

Population - entire collection of objects or individuals about which information is desired.

→ easier to take a sample

- ◆ **Sample** - part of the population that is selected for analysis
- ◆ **Watch out for:**

- Limited sample size that might not be representative of population
- ◆ **Simple Random Sampling-**

Every possible sample of a certain size has the same chance of being selected

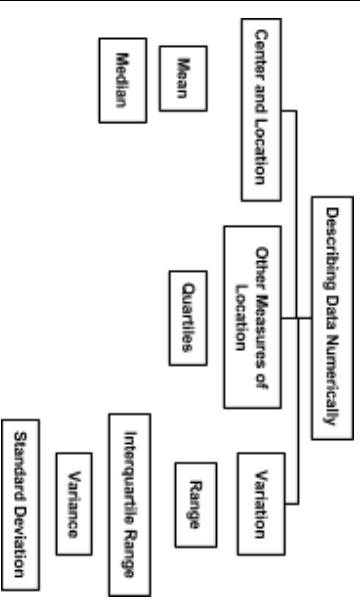
Observational Study - there can always be lurking variables affecting results

- i.e., strong positive association between shoe size and intelligence for boys
- ****Should never show causation**

Experimental Study - lurking variables can be controlled; can give good evidence for causation

Descriptive Statistics Part I

→ Summary Measures



→ **Mean** - arithmetic average of data values

- ◆ ****Highly susceptible to extreme values (outliers).**

Goes towards extreme values

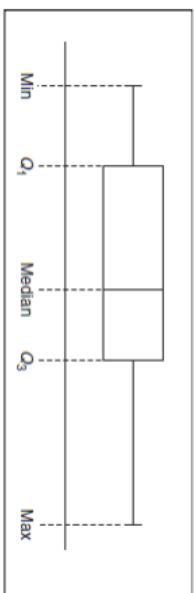
- ◆ Mean could never be larger or smaller than max/min value but could be the max/min value
- ◆ s_x^2 gets rid of the negative values
- ◆ units are squared

→ **Median** - in an ordered array, the median is the middle number

- ◆ ****Not affected by extreme values**

→ **Quartiles** - split the ranked data into 4 equal groups

- ◆ **Box and Whisker Plot**



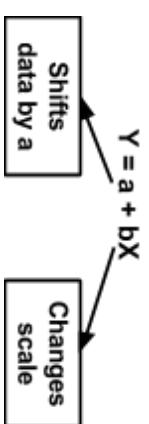
→ **Standard Deviation** - shows variation about the mean

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

- ◆ highly affected by outliers
- ◆ has same units as original data
- ◆ finance = horrible measure of risk (trampoline example)

Descriptive Statistics Part II

Linear Transformations



→ **Range** = $X_{\text{maximum}} - X_{\text{minimum}}$

- ◆ **Disadvantages:** ignores the way in which data are distributed; sensitive to outliers

- **Interquartile Range (IQR)** = $3rd quartile - 1st quartile$
- ◆ Not used that much
- ◆ Not affected by outliers

- Linear transformations change the center and spread of data
- $V\text{ar}(a+bX) = b^2 V\text{ar}(X)$
- $\text{Average}(a+bX) = a+b/\text{Average}(X)$

→ Effects of Linear Transformations:

- ◆ $mean_{new} = a + b * mean$
- ◆ $median_{new} = a + b * median$
- ◆ $stddev_{new} = |b| * stddev$
- ◆ $IQR_{new} = |b| * IQR$

→ **Z-score** - new data set will have mean 0 and variance 1

$$z = \frac{X - \bar{X}}{S}$$

Empirical Rule

→ Only for mound-shaped data

Approx. 95% of data is in the interval:

$$(\bar{X} - 2S_x, \bar{X} + 2S_x) = \bar{X} + / - 2S_x$$

→ only use if you just have mean and std. dev.

