

Clustering

* taken from Ron Yarba's slides

- Q. Cluster NBA players into offensive Roles or Archetypes based on how they play.

- * First, need variables that describe how they play.
- * Todd Whitehead: try using Synergy Play Type data!

```
[1] "FGA_freq_Spotup"           "FGA_freq_PRRollMan"          "FGA_freq_Postup"  
[4] "FGA_freq_Cut"              "FGA_freq_OffRebound"         "FGA_freq_PRBallHandler"  
[7] "FGA_freq_Isolation"        "FGA_freq_OffScreenOrHandoff"
```

- * Then, with this data, how to cluster players??

What is **unsupervised learning**?

We have p variables for n observations x_1, \dots, x_n , and for observation i :

$$x_{i1}, x_{i2}, \dots, x_{ip}$$

- *unsupervised*: none of the variables are **response** variables, i.e., there are no labeled data

Think of unsupervised learning as **an extension of EDA...**

- ⇒ **there is no unique right answer!**

What is clustering (aka cluster analysis)?

ISLR 10.3:

very broad set of techniques for finding subgroups, or clusters, in a dataset

- observations **within** clusters are **more similar** to each other,
- observations **in different** clusters are **more different** from each other

How do we define **distance / dissimilarity** between observations?

- e.g. **Euclidean distance** between observations i and j

$$d(x_i, x_j) = \sqrt{(x_{i1} - x_{j1})^2 + \dots + (x_{ip} - x_{jp})^2}$$

Units matter!

- one variable may *dominate* others when computing Euclidean distance because its range is much larger
- can standardize each variable / column of dataset to have mean 0 and standard deviation 1 with **scale()**
- **but we may value the separation in that variable!** (so just be careful...)

What's the clustering objective?

- C_1, \dots, C_K are sets containing indices of observations in each of the K clusters
 - if observation i is in cluster k , then $i \in C_k$
- We want to minimize the **within-cluster variation** $W(C_k)$ for each cluster C_k and solve:

$$\underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K W(C_k) \right\}$$

- Can define using the **squared Euclidean distance** ($|C_k| = n_k = \# \text{ observations in cluster } k$)
$$W(C_k) = \frac{1}{|C_k|} \sum_{i,j \in C_k} d(x_i, x_j)^2$$
- Commonly referred to as the within-cluster sum of squares (WSS)

Lloyd's algorithm

- 1) Choose K random centers, aka **centroids**
- 2) Assign each observation closest center (using Euclidean distance)
- 3) Repeat until cluster assignment stop changing:
 - Compute new centroids as the averages of the updated groups
 - Reassign each observations to closest center

Converges to a local optimum, not the global

Results will change from run to run (set the seed!)

Takes K as an input!

things to check when clustering

- do the units of the variables make sense?
- Should you standardize the variables?
- is "nstart", the number of starting random configurations, large enough?

So, how do we choose the number of clusters?!



There is no universally accepted way to conclude that a particular choice of K is optimal!

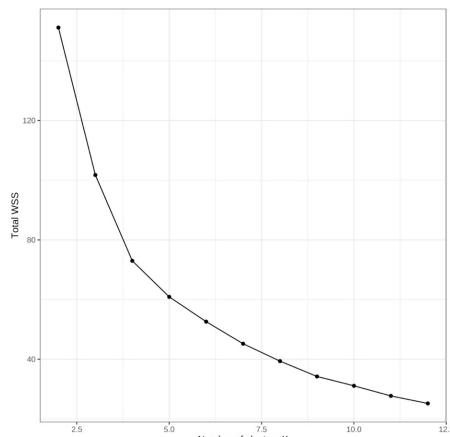
Popular heuristic: elbow plot (use with caution)

Choose K where marginal improvements is low at the bend (hence the elbow)

This is just a guideline and should not dictate your choice of K !

Gap statistic is a popular choice (see `clusGap` function in `cluster` package)

Next Tuesday: model-based approach to choosing the number of clusters!



Better alternative to `nstart`: K-means++

Pick a random observation to be the center c_1 of the first cluster C_1

- This initializes a set $\text{Centers} = \{c_1\}$

Then for each remaining cluster $c^* \in 2, \dots, K$:

- For each observation (that is not a center), compute $D(x_i) = \min_{c \in \text{Centers}} d(x_i, c)$
 - Distance between observation and its closest center $c \in \text{Centers}$
- Randomly pick a point x_i with probability: $p_i = \frac{D^2(x_i)}{\sum_{j=1}^n D^2(x_j)}$
- As distance to closest center increases \Rightarrow probability of selection increases
- Call this randomly selected observation c^* , update $\text{Centers} = \text{Centers} \cup c^*$
 - Same as `centers = c(centers, c_new)`

