

PCA (Principal Component Analysis)

- Dimension reduction (denoise)
- Capture group structure
- Data visualization (high dimension)

Case Study : ASVAB tests and AFQT

- Screening tests to join army
- ASVAB : 10 tests Word, Parag, math, arith, science, Coding ...

AFQT : Word, Parag, math , arith

Goal : 1. How to summarize the tests to reveal one's intelligence or some abilities?

2. How to display 10 dimensional data
3. How AFQT is formed ?

0.

A toy example:

Case 1:

Two sets of tests, x_1, x_2

Goal: to see how students do?

x_1 : spread out

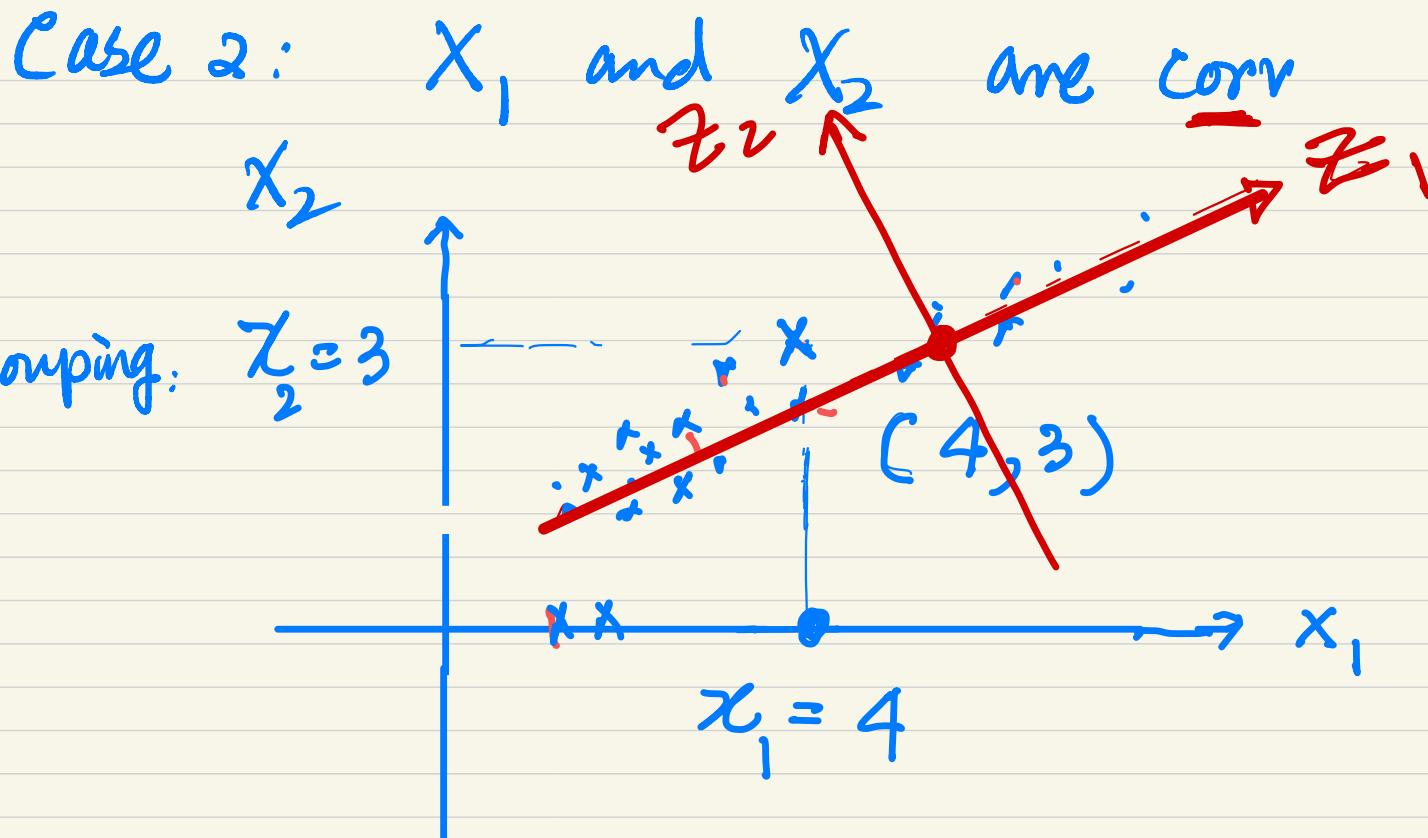
x_2 : Every one got ≈ 30
too hard!



x_2 provides no information: throw it away

only uses x_1 !

x_1 : Principal Component



option 1: Keep both $X_1 + X_2$

$$PC1: Z_1 = \underbrace{\varphi_{11}X_1 + \varphi_{21}X_2}_{\text{1: Use } X_1 \text{ alone}}$$

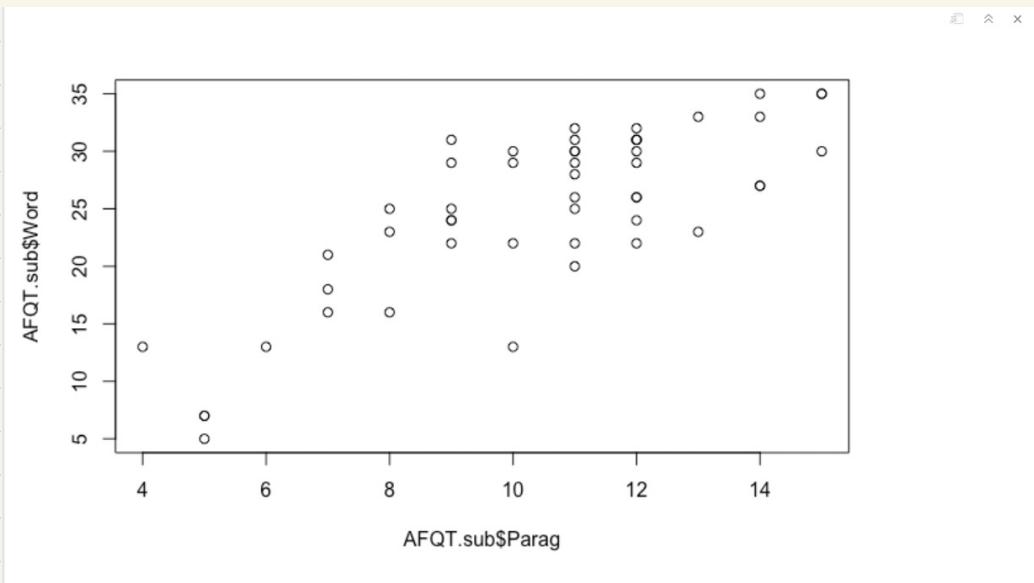
$$PC2: Z_2 = \varphi_{12}X_1 + \varphi_{22}X_2 \quad 3: \text{take weighted sum : } Z_1$$

