DSE210 HW5 - Ryan Inghilterra

Ryan Inghilterra
03/15/2018
DSE 210

## HW 5 - Worksheet 10 - PCA and SVD

① set of vectors $\begin{pmatrix} 3 \\ 4 \\ 0 \end{pmatrix}$, $\begin{pmatrix} 4 \\ -3 \\ 0 \end{pmatrix}$, $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$ an orthonormal basis of $\mathbb{R}^3$?

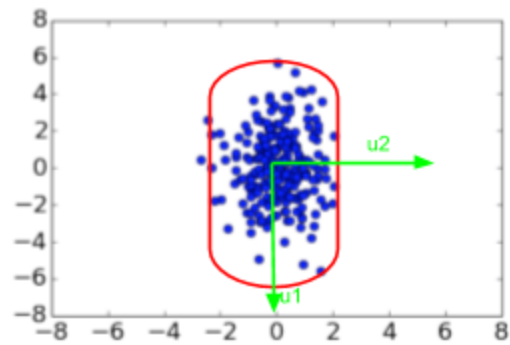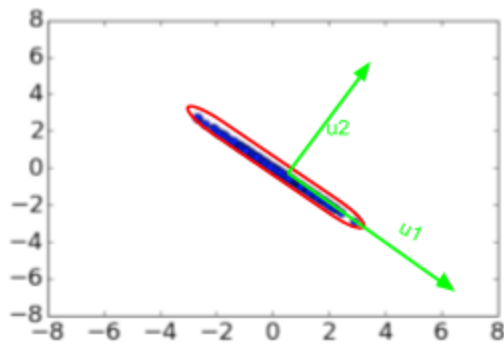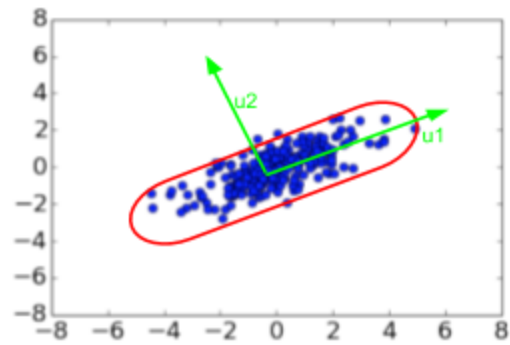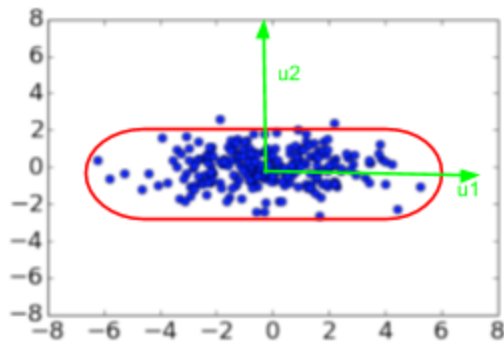to be orthonormal first off all must be unit vectors

so $\|u\| = 1$, we can see that for $V_1 = \begin{pmatrix} 3 \\ 4 \\ 0 \end{pmatrix}$, $\|V_1\| = \sqrt{3^2 + 4^2 + 0^2} = \sqrt{25}$

$= 5$

$\|V_1\| = 5$, so not a unit vector

therefore NO set of vectors not an orthonormal basis of $\mathbb{R}^3$

② SEE NEXT PAGE

2. The following four figures show different 2-dimensional data sets. In each case, make a rough sketch of an ellipsoidal contour of the covariance matrix and indicate the directions of the first and second eigenvectors (mark which is which).

③ Let $u_1, u_2 \in \mathbb{R}^d$ be two vectors, $\|u_1\| = \|u_2\| = 1$ and $u_1 \cdot u_2 = 0$

Define $U = \begin{pmatrix} \uparrow & \uparrow \\ u_1 & u_2 \\ \downarrow & \downarrow \end{pmatrix}$
$\qquad\qquad d \begin{pmatrix} \uparrow & \uparrow \\ u_1 & u_2 \\ \downarrow & \downarrow \end{pmatrix}$
$\qquad\qquad\qquad \underset{2}{x}$

(a) what are dimensions of following?

