

Dimension Reduction (PCA)

07-10-2023

What is the goal of dimension reduction?

We have p variables (columns) for n observations (rows) BUT which variables are **interesting**?

- **REFINED QUESTION:** can we “project” data to a lower dimension but keep maximal information?
- **SHARPER QUESTION:** is there another “basis”, which better expresses the information in our original data set?

Linear Algebra Interlude

- 2D vectors - magnitude (‘norm’) and direction (‘dot product’)
- Projection: length of the shadow of the given vector over another vector

$$Proj_w(v) = (v^T w^*) w^*$$

where:

$$w^* = \frac{w}{||w||}$$

and

$v^T w^*$ = degree of information preserved about v after projecting onto w

MATRICES can be thought of as:

- data
- functions (linear transformations)

EIGENVALUES AND EIGENVECTORS:

$$Au = \lambda u$$

λ = eigenvalue

u = eigenvector

Importance:

- Eigenvectors basically stay invariant to rotation after being acted on by A – “holding ground after being acted on by A ”

PCA

GOAL: can we find p new directions that preserves:

- linearity
- maximizes variance explained
- are orthogonal

Let u be the vector that preserves the most information from the data:

$$\max \sum_{i=1}^p (x_i^T u)^2$$

s.t. $u^T u = 1$, or equivalently: $u^T u - 1 = 0$

Then to find the other principal components:

$$\max \sum_{i=1}^p (x_i^T u)^2$$

s.t. $u_2^T u_2 = 1$ AND $u_1 \perp u_2$

- PCA explores the covariance between variables and combines variables into a smaller set of uncorrelated variables called principal components (PCs)
- The first principal component is the linear combination of the p variables that has the **largest variance**. The amount of variability captured goes in descending order.

Singular Value Decomposition (SVD)

X is the covariance matrix

$$X = UDV^T$$

* Matrices U and V contain the left and right singular vectors of scaled matrix X

- D is the diagonal matrix of the singular values
- SVD simplifies matrix-vector multiplication as rotate, scale, and rotate again
- V is called the loading matrix

