

Biost 518 / Biost 515

Applied Biostatistics II / Biostatistics II



Scott S. Emerson, M.D., Ph.D.
Professor of Biostatistics
University of Washington

Review of Correlation
(From Lecture 7, Biost 517, Fall 2012)

January 12, 2015

Correlation



Correlation Coefficient



- A measure of the tendency of the largest measurements for one variable to be associated with the largest measurements of the other variable
 - Dimensionless
 - The sample correlation r estimates the population correlation ρ (rho)

Pearson's Correlation Coefficient



- Definition of correlation between X and Y:

$$r = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

$$= \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sqrt{\sum_{i=1}^n X_i^2 - n\bar{X}^2} \sqrt{\sum_{i=1}^n Y_i^2 - n\bar{Y}^2}}$$

Possible Values of r

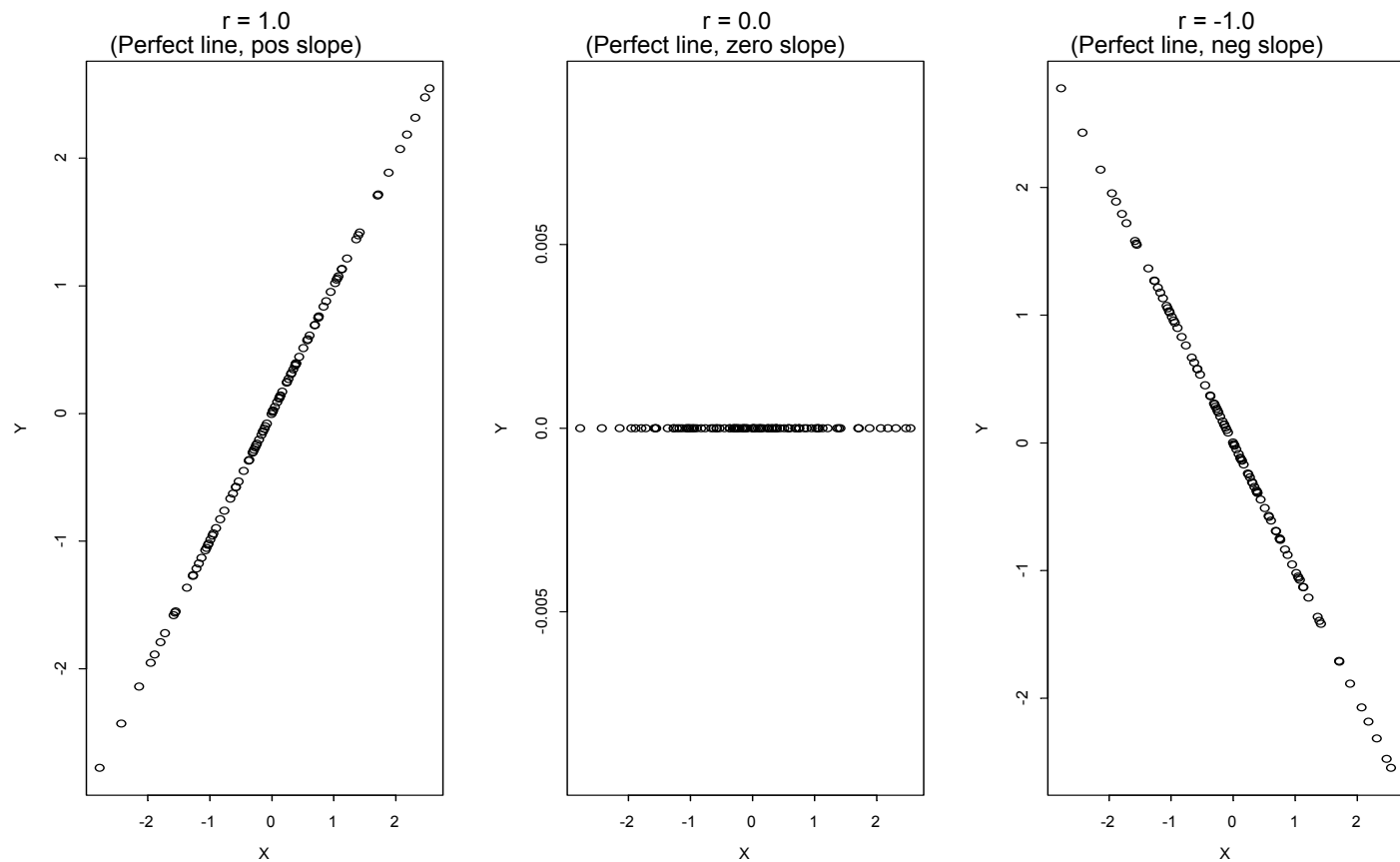


- Range of r : $-1 \leq r \leq 1$
 - $r = 1$: perfect positive correlation
 - a graph of X vs Y will be a straight line with positive slope
 - $r = -1$: perfect negative correlation
 - a graph of X vs Y will be a straight line with negative slope
 - $r = 0$: no correlation

