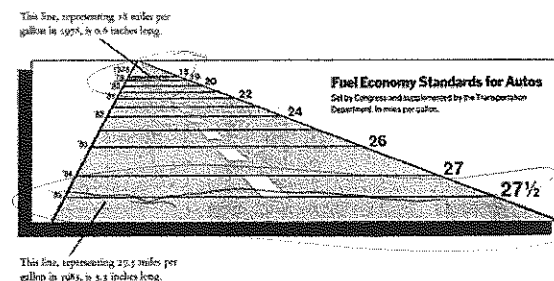
5/8

Optimal: $LF = 1$ Distortion = "graph doesn't match the data"

increased effect: $LF > 1$

decreased effect: $LF < 1$

Visual representation of data is inconsistent with numerical representation



size of "effect" in graph

Lie factor = size of "effect" in data

"effect" = change in amount of some feature/variable

→ in data

Tufte suggests optimizing the Lie Factor:

$$\text{Actual \% increase: } \frac{127.5 - 18}{18} \times 100\% \approx 52.8\%$$

$$\text{Graphical \% increase: } \frac{5.3 \text{ in} - 0.6 \text{ in}}{0.6 \text{ in}} \times 100\% \approx 783\%$$

$$LF = \frac{783\%}{52.8\%} \approx 14.83$$

"Decorating" and Data-Ink

Graphics should not draw the viewer's attention away from the data.
Extras get in the way.

Note: Decoration does not refer to appropriate graph labeling.
Labels should always be clear, detailed, and thorough.
Label key parts of the data. Add text explanations if necessary.

Data Ink should primarily present information about the data:
the non-erasable, non-redundant core of a graphic

Tufte suggests using the *data-ink ratio*:

$$DI = \frac{\text{data ink used to represent data}}{\text{total ink in graphic}} \times 100\%$$

% of ink devoted to non-redundant information in graph

7/8

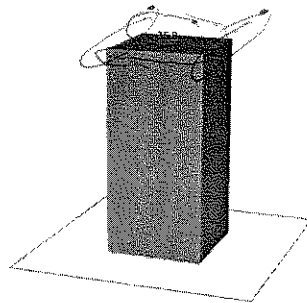
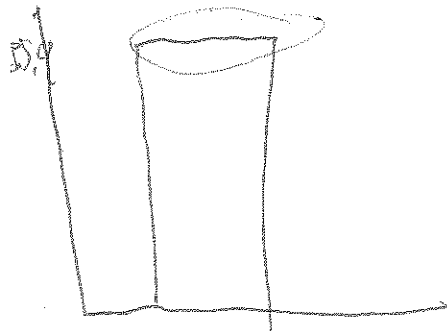
Ideally, maximize DI ratio (max of 1)

"Decorating" / Data-Ink

Two ways to increase the proportion of data-ink:

Remove non-data-ink:

Remove redundant data-ink:



- 1) height of front face
- 2) - - - - side face (right)
- 3) - - - - back line
- 4) - - - - left line
- 5) 35.9
- 6) area of front face
- 7) height of location of 35.9

Today: Grammar of Graphics
1-D Categorical
Friday: ggplot2, 1-D Categorical

January 25, 2017

ggplot2: Based on "The Grammar of Graphics" (Wilkinson, 2005)

Each plot can be broken down into core components.
Wilkinson defines core components in book.

Hadley Wickham puts this into practice in R via ggplot2.

1/8
ggplot(cars93) + geom_point(aes(x = fuel_tank_capacity, y = MPG.city, color = Type))

R Package ggplot2 – Hadley Wickham

Core components of a plot:

1. **data**: in ggplot2, data must be stored as an R data frame
2. **coordinate system**: describes 2-D space that data is projected onto
e.g., Cartesian coordinates, polar coordinates, map projections, ...
3. **geometries**: describe type of geometric objects that represent data
e.g., points, lines, polygons, ...
4. **aesthetics**: describe visual characteristics that represent data
e.g., for example, position, size, color, shape, transparency, fill
5. **scales**: for each aesthetic, describe how it is converted into values
that are displayed on the actual graph
e.g., log scales, color scales, size scales, date scales, ...
6. **stats**: describe statistical transformations that help summarize data
e.g., counts, means, medians, regression lines, ...
7. **facets**: describe how data is split into subsets and displayed as
multiple small graphs (particularly important for categorical data!)

map pieces of
the data to
pieces of
the graph

of obs
in data
↑

1-D Categorical Data

Recall: Data can be categorical or continuous

Categorical data can be ordered or unordered / nominal

1-D Categorical Data Structure:

for any obs. in this vector,
that obs. can take one of the following categories

How could we summarize this data?
What information would you report?