Chebyshev's Rule

→ Use for any set of data and for any number k , greater than 1 (1.2, 1.3, etc.)

$$\rightarrow 1 - \frac{1}{k^2}$$

→ (Ex) for $k=2$ (2 standard deviations), 75% of data falls within 2 standard deviations

Detecting Outliers

→ Classic Outlier Detection

- ◆ doesn't always work
- ◆ $|z| = \left| \frac{X - \bar{X}}{S} \right| \geq 2$

The Boxplot Rule

- ◆ Value X is an outlier if:

$$X < Q1 - 1.5(Q3 - Q1)$$

Or

$$X > Q3 + 1.5(Q3 - Q1)$$

→ Skewness
measures the degree of asymmetry exhibited by data

- ◆ Correlation doesn't imply causation
- ◆ The correlation of a variable with itself is one
- ◆ negative values = skewed left
- ◆ positive values = skewed right
- ◆ if $|skewness| < 0.8$ = don't need to transform data

Combining Data Sets

$$\rightarrow \text{Mean } (\bar{Z}) = \bar{\bar{Z}} = a\bar{X} + b\bar{Y}$$

$$\text{Var } (\bar{Z}) = S_{\bar{Z}}^2 = a^2 Var(X) + b^2 Var(Y) + 2abCov(X, Y)$$

Measurements of Association

→ Covariance

- ◆ Covariance > 0 = larger x, larger y
- ◆ Covariance < 0 = larger x, smaller y
- ◆ $S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$
- ◆ Units = Units of x . Units of y
- ◆ Covariance is only +, -, or 0 (can be any number)

Portfolios

→ Return on a portfolio:

$$R_p = w_A \bar{R}_A + w_B \bar{R}_B$$

- ◆ weights add up to 1
- ◆ return = mean
- ◆ risk = std. deviation

→ Variance of return of portfolio

$$S_p^2 = w_A^2 S_A^2 + w_B^2 S_B^2 + 2w_A w_B (S_{A,B})$$

- ◆ $r_{xy} = \frac{covariance_{xy}}{(std.dev.x)(std.dev.y)}$
- ◆ correlation is between -1 and 1
- ◆ Sign: direction of relationship
- ◆ Absolute value: strength of relationship (-0.6 is stronger relationship than +0.4)

Probability

- measure of uncertainty
- all outcomes have to be exhaustive (all options possible) and mutually exclusive (no 2 outcomes can occur at the same time)

Magnitude of r	Interpretation
.00-.20	Very weak
.20-.40	Weak to moderate
.40-.60	Medium to substantial
.60-.80	Very strong
.80-1.00	Extremely strong

Probability Rules

- Probabilities range from $0 \leq Prob(A) \leq 1$
- The probabilities of all outcomes must add up to 1
- The complement rule = A happens or A doesn't happen

$$P(\bar{A}) = 1 - P(A)$$

$$P(A) + P(\bar{A}) = 1$$

4. Addition Rule:

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

Contingency/Joint Table

- To go from contingency to joint table, divide by total # of counts
- everything inside table adds up to 1

		B	\bar{B}
A	P(A and B)	P(A and \bar{B}) $= P(\bar{B} A)P(A B)$	
	P(\bar{A} and B) $= P(\bar{A} B)P(B \bar{A})$	P(\bar{A} and \bar{B}) $= P(\bar{A} \bar{B})P(\bar{B} \bar{A})$	

Called the rule of total probability

$$P(A) = P(A \text{ and } B) + P(A \text{ and } \bar{B})$$

$$= P(A|B)P(B) + P(A|\bar{B})P(\bar{B})$$

Called the rule of total probability

Conditional Probability

$$\rightarrow P(A|B)$$

$$\rightarrow P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$

- Given event B has happened, what is the probability event A will happen?
- Look out for: "given", "if"

Independence

- Independent if:

$$P(A|B) = P(A) \text{ or } P(B|A) = P(B)$$

- If probabilities change, then A and B are dependent
- **hard to prove independence, need to check every value

Multiplication Rules

- If A and B are INDEPENDENT:

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

→ Another way to find joint probability:

$$P(A \text{ and } B) = P(A|B) \cdot P(B)$$

$$EMV = X_1(P_1) + X_2(P_2) + \dots + X_n(P_n)$$

2 x 2 Table

Decision Tree Analysis

- square = your choice
- circle = uncertain events

Discrete Random Variables

- $\rightarrow P_X(x) = P(X = x)$
- $\rightarrow \mu_x = E(x) = \sum x_i P(X = x_i)$
- $\rightarrow Example: (2)(0.1) + (3)(0.5) = 1.7$