$\Rightarrow PCA$

- How to do it?
- How much information lost?

1. Simple case : $x_1 = \text{Parag} \Rightarrow \bar{x}_1 = 10.4, s_1 = 2.75$
 $x_2 = \text{Word} \qquad \bar{x}_2 = 24.9, s_2 = 7.41$

$n = 50$



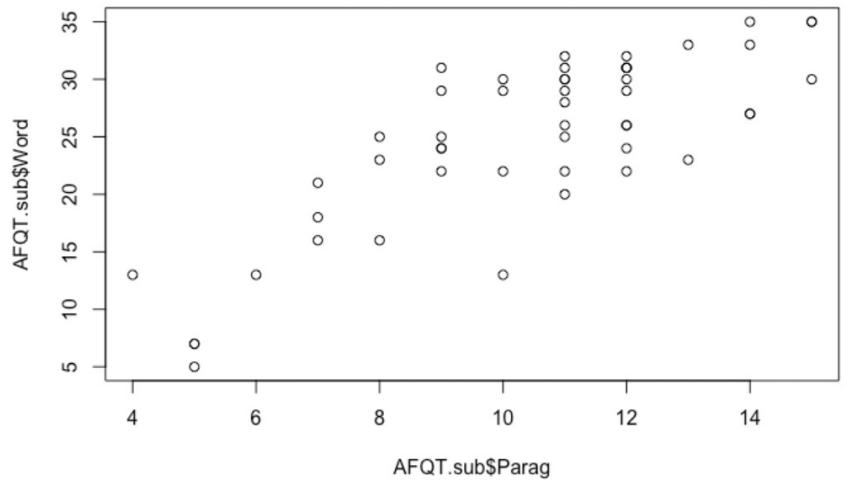
Goal : Find a new score

$$z_1 = \varphi_{11} \cdot x_1 + \varphi_{21} \cdot x_2, \varphi_{11}, \varphi_{21} \text{ unknown}$$

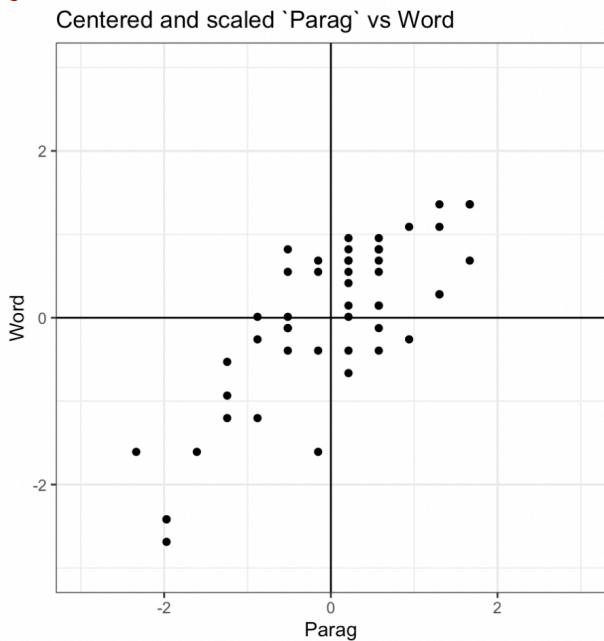
z_1 : Informativ

φ : Informativ ?