Straight Line Relationships



- Pearson's correlation coefficient with linear data



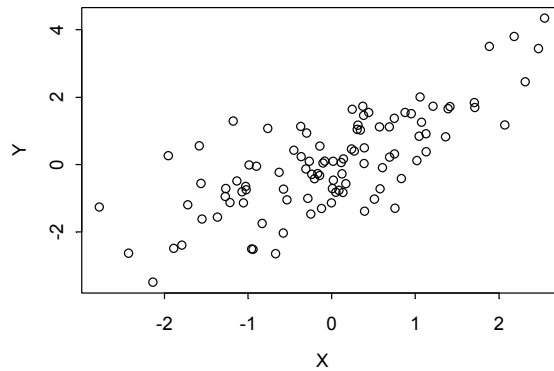
6

Linear Trends in Data

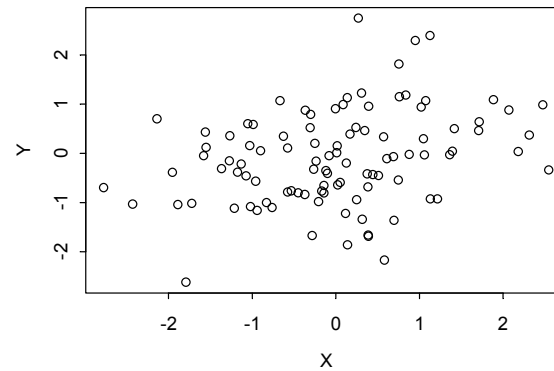


- Pearson's correlation coefficient with variable data

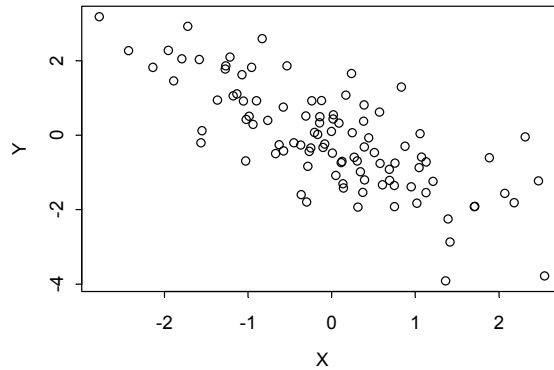
$r = 0.75$



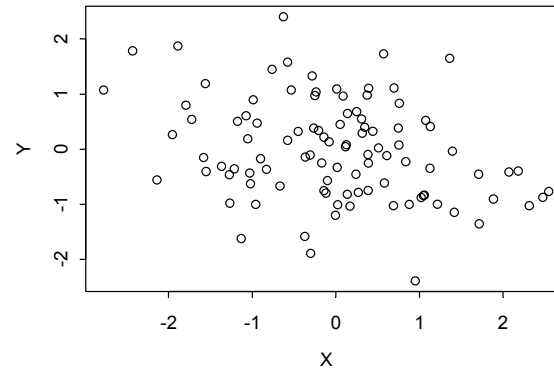
$r = 0.30$



$r = -0.75$



$r = -0.30$



Correlation and Independence



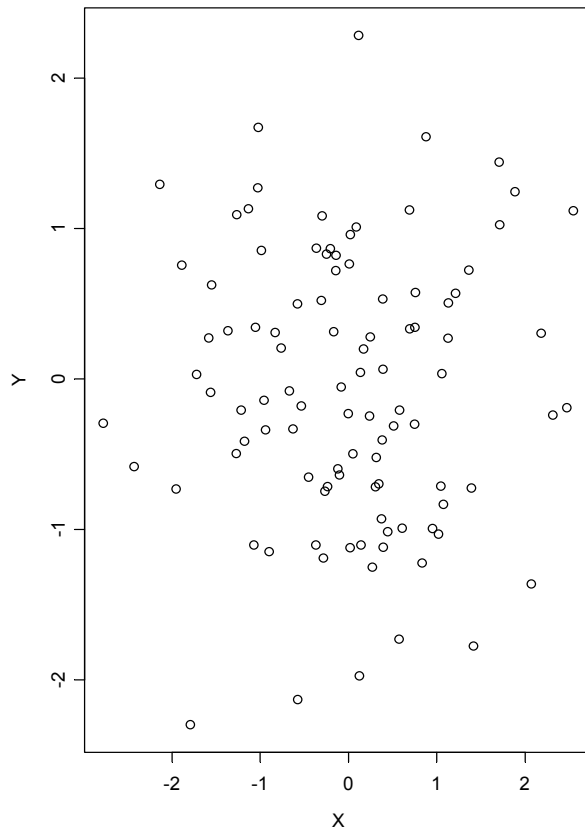
- Independent variables will have $\rho = 0$
 - (and r tending to be close to 0)
- However, uncorrelated variables are not necessarily independent
 - Correlation measures linear trend in the mean of one variable in groups defined by the other