Expectation

- $\rightarrow \sigma^2 = E(x^2) - \mu_x^2$
- $\rightarrow Example: (2)^2(0.1) + (3)^2(0.5) - (1.7)^2 = 2.01$

Rules for Expectation and Variance

- $\rightarrow \mu_s = E(s) = a + b\mu_x$
- $\rightarrow Var(s) = b^2 \cdot \sigma^2$

Jointly Distributed Discrete Random Variables

- $\rightarrow Independent \text{ if: }$

$$P_{x,y}(X = x \text{ and } Y = y) = P_x(x) \cdot P_y(y)$$

→ Expected Value Solution =

$$\boxed{\begin{aligned} \text{Example: } EV(\text{Average factory}) &= 90(.3) + 120(.5) + (-30)(.2) \\ &= 81 \end{aligned}}$$

→ Combining Random Variables

- ♦ If X and Y are independent:

$$E(X + Y) = E(X) + E(Y)$$

$$Var(X + Y) = Var(X) + Var(Y)$$

♦ If X and Y are dependent:

$$E(X + Y) = E(X) + E(Y)$$

$$Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$$

- 3.) At least one success
- $$P(\text{at least 1 success}) = 1 - (1 - p)^n$$
- 4.) At least one failure
- $$P(\text{at least 1 failure}) = 1 - p^n$$
- 5.) Binomial Distribution Formula for x=exact value

$$Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$$

→ Covariance:

$$\begin{aligned} \text{Cov}(X, Y) &= E(XY) - E(X)E(Y) \\ \rightarrow \text{If } X \text{ and } Y \text{ are independent, } \text{Cov}(X, Y) &= 0 \end{aligned}$$

| Calculate the Covariance

- We will use the formula $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$
 - For a die $E(X) = E(Y) = 3.5$
 - We need to find $E(XY)$
- | Probability | X | Y | XY | Prob × XY |
|-------------|---|---|----|------------|
| 1/6 | 1 | 6 | 6 | 6/6 = 1 |
| 1/6 | 2 | 5 | 10 | 10/6 = 5/3 |
| 1/6 | 3 | 4 | 12 | 12/6 = 2 |
| 1/6 | 4 | 3 | 12 | 12/6 = 2 |
| 1/6 | 5 | 2 | 10 | 10/6 = 5/3 |
| 1/6 | 6 | 1 | 6 | 6/6 = 1 |
- $$\text{E}[XY] = \text{sum } q_i \cdot y_i = 9/3 = 3$$
- So $\text{Cov}(X, Y) = 9.33 - (3.5)(3.5) = -2.91$
 - The covariance is negative because larger values of X are associated with smaller values of Y.

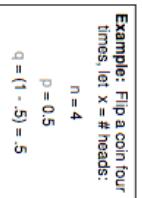
Binomial Distribution Formula

$$P(X=x) = \frac{n!}{x!(n-x)!} p^x q^{n-x}$$

x = number of 'successes' in sample,
(x = 0, 1, 2, ..., n)

p = probability of "success" per trial
q = probability of "failure" = (1 - p)

n = number of trials (sample size)



6.) Mean (Expectation)

$$\mu = E(X) = np$$

7.) Variance and Standard Dev.

$$\sigma^2 = npq$$

$$\sigma = \sqrt{npq}$$

$$q = 1 - p$$

Binomial Example

- 3) During the semester a professor cycles to school on 5 days of the week. On any given day, the probability that he arrives at school after 9am is 0.1. For a period of 4 weeks (20 days), calculate the probability that he arrives after 9am

- b) On at least 1 day but no more than 3 days

- doing something n times
- only 2 outcomes: success or failure
- trials are independent of each other
- probability remains constant

1.) All Failures

$$P(\text{all failures}) = (1 - p)^n$$

Continuous Probability Distributions

- the probability that a continuous random variable X will assume any particular value is 0

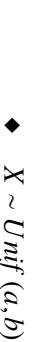
→ Density Curves

- ♦ Area under the curve is the probability that any range of values will occur.
- ♦ Total area = 1

- It is described by the function:

$$f(x) = \frac{1}{b-a}, \text{ where } a \leq x \leq b$$

$$\bullet X \sim \text{Unif}(a, b)$$



Uniform Example

- 5) Suppose the number of donuts a nine-year old child eats per month is uniformly distributed from 0.5 to 4 donuts, inclusive

