



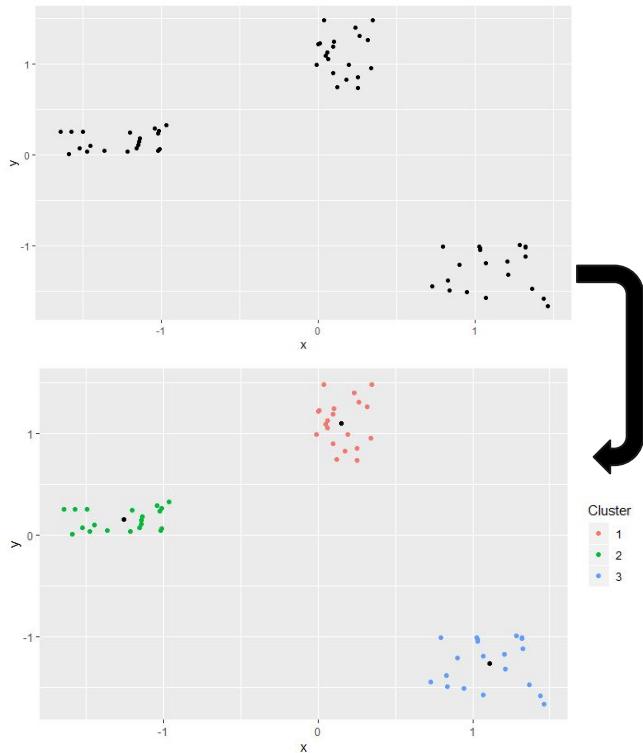
# Clustering With PCA and K-means

Brandon Chan

---

# What is clustering?

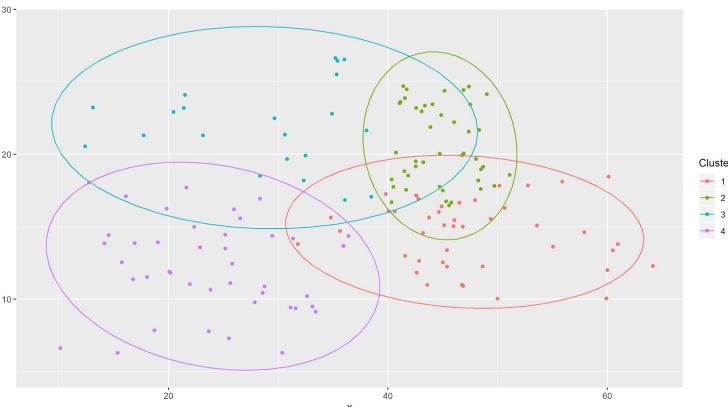
- Clustering is when you group similar data points together
- One group of similar data points makes up a cluster
- One key characteristic of clustering is that the groups/clusters are not predefined



# Why cluster?

- Clustering is a great way to get the bigger picture of a data set
- Clustering is especially useful for understanding big data sets with lots of variables
- Also, by grouping together similar data points we can:
  - Find relationships among variables
  - Spot outliers
  - Uncover underlying patterns in the data

	name	short	gp	pos	pttFrequency	phs	lgs	fpm	pp	psyWeariness	psyPerTeam	psyPerTeamS	psyPerPurc	psyPerGame	psyPerGame	psyPerGameS	lgsPerGame	lgsPerGameS	lgsPerPurc	lgsTotal	pttFG	pttTS%	pttTDow
ATL	Atlanta	ATL	81	207	0.32773209	170	179	81	0.32121004	104	301	2.03333368	2.00776423	2.00776423	0.73000042	1.45079102	1.45079102	118	0.44701012	0.4444			
SAC	Sacramento	SAC	82	206	0.16751702	345	274	118	1.165549	428	59	1.80701010	2.01751757	2.01751757	1.16145541	1.16145541	1.16145541	1.92059227	158	0.43335766	0.803		
ATL	Atlanta	ATL	80	148	0.09800540	175	152	83	1.182450	445	60	1.80000000	2.16750000	2.16750000	0.73000000	0.73000000	0.73000000	70	0.44999997	0.379			
TOR	Toronto	TOR	65	159	0.16989101	180	152	51	1.165549	428	66	2.00000000	2.01750000	2.01750000	0.73000000	2.01750000	2.01750000	59	0.44999997	0.379			
CLE	Cleveland	CLE	72	172	0.16989101	64	60	43	1.165549	428	60	2.01750000	2.01750000	2.01750000	0.73000000	1.98999952	1.98999952	62	0.44999974	0.342			
CLE	Cleveland	CLE	14	18	0.12050103	21	15	7	1.121051	491	14	1.82000000	2.01750000	2.01750000	0.73000000	1.97142327	1.97142327	8	0.44999997	0.323			
TOR	Toronto	TOR	60	253	0.20698803	269	226	101	1.142329	424	71	2.00000000	2.01750000	2.01750000	0.73000000	1.96200000	1.96200000	125	0.44999995	0.330			
DET	Detroit	DET	72	192	0.09709802	165	178	83	0.947917	207	281	2.09889987	2.02777775	2.02777775	0.70700000	1.97722322	1.97722322	115	0.39332828	0.0491			
CHI	Chicago	CHI	22	46	0.14713800	45	43	17	0.970516	244	263	2.09009009	2.04645453	1.98454845	0.77027272	1.81781916	1.81781916	26	0.35850484	0.323			
LAL	Lakers	LAL	47	123	0.21317303	68	118	33	0.970454	90	405	2.07021238	2.08106988	2.08106988	0.70217686	1.80810644	1.80810644	65	0.37991021	0.2415			
F	GSW	GSW	73	174	0.20237765	169	162	55	0.970454	90	263	2.08910000	2.08910000	2.08910000	0.70217686	1.80810644	1.80810644	52	0.34320000	0.2350			
DET	Detroit	DET	23	49	0.20332000	47	45	15	0.970454	90	250	2.08910000	2.08910000	2.08910000	0.70217686	1.80810644	1.80810644	52	0.34320000	0.2350			
DET	Detroit	DET	77	148	0.16277300	273	216	80	1.000000	388	156	1.22700003	1.56544543	2.02191001	1.31310103	1.82278762	1.82278762	125	0.41200000	0.3851			
TOR	Toronto	TOR	67	160	0.39711899	173	176	80	0.986173	156	340	2.08957071	2.05208957	2.05208957	0.99532239	1.75154219	1.75154219	116	0.34090009	0.471			
LAL	Lakers	LAL	67	168	0.30950000	162	158	58	0.986173	219	277	2.05208906	2.47179145	2.35820000	1.52288000	1.52288000	1.52288000	102	0.34040038	0.323			
WAS	Washington	WAS	62	266	0.15203400	309	241	110	1.118105	425	60	2.04900004	2.76502688	2.90902409	1.34160341	1.97763000	1.97763000	131	0.45614564	0.355			
NOR	New Orleans	NOR	18	38	0.31400000	45	35	18	1.142329	447	45	2.01111111	2.00000000	1.64444444	1.00000000	0.64444444	0.64444444	17	0.41420071	0.3230			
MEM	Memphis	MEM	36	41	0.20907000	58	52	19	0.905020	208	287	1.00000000	1.00000000	1.00000000	0.64444444	0.97000000	0.97000000	33	0.34320000	0.2417			
IND	Indiana	IND	24	56	0.20907000	240	24	13	0.905020	68	437	1.05714248	1.42027148	1.42027148	0.20000000	1.24282771	1.24282771	47	0.31999987	0.2418			
CHI	Chicago	CHI	81	162	0.09600000	206	148	48	1.171051	478	18	2.00000000	2.04220000	1.88000000	0.83900017	1.00000000	0.00000000	81	0.44917064	0.347			
POR	Portland	POR	70	214	0.14800000	233	202	64	1.080700	381	114	2.07140048	2.08571400	1.20000000	1.60717429	1.10	0.41814548	0.351					
ATL	Atlanta	ATL	45	101	0.24848101	94	69	30	0.905020	100	356	2.04666666	2.00000000	1.97777775	0.86666667	1.31111101	0.59	0.33707965	0.0477				
TOR	Toronto	TOR	60	158	0.11000000	181	146	66	1.145170	427	66	2.00000000	2.00000000	1.00000000	1.32200000	1.32200000	80	0.43200079	0.371				
DAL	Dallas	DAL	28	68	0.42000000	101	79	36	1.147700	420	65	2.05464278	2.04270000	2.04270000	1.24179511	1.40270000	1.40270000	45	0.45990000	0.534			
LAC	Lakers	LAC	4	6	0.42000000	6	2	2	0	0	2	2.05464278	2.04270000	2.04270000	1.24179511	1.40270000	1.40270000	4	0.45990000	0.534			
NYC	New York	NYC	73	204	0.35100000	241	191	76	0.920471	311	184	2.07400055	2.00000000	1.64444444	1.32700014	1.00000010	118	0.36170114	0.2238				
WAS	Washington	WAS	77	208	0.32000000	310	256	103	1.072000	385	130	2.07334075	2.03874003	2.02471032	1.37100014	1.00000000	155	0.41334775	0.3338				
WAS	Washington	WAS	80	210	0.20500000	233	187	79	1.119051	399	95	2.02000000	2.07700000	0.98700000	1.32700000	1.00000000	108	0.42400098	0.3664				
MIN	Minnesota	MIN	65	172	0.53700000	164	150	50	0.954000	212	283	2.04913005	2.02070002	2.02070002	1.09200077	1.75500054	1.75500054	100	0.33533303	0.480			
MDM	Memphis	MDM	44	162	0.33100000	135	141	48	0.882000	163	260	3.47702773	3.36618162	3.20404468	1.00000000	2.11500000	2.11500000	92	0.40420053	0.468			



