

Lesson Plan 3/6

~

Tuesday, March 6, 2018 5:26 AM

Admin:

- Midterm grading delay... coming soon
- HW2 grades posted tomorrow

HDSC:

- The amazing life and contributions of John W. Tukey: <https://www.stat.berkeley.edu/~brill/Papers/life.pdf>

Project 1:

- Questions/comments?

EDA (cont'd)

- Range
- Central Tendency
 - o Why does the mean matter?
- Spread
 - o Variance
 - o IQR
- Outliers
- Skew
- Covariance and Correlation

Exploratory Data Analysis

- Central Tendency
- Spread

- Skew
- Outlier

Central Tendency

- median
- mean (average)

Notation:

n entities
p attributes

Single attribute

\hat{x}_i : i -th observation $i \in \{1, n\}$

$x_{(1)}$: rank smallest value

$x_{(n)}$: largest value

$\rightarrow x_{(\frac{n}{2})}$: median (half of the values
 $\leq x_{(\frac{n}{2})}$)

Sample mean:

$$\rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

* Mean is the optimal value
of a distance optimization

problem

Intuition: central tendency: a value "close" to a lot of the data.

① Define close:

$$(x_i - \mu)^2 \quad \text{(squared distance)}$$

② "a lot of the data"

$$RSS(\mu) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Residual sum of squares

* The sample mean \bar{x} minimizes RSS(μ) over all values of μ

$$\bar{x} = \arg \min_{\mu} \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

→ Need to show

$$\frac{1}{n} \sum_{i=1}^n x_i \text{ minimum} \Rightarrow \text{RSS}(\mu)$$

Find minimum + cst

First derivative

$$f'(x_0) = 0 \Rightarrow$$

$$x_0 = \arg \min_{x^1} f(x)$$

$$\frac{d}{d\mu} \text{RSS}(\mu) = \frac{d}{d\mu} \frac{1}{2n} \sum (x_i - \mu)^2$$

$$= \frac{1}{2n} \sum_i \frac{d}{d\mu} (x_i - \mu)^2$$

(sum rule)

$$= \frac{1}{2n} \sum_i 2(x_i - \mu) \frac{d}{d\mu} (x_i - \mu)$$

(exponent * chain rule)

$$= \frac{1}{2n} \sum_i 2(x_i - \mu) * (-1)$$

$$= \frac{1}{n} \sum \mu - \frac{1}{n} \sum x_i$$

$$= \mu - \frac{1}{n} \sum x_i$$

$$\frac{d}{d\mu} \text{RSS}(\mu) = 0 \implies \mu - \frac{1}{n} \sum x_i = 0 \implies$$

$$\hat{\mu} = \frac{1}{n} \sum x_i = \bar{x}$$

Estimates are minimizers of
optimization problems over
data

* median(x) minimizes

$$SAD(m) = \sum_{i=1}^n |x_i - m|$$

Sample Variance:

$$\text{var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = RSS(\bar{x})^*$$

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

standard deviation

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

$$\text{Var}(x) = \sqrt{n} \cdot \text{Cov}(x)$$

Rank-based spread:

IQR: inter-quartile range

$$\underbrace{[x_{(\frac{n}{4})}, x_{(\frac{3n}{4})}]}_{\text{half the data within this range}} \quad (\text{half the data within this range})$$

Outliers:

values that are unusually "far away" from center

$$\text{outliers}_{\leq 1}(x) = \{x_i : |(x_i - \bar{x})| > K * \text{sd}(x)\}$$

$$K = 1.5$$

$$K = 3$$

$$\text{outliers}_{IQR}(x) = \left\{ x_i : \begin{array}{l} x_i < X_{\left(\frac{1}{4}\right)} - K * \text{IQR}(x) \\ \text{or} \\ x_i > X_{\left(\frac{3n}{4}\right)} + K * \text{IQR}(x) \end{array} \right\}$$

Skew

Is the bigger spread below us.
above central tendency

$$X_{\left(\frac{n}{2}\right)} - X_{\left(\frac{n}{4}\right)}$$

lower spread

$$X_{\left(\frac{3n}{4}\right)} - X_{\left(\frac{n}{2}\right)}$$

higher spread

Central: median mean $\begin{cases} \text{optimality} \\ \text{of these} \\ \text{estimates} \end{cases}$

> spread: variance
standard deviations

IQR

skew : IQR

outliers: sd rule
IQR rule



Covariance

Correlation

|

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

... 1 1

Unit-less

$$\text{cor}(x, y) = \frac{\text{cov}(x, y)}{\text{sd}(x) \text{sd}(y)}$$

[-1, 1]