- a) Find the probability that a randomly selected nine-year old child eats more than two donuts in a month.

$$\begin{aligned} X &\sim \text{Unif}(a, b) & f(x) &= \frac{1}{b-a}, \text{ where } a \leq x \leq b \\ X &\sim \text{Unif}(0.5, 4) & f(x) &= \frac{1}{3.5}, \text{ where } 0.5 \leq x \leq 4 \end{aligned}$$

$$\text{Probability} = \text{Area} = \text{Width} \times \text{Height}$$

$$\text{Probability} = 2 \cdot \frac{1}{3.5}$$

(Example cont'd next page)

$$\begin{aligned} P(x = 1) &= \frac{20!}{(20-1)!} \cdot (0.1)^1 \cdot (0.9)^{19} = 0.27017034353 \\ P(x = 2) &= \frac{20!}{(20-2)!} \cdot (0.1)^2 \cdot (0.9)^{18} = 0.28517980706 \\ P(x = 3) &= \frac{20!}{(20-3)!} \cdot (0.1)^3 \cdot (0.9)^{17} = 0.1901987138 = 0.745470022 \end{aligned}$$

	<pre> . summarize Variable Obs Mean Std. Dev. Min Max before 18 8.888889 .6733995 8 10 after 18 7.666667 1.484214 5 10 diff 18 -1.222222 1.388594 -1 3 </pre>				
One-sample t test					
<pre> x Obs Mean Std. Err. Std. Dev. [95% Conf. Interval] mean = mean(x) Ho: mean = 0 Ha: mean < 0 Pr(T < t) = 0.9995 Pr(T > t) = 0.0010 Pr(T > t) = 0.0005 </pre>					

	<p>→ Interpretation of slope - for each additional x value (e.g. mile on odometer), the y value decreases/ increases by an average of b_1 value</p> <p>→ Interpretation of y-intercept - plug in 0 for x and the value you get for \hat{y} is the y-intercept (e.g. $y=3.25-0.0614 \times \text{SkippedClass}$, a student who skips no classes has a gpa of 3.25.)</p> <p>→ danger of extrapolation - if an x value is outside of our data set, we can't confidently predict the fitted y value</p>
<p>Simple Linear Regression</p> <ul style="list-style-type: none"> → used to predict the value of one variable (dependent variable) on the basis of other variables (independent variables) → $\hat{Y} = b_0 + b_1 X$ → Residual: $e = Y - \hat{Y}_{\text{fitted}}$ → Fitting error: ◆ $e_i = Y_i - \hat{Y}_i = Y_i - b_0 - b_1 X_i$ ◆ e is the part of Y not related to X → Values of b_0 and b_1 which minimize the residual sum of squares are: 	

	<p>(slope) $b_1 = r \frac{S_y}{S_x}$</p> <p>$b_0 = \bar{Y} - b_1 \bar{X}$</p>																								
<p>Properties of the Residuals and Fitted Values</p> <ol style="list-style-type: none"> 1. Mean of the residuals = 0; Sum of the residuals = 0 2. Mean of original values is the same as mean of fitted values $\bar{Y} = \hat{\bar{Y}}$ 																									
<p>We have the decomposition of our observation</p> $Y = \hat{Y} + e$ <p style="text-align: center;"><small><i>y that's related to x</i></small> <small><i>y that's unrelated to x</i></small> <small><i>y that's leftover</i></small></p>																									
	<p>→ Good fit: if SSR is big, SEE is small</p> <p>→ $\text{SST}=\text{SSR}$, perfect fit</p> <p>→ R^2: coefficient of determination</p> $R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$ <p>→ R^2 is between 0 and 1, the closer R^2 is to 1, the better the fit</p> <p>→ Interpretation of R^2: (e.g. 65% of the variation in the selling price is explained by the variation in odometer reading. The rest 35% remains unexplained by this model)</p> <p>→ ** R^2 doesn't indicate whether model is adequate**</p> <p>→ As you add more X's to model, R^2 goes up</p> <p>→ Guide to finding SSR, SSE, SST</p>																								
	<p>Analysis of Variance</p> <table border="1"> <thead> <tr> <th>SOURCE</th> <th>DF</th> <th>SS</th> <th>MS</th> <th>SSR/k</th> <th>SSE/(n-k-1)</th> </tr> </thead> <tbody> <tr> <td>Regression</td> <td>k</td> <td>SSR</td> <td>SSR/k</td> <td></td> <td></td> </tr> <tr> <td>Error</td> <td>n-k-1</td> <td>SSE</td> <td>SSE/(n-k-1)</td> <td></td> <td></td> </tr> <tr> <td>Total</td> <td>n-1</td> <td>SST</td> <td></td> <td></td> <td></td> </tr> </tbody> </table>	SOURCE	DF	SS	MS	SSR/k	SSE/(n-k-1)	Regression	k	SSR	SSR/k			Error	n-k-1	SSE	SSE/(n-k-1)			Total	n-1	SST			
SOURCE	DF	SS	MS	SSR/k	SSE/(n-k-1)																				
Regression	k	SSR	SSR/k																						
Error	n-k-1	SSE	SSE/(n-k-1)																						
Total	n-1	SST																							