- $U$ — $\boxed{d \times 2}$
- $U^T$ — $\boxed{2 \times d}$
- $UU^T$ — $\boxed{d \times d}$
- $u_1 u_1^T$ — $\boxed{d \times d}$

(b) what are the differences between following four mappings?

1. $x \mapsto (u_1 \cdot x, \, u_2 \cdot x)$
2. $x \mapsto (u_1 \cdot x) u_1 + (u_2 \cdot x) u_2$
3. $x \to V^T x$
4. $x \to UU^T x$

mapping 1 and 3 are equivalent, both are projections
of $x$, from $x \in \mathbb{R}^d$ to $\mathbb{R}^2 \in \mathbb{R}^2$ (d to K dimensions)
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{K=2}{\uparrow}$

mapping 2 and 4 are equivalent both are reconstructions
$x \in \mathbb{R}^d$ projected to $\mathbb{R}^2$ reconstructed back onto $\mathbb{R}^d$
$\qquad\qquad\qquad\qquad \underset{K=2}{}$

Ryan Inghilterra
DSE #10
03/15/2018

# HW - Worksheet 11 - Sampling

① 9 red marbles, 1 blue marble. 900 rand. draws w/ replacement.

$n = 900$    $Pr(red) = 0.9 = p$

$N(n \cdot p, \; n \cdot p (1-p))$    since observational number

$N(900 \cdot (0.1), \; 900 \cdot (0.9) \cdot (0.1))$

$= \boxed{N(810, 81)}$

② in world say 1% people are left handed. sample 200 people at random.
Give 99% confidence interval for # of them left handed

$N = 200$    Find STD, $var = np(1-p)$    $Pr(L) = 0.01 = p$

$STD = \sqrt{200(0.01)(0.99)}$    mean $M = n \cdot p$

$= \sqrt{1.98} \approx 1.41$    $M = 200 \cdot 0.01$

multiply STD by 3 since 99%    $M = 2$

99% confidence interval left handed $= 2 \pm \sqrt{1.98}(3)$

$\boxed{= 2 \pm 4.22}$

③ 20 wedges numbered 1-20. Half red, half black, 100 darts thrown

$n = 100$    $Pr(R) = 0.5$    $Pr(B) = 0.5$    $E(X_i) = 1/20$

(a) $X_i$ # darts fall in wedge 1. What is $E(X_i)$ and $Var(X_i)$

$E(X_i) = \left(\frac{1}{20}\right)(100) = \boxed{5 = E(X_i)}$

$\quad\quad\quad\quad\;\; p \quad\quad n$

$var(X_i) = n \cdot p (1-p) = (100)(.05)(.95) = \boxed{4.75 = Var(X_i)}$

(b) upper bound on $X_i$, 95% confidence.    (normal approx)

$M +- 2 \cdot std \to 5 \pm 2 \cdot \sqrt{4.75} \to \boxed{5 \pm 2.18}$

c) $Z_r = $ # wedges red, $Z_b = $ # wedges black, $Z = |Z_r - Z_b|$

$Y_i = \begin{cases} 1 & \text{if dart on red} \\ -1 & \text{if dart on black} \end{cases}$, $Z_r - Z_b = Y_1 + Y_2 \dots Y_{100}$    want 99% confidence interval

(c) $\boxed{E(Y_i) = 0}$  $\boxed{var(Y_i) = 1}$    d.) central limit theorem $Z_r - Z_b$ (approx normal)

(e) 99% conf. interval →    $E(Y) = 0$  $var(Y) = n$, $std(Y) = \sqrt{n} = \sqrt{100} = 10$

$M \pm 3 \cdot std = 0 \pm 3(10)$  $\boxed{0 \pm 30}$    $\boxed{M = 0, \; std = 10, \; N(0,100)}$

④ colorblind appear 1% people. How large sample for prob. of it containing
at least 1 colorblind person to be at least 95%.