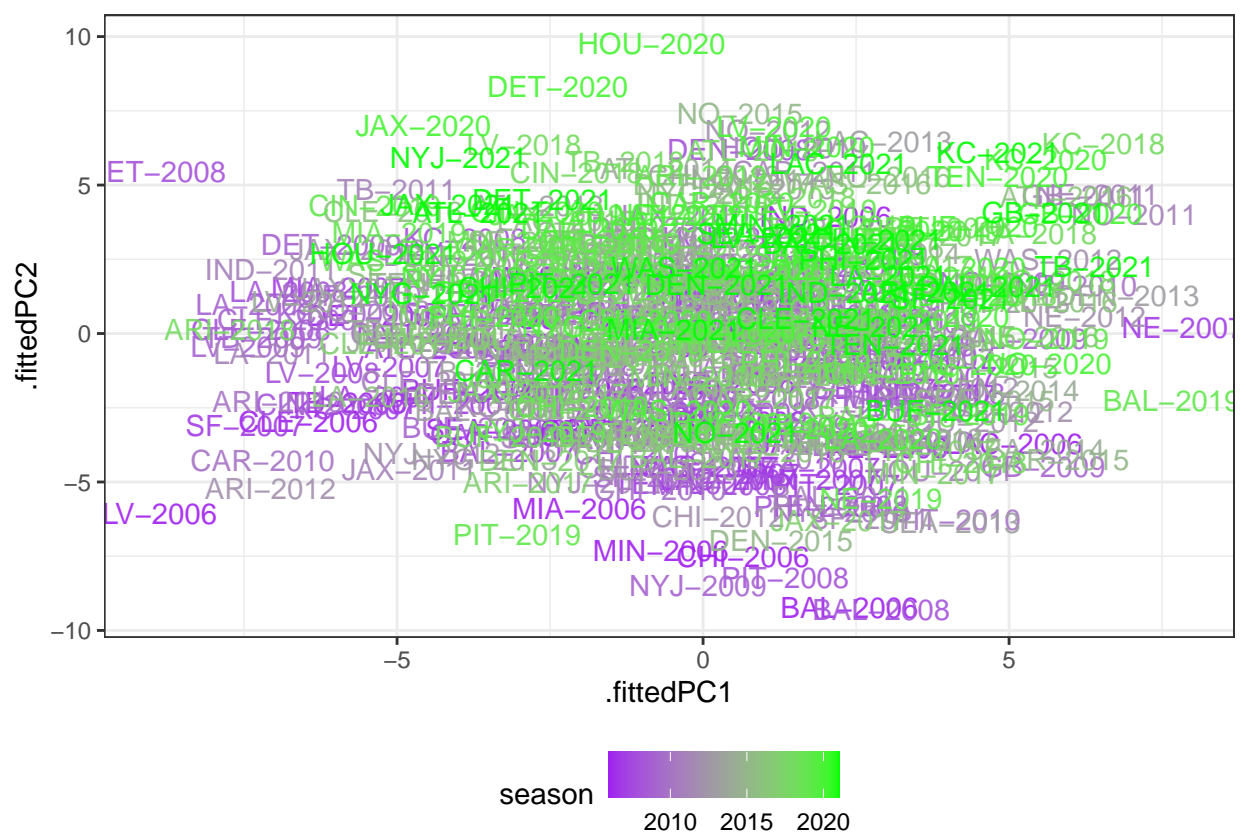
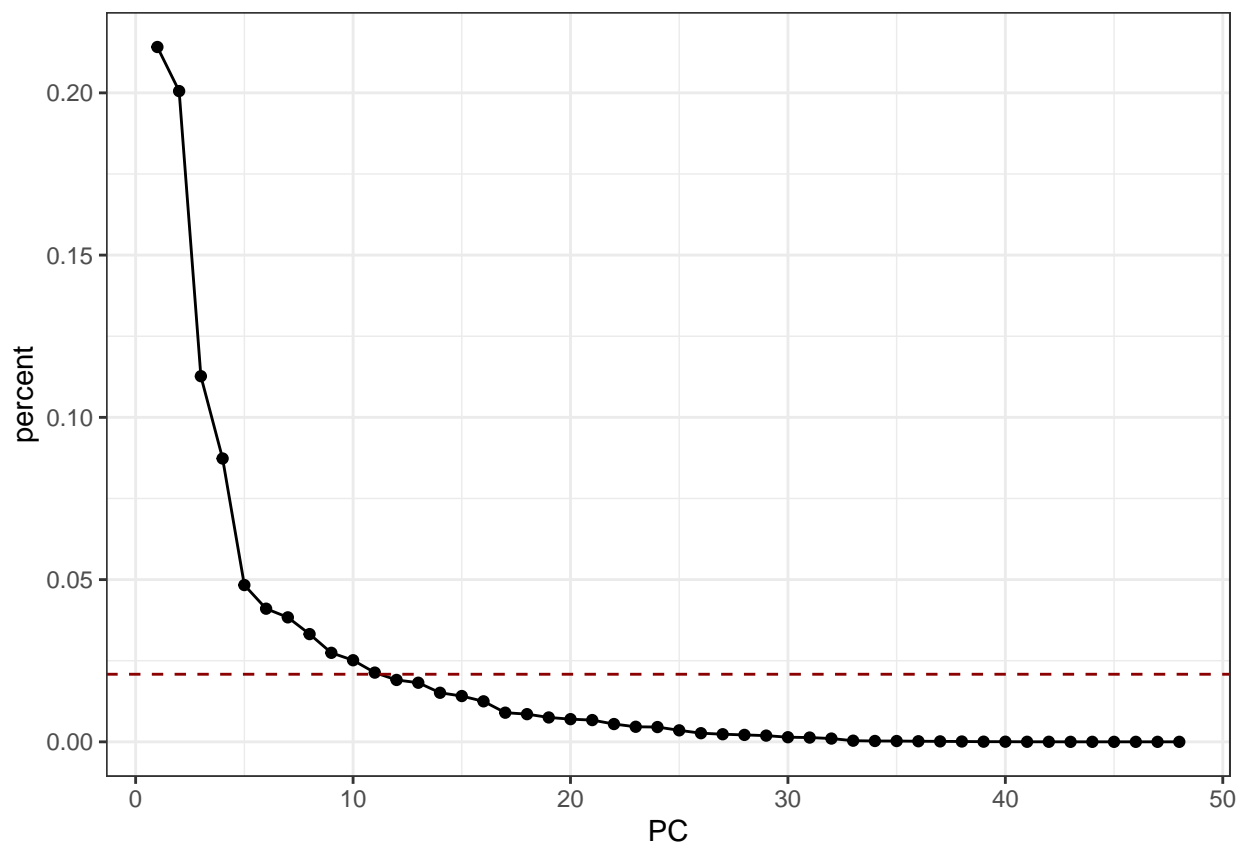
$Z = XV$ is the PC matrix