 - It is possible that a nonlinear association exists between two variables, and that the first order trend is a zero slope

Uncorrelated Variables

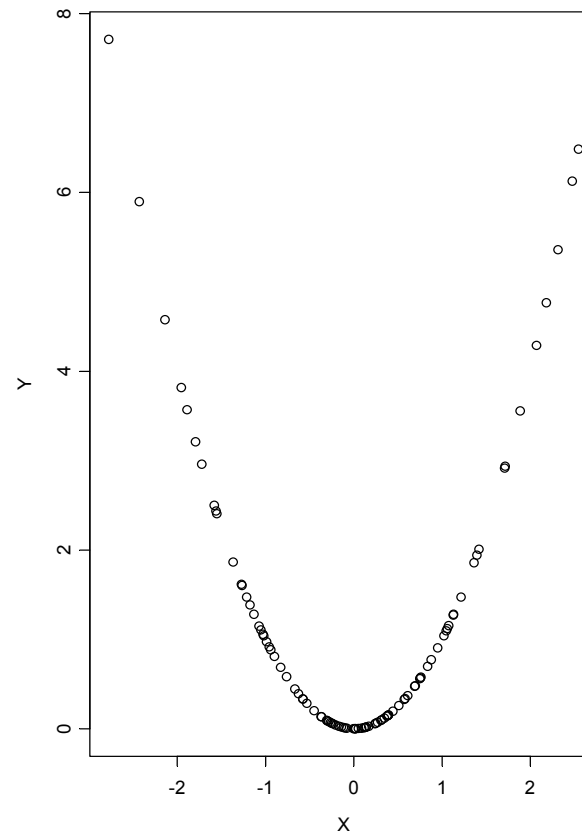


- No linear trend between the variables

$r = 0.0$ (Independence)



$r = 0.0$ (Association, but no linear trend)



Software Commands



- Stata
 - `correlate varlist`
 - Correlation of all pairs of variables
 - Missing data deleted on a casewise basis
 - `pwcorr varlist`
 - Correlation of all pairs of variables
 - Missing data deleted on a pairwise basis
- R
 - `FEV <- read.table("fev.txt", header=TRUE)`
 - `correlate(FEV)`
 - Correlation of all pairs of variables
 - Missing data deleted on a pairwise basis by default

Ex: Correlation in FEV Data

.....