Original



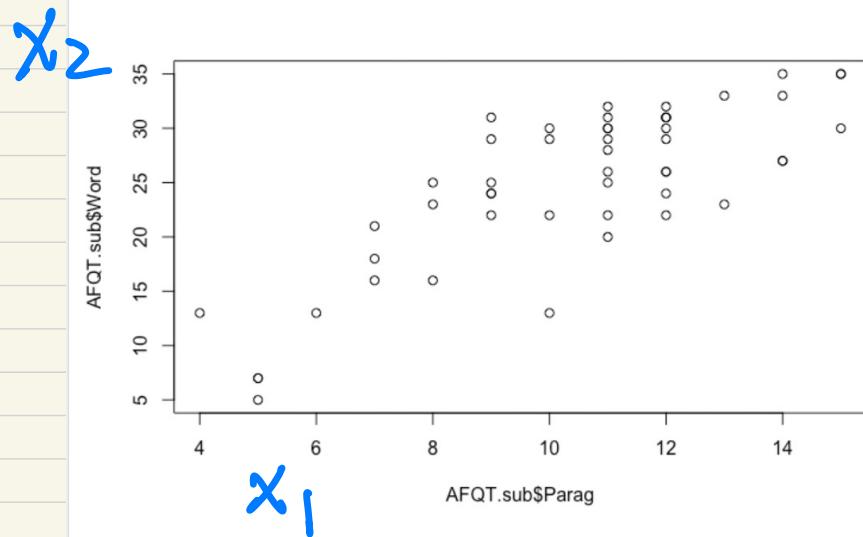
Centred & Scaled



First: Center and scale each

For simplicity consider

$n = 50$ People



Original

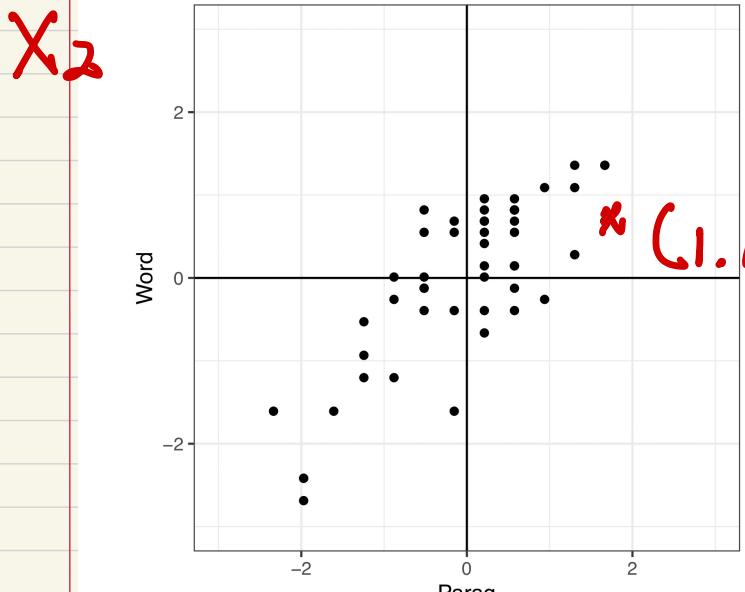
Parag : Ave = 10.4 $SD = 2.75$

Word : Ave = 24.9 $SD = 7.41$

Centered & Scaled

$$x_1 = \text{Parag-Center-Scaled} = \frac{\text{Parag} - 10.4}{2.75}$$

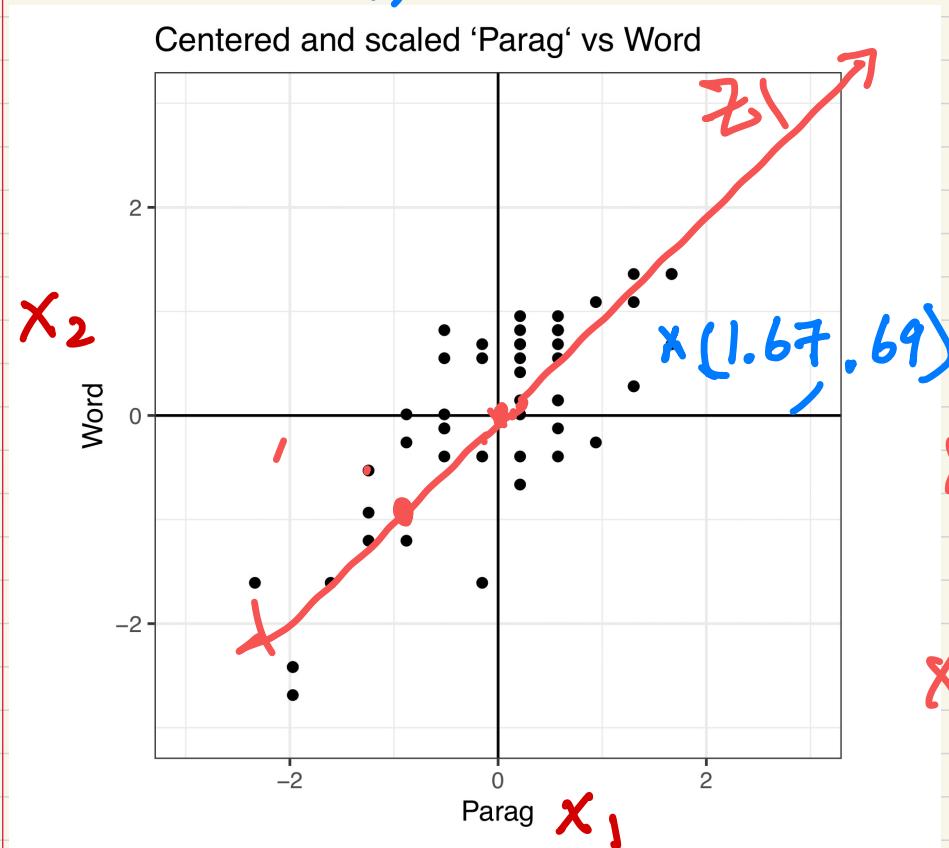
$$x_2 = \text{Word-Center-scaled} = \frac{\text{Word} - 24.9}{7.41}$$



$(1.67, .69)$

$\alpha: (1.67, .69) ?$

First: Center and scale each
 x_1, x_2



For simplicity consider

$n = 50$ People

Original

Parag : Ave = 10.4 SD = 2.75

Word : Ave = 24.9 SD = 7.41

Centered & Scaled

$$x_1 = \text{Parag-Center-Scaled} = \frac{\text{Parag} - 10.4}{2.75}$$

$$x_2 = \text{Word-Center-scaled} = \frac{\text{Word} - 24.9}{7.41}$$

a new score
 sits on the line

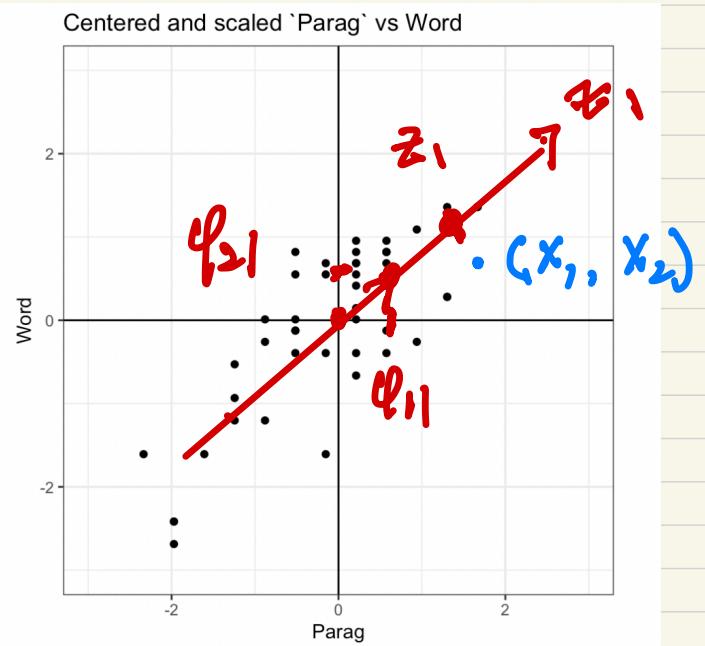
$$\text{PC1: } z_1 = \varphi_{11}x_1 + \varphi_{21}x_2$$