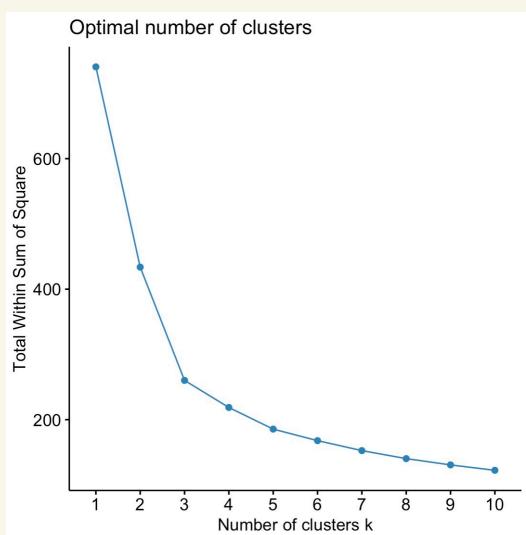
Then run K-means using these Centers as the starting points

Time to Cluster NBA players into offensive Roles by Synergy FGA play type frequencies

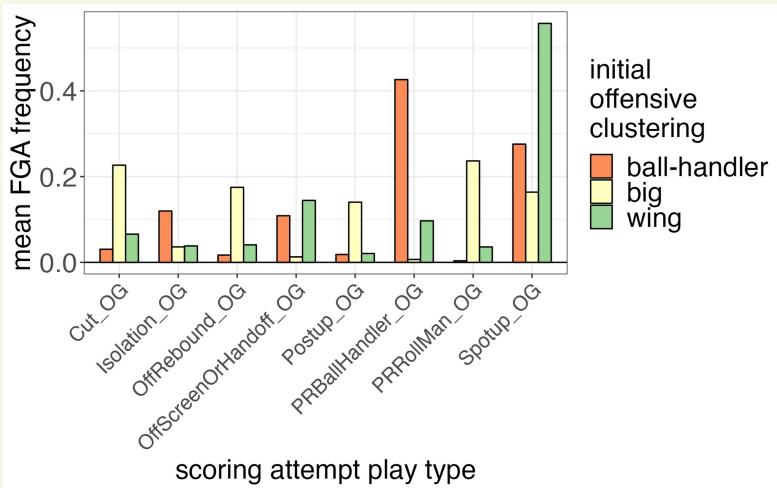
```
### Hard Clustering: K means
```

```
set.seed(387397)
```

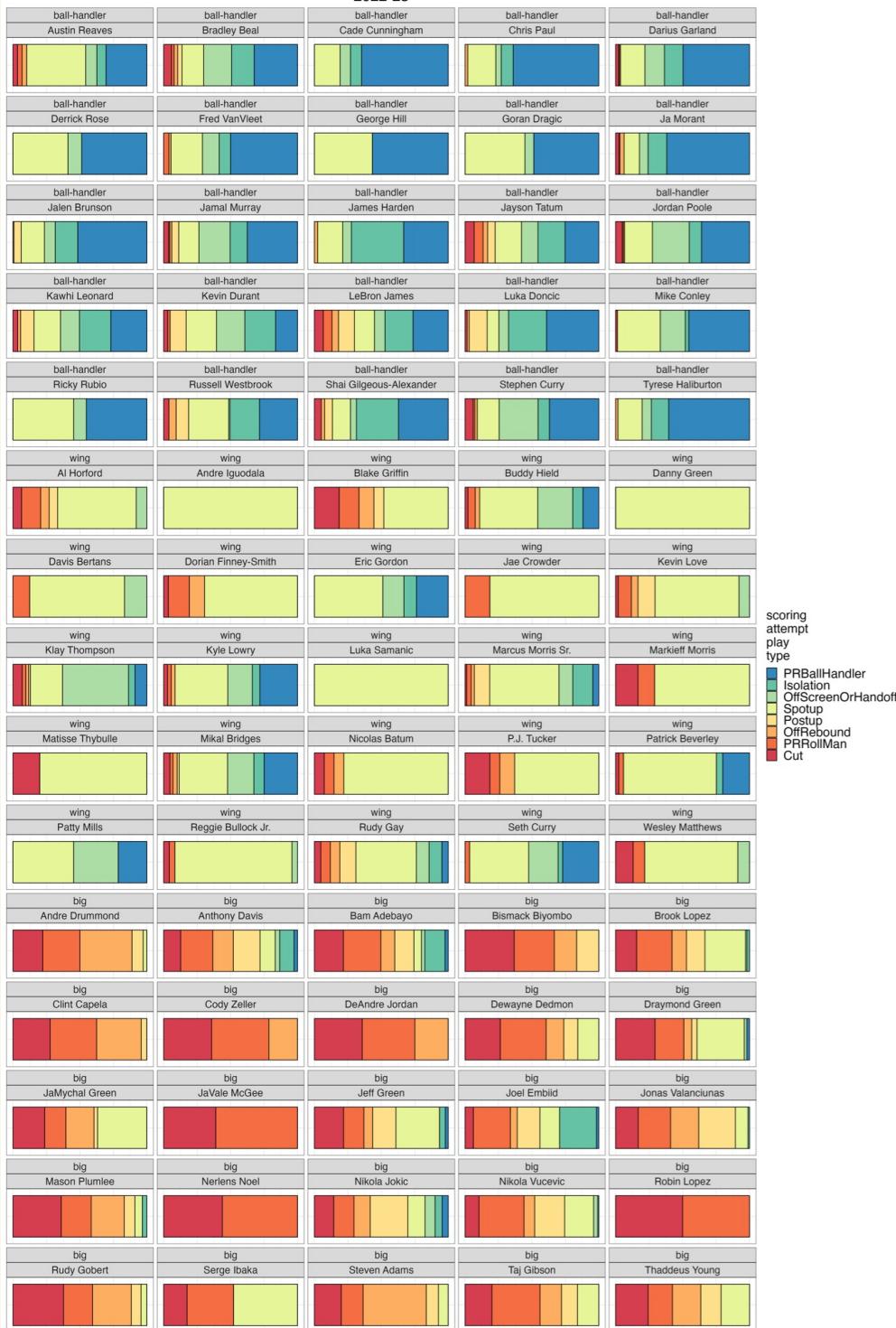
```
km1 <- flexclust::kcca(mat_train, K1)  
print(km1)
```