→ "counts" $\{C_1, C_2, \dots, C_K\}$
frequencies / proportions / percentages
in each category

3/8

- range of Percentages
- smallest / least frequent category?
- largest / most frequent category?
- "outliers", es. category w/ only one obs

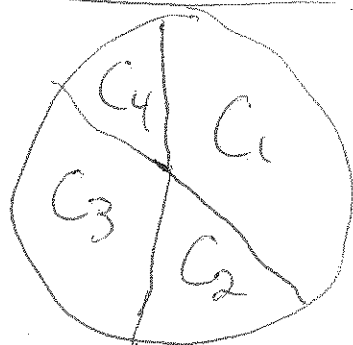
1-D Categorical Data

To show the differences among the categories, need to use area plots:

Differences in area correspond to differences
in category frequency (count, prop. perc.)
Each area of graph corresponds to a category

Examples of area plots?

Pie chart



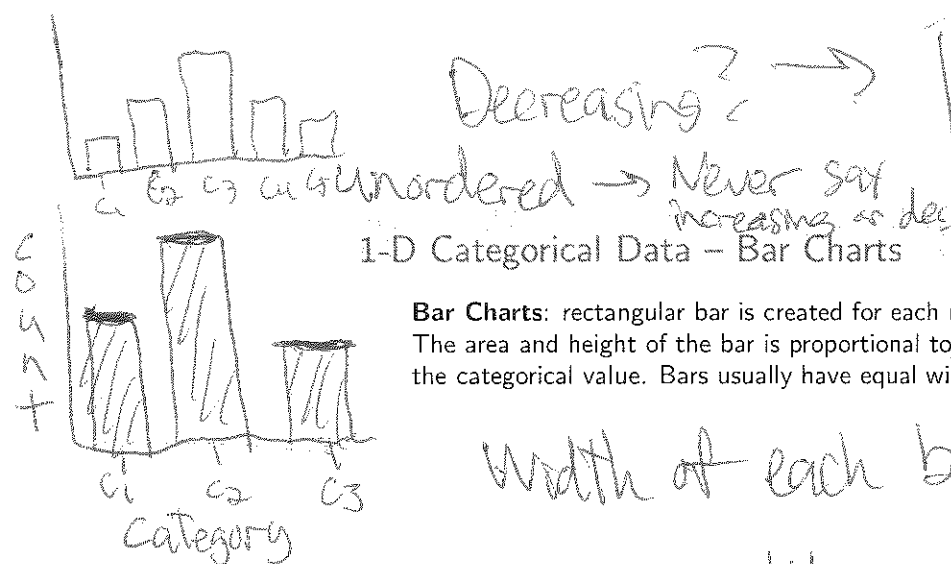
bar chart



Spine charts

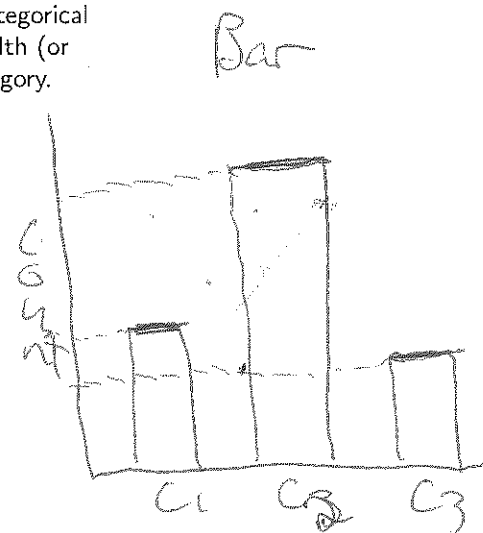
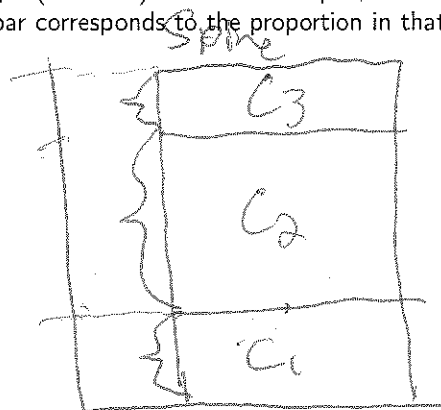
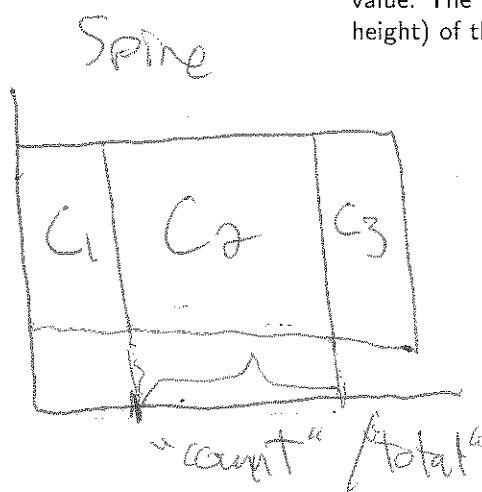
rose diagrams

4/8



1-D Categorical Data - Spine Charts

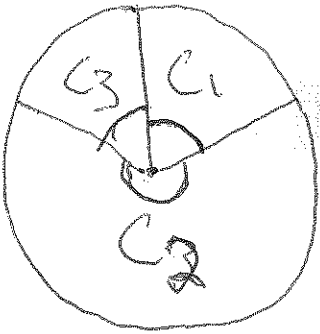
Spine Charts: rectangular bar is created for each unique categorical value. The height (or width) of all bars is equal, and the width (or height) of the bar corresponds to the proportion in that category.