Cluster

1

2

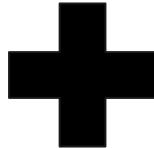
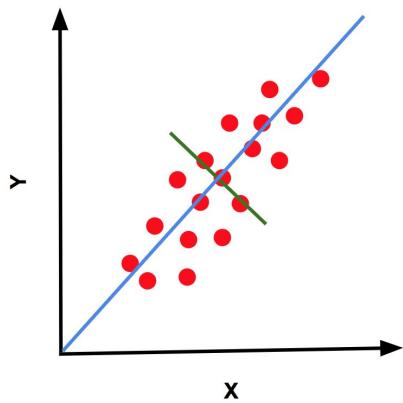
3

4

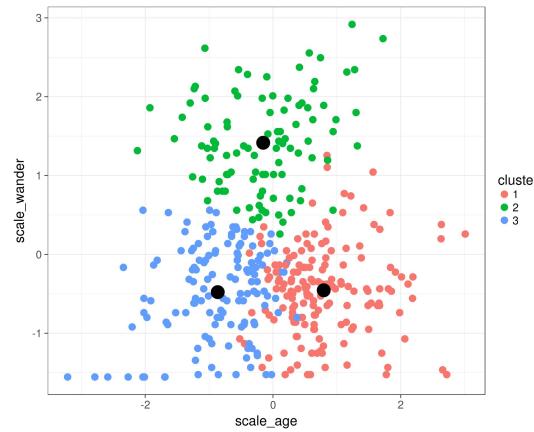
---

# How?

## Principal Component Analysis



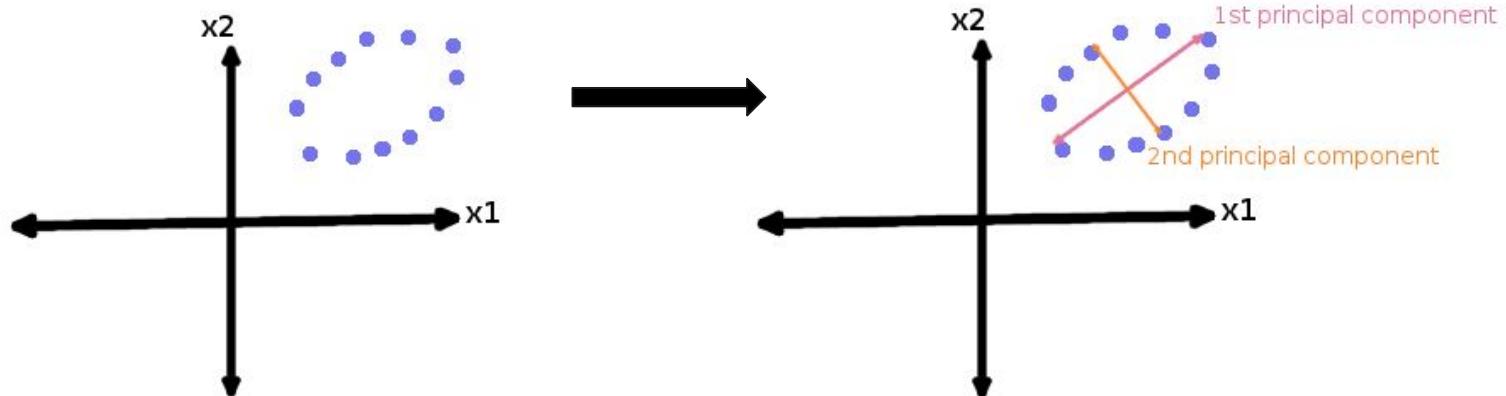
## K-means clustering



---

# Principal Component Analysis

- Principal Component Analysis (PCA) is a transformation of the original data
- PCA plots the original data points on a new coordinate system such that the variation among the data points is maximized



---

## Some Things to Note

- Each axis in the new coordinate system is known as a Principal Component
- Each Principal Component represents one way that the data points can be transformed
- Each transformation results in a different amount of variation among the data points
- The Principal Components are ordered from the highest amount of variation to the least amount of variation

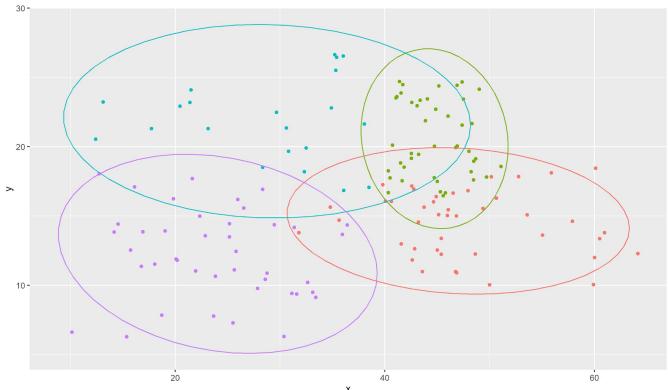
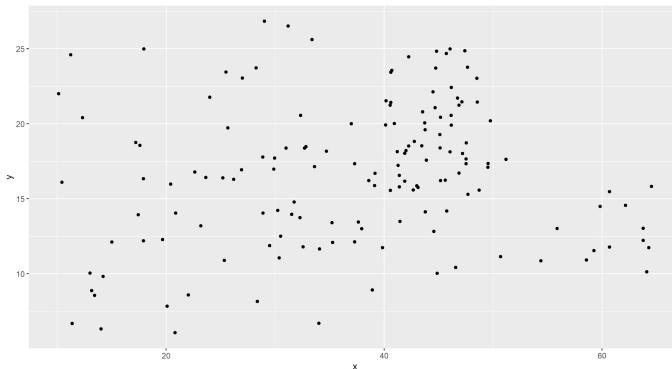


**In summary, PCA transforms the points so that  
there is the MOST variation in the LEAST  
number of axes.**

---

# K-means Clustering

- Simply put, here's how k-means clustering works:
  - The user decides how many clusters they want to have
  - The k-means algorithm decides where the centers of those clusters should be
  - Every data point is assigned to the nearest cluster



