

Lesson Plan 3/8

Thursday, March 8, 2018 3:11 PM

~

Admin

- Grades are coming, we swear...

HDSC

- Visual explanation of stats and probability: <http://students.brown.edu/seeing-theory/>

Project 1

- Questions/comments?
- Operating on all pairs of rows from two tables

Data transformations

Handling missing data

Example:

- dates : difference in days

$$d = |x_i - x_j|$$

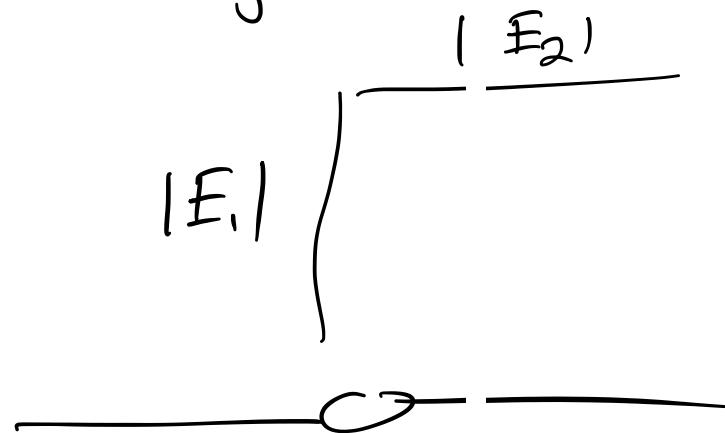
$$s = e^{-d}$$

- cat similarity $\sim \text{if } v_1 = v_2$

$$S(v_1, v_2) = \begin{cases} 1 & v_1 = v_2 \\ 0 & v_1 \neq v_2 \end{cases}$$

Similarity matrix

$$S_{ij} : \text{similarity}(i, j)$$



- ① Transform so that it satisfies

1 0 , 1

- central tendency
- spread

Can we perform transformation

- C.E. = 0
- spread = 1

Removes units

\rightarrow Standardization

Given data x_1, \dots, x_n

$$N = \bar{x} + 1 \text{ std}$$

$$z_i = \frac{x_i - \bar{x}}{sd(x)} \quad \begin{matrix} \leftarrow \text{data units} \\ \in \text{data units} \end{matrix}$$

Q1) \bar{z} :

$sd(\bar{z})$:

$$\bar{z} = \frac{1}{n} \sum_i z_i \quad (\text{def'n})$$

$$= \frac{1}{n} \sum_i \left(\frac{x_i - \bar{x}}{sd(x)} \right) \quad (\text{plug-in})$$

$$= \frac{1}{n} \sum_i \frac{x_i}{sd(x)} - \frac{1}{n} \sum_i \left(\frac{\bar{x}}{sd(x)} \right)$$

$$\begin{aligned}
 &= \frac{1}{sd(x)} \left[\frac{1}{n} \sum_i x_i - \frac{1}{n} \sum_i \bar{x} \right] \\
 &= \frac{1}{sd(x)} \left\{ \bar{x} - \frac{1}{n} \cdot (n\bar{x}) \right\} \\
 &= 0
 \end{aligned}$$

$$sd(z) = 1$$

Centering : $z_i = (x_i - \bar{x})$ \bar{z} ? $sd(z)$?

Scaling : $z_i = \frac{x_i}{sd(x)}$ \bar{z} ? $sd(z)$?



$\longleftarrow \cup$

Categorical \rightarrow numeric

One-hot-encoding

\rightarrow Given categorical variable x
can take values $\{v_1, v_2, \dots, v_m\}$

$$\Rightarrow \begin{array}{c} v_1 \quad v_2 \quad v_3 \quad \dots \quad v_m \\ \hline | \quad | \quad | \quad \dots \quad | \\ v_{i1} = \begin{cases} 1 & \text{if } x_i = v_1 \\ 0 & \text{otherwise} \end{cases} \end{array}$$

$$v_{ij} = \begin{cases} 0 & \text{if } x_i = v_j \\ 1 & \text{if } x_i \neq v_j \\ 0 & \text{otherwise} \end{cases}$$

— C —

Reducing data skew with logarithmic transform

$$z_i = \log(x_i)$$

Only works
if $x_i > c >$



shifted log

\rightarrow $x_i - \min(x)$

$$z_i = \log(x_i - \min(x))$$

signed shifted log transform

$$z_i = \text{sign}(x_i) * \log(|x_i| + 1)$$

$$\hookrightarrow \text{sign}(x_i) = \begin{cases} +1 & \text{if } x_i > 0 \\ 0 & \text{if } x_i = 0 \\ -1 & \text{if } x_i < 0 \end{cases}$$

→

Handling missing data

→ missing @ random

vs. " systematically

→ imputation & its effect
on central tendency &
spread