

Today: More on 2-D Continuous Data

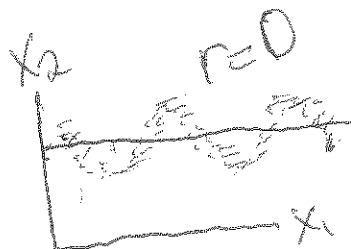
Sam Ventura
36-315

Today: More on 2-D and 3-D Continuous Data

Department of Statistics
Carnegie Mellon University

October 10, 2016

correlation \neq causation



Correlation
 $-1 \leq r \leq 1$

$r \approx -1$: strong negative linear relationship
 $r \approx 0$: no linear relationship
 $r \approx 1$: strong positive linear relationship

2-D Continuous Data - Goals \rightarrow relationship between the two vars?

Data structure

Dataset: n rows, \rightarrow # obs.
2 columns

$\{X_1, X_2\} \in \mathbb{R}^2$

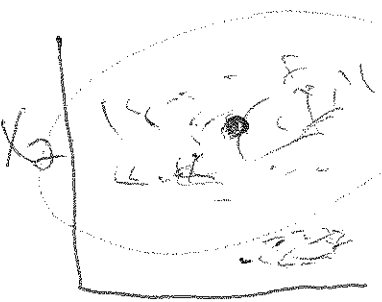
\rightarrow each variable is continuous

Examining variables individually: histograms, boxplots, ruggs, density plots, violin plots
 - center (mean, median), spread (variance, sd, range, IQR)
 skew, modality

Examining variables together/jointly: "Is there a relationship?"
 "what is the relationship?"
linear relationship: $y = mx + b$, $X_2 = \hat{\beta}_0 + \hat{\beta}_1 X_1$
nonlinear relationships: $y = f(x)$, $X_2 = f(X_1)$
 \rightarrow any non-linear fun

2-D Continuous Data - Scatterplots

scatterplots: 2-D plot of X_1 and $X_2 \rightarrow$ mapped to y aesthetic
 \hookrightarrow in the cartesian plane \hookrightarrow mapped to x aesthetic



what information do we get out of scatterplots?

\hookrightarrow is there a relationship? how strong?
what direction is the linear relationship?

\hookrightarrow Trends (linear or nonlinear)

\hookrightarrow joint distribution of X_1 and X_2

\hookrightarrow group structure / clusters of obs

\hookrightarrow similar to modality in 1-D cont.

2-D Continuous Data - Trends

Linear trends: linear regression

uses "least squares" to find the "best" linear fit between X_1 and X_2 . Trying to estimate β_0 and β_1

in this equation: $X_2 = \beta_0 + \beta_1 X_1$ ($y = mx + b$)

short version: minimizes sum of squared errors from points to the regression line

Non-linear trends: "smoothing" splines, loess regression / smoothing

\hookrightarrow fit non-linear / non-parametric, "local" regression to X_1, X_2

Uses a "moving window" to estimate the best fit at

each part of the data potential issues: pick bandwidth overfitting, especially to outliers