---

# The Data

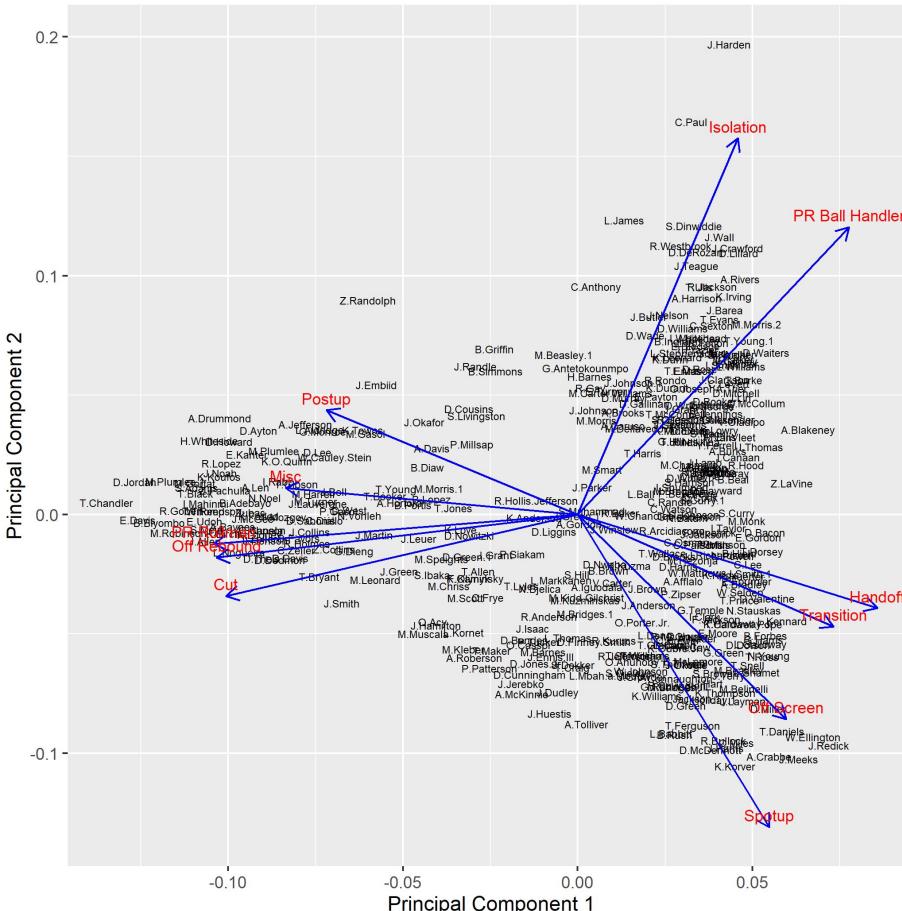
One data point - One NBA Player who played at least 1000 minutes in the past 3 seasons

Variables - The frequencies at which they took part in the following plays:

- Rollman in the pick and roll
- Off rebound
- Cut
- Postup
- Isolation
- Spotup
- Off screen
- Transition
- Ball handler in the pick and roll
- Handoff
- Misc

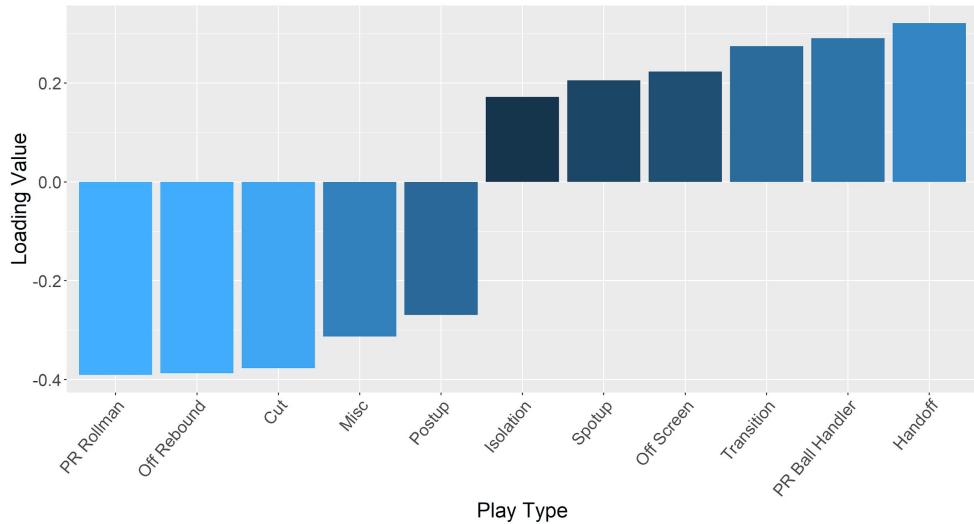


# Biplot





# PC1 Loadings

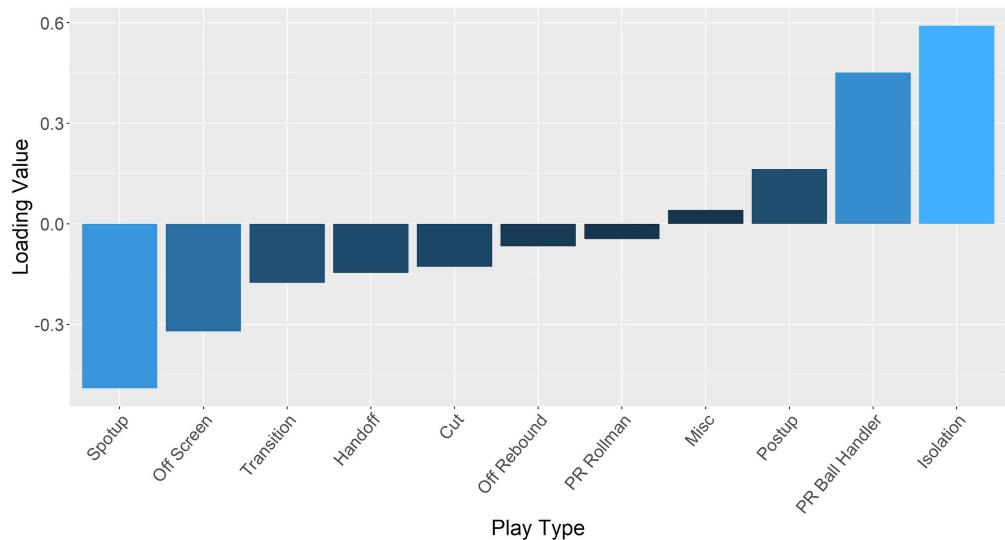


- High Handoff, PB Ball Handler, and Transition
- Low PR Rollman, Off Rebound, and Cut
- Separates big men from guards and forwards





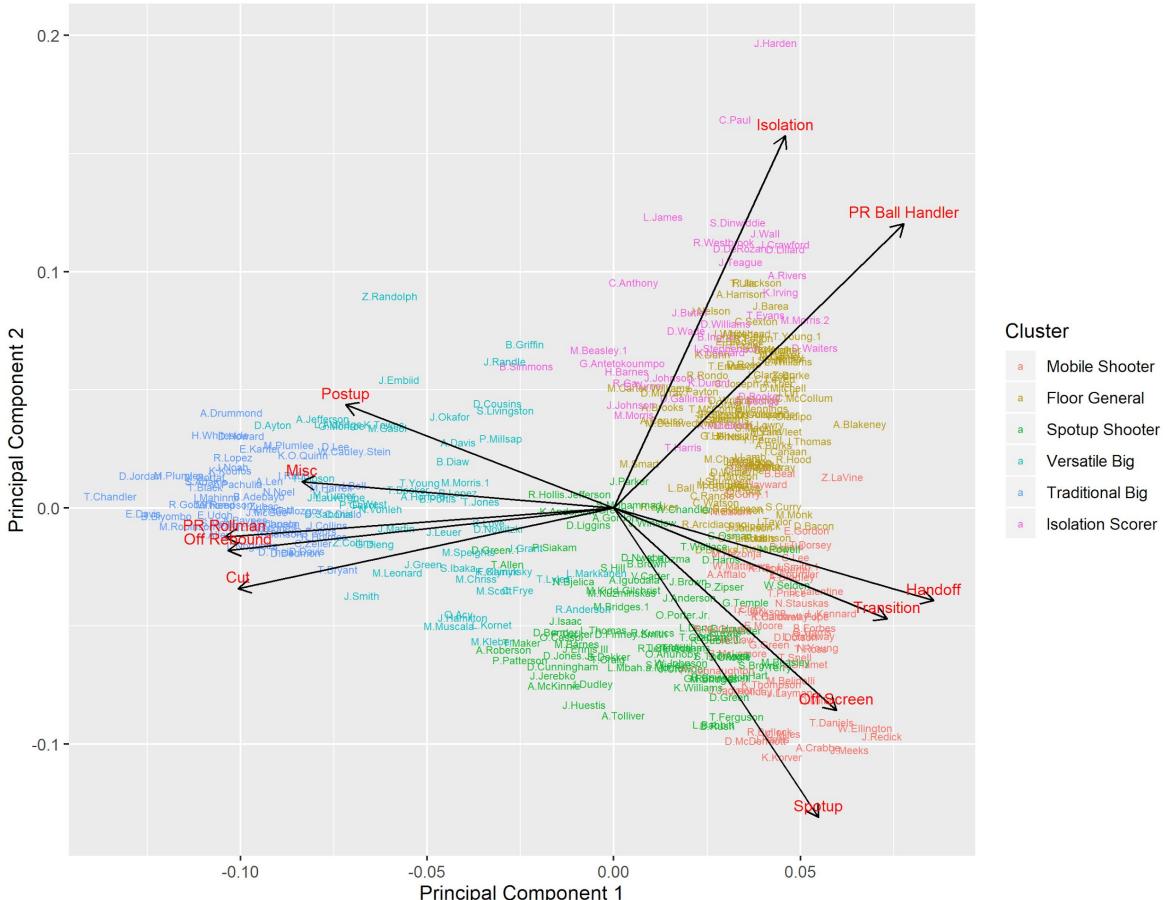
# PC2 Loadings



- High Isolation, and PR Ball Handler
- Low Spotup and Off Screen
- Separates ball handlers from shooters

