# Ch. 3: Correlation & Regression

- Exploring relationships between 2 variables
- Scatterplots
- Correlation
- Linear Regression
- Other Correlation Coefficients

# Bivariate relationships

- "is factor A related to factor B"?
- Methods of analysis:
  - Anecdotal / Clinical -- often forms the basis for further systematic research & data collection
  - Numerically -- check values & % at extremes
  - Visually -- scatterplots
    - easy to see relationships and problems w/data
    - hard to prove / test
  - Statistically -- correlation & regression
    - hard to detect problems w/data
    - easy to test hypothesis

# Anecdotal / Clinical

- Many interesting findings in psychology first originate from non-scientific approaches
- "Intuition" that something is related through experiencing multiple situations
- Pattern recognition
- Human brains are both excellent and terrible pattern recognizers
- Problems -- faulty memory, confirmation biases, prejudice, etc…
- First step after a "gut" feeling is to begin collecting data.

# Simple numerical analysis

- Simplify the situation by using Categorical variables (or reducing Continuous variables to Categorical variables)
- Use extreme cases to maximize effect
- Compute percentages in a 2x2 matrix
- Do the results suggest an effect?

- Compute Chi-square statistic to judge significance

# Example

- "I think there is brain dysfunction in HIV disease" as measured by neuropsychological testing
- Medical status: control vs. HIV+ symptomatic
- NP test results: normal vs. impaired

|  |  | Medical Status | |
|---|---|---|---|
|  |  | Control | HIV+ |
| NP Status | Normal | 85% | 52% |
|  | Impaired | 15% | 48% |

# Issues

- Pro: Results are very easy to understand from a human or clinical point of view.
- Con: dividing continuous variables into two values throws away a lot of data and statistical power

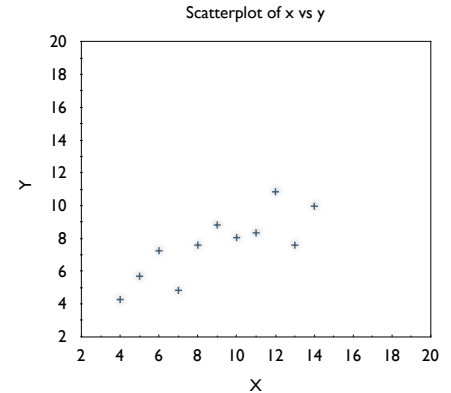- Graphical and Statistical methods should be used as well.

# Scatterplots

- Graph two variables in relation to each other on two-dimensional X,Y axis
- Easy to see relations between data
- Easy to see problems with data

- Hard to prove or determine if an apparent relationship is "significant"
- Hard to interpret data clinically or in "common sense" terms

---

# Scatterplots

| x | y |
|------|-------|
| 10.0 | 8.04 |
| 8.0 | 7.58 |
| 13.0 | 7.58 |
| 9.0 | 8.81 |
| 11.0 | 8.33 |
| 14.0 | 9.96 |
| 6.0 | 7.24 |
| 4.0 | 4.26 |
| 12.0 | 10.84 |
| 7.0 | 4.82 |
| 5.0 | 5.68 |



Scatterplot of x vs y

---

# Linear Regression

- Assume that two variables are related, and that this relationship is <u>linear</u> -- model the data by a simple straight line for the data.
- For any given data set, we pick the line that best "fits" our data
- Similar terms: linear regression, fitting a line, finding the trend, creating a trendline, best fit line, etc.
- Residuals = difference between prediction and actual value
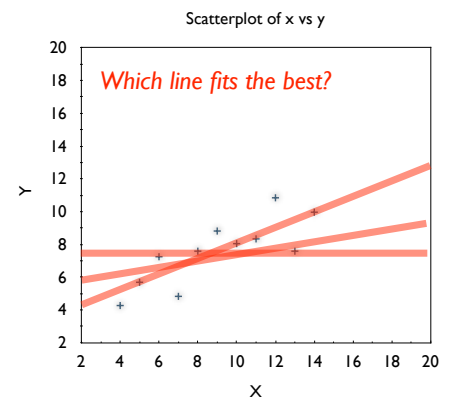- Linear Regression minimizes the square of the residuals, often called "Ordinary Least Squares"

---

# Linear Regression

*Equation:*
$y = 3.0 + 0.5x$

*Correlation*
$r_{x,y} = 0.816$



Scatterplot of x vs y

*Which line fits the best?*

---

# Anscombe's Quartet I

| x | y |
|------|-------|
| 10.0 | 8.04 |
| 8.0 | 7.58 |
| 13.0 | 7.58 |
| 9.0 | 8.81 |
| 11.0 | 8.33 |
| 14.0 | 9.96 |
| 6.0 | 7.24 |
| 4.0 | 4.26 |
| 12.0 | 10.84 |
| 7.0 | 4.82 |
| 5.0 | 5.68 |