Spine: very hard to visually determine category counts

Bars: very easy

Issues: - difficult to compare areas/angles
 - Bad w/ lots of categories
 - No total # of obs; ~~hard~~ hard to determine exact counts/proportions



1-D Categorical Data - Pie Charts

Pie Charts: circle divided up into sections ("pie slices") such that the area of each section is proportional to the number of observations with each unique categorical value.

Area of a circle: $A = \pi r^2$

Area of a slice: $A_{\text{sector}} = \pi r^2 \cdot \frac{\theta}{360}$

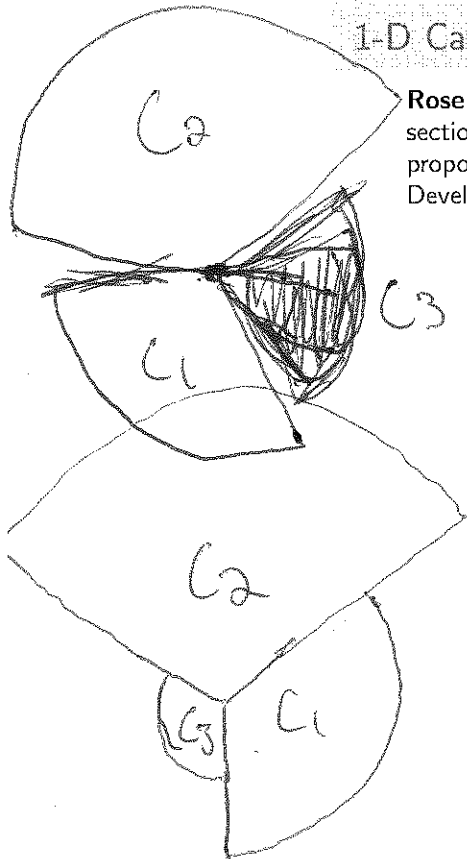
θ = angle in degrees (0 to 360)

$r \propto 1$ (nothing)

$A \propto$ "counts" in each category } or proportions
 $\theta \propto$ "counts" in each category }

1-D Categorical Data - Rose Diagrams

Rose Diagrams: circle sections are created for each category. All sections have the same width/arc/angle. The radius of each section is proportional to the category frequency. Sections are called "petals". Developed by Florence Nightingale (example posted to Blackboard).