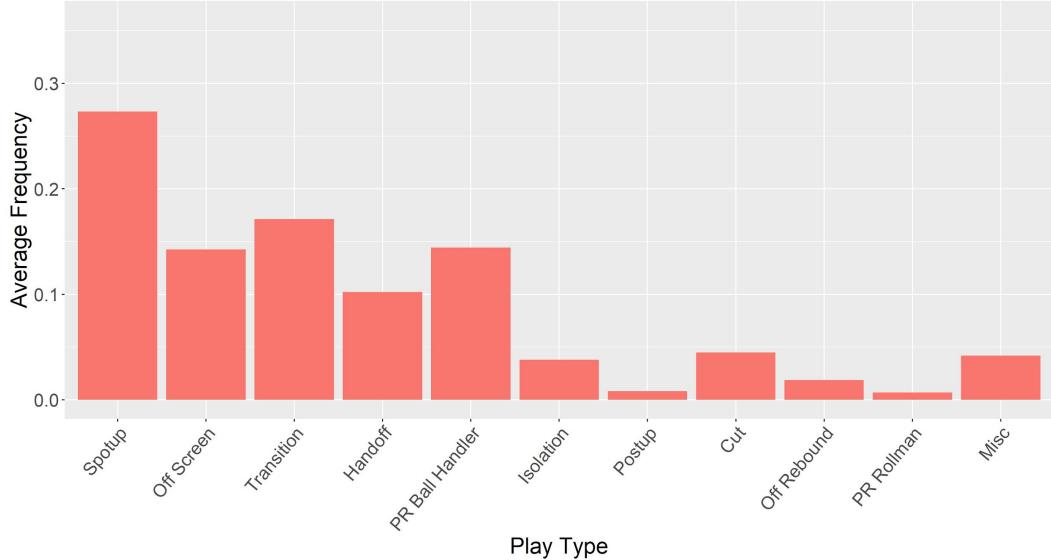


# Clusters



---

# Mobile Shooters

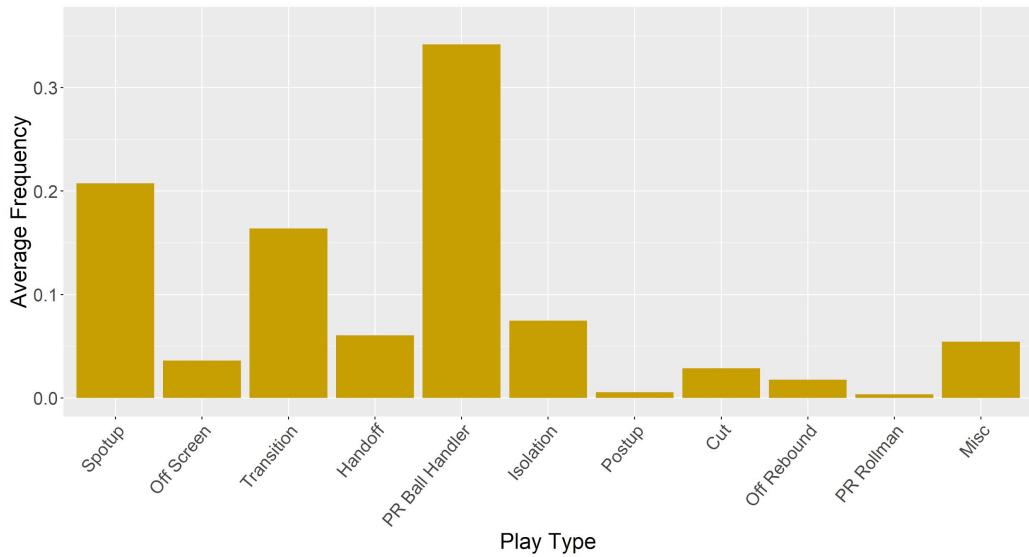


- High Spotup, Off screen, Transition
- Low PR Rollman, Postup, Off Rebound
- Key members: Buddy Hield, JJ Redick, Klay Thompson





# Floor Generals

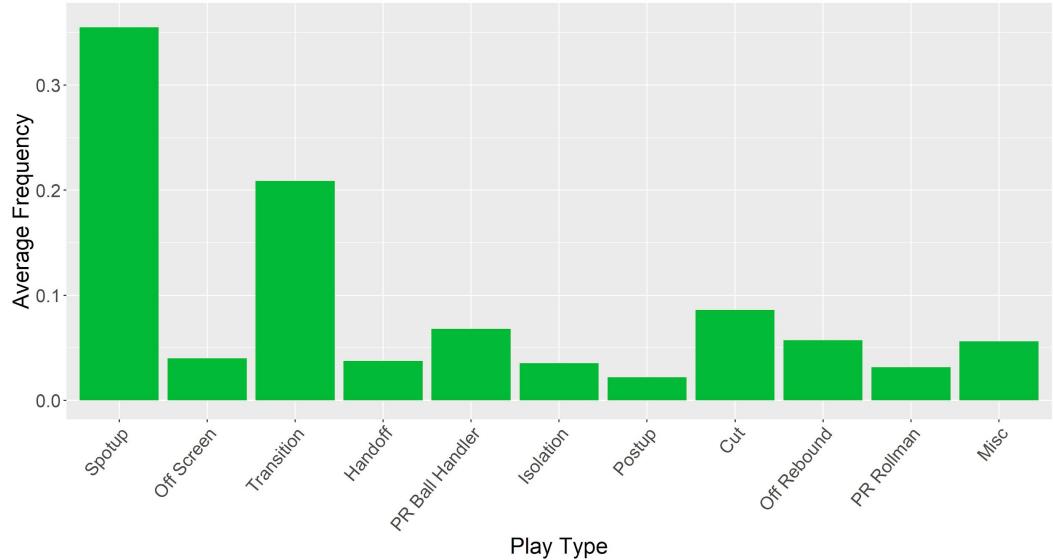


- High PR Ball Handler, Spotup, Transition
- Low PR Rollman, Postup, Off Rebound
- Key Members: Steph Curry, Lonzo Ball, Rajon Rondo