$M \pm 2 \cdot std$   since 95% conf. interval   $M = (0.1 \cdot N)$  $var = (0.1) n (0.99)$

think in terms of prob. nobody colorblind less than 5%

$$.99^N < .05^1 \qquad \boxed{N > 298.073}$$

$N \ln(.99) < \ln(.05)$   (sample size need to be 299 or greater)

⑦ tolerate chance errors of 1% or so of estimate. Use sample size
100, 250, 10,000?  Aux. source of info suggests population % range 20% to 40%

want estimate to be $M \pm .01$
$$\sigma = \sqrt{\frac{p \cdot (1-p)}{n}} \qquad \text{so } p = .40 \text{ since will}$$
be worse then

for $N = 100$  $\sigma = \sqrt{\frac{(0.4)(0.6)}{100}} = .05 \leftarrow$ too high!

for $N = 250$  $\sigma = \sqrt{\frac{(0.4)(0.6)}{250}} = .03 \leftarrow$ too high!   $M = p = 0.4$   $P = 0.2$

for $N = 10,000$  $\sigma = \sqrt{\frac{(0.4)(0.6)}{10,000}} = .0048 \leftarrow$ good ✓

Since we want 95% conf. interval that is

$M \pm 0.01$     so where   $2 \cdot \sigma \leq 0.01$

only for $N = 10,000$, is $2\sigma < 0.01$

$2(.0048) = .0096 < .01$ ✓

so should use sample size of $\boxed{10000}$

⑧ city there are 100,000 people age 18 to 24. random sample 500 people
194 enrolled in college. Estimate % of all people 18-24 in city in college
give 95.5% confidence interval.

$n = 500$   $194/500 = \hat{p}$   (fraction)

want $M \pm 2\sigma$   $\frac{1}{9.388}$   $M = \hat{p} = 0.388$  $\sigma = \sqrt{\frac{p(1-p)}{n}}$

$$= \sqrt{\frac{(.388)(1-.388)}{500}}$$

so 95.5 conf. interval is
(0.388 ± 0.02)

$\sigma = 0.02$

Ryan Inghilterra
PSE 210
03/15/2018

## HW 5 Worksheet 11 - Sampling. Continued

⑨ $n = 1000$ what sample size should they use?

(SAME sample size of 1000)

what matters is the sample size, not overall population size.

⑩ $n = 100$  $\hat{\mu} = 307$  $\hat{\sigma} = 30$  $\sigma = \dfrac{\hat{\sigma}}{\sqrt{n}} = \dfrac{30}{\sqrt{1000}} = 0.949 = \sigma$

$\boxed{M = 307 \quad \sigma = 0.95}$

we dont know, so use formula w/ sample std $\hat{\sigma}$

Ryan Inghilterra
DSE 210
03/15/2018

(HW5)
## Worksheet 12 - Hypothesis Testing

① 1990 US, 2.1 million deaths from all causes, compared to 1.7 million in 1960. 25% increase, show public health got worse over 1960-1990?

[NO], not necessarily, all causes can include causes other than health, also not taking into account population increase.

② smoking study

ⓐ controlled or observational? (observational)

ⓑ don't want to confound different biases between age groups

ⓒ that conclusion is not necessarily true and can be misleading. People could have stopped smoking b/c of health issues. And those who were healthy had no reason to stop smoking.

③ oral contraceptive study

ⓐ controlled or observational? (observational)

ⓑ want to control for confounding factors. Age, education levels, and marital status can have a big influence on the results of this experiment.