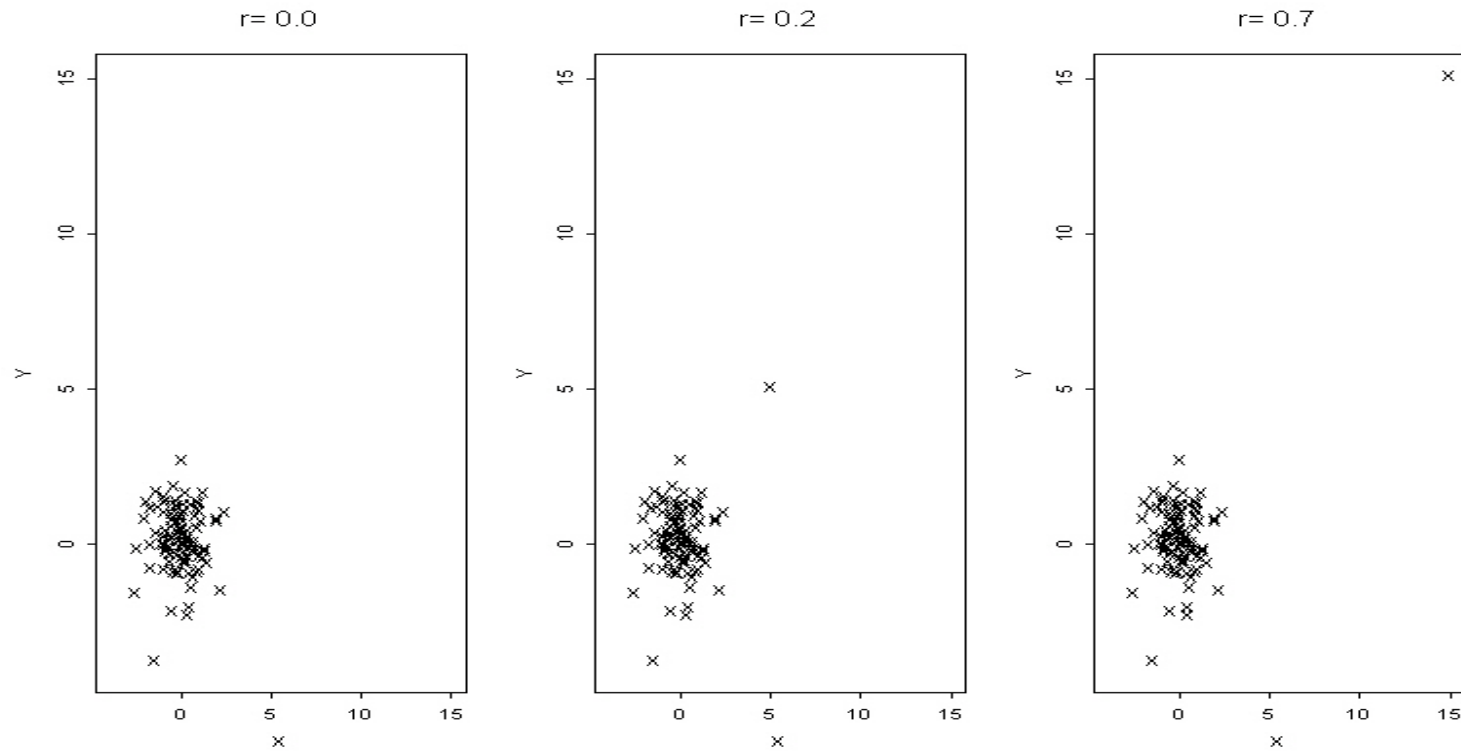
```
. corr subjid age fev height sex smoke
(obs=654)
-----+-----
```

	subjid	age	fev	height	sex	smoke
subjid	1.0000					
age	-0.0112	1.0000				
fev	-0.0147	0.7565	1.0000			
height	-0.0317	0.7919	0.8681	1.0000		
sex	0.0407	-0.0291	-0.2084	-0.1590	1.0000	
smoke	-0.0601	-0.4043	-0.2454	-0.2804	-0.0756	1.0000

- Some of these correlations don't make much sense
 - subjid is a nominal variable
 - sex, smoke are binary variables

Effect of Outliers on r

- Pearson's correlation coefficient can be greatly affected by outliers



Spearman's Rank Correlation



- To decrease the influence of outliers, Spearman's rank correlation coefficient computes the correlation of the ranks of the data
- In the previous example, the rank correlation is always the same: approximately 0.07

Software: Spearman's Correlation



- Stata: `"spearman var1 var2"`
 - Correlation of one pair of variables
 - Cases with missing data for either variable are deleted, and then ranks are computed
- R:
 - `FEV <- read.table("fev.txt", header=TRUE)`
 - `correlate(FEV, method="spearman")`
 - Correlation of all pairs of variables
 - Missing data deleted on a pairwise basis by default

