

Admin

- Hw3 due 9pm
- Exam #1 [2/14]
- mini projects
- Hw graded!
- 2/10 signups
- 2/23 slides due
- 2/24, 2/25 prez

Agenda

1. normalizing data
2. Variance, Std dev
3. Python

1. Normalizing Data

↳ tuesday: correlation
 ↳ linear regression } \overline{x} (ind) \overline{y} (dep)

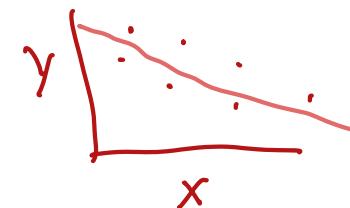
Conc: adm rate, tuition highly correlated

Another ~~BB~~ in time series data

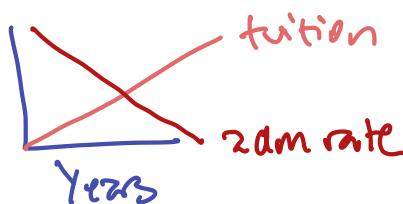
- sequence of observations assoc. with points in time
 - ↳ successive, equally spaced

- time series analysis: extracting predictions and filling in gaps

- given x , what is y ?
- given 2 new x , what would y be?



- ↳ cold end up here if values up and down
- ↳ but both are consistently ↑ or ↓ so they may also correlate with time

Expectation

time: ind variable
 tuition, adm: dep variable

↳ But! per year: tuition & features
 adm rates

2014 vs. 2015?
 2020 vs. other year?
 2008 vs. 2024
 or
 2020 vs. 2024
 {Euclidean distance}!

Euclidean distance assumes
 that: unit change in any
 direction is equally significant

(ex) 2014 vs. 2015

adm	32.22	28.48
tuition	60290	62800

$$\delta(2014, 2015) = \sqrt{(z-z')^2 + (t-t')^2} = 2510$$

2014 vs 2015

Drop in adm rate?

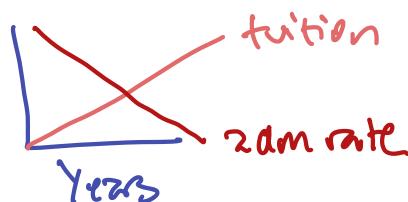
$$\begin{aligned} \text{adm rate} &= 18 \\ \text{tuition} &= 62800 \end{aligned}$$

$$\sqrt{(32.22 - 18)^2 + (60290 - 62800)^2}$$

2510.04

2014 vs. 2016

expectation



problem #1: one feature has huge impact over the other

problem #2: communication



To solve:

tuition: 50,000 to 78,000

adm : 8 to 33

Solution → normalization!

$$\cdot \max = 1 \quad \cdot \min = 0$$

→ smallest val in population: 0

largest val in pop: 1

for every $x_1, \dots, x_{m_{2X}}, x_{\min}$ everything else scaled

$$x_{\text{scaled}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

(ex) pop = 1, 2, 3, 4, 5

min 1
max 5

↳ Feature scaling

$$\frac{3-1}{5-1} = \frac{2}{4} = \text{_____}$$

min(max normalization):
↳ 0 to 1 range

downside: no units
 |
 | suppress athletes

2. Varianz + Std Deviation

↳ How different from the mean?

The figure consists of two horizontal red lines. The left line has six tick marks above it and one tick mark below it, with the word "low var" written below it. The right line has four tick marks above it and three tick marks below it, with the word "high var" written below it.

- Variance average squared difference from mean

62

Used for: Comparing 2 datasets
(end up w/bigger #'s)

- std dev very difference from mean

6

used for one dataset

Used for - one dataset

(keeps units you started with)

$$\sigma^2 = \frac{1}{N-1} \sum (x_i - \bar{x})^2$$

$$6 = \sqrt{6^2}$$

(basically a mean, but we divide by $N-1$)

(Bessel's, to correct for not having every data point)

$$\text{lst2} = \{-12, 30, 50, -55, -75, 80\}$$

$$s^2(\text{lst2}) \quad \bar{x} = 3$$

3708

$$\text{lst1} = \{-1, 3, 5, -3, 6, 8\}$$

$$s^2(\text{lst1}) \quad \bar{x} = 3$$

18

Variance: lst2 varies \approx lot more than lst1

$$s^2(\text{lst2}) = 60.89$$

$$s^2(\text{lst1}) = 4.24$$

↳ on avg, how far from mean

↳ on avg, how far from mean

In Python: `{statistics.variance}`

Outlier: > 2 std devs from mean