---

# Spotup Shooters

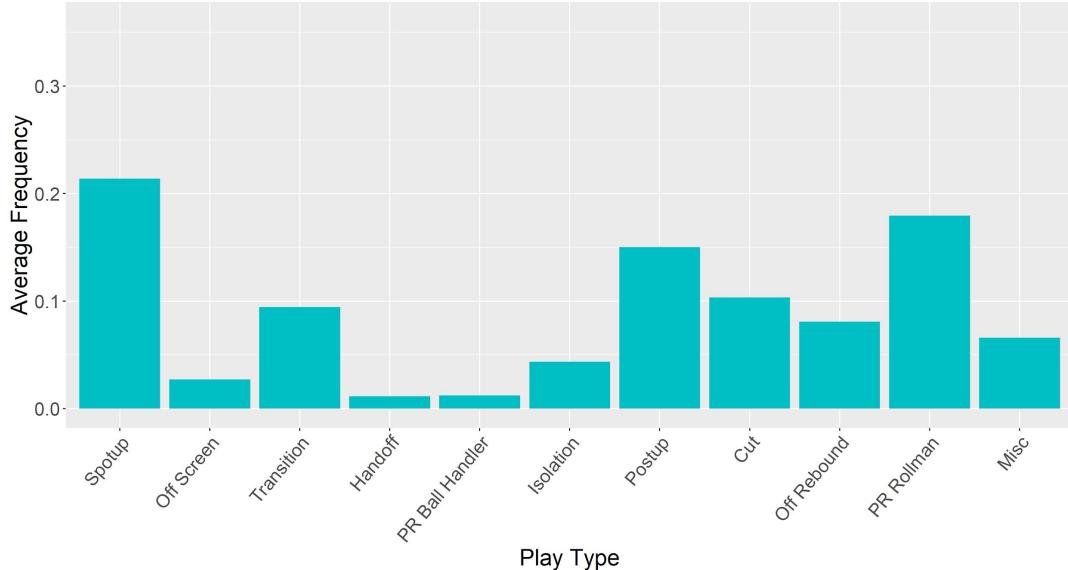


- High Spotup and Transition
- Low Off Screen and Postup
- Key members: Otto Porter Jr., Anthony Tolliver, Nemanja Bjelica



---

# Versatile Bigs

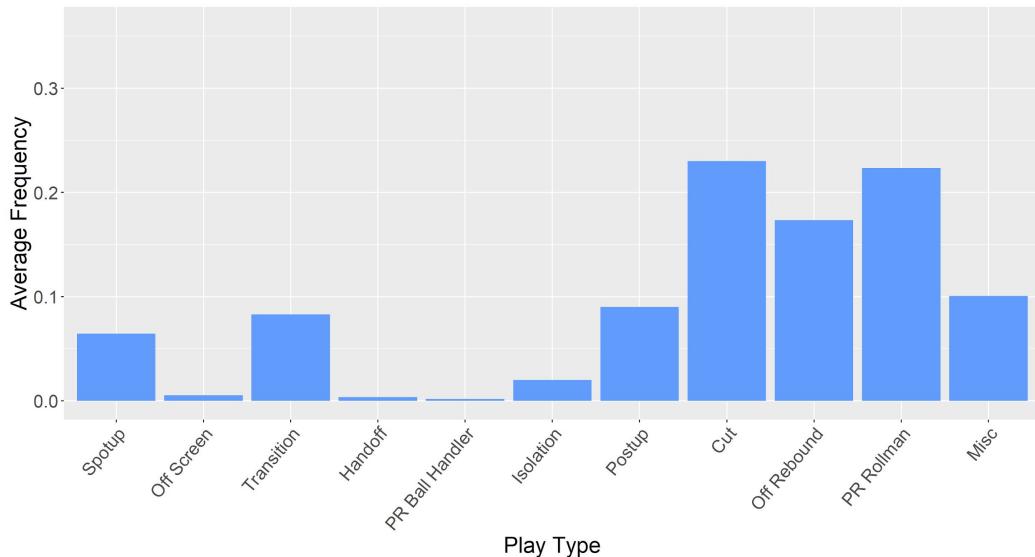


- High Spotup, PR Rollman, Postup
- Low Handoff, PR Ball Handler, Off Screen
- Key members: Demarcus Cousins, Anthony Davis, Blake Griffin





# Traditional Bigs

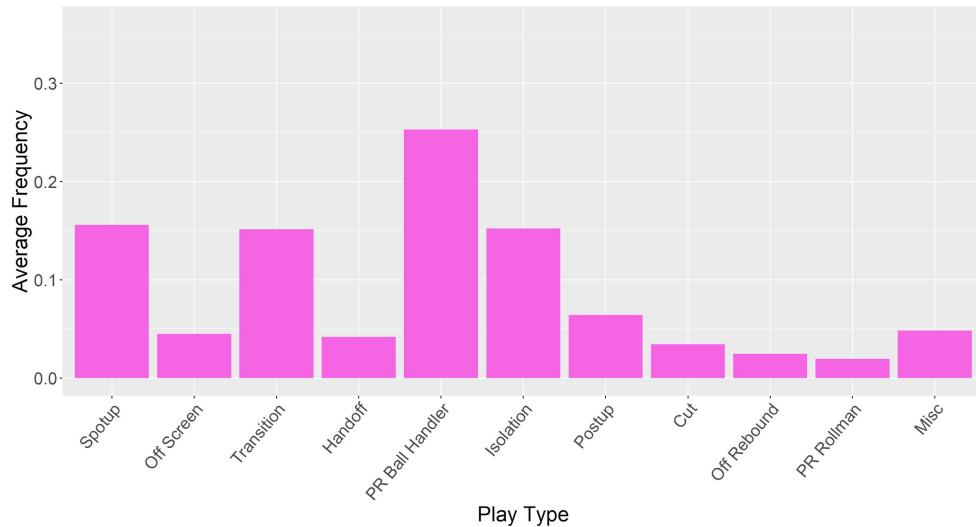


- High Cut, PR Rollman, Off Rebound
- Low PR Ball Handler, Handoff, Off Screen
- Key members: Rudy Gobert, DeAndre Jordan, Andre Drummond



---

# Isolation Scorers



- High PR Ball Handler, Isolation, Transition
- Low PR Rollman, Off Rebound, Cut
- Key members: James Harden, Lebron James, Russell Westbrook



---

# Main Takeaways

- The best way to differentiate NBA players is by separating big men from non-big men
- The second-best way to differentiate NBA players is by separating ball handlers from non-ball handlers.

## Big Men

Postups  
Being the rollman  
Grabbing rebounds  
Cuts  
Miscellaneous

## Ball Handlers

Isolation  
Running the pick and roll

## Shooters

Spotting up  
Running in Transition  
Handoffs  
Running off screens



---

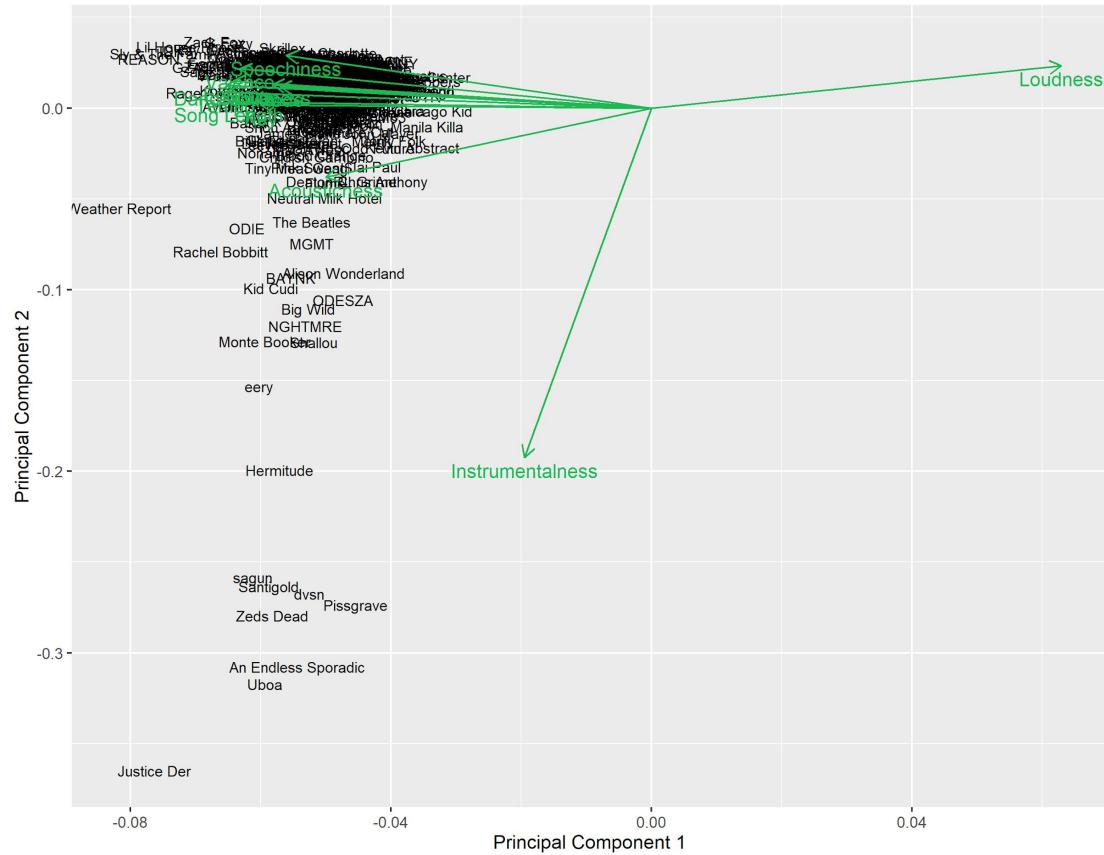
# The Data

One data point - One artist that I have listened to in the past 4 months

Variables - the average of these attributes from the songs I listened to from each artist:

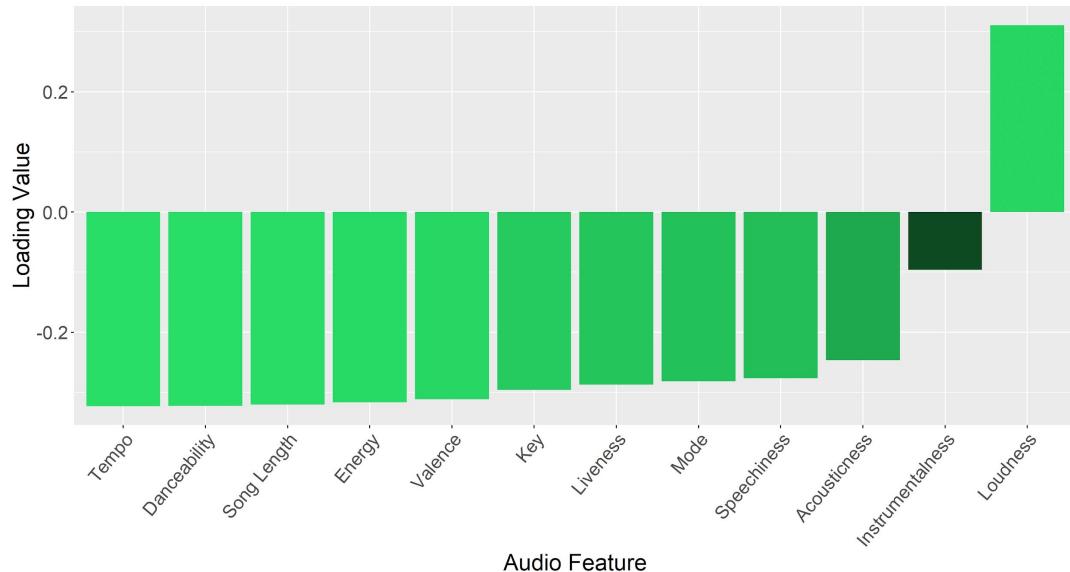
- Tempo
- Danceability
- Song Length
- Energy
- Valence
- Key
- Liveness
- Mode
- Speechiness
- Acousticness
- Instrumentalness
- Loudness







# PC1 Loadings

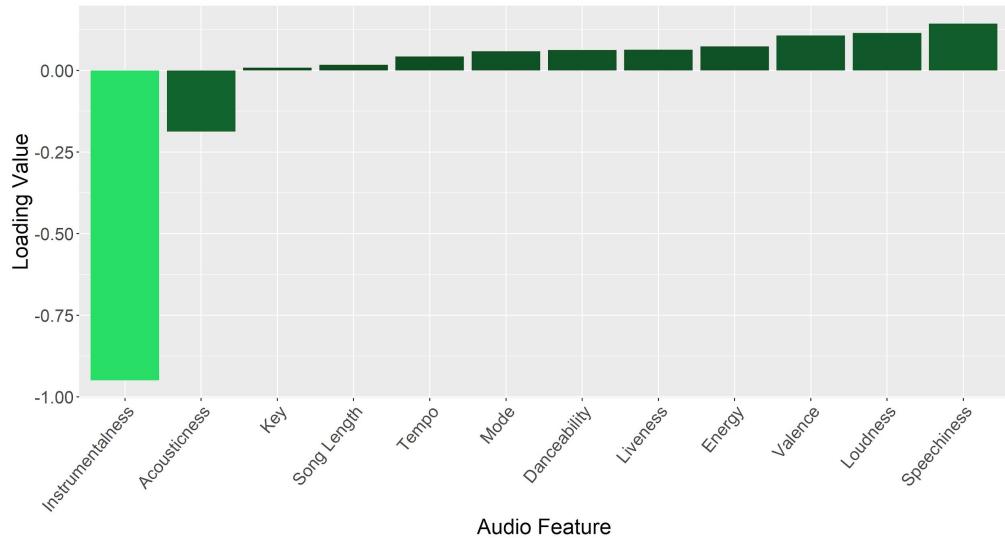


- High Loudness
- Low everything else
- Separates loud artists from not loud artists





## PC2 Loadings

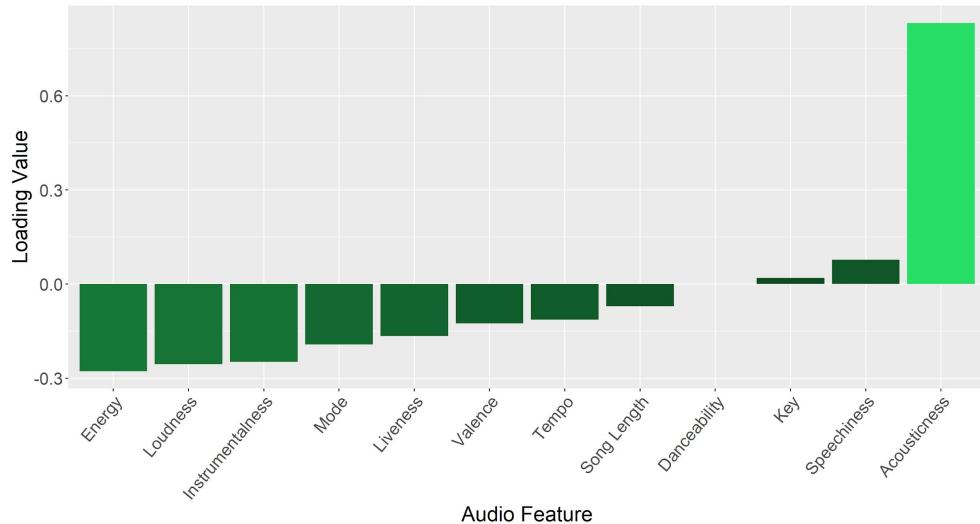


- High most variables
- Very low Instrumentalness
- Separates artists who play instrumentals and artists that use vocals





