

## Introduction

Motivation: Data Science has potential in inspecting sports data and helping sports teams make better-informed decisions.

Aim: To find structure in Purdue Women's Soccer training data.  
22 players, 48 drills, 9 features.

Approach: Unsupervised learning problem: clustering, while extracting most relevant features

## Methods

### Data Preparation

Extracted players: 22 observations, 9 features.

Extracted drills: 48 observations, 9 features.

Features scaled according to median using formula:

$$\frac{\text{value} * 100}{\text{median}(\text{feature})}$$

Players before scaling									Players after scaling																
1181201	68.1037	5.99838	102.6986	5.849373	13.54287	9.19158	106.4908	118.9789	118.8767	58.79588	105.8014	48.45585	115.1288	121.1431	106.8816	123.8118	112.4286	118.1841	88.18287	89.8718	72.89134	81.04152	120.9495	88.1141	105.2434

### Feature Reduction

Feature selection:

Correlation matrix using Pearson's method:

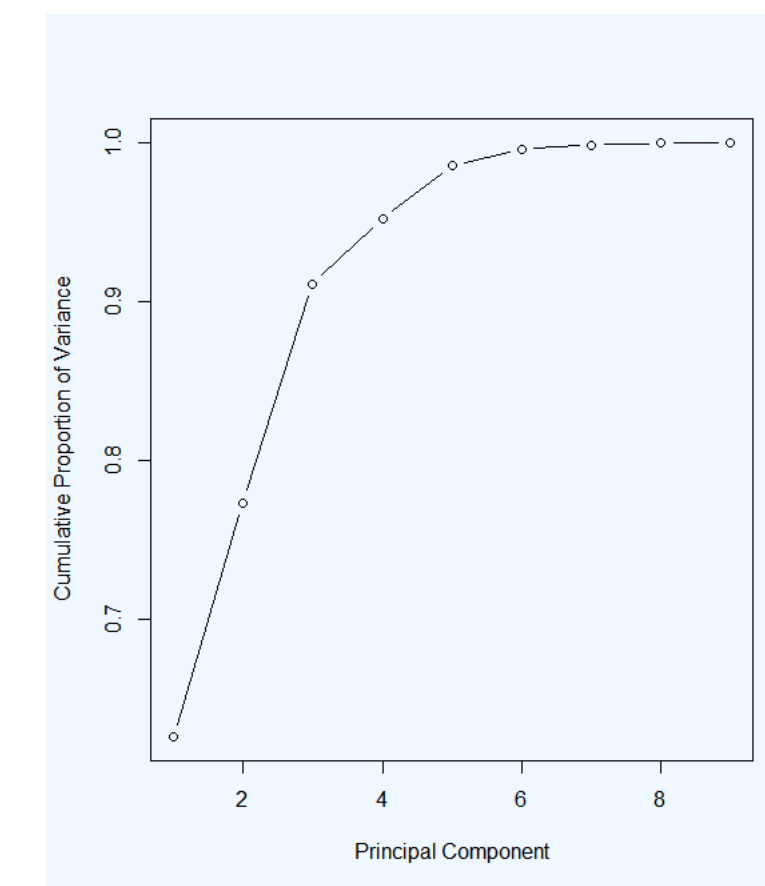
$$\frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{(n\Sigma x^2 - (\Sigma x)^2)(n\Sigma y^2 - (\Sigma y)^2)}}$$

	Distance Total	Distance Per Min	High Speed Running	HML Dis-tance	Sprints	Accelerations	Decelerations	HML Dis-tance	Average Heart Rate
Distance Total	1	0.8186695	0.4207638	0.8683186	0.4559951	0.4621517	0.6832496	0.6167345	0.2669202

### Feature transformation:

Principal components analysis: Find mutually orthogonal axes with maximum variance.  
First 6 principal components account for 99.6% of variance.

Principal Component #	Proportion of variance
1	62.6414207
2	14.71655923
3	13.77081087
4	4.064271592
5	3.373364063
6	1.051230178
7	0.24267435
8	0.101564035
9	0.038104976



### Finding similarities in results

K! possibilities of matching k clusters from two clustering results.

Clusters from result #1: {A,B,C,D}  
Clusters from result #2: {a,b,c,d}

Hamming distance used to find best possibility.