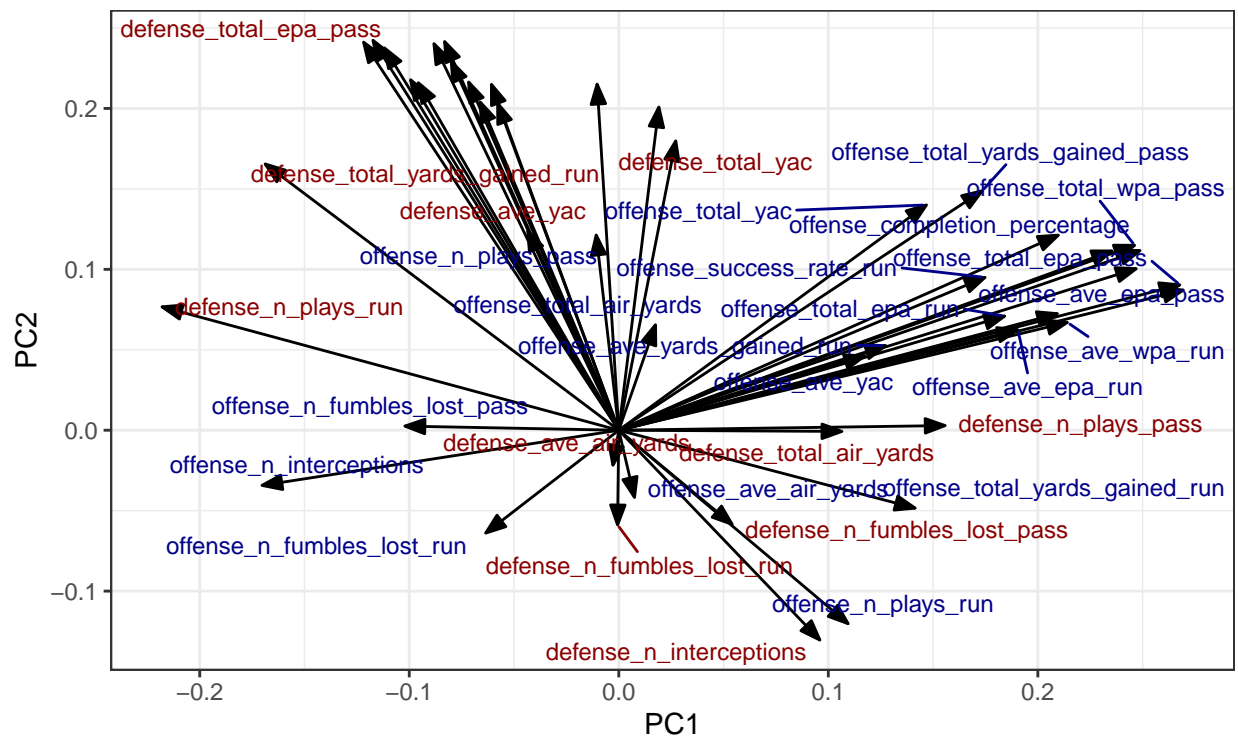
Eigenvalue Decomposition

- V are **eigenvectors** of $X^T X$
- U are the **eigenvectors** of XX^T
- The singular values (diagonal of D) are square roots of the **eigenvalues** of $X^T X$ or XX^T
- Meaning that $Z = UD$