$(\varphi_{11}, \varphi_{21})$: Loading

$$\max_{\varphi_{11}, \varphi_{21}} \text{Var}(z_1) = \max_{\varphi_{11}, \varphi_{21}} \text{Var}(\varphi_{11}x_1 + \varphi_{21}x_2)$$

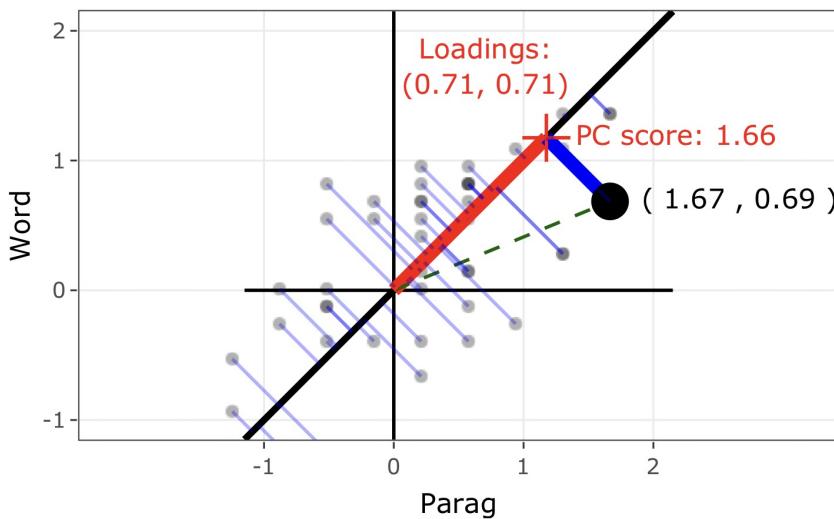
$$\varphi_{11}^2 + \varphi_{21}^2 = 1$$

$$z_1 = \varphi_{11} x_1 + \varphi_{21} x_2$$



Comments : Interpretations

Principle Component 1 of 'Parag' vs 'Word' (



In the above plot we want to demonstrate the following beautiful geometric interpretation of PCA.

Fact 1: A line which minimizes the total squared distance must go through the origin (or sample means)

Fact 2: By the Pythagorean theorem, for any point:

$$\text{PC score}^2 + \text{Perpendicular distance}^2 = \text{Distance to origin}^2$$

Summary :- PC1 : linear combination of x_1, x_2

Interpretation I: $\text{Var}(z_1)$ is the largest

Interpretation II: The total perpendicular distance² is the smallest

Solution: $\varphi_{11} = .71, \varphi_{21} = .71$

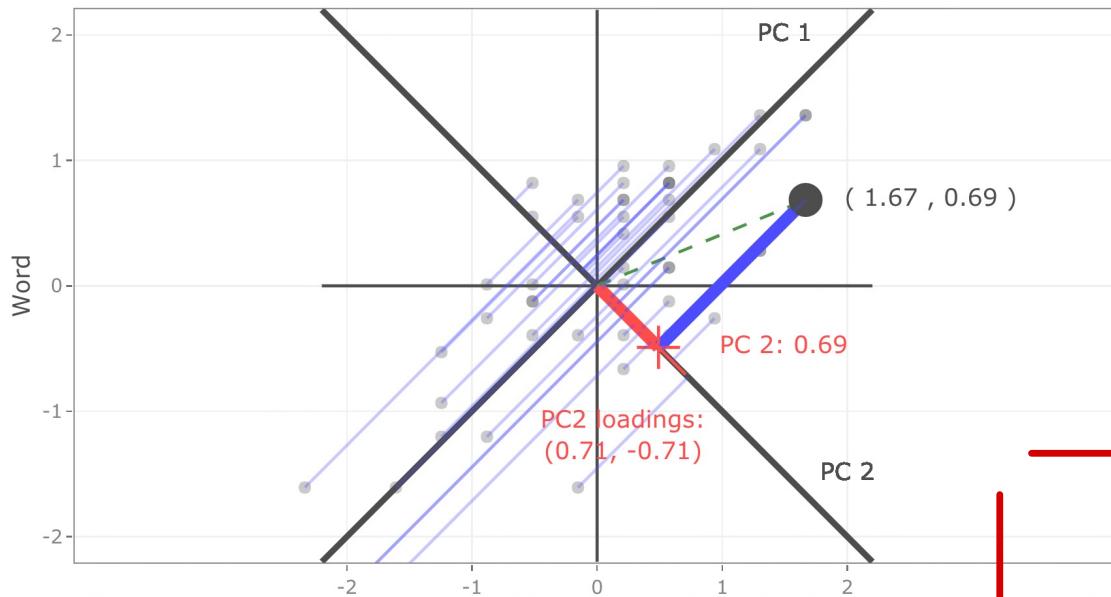
$$\text{PC1: } z_1 = .71x_1 + .71x_2$$

$$\text{ex. } (x_1, x_2) = (1.67, 0.69)$$

$$\begin{aligned}\text{PC1: } z_1 &= .71 \times 1.67 + .71 \times 0.69 \\ &= 1.66\end{aligned}$$

More PC's

Second Principle Component of 'Parag' vs 'Word' (centered, scaled)



Given PG1, look
for PC2:

$$Z_2 = \varphi_{12} X_1 + \varphi_{22} X_2$$

$$\max V_{\text{var}}(Z_2)$$

$$\varphi_{12}^2 + \varphi_{22}^2 = 1$$

$$(\varphi_{11}, \varphi_{21}) \perp (\varphi_{12}, \varphi_{22})$$

Solution Loadings

$$\text{PC2 : } \varphi_{12} = .71 \quad \varphi_{22} = -.71$$

$$\text{PC2 : } Z_2 = .71 X_1 - .71 X_2$$

$$(X_1, X_2) = (1.67, .69) \quad \text{PC2} = Z_2 = .71 \times 1.67 - .71 \times .69 \\ = .7$$

How to get PC1 ?