ⓒ study only showed there exists a relationship between the pill and cervical cancer. Did not prove it is a causal relationship

⑥ 10,000 tossings, coin comes up heads 5,400 times. Conclude coin is biased

ⓐ (null hypothesis: coin is unbiased)

ⓑ compute z-statistic and p-value

observed $= 5400$   expected $= 5000$

$$std = \sqrt{n \cdot p(1-p)} \quad p = 0.5$$

$$std = \sqrt{(10000)(.5)(.5)} \quad n = 10,000$$

$$std = \sqrt{2500} = 50$$

$$z\text{-stat} = \frac{observed - expected}{std}$$

$$= \frac{5400 - 5000}{50} = \frac{400}{50} = 8 = z\text{-stat} \quad (p\text{-value is} < 0.0001)$$

ⓒ strong evidence against the null, reject at 0.05 significance level (conclude coin is biased)

(7) die is rolled 100 times. total # of spots = 368 instead of expected 350.
Can this be explained as chance variation or die loaded?

null hypothesis: die is fair    mean of single dice roll = 3.5

std (dice roll) = $\frac{1}{6}((1-3.5)^2 + (2-3.5)^2 + \ldots (6-3.5)^2} = \sqrt{35/12} = 1.71$

divide std 1.71 by √sample size n , $\frac{1.71}{\sqrt{100}} = 0.171$

mean of sample dice roll = 3.68

Z-stat = $\frac{(3.68 - 3.5)}{0.171} = \frac{0.18}{0.171} = 1.05$       p-value is 0.146859

result is not significant at $p < 0.05$
therefore we cannot reject the null, we accept the null.
( Conclusion: die are not loaded )


(8) other things being equal, which is better for null hypothesis:
a higher p-value or lower p-value?
( a higher p-value ), higher p-value more likely
to accept the null. since p-value is the probability
of observing the experiment results under the null.

Ryan Inghilterra
DSE 210
03/15/2018

## (HW 5) - Worksheet 12 - Hypothesis Testing Cont.

⑨ drug abuse survey, each year 1985 and 1992, 700 random people sampled

ⓐ age 18 to 25, percentage of marijuana users dropped from 21.9%
to 11.0%. Real or chance variation? $n = 700$

observed difference = $.219 - .11 = .109$

$$\sigma_{1985} = \sqrt{\frac{(.219)(.781)}{700}} = .01563 \qquad \sigma_{1992} = \sqrt{\frac{(.11)(.89)}{700}} = .01182$$

null hypothesis: $X_{1985} = X_{1992}$ (drop is from chance variation)

calculate $\sigma = \sqrt{\sigma_{1985}^2 + \sigma_{1992}^2}$   $\sigma = \sqrt{(.01563)^2 + (.01182)^2} = .0196$

so   Z-stat = $\dfrac{.109}{.0196} = 5.56$, reject the null.

Conclusion: Drop in marijuana use from 1985 to 1992 is real!

ⓑ cigarette 36.9 to 31.9, difference real or chance in variation?

observed difference = $.369 - .319 = .05$

null hypothesis: $X_{1985} = X_{1992}$ (drop is from chance variation)

$$\sigma_{1985} = \sqrt{\frac{(.369)(.631)}{700}} = .01823 \qquad \sigma_{1992} = \sqrt{\frac{(.319)(.681)}{700}} = .01761$$

calculate $\sigma = \sqrt{(.01823)^2 + (.01761)^2} = .02534$

so   Z-stat = $\dfrac{.05}{.02534} = 1.97 \rightarrow$ p-value is .024419

we are testing at 95% confidence interval, result is significant
at $p < 0.05$ level
therefore reject null

Conclusion: Drop in cigarette use from 1985 to 1992 is real!

(10) 1000 freshman sample. Average # hours = 12.2 worked
std = 10.5 for public school
- other survey at private university Average # hours = 9.2 std = 9.9
difference between two averages due to chance?

hypothesis

null: means of the two distributions (public vs private universities)
are the same. Differences due to chance

$\theta = \sqrt{\theta_1^2 + \theta_2^2} \longrightarrow \qquad \theta = \sqrt{(10.5)^2 + (9.9)^2} \qquad = 14.43$

z stat = $\frac{12.2 - 9.2}{14.43}$ = $\frac{3}{14.43}$ = 0.207 $\qquad$ p-value $\approx$ 0.418 oc5

↑ high p-value!

result is not significant at $p < 0.05$

therefore we cannot reject the null, we accept the null.

(Conclusion: the differences between the two averages are
due to chance)