A	a	B	b	C	c	D	d
A	a	B	b	C	c	D	c
A	a	B	c	C	b	D	d
A	a	B	c	C	d	D	d
A	a	B	d	C	b	D	c
A	a	B	d	C	c	D	b
A	b	B	a	C	c	D	c
A	b	B	a	C	d	D	c
A	b	B	c	C	a	D	d
A	b	B	c	C	d	D	a
A	b	B	d	C	a	D	c
A	b	B	d	C	c	D	a
A	c	B	a	C	b	D	d
A	c	B	a	C	d	D	b
A	c	B	b	C	a	D	d
A	c	B	b	C	d	D	a
A	c	B	d	C	a	D	b
A	c	B	d	C	b	D	a
A	d	B	a	C	c	D	b
A	d	B	a	C	d	D	a
A	d	B	b	C	a	D	c
A	d	B	b	C	c	D	a
A	d	B	c	C	a	D	b
A	d	B	c	C	b	D	a

### Clustering

#### K-means:

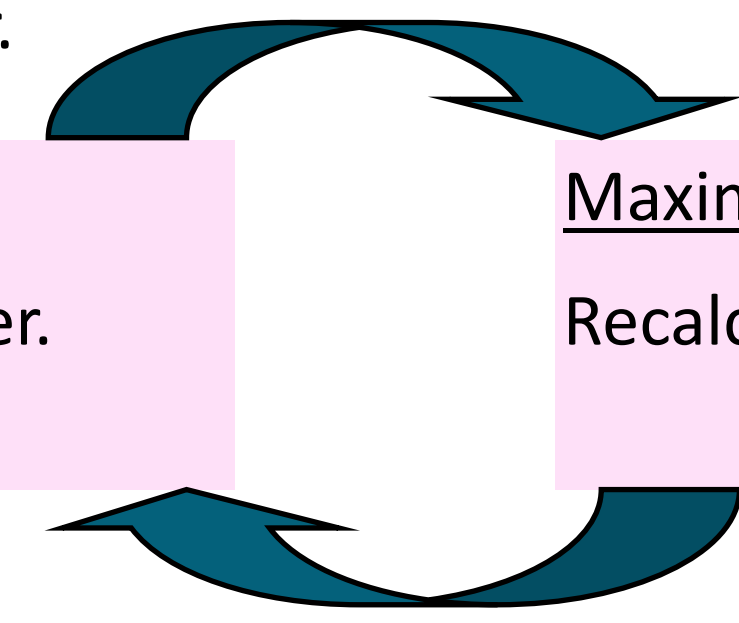
Initialization: Each point is a center.

#### Expectation

Assign each point to closest center.

#### Maximization

Recalculate centers as means of clusters.



Terminating Condition: Centers converge.

Challenge: Random initialization of centers yielded different results at each execution.

Solution: Use centers for which total within-cluster sum-of-squares is minimum.