$\theta \propto 1$

$r \propto$ "counts"

$A \propto (\text{"counts"})^2$

OR $r \propto \sqrt{\text{count}}$

$A \propto \text{count}$

↓

Issue = distortion

Issue = requires extra data manipulation



$P(X=C_1) = 1/4$, $P(X=C_2) = 1/2$, $P(X=C_3) = 1/4$
 or pie, rose, spine

What Does a Bar Chart Show?

Marginal Distribution:

Categorical variable X has K categories

"Distribution" of a categorical variable
 interested in $P(X=C_k)$ for each k
 → "truth" or "true distribution"

Empirical Distribution:

→ "what did we observe"
 "observed distribution"

"our best estimate of the true (marginal) distribution given the information/data that we have"

of the variable that we plotted

$$\begin{aligned}\hat{P}(X=C_1) &= 0.26 \\ \hat{P}(X=C_2) &= 0.48 \\ \hat{P}(X=C_3) &= 0.26\end{aligned}$$

we want to visualize uncertainty

Bar Charts: Counts vs. Proportions

Counts/Frequencies of each category:

→ no sample size

Proportions:

Advantage: we get an idea of sample size

But, we do get interesting statistical info

estimates of $P(X=C_j) = \hat{P}_j = \frac{\text{\# of obs. in category } j}{\text{total \# of obs.}}$

standard error on \hat{P}_j : $se(\hat{P}_j) = \sqrt{\frac{\hat{P}_j(1-\hat{P}_j)}{N}}$ $N \rightarrow$ sample size

Confidence interval around \hat{P}_j : $\hat{P}_j \pm Z_{\alpha/2} \cdot \sqrt{\frac{\hat{P}_j(1-\hat{P}_j)}{N}}$

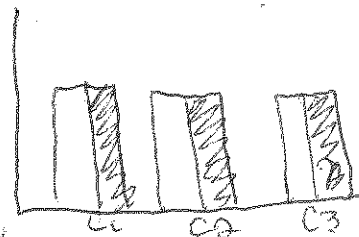
$\alpha = 0.99 \rightarrow 99\% \text{ CI}$ $\alpha = 0.95 \rightarrow 95\% \text{ CI}$

Chi-Square Test for Independence

Chi-squared test: Statistical test used to determine whether there is a significant difference between the expected frequencies and the observed frequencies in one or more categories (of a categorical variable).

2-D Categorical: Used to test differences in the conditional distributions (more on this Wednesday)

1-D Categorical: assume we have K categories

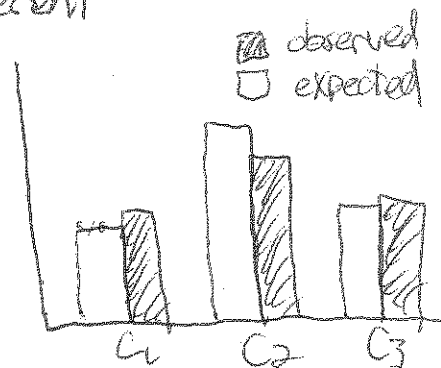


Setup #1

$H_0: P_1 = P_2 = P_3 = \dots = P_K \rightarrow$ "all proportions are equal"
 H_a : at least one of the P_j 's is different

Setup #2

$H_0: P_1 = P_1^*, P_2 = P_2^*, \dots, P_K = P_K^*$
 H_a : at least one is not the same as expected



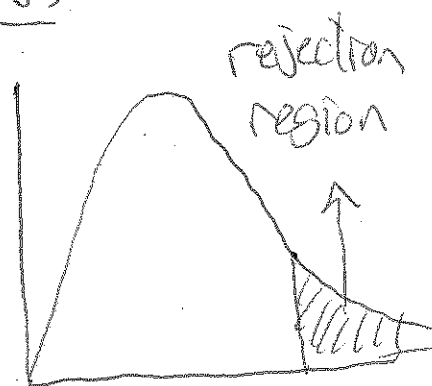
$P_j^* =$ expected proportion in Category j