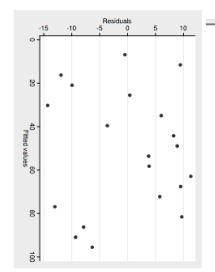
Regression Diagnostics

Standardize Residuals

$$r_i = \frac{e_i}{s_e} \approx \frac{\varepsilon_i}{\sigma} \sim N(0,1)$$

Check Model Assumptions

→ Plot residuals versus Yhat



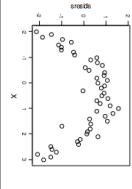
This is the way a residual plot looks when the model fits the data:
No obvious pattern!!!!
resids unrelated to X!!!!!!

- Outliers
- ◆ Regression likes to move towards outliers (shows up as R^2 being really high)
- ◆ want to remove outlier that is extreme in both x and y

• Outtest

- Nonlinearity (ovtest)
- ◆ Plotting residuals vs. fitted values will show a relationship if data is nonlinear (R^2 also high)

As a diagnostic, we plot the standardized residuals versus X:



The nonlinearity is even more evident in the residual plot!! What is wrong with fitting a linear regression to this data?

- ◆ Log transformation - accommodates non-linearity, reduces right skewness in the Y, eliminates heteroskedasticity
- ◆ **Only take log of X variable

so that we can compare models.
Can't compare models if you take log of Y.

• Transformations cheatsheet

→ Normality (sktest)

- ◆ H_0 : data = normality
- ◆ H_a : data \neq normality

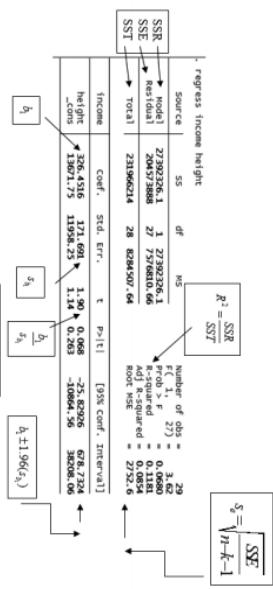
- ◆ don't want to reject the null hypothesis. P-value should be big

• ovtest

Ramsey RESET test using powers of the fitted values of y
Ho: model has no omitted variables
 $F(3, 6044) = 158.43$
Prob > F = 0.0000

Summary of Regression Output

Guide to Regression Output



$$S_{\bar{x}} = \sqrt{\frac{S_x^2}{n-k-1}}$$

- ◆ Homoskedastic: band around the values
- ◆ Heteroskedastic: as x goes up, the noise goes up (no more band, fan-shaped)
- ◆ If heteroskedastic, fix it by logging the Y variable
- ◆ If heteroskedastic, fix it by making standard errors robust

→ Multicollinearity

- ◆ when x variables are highly correlated with each other.
- ◆ $R^2 > 0.9$
- ◆ pairwise correlation > 0.9
- ◆ correlate all x variables, include y variable, drop the x variable that is less correlated to y

→ Homoskedasticity (hettest)

- ◆ H_0 : data = homoskedasticity
- ◆ H_a : data \neq homoskedasticity

```
sktest r45
Skewness/Kurtosis tests for normality joint=
Variable   Pr(Skewness)  Pr(Kurtosis) adj chisq(df)  Pr>chisq
res        0.869      0.046      4.25      0.1195
```