#### Spectral Clustering:

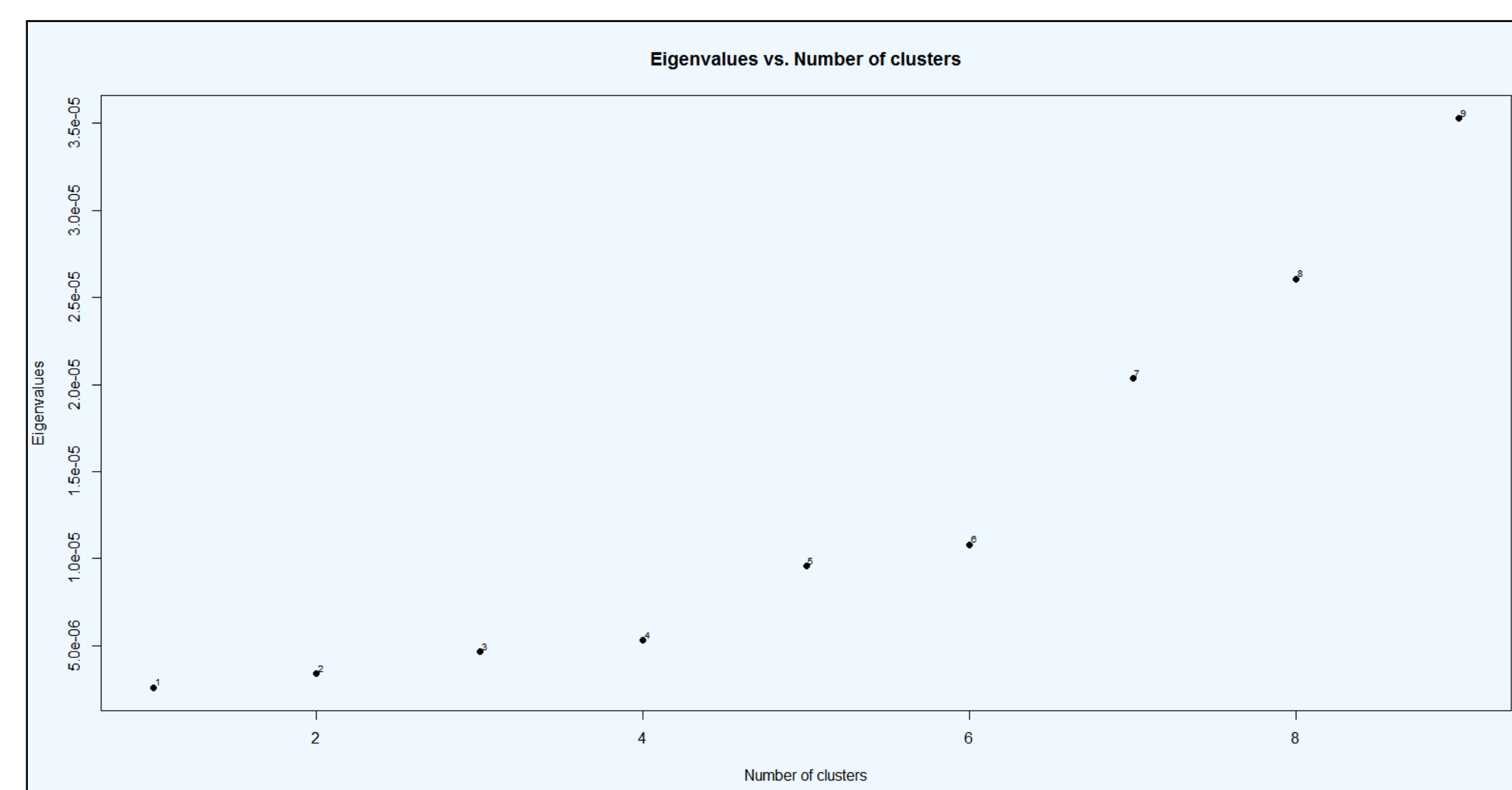
Gaussian Kernel Similarity Function

$$W_{ij} = e^{-\frac{|x_i - x_j|^2}{2\sigma^2}}$$

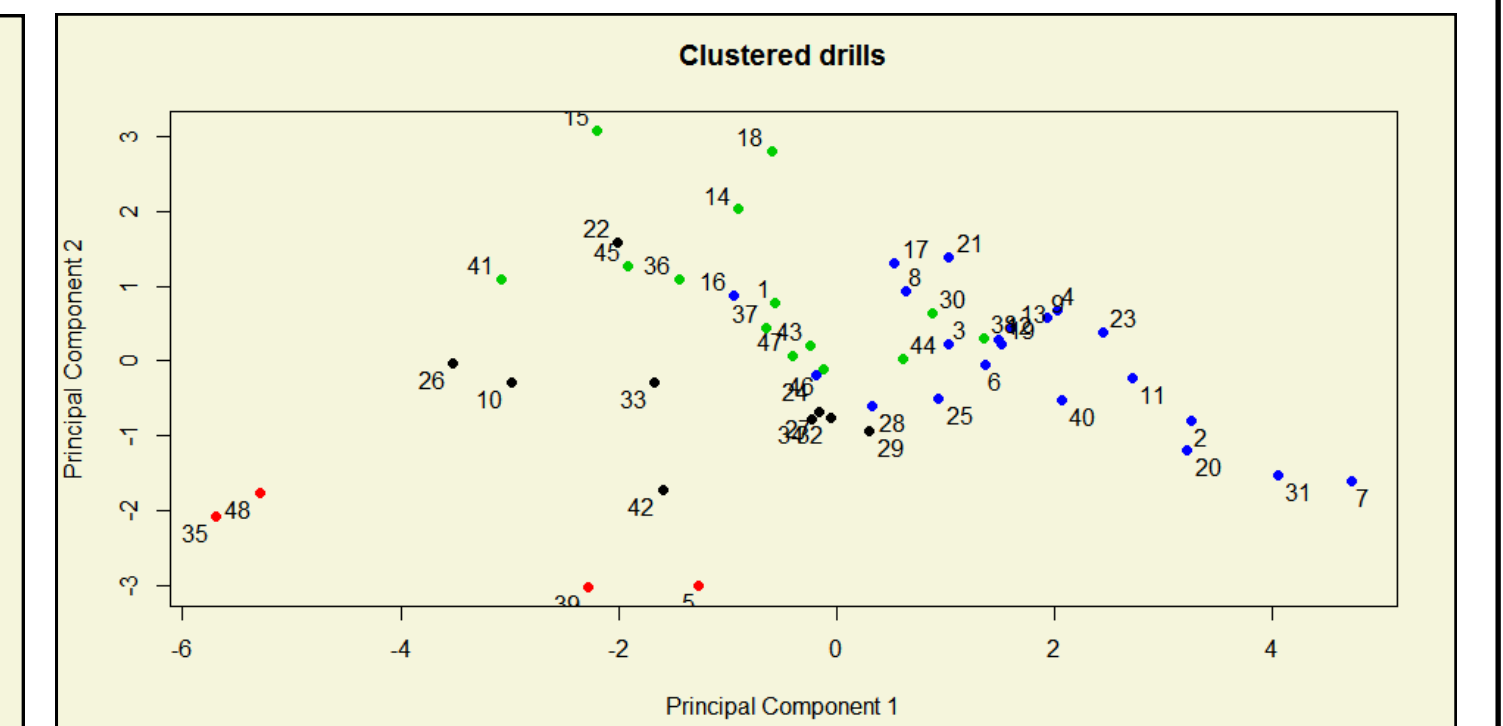
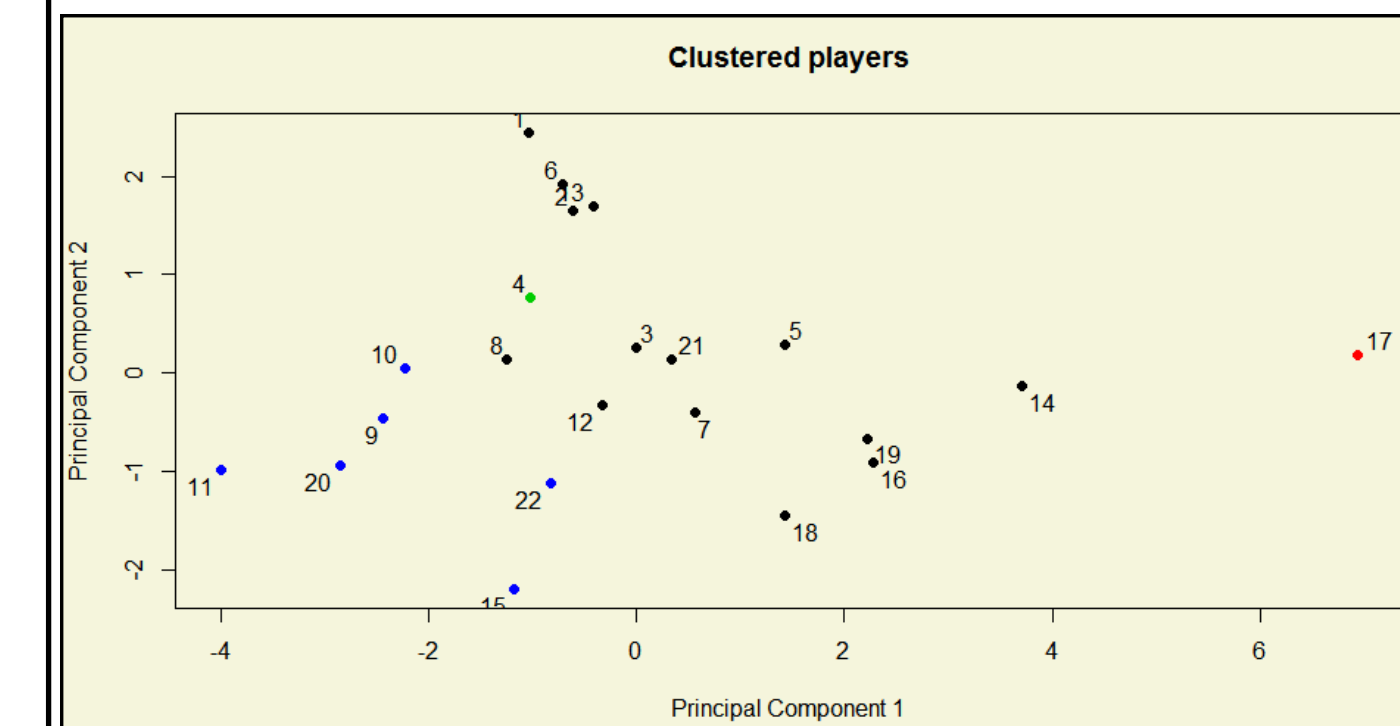
Largest k eigenvalues of normalized laplacian of affinity matrix are subjected to k-means.

Eigengap method: Optimum k is given by index of first encounter of big jump in eigenvalues.

The appropriate sigma was found for which k=4.



## Result



Combined results from k-means and spectral clustering of both the six selected features as well as the first six principal components.

## Conclusion and Future Work

Conclusion: Fairly consistent clusters were discovered within both drills and players.

#### Next steps:

Short term: Explore ways to make affinity matrix for spectral clustering more accurate.

Long term: Perform semi-supervised learning with future data.

## References

[1] Ng, A. Y., Jordan, M. I., & Weiss, Y. (2002). On spectral clustering: Analysis and an algorithm. *Advances in Neural Information Processing Systems, Advances in Neural Information Processing Systems, 2002*.

[2] Shlens, J. (2014). A Tutorial on Principal Component Analysis.

[3] R Core Team (2016). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

## Acknowledgement

Purdue SURF Program: Sponsorship  
Jampani Dwaraknath Reddy: Fruitful discussions