Ex: Correlation in PSA Data



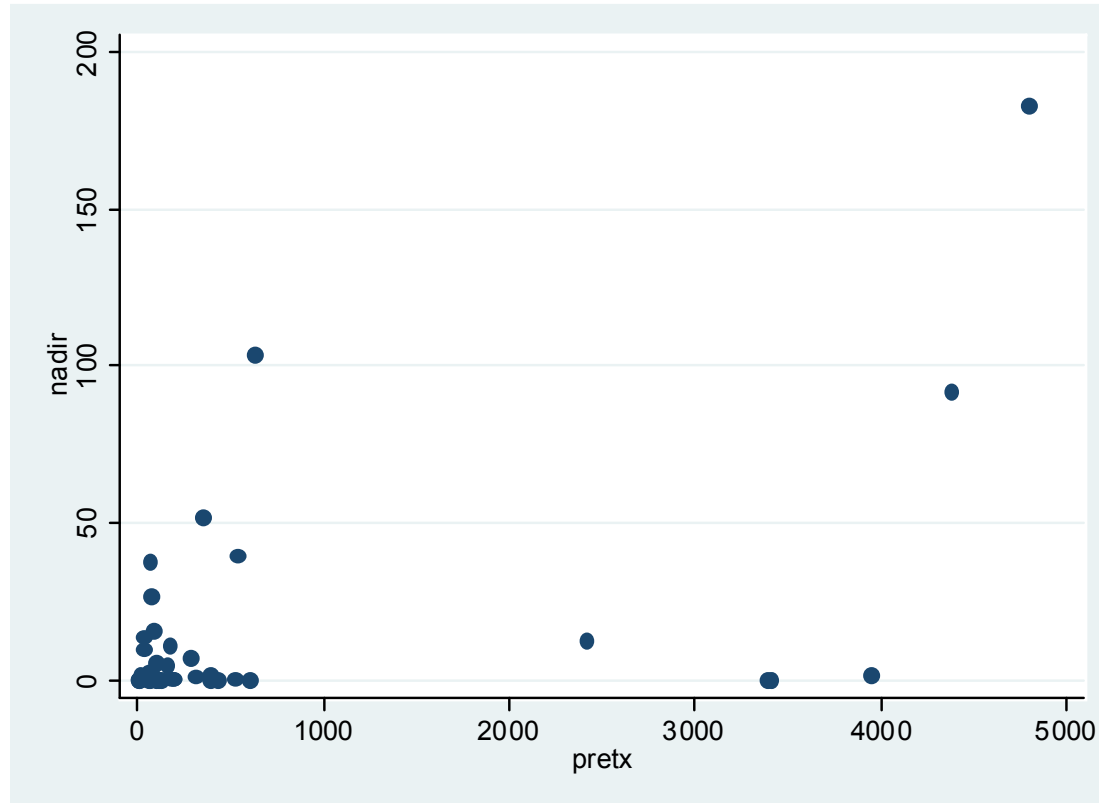
```
corr nadir pretx
(obs=43)
      |      nadir      pretx
-----+-----
nadir|      1.0000
pretx|      0.5371      1.0000
```

```
spearman nadir pretx
Number of obs =      43
Spearman's rho =      0.1489
```