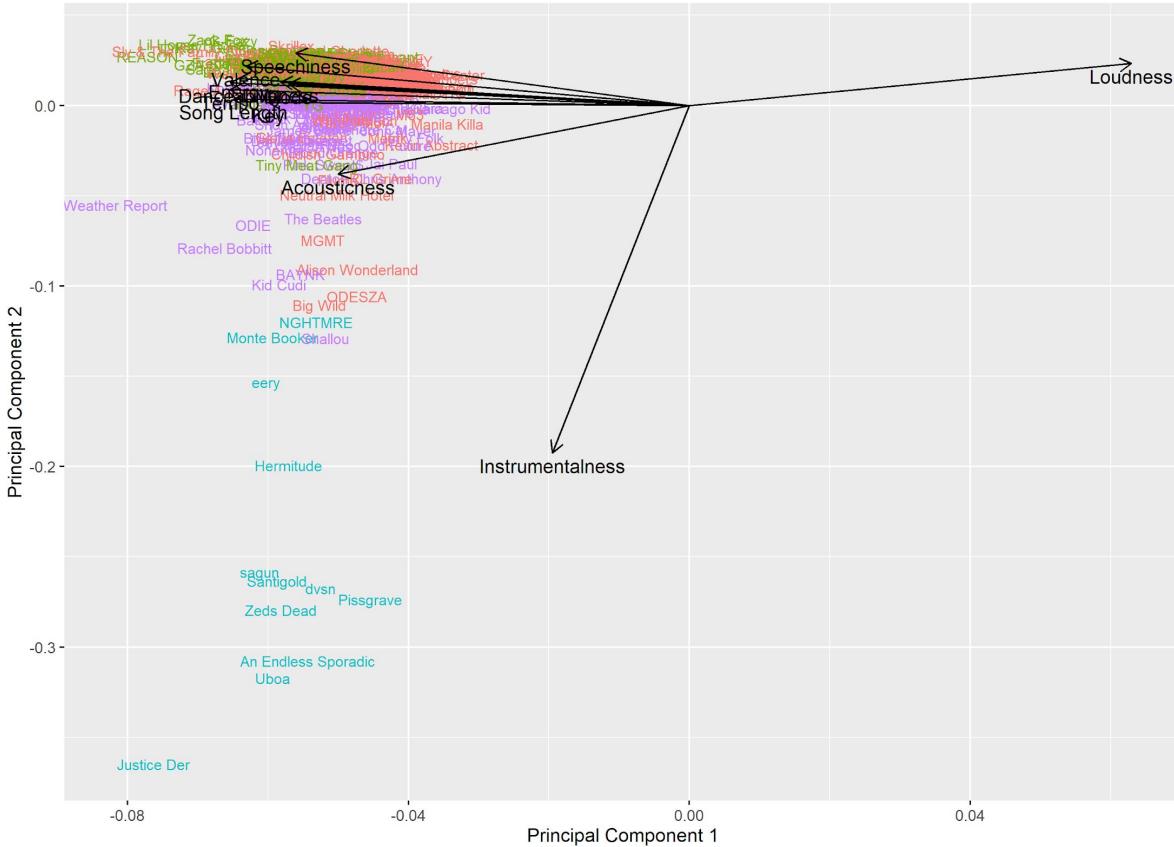
# PC3 Loadings



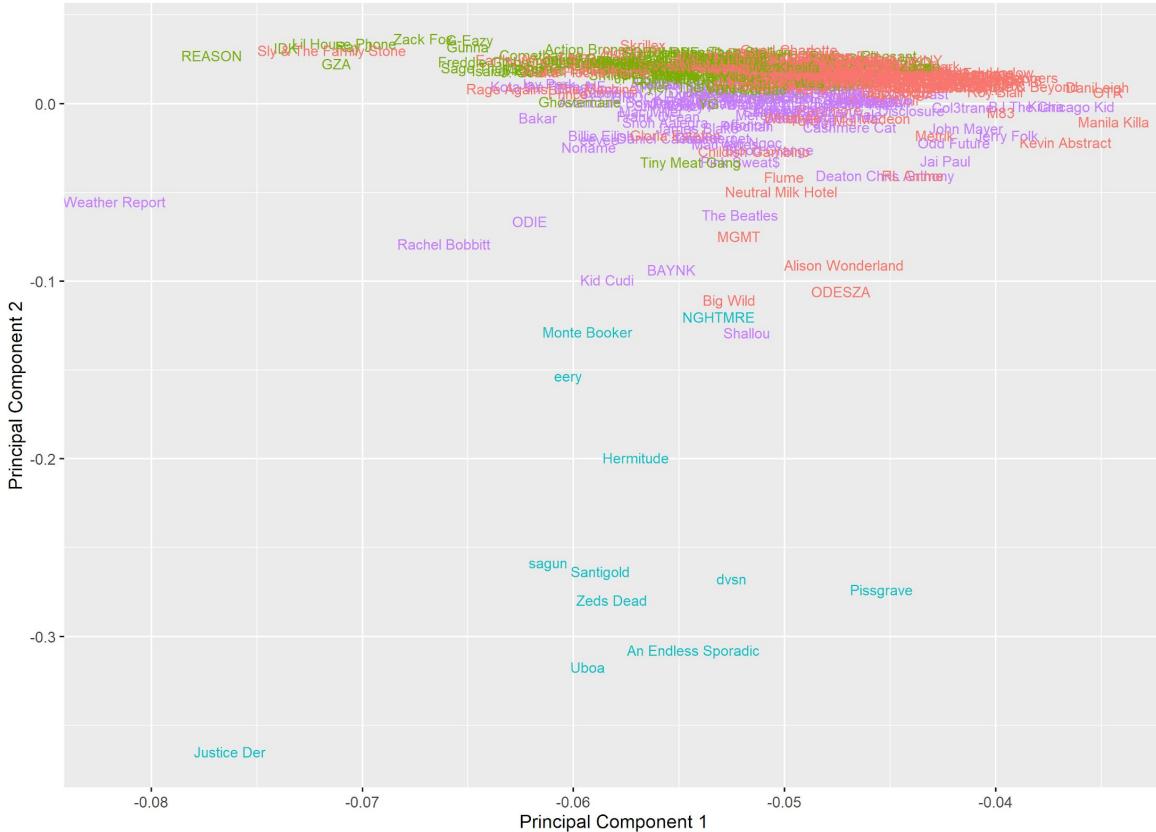
- High Acousticness
- Low everything else
- Separates acoustic instrumentation artists from artists that use other instruments



# Clusters

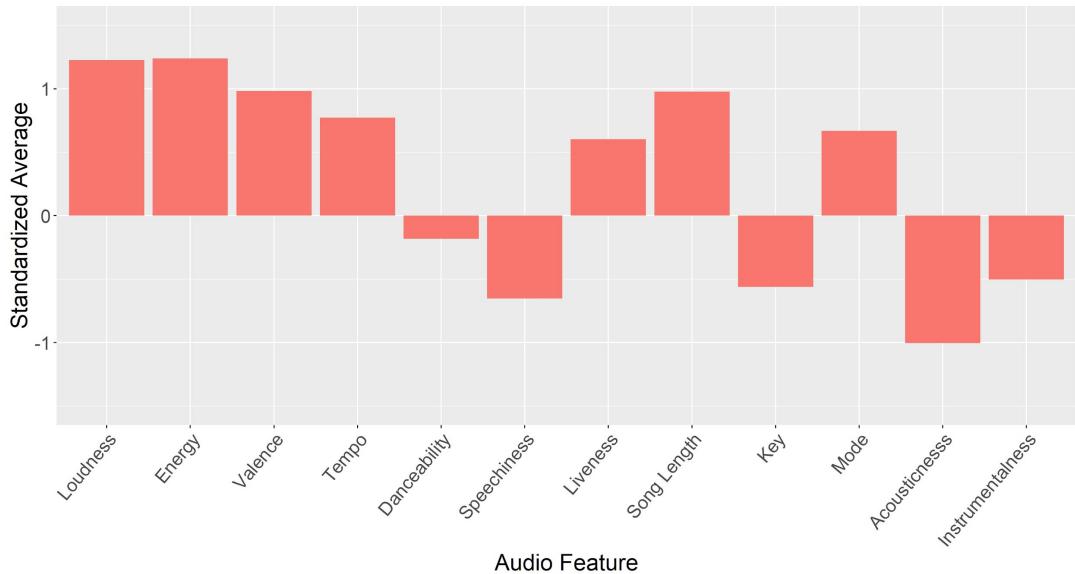


# Clusters





# Turn Up

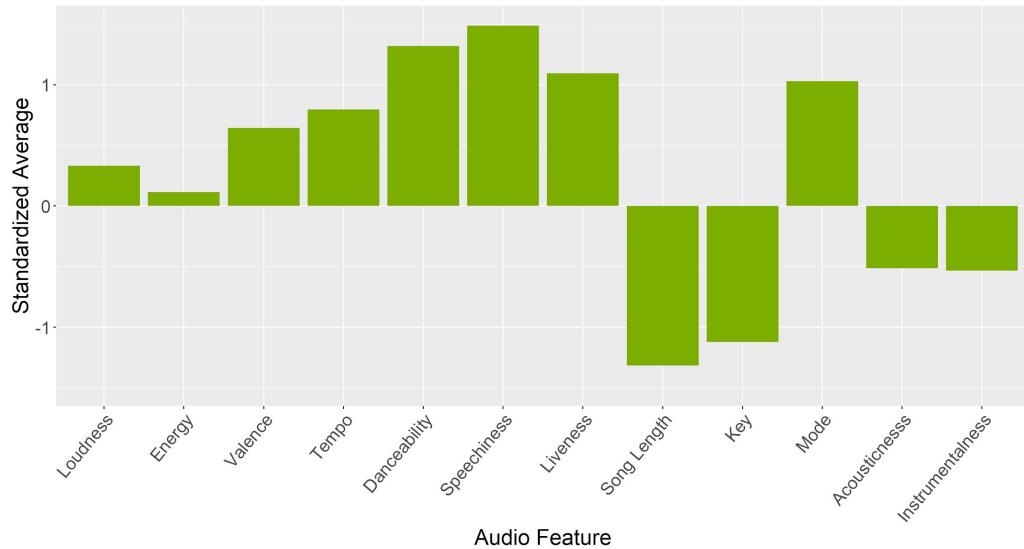


- High Loudness, Energy, Valence
- Low Acousticness, Speechiness, Instrumentalness
- Key members: Taylor Swift, Rage Against the Machine, Flume