Computing and Interpreting the Chi-Square Test

Test Statistic:
$$\chi^2 = \sum_{j=1}^K \frac{(O_j - E_j)^2}{E_j}$$

$O_j =$ observed # in each category j

$E_j =$ expected #



Interpretation:

Large χ^2 statistic means what we observed is very different from what we expected, so we will reject the null hypothesis (H_0) and conclude that there is evidence that the null is not true at some level α .

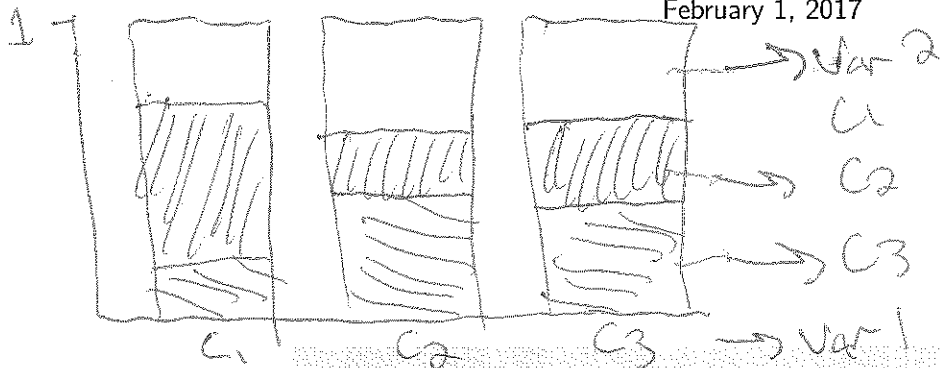
Today: Marginal and Conditional Distributions
2-D Categorical Data
Contingency Tables and Mosaic Plots
Friday: 2-D Categorical

Sam Ventura
36-315

Department of Statistics
Carnegie Mellon University

February 1, 2017

Proportional bar chart



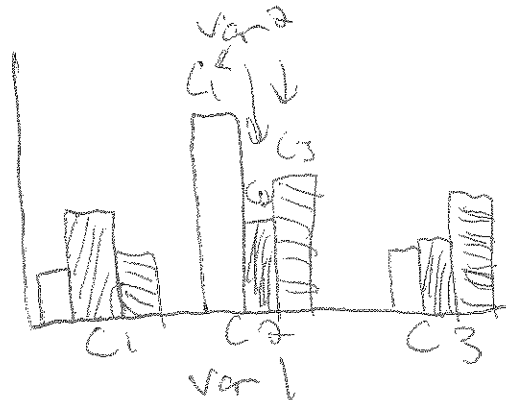
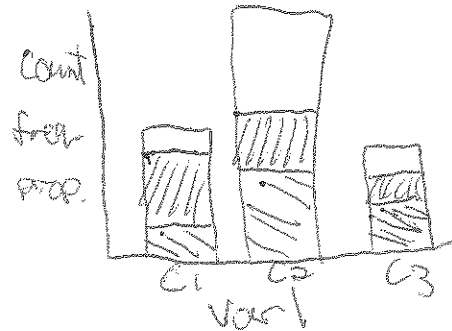
1/7

2-D Categorical: Stacked and Side-by-Side Bar Charts

Stacked bar chart

2-D Bar Charts:

- ☐ Var 2, C1
- ☐ Var 2, C2
- ☐ Var 2, C3



What does a 2-D bar chart show?:

stacked

- Marginal of V_1 : easy 😊
- conditional distributions of V_2 given V_1 : 😞

side-by-side

- Marg. V_1 : OK 😊
- cond. V_2/V_1 : 😊

proportional

- Marg. of V_1 : 😞
- cond. of V_2/V_1 : 😊

2/7

Contingency Tables \rightarrow 2 variables

Contingency Tables Give Us:

Empirical {
 Marginal distribution (of each variable individually)
 conditional distns (of $V_1 | V_2$ or $V_2 | V_1$)
 Joint distn (of both V_1 and V_2 simultaneously)

Two Categorical Variables:

Var 1: K_1 categories; Var 2: K_2 categories

What Are We Interested In?

n = total # of obs
 $n_{i\cdot}$ = total # in row i
 $n_{\cdot j}$ = total # in column j

we know

• Observed counts in each cell

$\rightarrow O_{ij}$ = # obs in cat i of V_1 and cat j of V_2

we want to know: Expected counts in each cell $\rightarrow E_{ij}$

Contingency Tables and Marginal/Conditional Distributions

Contingency Tables and Marginal/conditional Distributions

Variable 1
 \downarrow

| | | Variable 2 | | | | |
|-----------|---------------|---------------|----------|-----------------|----------------|--|
| | | C_1^* | C_2^* | ... | $C_{K_2}^*$ | |
| C_1 | O_{11} | O_{12} | ... | O_{1K_2} | $n_{1\cdot}$ | \rightarrow Marginal distn of Var 1 |
| C_2 | O_{21} | O_{22} | ... | O_{2K_2} | $n_{2\cdot}$ | |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | |
| C_{K_1} | O_{K_11} | O_{K_12} | ... | $O_{K_1K_2}$ | $n_{K_1\cdot}$ | \rightarrow Conditional distn of V_2 given $V_1 = C_2$ |
| | $n_{\cdot 1}$ | $n_{\cdot 2}$ | ... | $n_{\cdot K_2}$ | n | |

\rightarrow Marginal distn of Var 2

\rightarrow cond. distn
of V_1 give
 $V_2 = C_2^*$

Today: 2-D Categorical Data Independence and Mosaic Plots

Sam Ventura
36-315

Department of Statistics
Carnegie Mellon University

February 6, 2017

1/4

for two events, A and B

Recall: Independence Rules from Probability

$$P(A) = P(A|B) \rightarrow \text{marginal dist'n of Var-1} = \text{conditional dist'n of Var-1 given Var-2}$$

$$P(B) = P(B|A) \rightarrow \text{marginal dist'n of Var-2} = \text{conditional dist'n of Var-2 given Var-1}$$

$$P(A \cap B) = P(A)P(B) \rightarrow \text{joint dist'n of Var-1, Var-2} = (\text{marginal of Var-1}) \times (\text{marginal of Var-2})$$

Can input contingency tables into chi-square tests for independence

E.g. `chisq.test(table(var1, var2))`

More on this in Lab 04

contingency table in R
test for independence of two categorical variables

2/4

i = category i of Var 1
 j = category j of Var 2

E_{ij} assumes
 Var 1 \perp Var 2
 \downarrow
 independent

Pearson Residuals

Pearson Residuals: Scaled difference between observed/expected

$$r_{ij} = \frac{O_{ij} - E_{ij}}{\sqrt{E_{ij}}} = \frac{\text{Observed} - \text{Expected}}{\sqrt{\text{Expected}}}$$

$r_{ij} \stackrel{u}{=} 0$: we observed roughly as many obs. in categories i and j as we would expect of Var 1 \perp Var 2

$r_{ij} > 0$: "too many" observed \geq assuming

$r_{ij} < 0$: "too few" observed \leq Var 1 \perp Var 2

Result: r_{ij} 's are asymptotically Normally distributed.

$$r_{ij} \sim N(\mu=0, \sigma^2=1)$$

Mosaic Plots

Mosaic Plots: Area plot for two categorical variables

Each cell in a contingency table gets a box

\hookrightarrow Area of each box \propto % of obs. in the corresponding cell of the table

\hookrightarrow width of each box \propto marginal distn of Var 1

\hookrightarrow heights of each box \propto conditional distn of Var 2 | Var 1

$|r_{ij}| > 2 \Rightarrow$ significant at the $\alpha=0.05$ level

$|r_{ij}| > 4 \Rightarrow$ significant at the $\alpha=0.01$ level

Can color the boxes by their differences from what was expected

Friday: Mosaic Plots in R