Ex: Nadir vs Pretreatment PSA

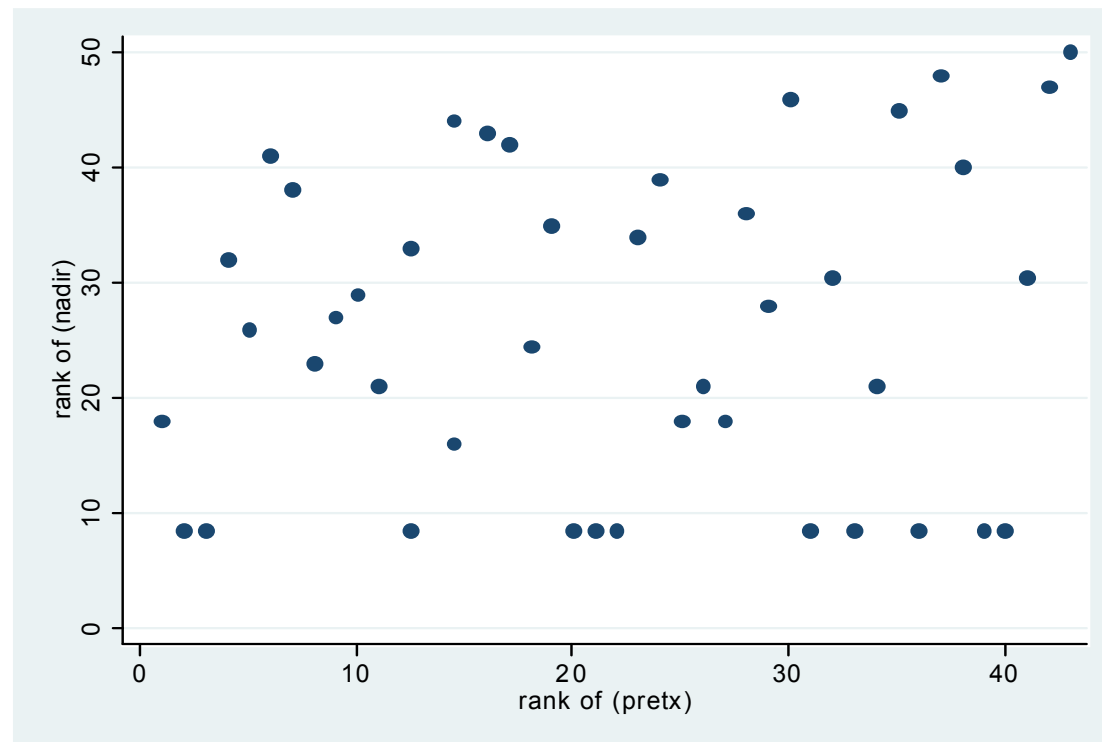


- Scatterplot of nadir versus pretx
– scatter nadir pretx



Ex: Nadir vs Pretx Ranks

- `egen rnknadir = rank(nadir)`
- `egen rnkpretx = rank(pretx)`
- `scatter rnknadir rnkpretx`



17

Ex: Spearman's Corr vs r



- Possible explanation for lower rank correlation with Spearman's
 - Perhaps outliers in distribution of nadir and/or pretx unduly inflate r
 - Perhaps transforming to ranks masks true linear association in skewed variables

Uses of Correlation



- By type of variable
 - Correlation is a mean, thus only makes sense when a mean does
 - Limited interpretability with categorical data
 - Of no scientific relevance with censored data
- By scientific question
 - Greatest relevance when looking for associations between variables
 - But not particularly generalizable across studies

Correlation and Regression



More Interpretable Formula for r



$$r \approx \beta \sqrt{\frac{\text{Var}(X)}{\beta^2 \text{Var}(X) + \text{Var}(Y | X = x)}}$$

β = (LS) slope between Y and X

$\text{Var}(X)$ = variance of X in sample

$\text{Var}(Y | X = x)$ = variance of Y in groups that
have same value of X

Properties of Correlation

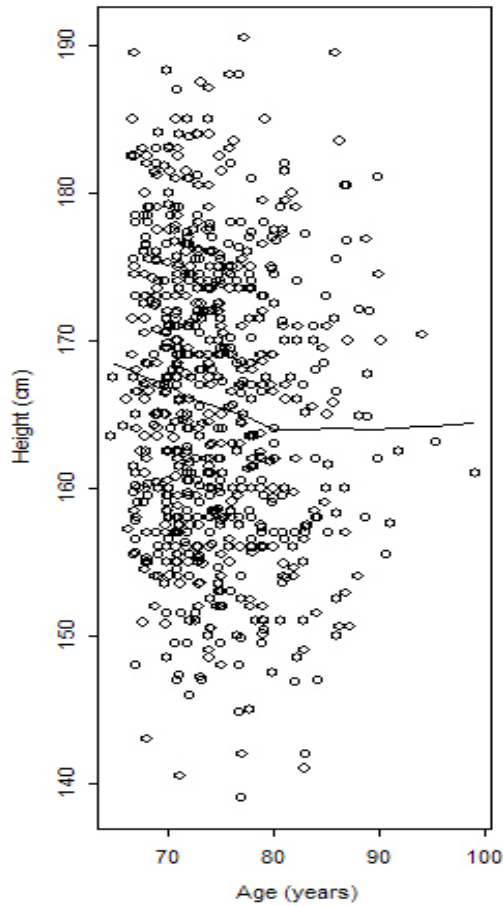


- Correlation tends to increase in absolute value as
 - The absolute value of the slope of the line increases
 - More negative the slope, the more negative the correlation
 - More positive the slope, the more positive the correlation
 - The variance of data decreases within groups that share a common value of X
 - The variance of X increases
 - (Sample size is unimportant in tendencies toward lower or higher correlation)