Scatterplot of x vs y

*Linear Regression*

---

# Anscombe's Quartet II

| x | y |
|------|------|
| 10.0 | 9.14 |
| 8.0 | 8.14 |
| 13.0 | 8.74 |
| 9.0 | 8.77 |
| 11.0 | 9.26 |
| 14.0 | 8.10 |
| 6.0 | 6.13 |
| 4.0 | 3.10 |
| 12.0 | 9.13 |
| 7.0 | 7.26 |
| 5.0 | 4.74 |



Scatterplot of x vs y

*Linear Regression*

## Anscombe's Quartet III

| x | y |
|---|---|
| 10.0 | 7.46 |
| 8.0 | 6.77 |
| 13.0 | 12.74 |
| 9.0 | 7.11 |
| 11.0 | 7.81 |
| 14.0 | 8.84 |
| 6.0 | 6.08 |
| 4.0 | 5.39 |
| 12.0 | 8.15 |
| 7.0 | 6.42 |
| 5.0 | 5.73 |

Scatterplot of x vs y

*Linear Regression*

---

## Anscombe's Quartet IV

| x | y |
|---|---|
| 8.0 | 6.58 |
| 8.0 | 5.76 |
| 8.0 | 7.71 |
| 8.0 | 8.84 |
| 8.0 | 8.47 |
| 8.0 | 7.04 |
| 8.0 | 5.52 |
| 19.0 | 12.50 |
| 8.0 | 5.56 |
| 8.0 | 7.91 |
| 8.0 | 6.89 |

Scatterplot of x vs y

*Linear Regression*

---

## Anscombe's Series 1-IV

---

## Anscombe's Quartet Summary

- The 4 series of data, though very different, have identical linear regression equations and identical correlations
- Each series has a Quantitative correlation, but it's clear (visually) that the relationships are Qualitatively different
- Each series should probably be handled differently, through techniques such as:
  - Trimmed Least Squares
  - Robust regression
- <u>Graph Your Data!</u>

---

## Residuals in Linear Regression

- X : dependent variable
- Y : independent variable
- Model:  predict Y from X
- Y' : (Y prime) = predicted Y
- $Y' = a + bX$
- Prediction is (usually) incorrect.  Difference between predicted (Y') and actual (Y) is called a "Residual"  = (Y - Y')
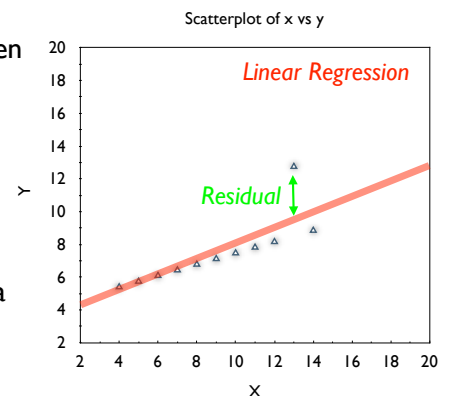- Calculation of best fit line minimizes the sum of the squared residuals  $\Sigma(Y-Y')^2$

---

## Residuals in Linear Regression

Residual is difference between actual Y and predicted Y' (Y - Y')

Graphically it is equal to how far away (vertically) a point is from the linear regression line

Scatterplot of x vs y

*Linear Regression*

*Residual*

# Correlation (r) Pearson's r

- Pearson's Product-Moment Correlation
- Measures the strength of the linear relationship between two variables
- Ranges between -1.0 and +1.0
- Is a special case of linear regression, when both X and Y have been turned into Z scores.
- r is transitive (correlation between X and Y is same as correlation between Y and X)
- $R^2$ = "explained variance" is the proportion of variation in the data explained by the model.
- $R^2$ ranges from 0 to 1.0  (0% to 100%)

# Correlation vs. Regression

|  | Linear Regression | Correlation |
|---|---|---|
| Scores | Raw | Z |
| Mean, Std Dev | sample means sample Std Dev | 0 1 |
| Equation | Y' = a + bX | Y' = r X |
| Slope | b = change in Y per change in X | r = correlation coefficient |
| Slope$^2$ | meaningless | $R^2$ = % variance explained |

# Other Correlation Coefficients

- Continuous (interval & ratio):  Pearson's r
- Ordinal (Ranked): A B C D…     1st, 2nd, 3rd...
  - Spearman's Rho: correlation between two ordinal / ranked variables.

- Dichotomous (yes/no, one/zero, T/F, Male/Female, Pass/Fail...)
  - True vs. Artificial?