R-func : prcomp()

COMPONENTS PC1 and PC2 .

```
```{r}
pc.parag.word <- prcomp(parag.word.scaled.centered)
names(pc.parag.word)
pc.parag.word$rotation # loadings
#pc.parag.word$x # PC scores
data.frame(parag.word, parag.word.scaled.centered,
pc.parag.word$x)[1:10,] #parag.word: original data
````
```

| | PC1 | PC2 |
|-------|-------|--------|
| Parag | 0.707 | 0.707 |
| Word | 0.707 | -0.707 |

Description: df [10 x 6]

Centered / Scaled

| | Parag | Word | Parag.1 | Word.1 | PC1 | PC2 |
|-----|-------|-------|---------|---------|--------|---------|
| | <int> | <int> | <dbl> | <dbl> | <dbl> | <dbl> |
| p1 | 9 | 22 | -0.517 | -0.3938 | -0.644 | -0.0868 |
| p2 | 9 | 25 | -0.517 | 0.0108 | -0.358 | -0.3729 |
| p3 | 10 | 29 | -0.153 | 0.5503 | 0.281 | -0.4972 |
| p4 | 15 | 35 | 1.666 | 1.3595 | 2.140 | 0.2169 |
| p5 | 14 | 33 | 1.302 | 1.0898 | 1.692 | 0.1504 |
| p6 | 11 | 26 | 0.211 | 0.1457 | 0.252 | 0.0462 |
| p7 | 11 | 30 | 0.211 | 0.6852 | 0.634 | -0.3353 |
| p8 | 9 | 24 | -0.517 | -0.1241 | -0.453 | -0.2776 |
| p9 | 12 | 26 | 0.575 | 0.1457 | 0.509 | 0.3035 |
| p10 | 12 | 24 | 0.575 | -0.1241 | 0.319 | 0.4942 |

Recover all PCI's.

$$\text{P1: } \text{PCI} = .707 \times (-.517) + .707 \times (-.394)$$
$$= .644$$

from the original data.

| skim_variable | numeric.mean | numeric.sd |
|---------------|--------------|------------|
| <chr> | <dbl> | <dbl> |
| Word | 24.9 | 7.41 |
| Parag | 10.4 | 2.75 |
| Math | 12.2 | 5.69 |
| Arith | 17.0 | 7.19 |

$$\text{PCI} = .707 \left(\frac{9 - 10.4}{2.75} \right) + .707 \left(\frac{22 - 24.9}{7.41} \right)$$
$$= \frac{-707}{2.75} \times 9 + \frac{707}{7.41} \times 22 - c$$

Weighted Sum of Parag & Word

Summary :

Parag word
Turn $(x_1, x_2) \rightarrow (z_1, z_2)$
PC₁, PC₂

- z_1, z_2 are linear comb. of x_1 and x_2
- $\text{Var}(z_1) \geq \text{Var}(z_2) \geq \begin{cases} \text{Var } x_1 \\ \text{Var } x_2 \end{cases}$
- $z_1 \perp z_2$. $\text{Cov}(z_1, z_2) = 0$
- No information lost using z_1, z_2
- May use z_1 alone if $\text{Var}(z_2)$ is much smaller

2. AFQT : Dim 4 : $x_1, x_2, x_3, x_4 \Rightarrow z_1, z_2, z_3, z_4$

We are ready to extend to

PC's for more variables.

AFQT : x_1, x_2, x_3, x_4 .

$$\text{PC}'_3 : z_1 = \varphi_{11}x_1 + \varphi_{21}x_2 + \varphi_{31}x_3 \\ + \varphi_{41}x_4$$

$$z_4 = \varphi_{14}x_1 + \varphi_{24}x_2$$

$$+ \varphi_{34}x_3 + \varphi_{44}x_4$$

R-function:

`prcomp()`

3.

Go to lecture

- 1) What does each PC score mean?
- 2) How many PCs should be used?
- 3) What interesting facts can be revealed by PCs?

I) Interpretations of PC's

a)

```
skim(data.AFQT) %>% select(skim_variable, numeric.mean, numeric.sd)
```

| skim_variable | numeric.mean | numeric.sd |
|---------------|--------------|------------|
| Word | 24.9 | 7.41 |
| Parag | 10.4 | 2.75 |
| Math | 12.2 | 5.69 |
| Arith | 17.0 | 7.19 |

default
↓ = FALSE

b) Prcomp(data.AFQT, scale = TRUE,

center = TRUE)
↑ default

rotation : loadings
x : PC scores.

```
```{r loading}
pc.4.loading <- pc.4$rotation #pc.4$x
knitr::kable(pc.4.loading)
````
```

| | PC1 | PC2 | PC3 | PC4 |
|-------|-------|--------|--------|--------|
| Word | 0.509 | -0.442 | 0.365 | -0.642 |
| Parag | 0.500 | -0.546 | -0.308 | 0.598 |
| Math | 0.496 | 0.483 | -0.652 | -0.310 |
| Arith | 0.494 | 0.523 | 0.589 | 0.368 |

All Centred, Scaled

PC1 = .5 Word + .5 Parag + .5 Math + .5 Arith

= .5(Word + Parag + Math + Arith)

= .5 Sum (4 tests) !!!

