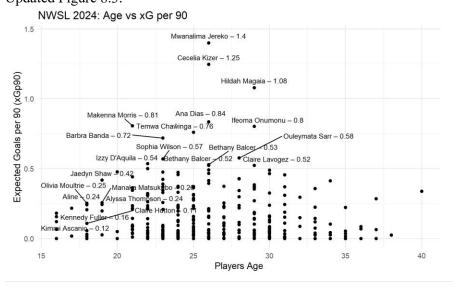
Module 7: Coding Assignment 2

#### Scenario 1: Data Collection and summarization

Similar to the nwsl\_player\_stats file within the ISAR package, I collected data for the 2024 NWSL season via FBref. Although this process was somewhat difficult as I had trouble finding the right table to scrape, as my initial code kept gathering squad related data instead of player, I was eventually able to output the 2024 NWSL player data and make a corresponding master excel file. To do this, instead of writing out a scrape, I simply copied in the entire 2024 NWSL player stats table into an excel spreadsheet, renamed the columns to be the exact same as that in the nwsl\_player\_stats spreadsheet already in ISAR for the 2022 NWSL season, and then begun subsequent analysis. Within this dataset contains 376 players across the 14 teams. The variables contained in this dataset that I will be using for the remainder of the assignment are player (players name), squad (the NWSL team), pos (players position), age (players age), min (total minutes players), starts (games started), xG (expected goals), xAG (expected assists), xGp90 (expected goals per 90 minutes), xAp90 (expected assists per 90), and xGxAp90 (sum of expected goals and assists per 90 minutes).

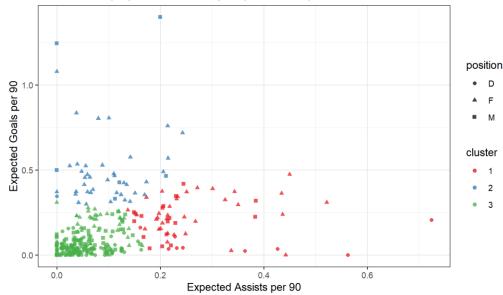
The EPL 2024 goalkeeping stats are already downloaded through ISAR and didn't require additional scraping. There are 44 goalkeepers in the dataset who played during the 2024/2025 Premier League season across all 20 teams. The variables used for this assignment include player (goalkeeper's name), nation (goalkeepers nationality),cPSxG/SoT (post-shot expected goals per shot on target), PSxG+/minus(Post-shot expected goals minus goals allowed), LaunchPCT (Launch completion percentage, Att (attempted passes), AvgLen (Average pass length), GK\_LaunchPCT (Goal kick launch percentage), and DefActions\_OPA\_90 (Number of defensive actions outside penalty area per 90 minutes).

**Scenario 2**: Clustering NWSL Players Updated Figure 8.3:



## Updated Figure 8.10:





Within the updated 8.10 figure, it looks like a high proportion of defenders are in cluster 3 which have low amounts of expected goals per 90 and assists per 90. There is then a higher proportion of forwards in group two, who have a higher expected goals per 90 versus expected assist per 90. Lastly, there are a lot of midfielders in cluster 1, as well as forwards, who have a higher rate of expected assists per 90 than expected goals per 90.

Table summarizing # of athletes assigned to each position and cluster

clusters_three <int></int>	<b>D</b> <int></int>	<b>F</b> <int></int>	<b>M</b> <int></int>
1	9	34	18
2	1	39	6
3	112	46	85

The table above confirms my prediction that cluster 3 is dominated by defenders. However, not all midfielders and forwards can score goals and assists, leading a solid proportion of them also in cluster 3. As forwards are the highest up the pitch, they dominate clusters 1 and 2, which contain players with higher rates of expected goals and assists per 90. Also, as cluster 1 is dominated by players with higher rates of expected assist per 90, there are 12 more midfielders in cluster 1 than cluster 2.

### Updated Figure 8.11:

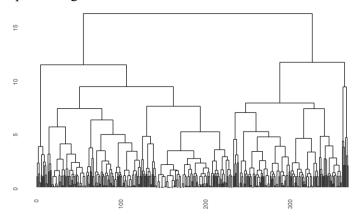


Table summarizing # of players assigned to each cluster:

cluster_k5 <int></int>	cluster_k4 <int></int>	count <int></int>	
1	1	129	
2	2	108	
3	2	111	
4	3	13	
5	4	8	

As shown in the table above, some clusters overlapped, leading to a reduction in assignments. This main change and overlap occurred during the reduction of the clusters from 5 to 4, resulting in the merging of clusters 2 and 3 (k=5) into a single cluster 2 (k=4). This merge combined 108 and 111 players from two separate k=5 clusters into a single k=4 cluster. The remaining clusters maintained a one-to-one ratio: k=5 cluster 1 became k=4 cluster 1, k=5 cluster 4 became k=4 cluster 3, and k=5 cluster 5 became k=4 cluster 4 (8 players). The only difference is that k=4 cluster 2 is now a combination of k=5 cluster 2 and 3. This consolidation reduced the total cluster count, with the most notable shift occurring between clusters 2 and 3.