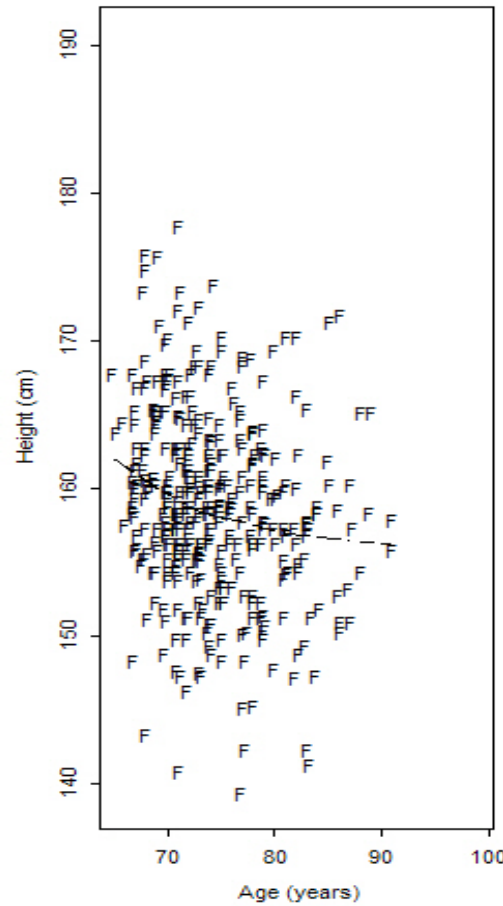
Ex: Height vs Age (by Sex)



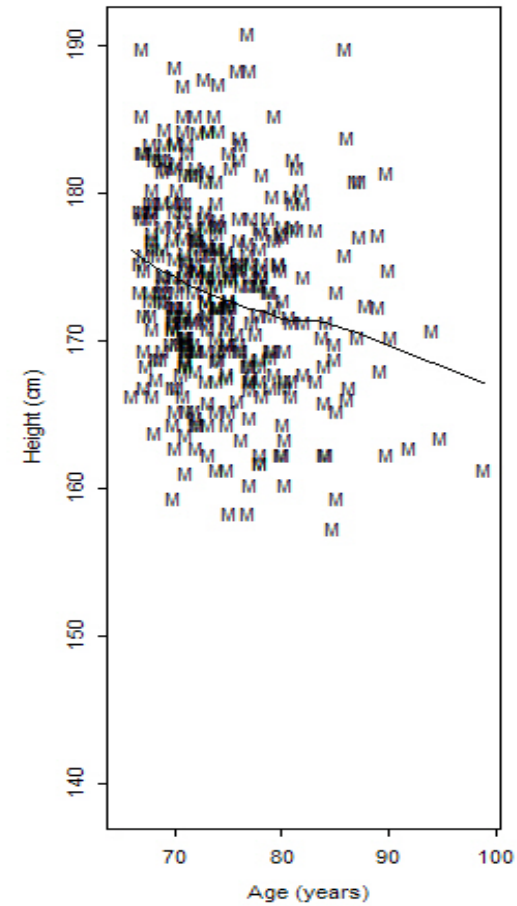
Both sexes: $r = -0.11$



Females: $r = -0.193$



Males: $r = -0.206$



Ex: Height vs Age (by Sex)