Example

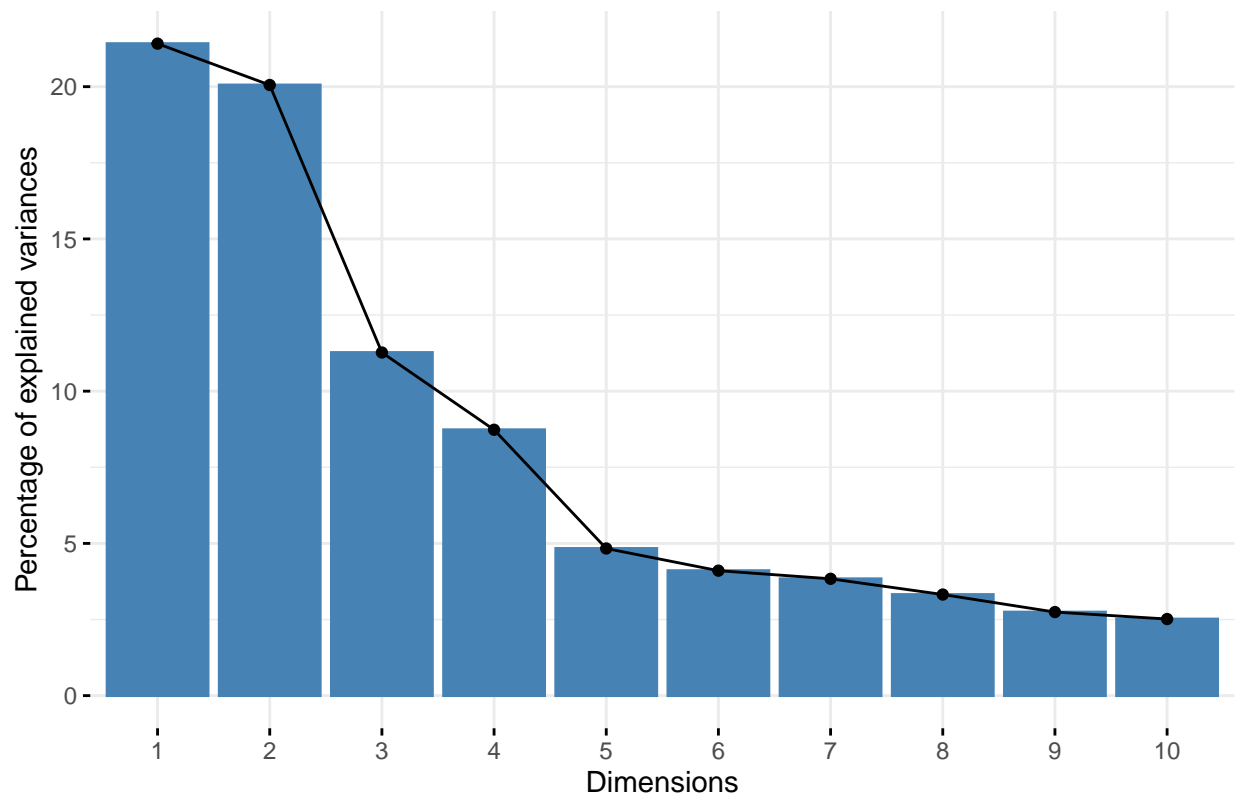
```
## Importance of components:
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  3.2060 3.1026 2.3257 2.04728 1.52301 1.40350 1.35714
## Proportion of Variance 0.2141 0.2006 0.1127 0.08732 0.04832 0.04104 0.03837
## Cumulative Proportion 0.2141 0.4147 0.5274 0.61468 0.66301 0.70405 0.74242
##          PC8      PC9      PC10     PC11     PC12     PC13     PC14
## Standard deviation  1.26250 1.14773 1.09881 1.01200 0.95689 0.93513 0.85233
## Proportion of Variance 0.03321 0.02744 0.02515 0.02134 0.01908 0.01822 0.01513
## Cumulative Proportion 0.77562 0.80307 0.82822 0.84956 0.86863 0.88685 0.90199
##          PC15     PC16     PC17     PC18     PC19     PC20     PC21
## Standard deviation  0.82315 0.77434 0.65692 0.64016 0.60076 0.5796 0.56756
## Proportion of Variance 0.01412 0.01249 0.00899 0.00854 0.00752 0.0070 0.00671
## Cumulative Proportion 0.91610 0.92859 0.93758 0.94612 0.95364 0.9606 0.96735
##          PC22     PC23     PC24     PC25     PC26     PC27     PC28
## Standard deviation  0.51349 0.47233 0.46768 0.41284 0.35810 0.33597 0.32018
## Proportion of Variance 0.00549 0.00465 0.00456 0.00355 0.00267 0.00235 0.00214
## Cumulative Proportion 0.97284 0.97749 0.98205 0.98560 0.98827 0.99062 0.99276
##          PC29     PC30     PC31     PC32     PC33     PC34     PC35
## Standard deviation  0.30357 0.26161 0.25289 0.22149 0.13146 0.11459 0.10964
## Proportion of Variance 0.00192 0.00143 0.00133 0.00102 0.00036 0.00027 0.00025
## Cumulative Proportion 0.99468 0.99610 0.99744 0.99846 0.99882 0.99909 0.99934
##          PC36     PC37     PC38     PC39     PC40     PC41     PC42
## Standard deviation  0.09672 0.08397 0.07385 0.05223 0.04814 0.03391 0.02901
## Proportion of Variance 0.00019 0.00015 0.00011 0.00006 0.00005 0.00002 0.00002
## Cumulative Proportion 0.99954 0.99968 0.99980 0.99985 0.99990 0.99993 0.99994
##          PC43     PC44     PC45     PC46     PC47     PC48
## Standard deviation  0.02562 0.02290 0.02213 0.02139 0.01718 0.01670
## Proportion of Variance 0.00001 0.00001 0.00001 0.00001 0.00001 0.00001
## Cumulative Proportion 0.99996 0.99997 0.99998 0.99999 0.99999 1.00000
```