## Scenario 3: Evaluating EPL Goalkeepers

### Recreate 8.9 step 1:

Table Summarizing Cluster Means for k=2

<b>cluster</b> <fctr></fctr>	PSxG_SoT <dbl></dbl>	PSxG_minus <dbl></dbl>	LaunchPct <dbl></dbl>	Att <dbl></dbl>	AvgLen <dbl></dbl>	GK_LaunchPct <dbl></dbl>	DefActions_OPA_90 <dbl></dbl>
1	0.3232143	0.425	32.87143	94.82143	31.48571	52.09643	1.110357
2	0.3033333	0.160	37.19333	499.00000	34.04667	62.44000	1.148667

### Table Summarizing Cluster Means for k=3

cluster <fctr></fctr>	PSxG_SoT <dbl></dbl>	PSxG_minus <dbl></dbl>	LaunchPct <dbl></dbl>	Att <dbl></dbl>	AvgLen <dbl></dbl>	GK_LaunchPct <dbl></dbl>	DefActions_OPA_90 <dbl></dbl>
1	0.3040000	1.6600000	48.60000	673.60000	38.66000	69.24000	0.894000
2	0.3009091	-0.3727273	30.96364	397.63636	31.64545	58.72727	1.368182
3	0.3248148	0.3740741	33.13704	88.81481	31.51481	51.96667	1.066667

### Percentage of Variation Explained:

k = 2:77.78 %

k = 3:89.26 %

Adding a third cluster captures roughly 90% of the variation, whereas including just two clusters captures on 78% of the variance in the goalkeepers data.

# Recreate 8.9 step 2:

Updated dataframe includes player name, nation, and k=2 and k=3 cluster

Player schr>	Nation <chr></chr>	cluster_k2 <int></int>	cluster_k3 <int></int>	
Nick Pope	eng ENG	1	2	
Alisson	br BRA	1	3	
randon Austin	eng ENG	i	3	
Itay Bayındır	tr TUR	1	3	
aniel Bentley	eng ENG	i	3	
lartin Dúbravka	sk SVK	1	3	
derson	br BRA	i	3	
ukasz Fabiański	pl POL	1	3	
raser Forster	eng ENG	1	,	
itezslav Jaros	cz CZE		3	
am Johnstone	eng ENG	1	,	
lip Jørgensen	dk DEN		,	
aoimhin Kelleher	ie IRL		3	
ntonín Kinský	cz CZE		3	
e Lumley	eng ENG	,	3	
e curriey lex McCarthy			3	
	eng ENG		3	
ijanet Muric	xk KVX br BRA		3	
eto obin Olsen	or BRA se SWE		3	
	de GER		3	
efan Ortega ex Palmer			3	
	eng ENG		3	
ákon Rafn Valdimarsson	is ISL	!	3	
son Steele	eng ENG		3	
kub Stolarczyk	pl POL	!	3	
ark Travers	ie IRL	!	3	
uglielmo Vicario	it ITA	1	3	
hristian Walton	eng ENG	!	3	
anny Ward	wls WAL	1	3	
ark Flekken	nl NED	2	1	
ean Henderson	eng ENG	2	1	
rdan Pickford	eng ENG	2	1	
aron Ramsdale	eng ENG	2	1	
atz Sels	be BEL	2	1	
phonse Areola	fr FRA	2	2	
pa Arrizabalaga	es ESP	2	2	
ads Hermansen	dk DEN	2	2	
ernd Leno	de GER	2	2	
miliano Martínez	ar ARG	2	2	

Recreate the second plot in Section 8.9 Step 4:

# Plot for k=2:

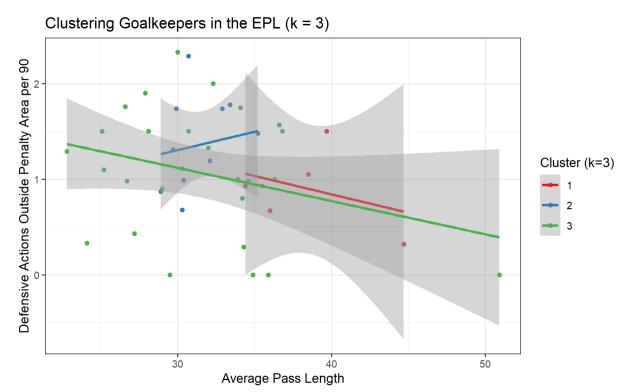
Clustering Goalkeepers in the EPL (k = 2)

Cluster (k=2)

Cluster (k=2)

Average Pass Length

Plot for k=3



Based on the two plots above clustering goalkeepers in the English Premier League for k = 2 and k = 3, we can identify different types of goalkeepers in the EPL. First off, the k = 2 chart groups keepers into two clusters: one with short passes and more defensive actions and a second with longer passes and fewer defensive actions. One expects a player with longer passes to be a traditional goalkeeper who boots it up the fielder versus a "sweeper-keeper" who can play in possession and come out of the box to intercept and make short passes. I find Manuel Neuer to be the perfect example for this cluster. The k=3 chart groups keepers into 3 different clusters. This time, there is a more neutral keeper with a medium pass length and about average defensive actions outside the penalty area per 90. Again, there is a sweeper keeper and a more traditional player