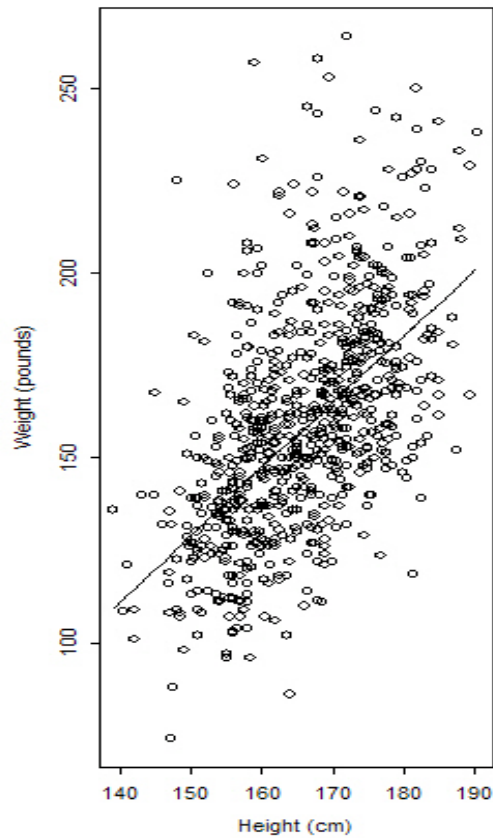


- Correlation between Height and Age
 - Males: $r = -0.206$; Females: $r = -0.193$
 - Combined: $r = -0.110$
- Less extreme r in combined sexes
 - Approximately same slope in each sex and overall
 - Approximately same variance of age in each sex and overall
 - Combined group has higher within group variance of height by age (due to sex effect)

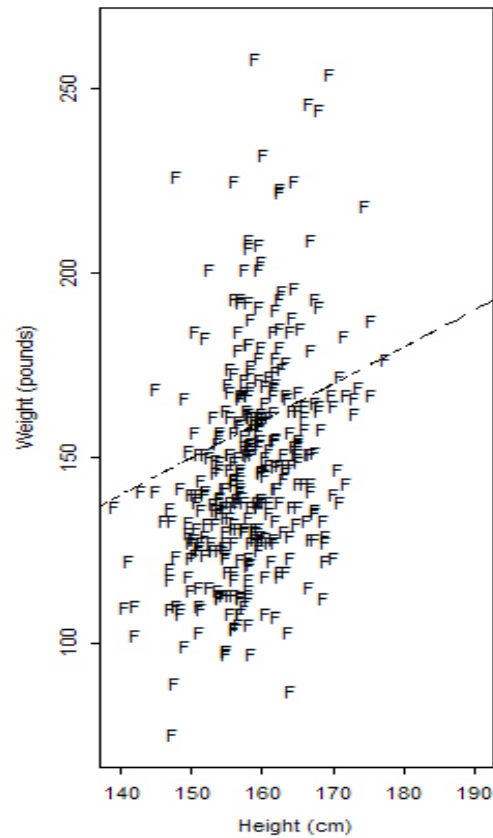
Ex: Weight vs Height (by Sex)



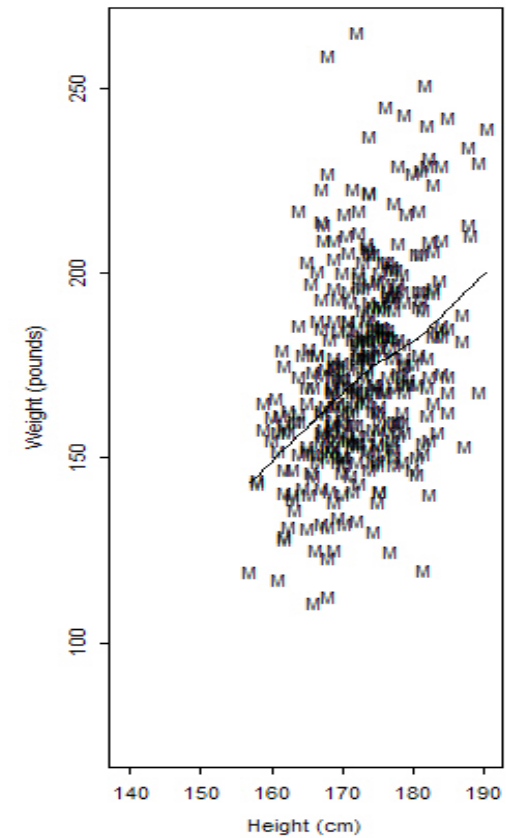
Both sexes: $r = 0.548$



Females: $r = 0.352$



Males: $r = 0.387$



Ex: Weight vs Height (by Sex)



- Correlation between Height and Weight
 - Males: $r = .387$; Females: $r = 0.352$
 - Combined: $r = 0.548$
- More extreme r in combined sexes
 - Approximately same slope in each sex and overall
 - Approximately same within group variance (by height) for each sex and overall
 - Combined group has higher variance of height

Scientific Relevance of r



- It should be noted that
 - the slope between X and Y is of scientific interest
 - the variance of $Y|X=x$ is partly of scientific interest, but it can be affected by restricting sampling to certain values of another variable
 - E.g., $\text{var}(\text{Height} | \text{Age})$ is less in males than when both sexes are included
 - the variance of X is often set by study design
 - This is often not of scientific interest