# Continuous vs. Dichotomous

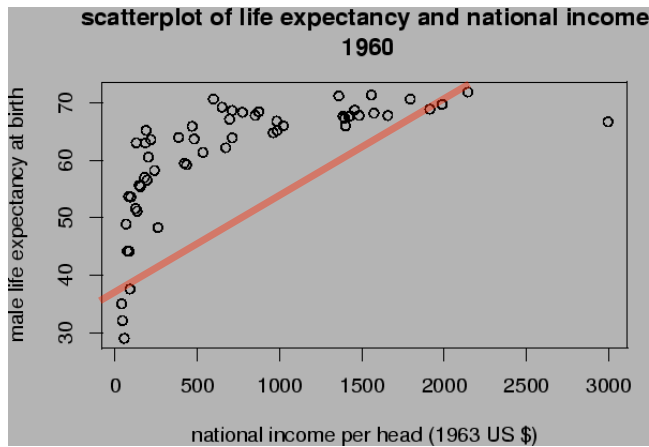| Type of X / Type of Y | Continuous | Artificial Dichotomous | True Dichotomous |
|---|---|---|---|
| Continuous | Pearson r | Biserial r | Point biserial r |
| Artificial Dichotomous | Biserial r | Tetrachroic r | Phi |
| True Dichotomous | Point biserial r | Phi | Phi |

# Correlation : Issues

- Technical / Calculation :
  - Non-normal distribution
  - Non-linear data and relationships
  - Outliers, data errors
  - Restricted Range
  - Shrinkage
- Interpretation:
  - Correlation =? Causation
  - Third variable explanations

# Non-linearity

- Linear regression & Correlation assume a linear relationship between X and Y
- Are most real-world relationships linear?
- Examples of non-linearity
- Solutions:
  - Intentionally restrict range of X
  - Rank variables then use Spearman's Rho
  - Transform variables (log, root, square, cube, etc.)
  - Use higher-order (polynomial) curve fitting, such as $Y = a + bX + cX^2 + dX^3$ ...
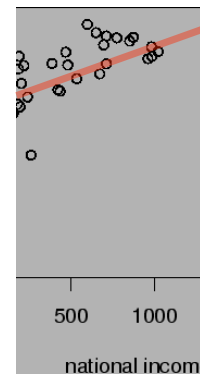
# Life expectancy / national income


scatterplot of life expectancy and national income 1960

# restrict range of X

# log transform X (or Y)


scatterplot of life expectancy and national income on the log scale, 1960

# Outliers & Data Errors?

# Correlation = Causation?

- A relationship (linear or otherwise) between X and Y tells us nothing about whether X causes Y
- Lack of correlation between X and Y does not mean that X doesn't cause Y

- Sleeping with your shoes on is correlated with waking up with a headache
- Ice cream sales are positively related to increase in drowning

# Shrinkage

- Least-squares regression attempts to fit the data set presented to it by reducing the observed residuals.
- This data set contains random errors.
- Thus, the parameters (equations) estimated for the linear regression line (and correlation coefficient) and residuals usually be higher than would be found in a separate data set.
- This reduction is called "Shrinkage"
- Cross-validation is best way to deal with it

# Cross Validation

- Step 1: With a given data set, compute the linear regression line that fits this data.
- Step 2: Apply this linear regression equation to a <u>different</u> data set.
- Step 3: Calculate the observed error in step 2. This is typically higher than seen in step 1, and a much better measure of fit.

- Note: sometimes you may artificially "create" two data sets by splitting a single data set in half.

# Hypothesis Testing

- All parameters (equations) we estimate from data have inherent error
- How do we know if a given estimate is correct?
- How big is the error likely to be (confidence intervals)?
- Inferential Statistics - covered later
  - Formulas to calculate probability, confidence intervals.
  - Higher N is better
  - "statistical significance" not the same as "clinical significance"

# Statistical and Clinical Significance

- These two terms are often confused and have very different meanings

- Statistical Significance: changes in DV are very unlikely to have been the result of random effects or chance. Often expressed as a P value ($p < .01$, or less than 1% chance to see these effects under H0)

- Clinical Significance : changes in DV are large enough to matter; the change was not trivial. If we accept H1, the conclusion is that H1's effect size is important.
  - depends on context. often evaluated in terms of cost/benefit or risk/benefit tradeoff

# Significance

- Example 1:
  - Two dice, Roll each once
  - Results: get a 3 and a 5
- Example B:
  - Two dice, Roll each 100 times
  - Results: Die A = 3.0, Die B = 3.10
- Example C:
  - Two dice, Roll each 100 times
  - Results: Die A = 3.0, Die B = 5.0