PC1: Weighted sum of 4 tests.

$$\begin{aligned}
 PC2 &\approx .5 (\text{Math} + \text{Arith} - (\text{Word} + \text{Parag})) \\
 &= .5 (\text{Sum}(\text{Math}) - \text{Sum}(\text{English})) \\
 &\quad !! !
 \end{aligned}$$

Remark 1: In the original scores

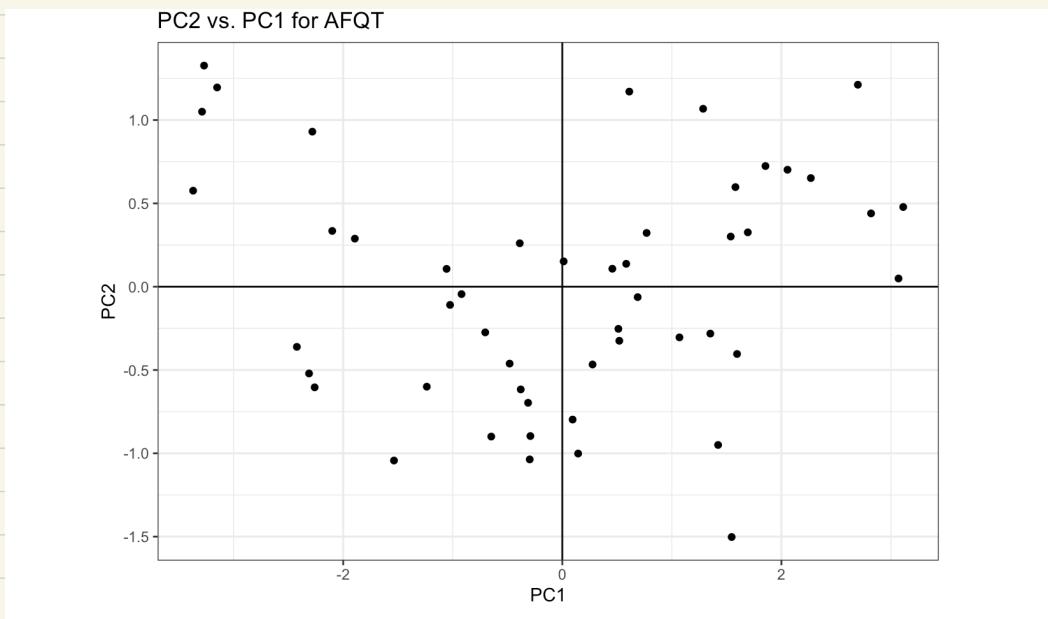
$$\begin{aligned}
 PC1 &= .5 \left(\left(\frac{\text{Math}}{5.69} + \frac{\text{Arith}}{7.19} \right. \right. \\
 &\quad \left. \left. - \left(\frac{\text{Word}}{7.41} + \frac{\text{Parag}}{2.75} \right) \right) + C_1 \right. \\
 &\quad \left. \uparrow \right)
 \end{aligned}$$

inverse prop to each sd !

Known
constant

Remark 2: Loadings are unique upto signs.

c) Plot (PC1, PC2) : informative



2) How many PC's should be used?

a) PVĒ : Proportion of Variance Explained

Var (PC) / Total Variances

```
```{r}
summary(pc.4)$importance #notice it is from summary()
```


	PC1	PC2	PC3	PC4
Standard deviation	1.750	0.699	0.5208	0.4238
Proportion of Variance	0.765	0.122	0.0678	0.0449
Cumulative Proportion	0.765	0.887	0.9551	1.0000

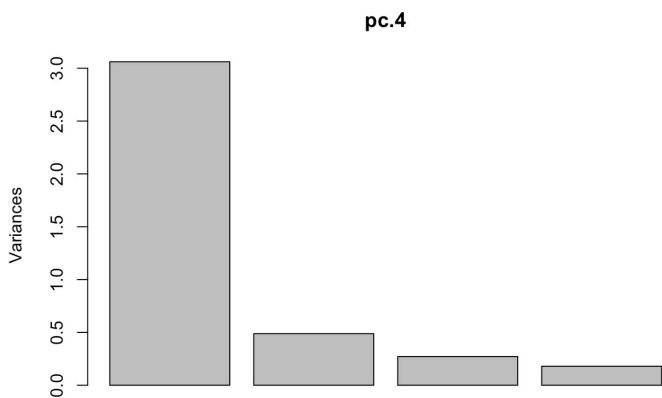

```

→ 3.06, .49, .27, .18

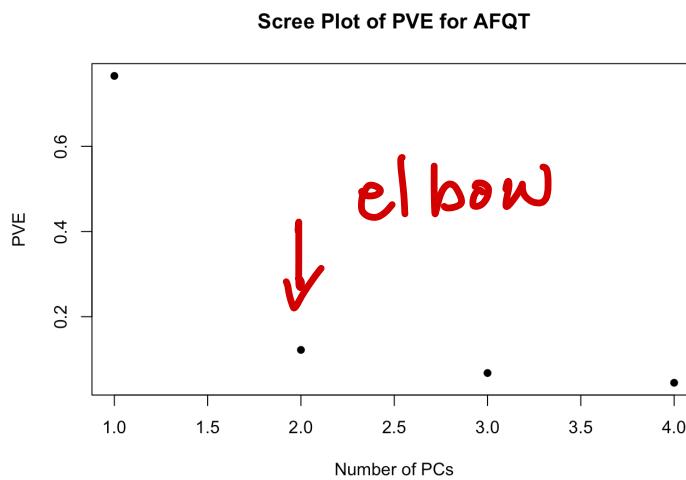
b) plot (pc.4)

Scree plots

```
plot(pc.4) # variances of each pc
```



```
plot(summary(pc.4)$importance[2, ], # PVE  
y lab="PVE",  
x lab="Number of PCs",  
pch = 16,  
main="Scree Plot of PVE for AFQT")
```