stat_type a defense a offense

Scree plot



A PCA plot showing the first two dimensions of variation (Dim1: 21.4%, Dim2: 20.1%) for 400 samples. The plot is divided into four quadrants by dashed lines at Dim1 = 0 and Dim2 = 0. The top-left quadrant (Dim1 < 0, Dim2 > 0) contains samples labeled 'F1' and 'F2'. The top-right quadrant (Dim1 > 0, Dim2 > 0) contains samples labeled 'F3' and 'F4'. The bottom-left quadrant (Dim1 < 0, Dim2 < 0) contains samples labeled 'M1' and 'M2'. The bottom-right quadrant (Dim1 > 0, Dim2 < 0) contains samples labeled 'M3' and 'M4'. The samples are numbered 1 through 400, with some numbers appearing in multiple locations, indicating they are part of different groups. The plot shows a clear separation between the four groups, with 'F1' and 'F2' clustered together, 'F3' and 'F4' clustered together, 'M1' and 'M2' clustered together, and 'M3' and 'M4' clustered together. The 'F' groups are generally located in the upper half of the plot, while the 'M' groups are generally located in the lower half. The 'F1' and 'F2' groups are located on the left side of the plot, while the 'F3' and 'F4' groups are located on the right side. The 'M1' and 'M2' groups are located on the left side of the plot, while the 'M3' and 'M4' groups are located on the right side.

Biplot displays both the space of observations and the space of variables