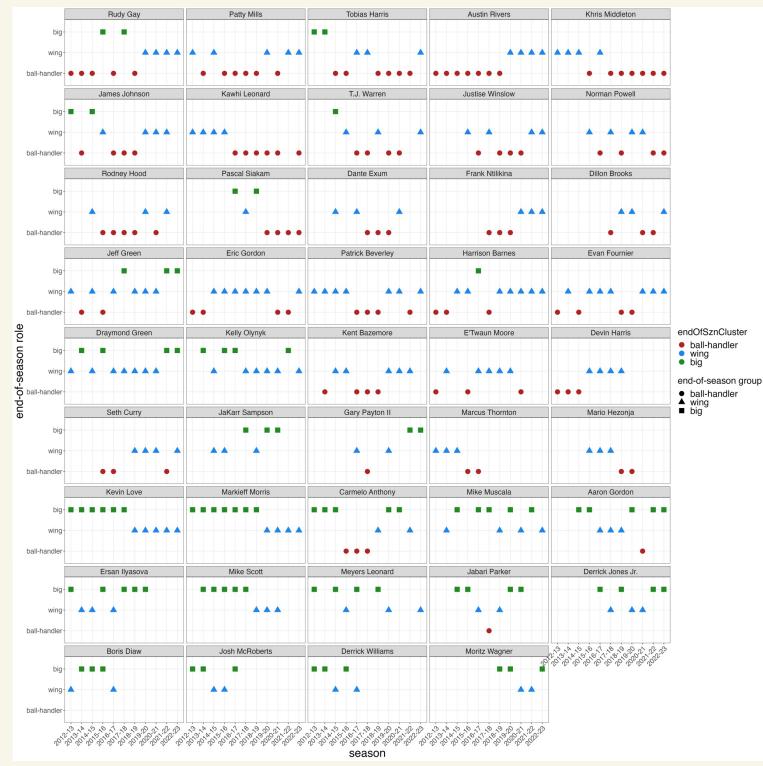
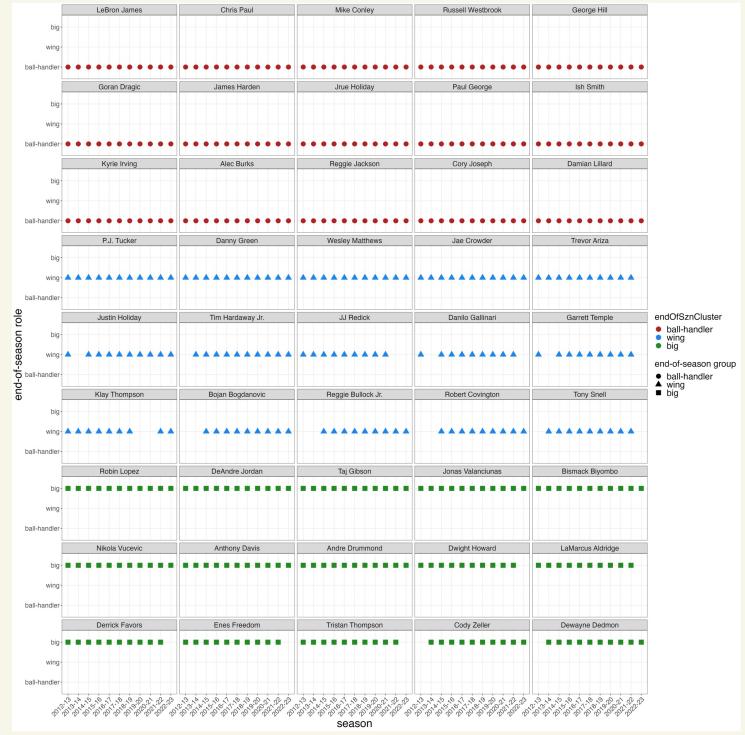


Elbow Plot
Began with K=3
clusters



2022-23





I further clustered the 3 groups into more refined Roles:

2022-23