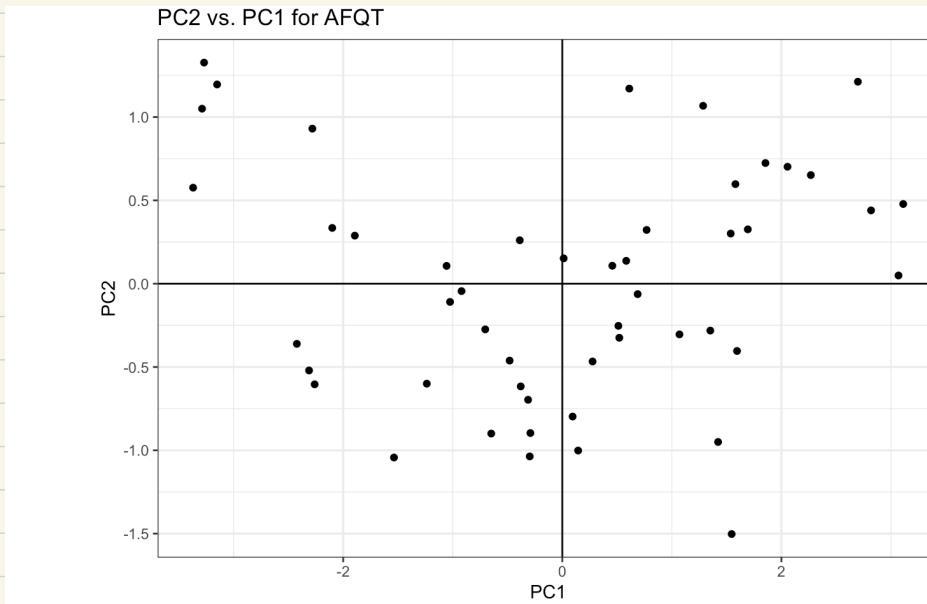


Elbow rule:

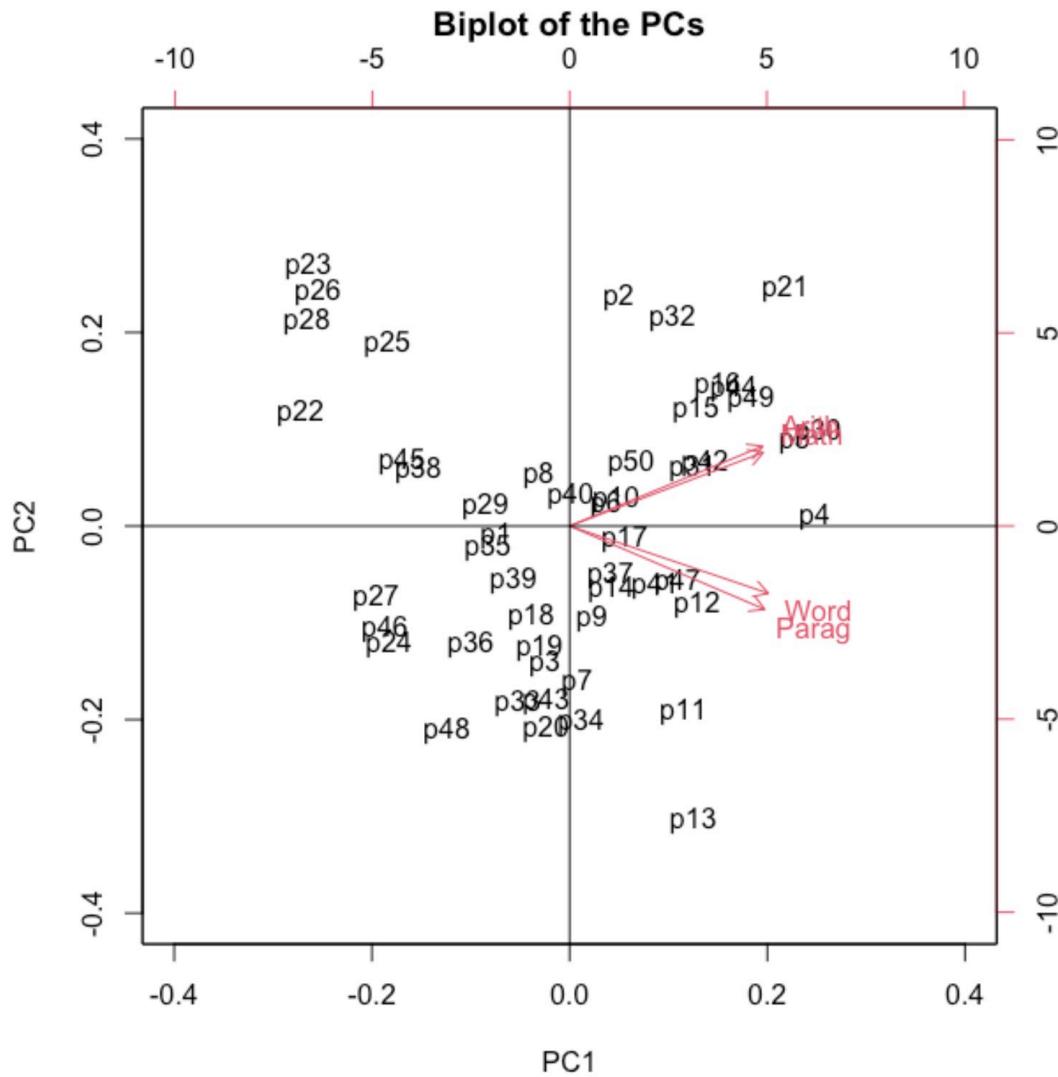
Point : The curve levels off afterwards

So we take two PC's

c) Plot (PC1, PC2) : see the variabilities of PC1, vs. PC2

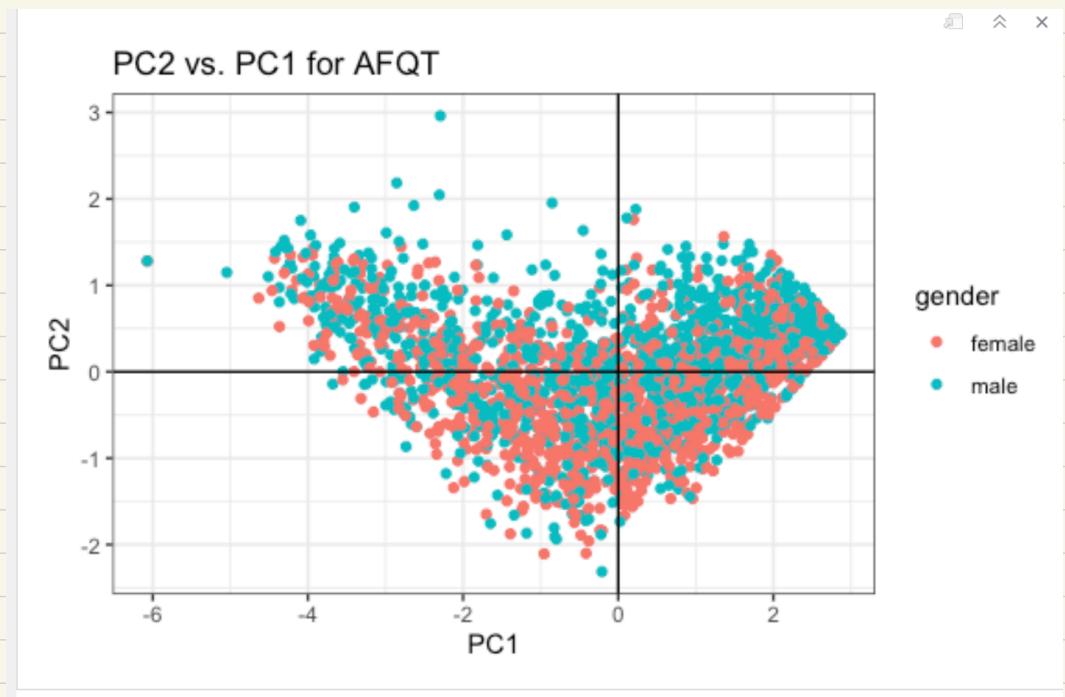


d) Biplot



3) a) Interesting findings?

Differences between males & females ??



b) How is AFQT created ??

4. a) How to find the loadings $\varphi_1 = \begin{pmatrix} \varphi_{11} \\ \varphi_{21} \end{pmatrix}$?

Facts: Let $\tilde{X} = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{50,1} & x_{50,2} \end{pmatrix}_{n \times 2}$

Then φ_1 is the leading eigen-vector of

$$\frac{(\tilde{X}^T \cdot \tilde{X})}{n-1}$$

Cov. matrix

i.e. $\left(\frac{\tilde{X}^T \cdot \tilde{X}}{n-1} \right) \varphi_1 = \lambda_1 \cdot \varphi_1$

↑
largest
eigenvalue

eigenvector

Prof:

$$\max_{\|\varphi_i\|^2=1} \text{Var}(z_i) = \varphi_i^T \cdot \frac{\bar{X}^T \cdot \bar{X}}{n-1} \cdot \varphi_i$$

$$\Rightarrow \frac{\bar{X}^T \cdot \bar{X}}{n-1} \cdot \varphi_i = \lambda_i \varphi_i$$



• λ_1 : largest eigen-value

$$\boxed{\text{Var}(z_1) = \lambda_1}$$

b) PCA secret: SVD (Singular Value Decomposition)

Take AFQT: $X_1 = \text{Wrd}$, $X_2 = \text{Pagr}$, $X_3 = \text{Math}$
 $X_4 = \text{Arith}$

Centered, Scaled

The data matrix

$$\underline{X} = \begin{pmatrix} X_1 & X_2 & X_3 & X_4 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots \end{pmatrix}$$

$n=50, p=4$

Fact 1: SVD: There exist v_1, \dots, v_4
 $u_1, u_2, u_3, u_4 : n \times 1 = 50 \times 1$, $d_1 > d_2 > d_3 > d_4$

$$\underline{X} = d_1 \cdot u_1 \cdot v_1^T + d_2 u_2 v_2^T + d_3 u_3 v_3^T + d_4 u_4 v_4^T$$

- u_1, \dots, u_4 Orthonormal
- v_1, \dots, v_4 :

Fact 2: PC's \Leftrightarrow SVD

- (v_1, v_2, v_3, v_4) : PC loadings.

- $PC_1 = d_1 u_1 \leftarrow \underline{\underline{X}} \cdot \underline{v}_1$
- $PC_2 = d_2 u_2 \leftarrow \underline{\underline{X}} \cdot \underline{v}_2$
- \vdots
- $PC_4 = d_4 u_4 \leftarrow \underline{\underline{X}} \cdot \underline{v}_4$

- $\lambda_1 = d_1^2 / n-1 = \text{Var}(PC_1)$
- \vdots
- $\lambda_4 = d_4^2 / n-1 = \text{Var}(PC_4)$

$\ddot{\cup}$

5. Recap PCA

Given a data with variables

x_1, x_2, \dots say $x_{p=10}$

- PC's are linear Comb. of $x_1, \dots x_{10}$

$$PC_1 = \varphi_{11}x_1 + \varphi_{21}x_2 + \dots + \varphi_{10,1}x_{10}$$

Properties

1. - $(\varphi_{11}, \varphi_{21}, \dots, \varphi_{10,1})$: loadings
 - $\text{Var}(PC_1)$ maximized
 - loadings / PC_1 unique up to the sign

- ϕ_{11} prop $\text{Cor}(PC_1, X_1)$
⋮

2. $PC_2, PC_3 \dots PC_{10}$ max #.

$$\text{Var}(PC_1) \geq \text{Var}(PC_2) \geq \dots \geq \text{Var}(PC_{10})$$

3. - All loadings are uncorr.

• All PC_i 's ?

4. How many PC_i 's ?

$$\text{PVE} = \frac{\text{Var}(PC_1)}{\text{Var}(PC_1) + \dots + \text{Var}(PC_{10})}$$

5. How to get PC loadings / PC's

- Prcomp()  Plot
biplot
- svd() : v loadings, $PC = d \cdot u$
- irlbac() much faster

more to come in
module - clustering