# Rap

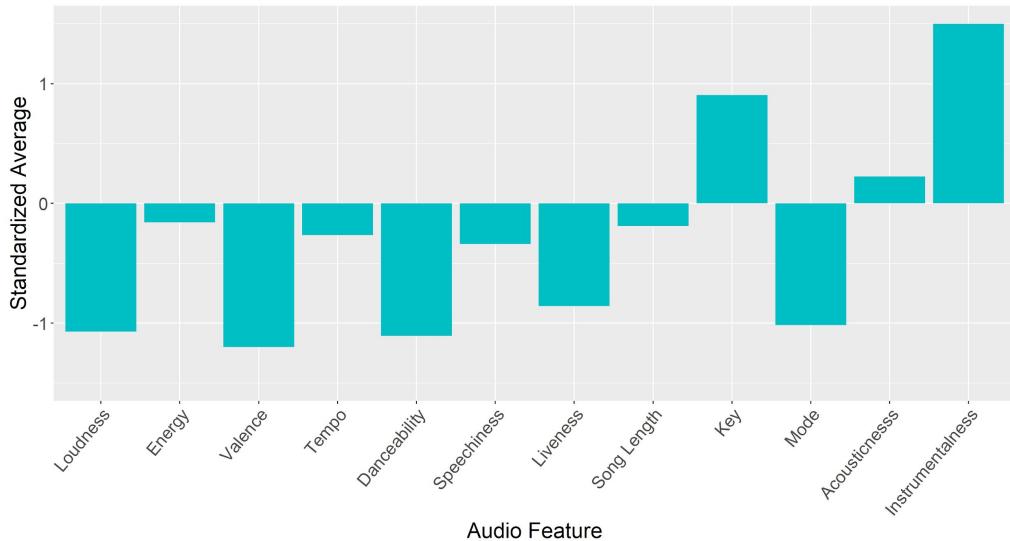


- High Speechiness, Danceability, Liveness
- Low Song Length and Key
- Key members: Kendrick Lamar, J. Cole, Wiz Khalifa





# Instrumental

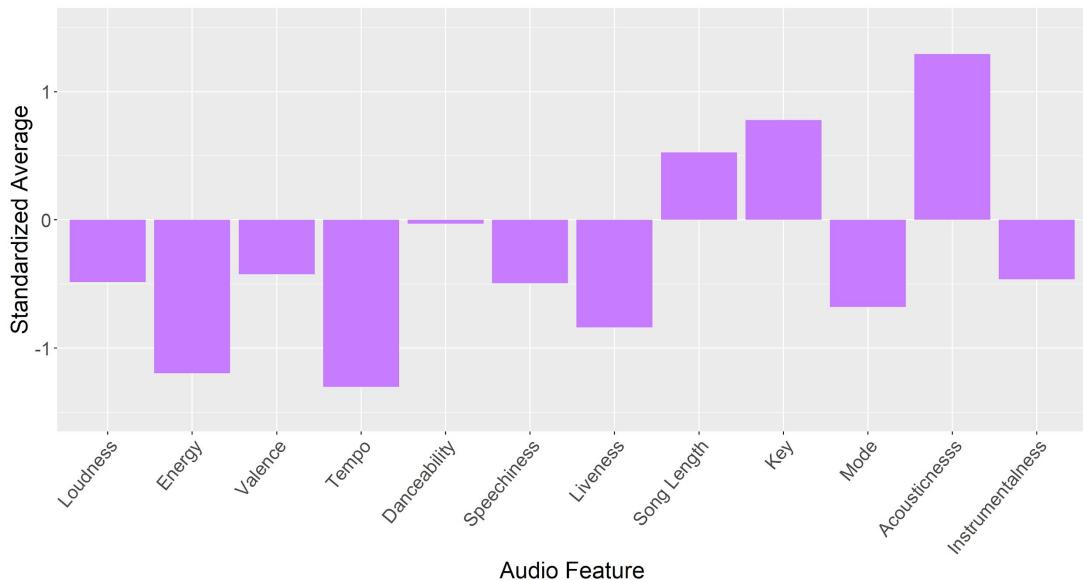


- Very high Instrumentalness and key
- Low everything else
- Key members: Justice, Derby, Hermitude





# Turn Down



- High Acousticness, Key, Song Length
- Low everything else
- Key members: Frank Ocean, Daniel Caesar, SZA



---

# Main Takeaways

- The best way to differentiate my favorite artists is by separating the loud artists from the quiet artists
- The second-best way to differentiate my favorite artists is by separating the instrumental artists from the artists that use vocals

Turn up

Loud  
Energetic  
Happy

Rap

Lyrical  
Danceable

Instrumental

Instrumental

Turn down

Acoustic  
Slow  
Soft



---

# The Data

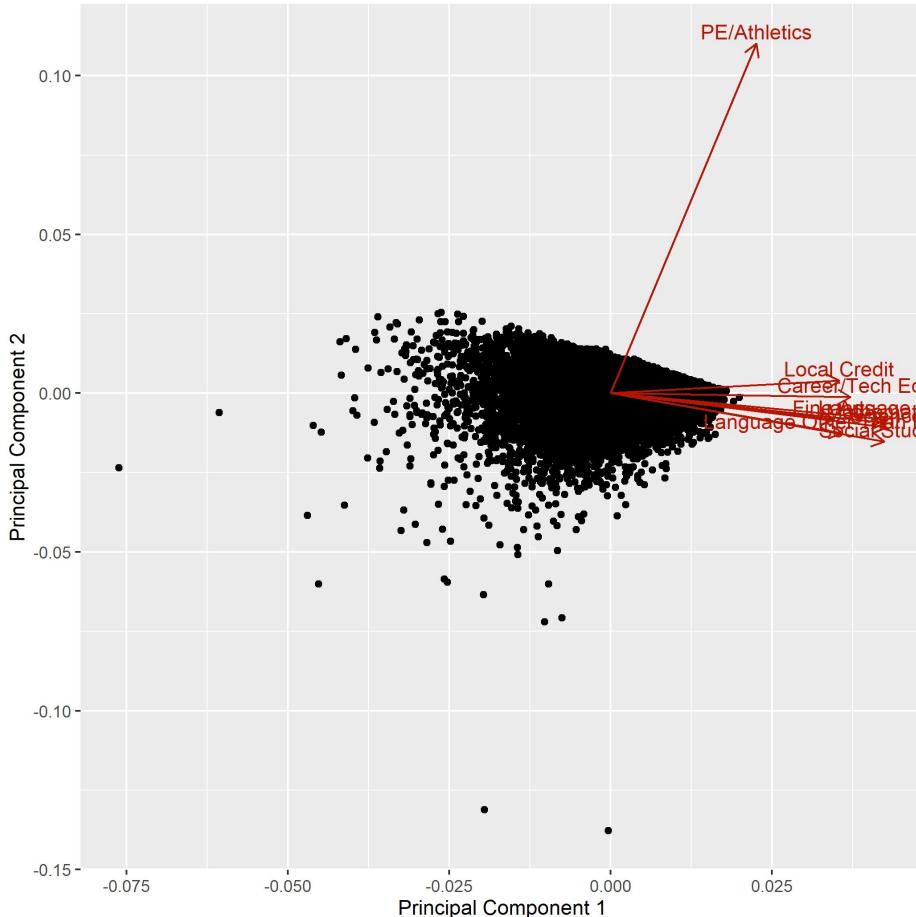
One data point - One student from the Aldine School District

Variables - A student's average grade in these subjects:

- Science
- Language Arts
- Mathematics
- Social Studies
- Career/Tech Ed
- Language Other Than English
- Local Credit
- Fine Arts
- PE/Athletics

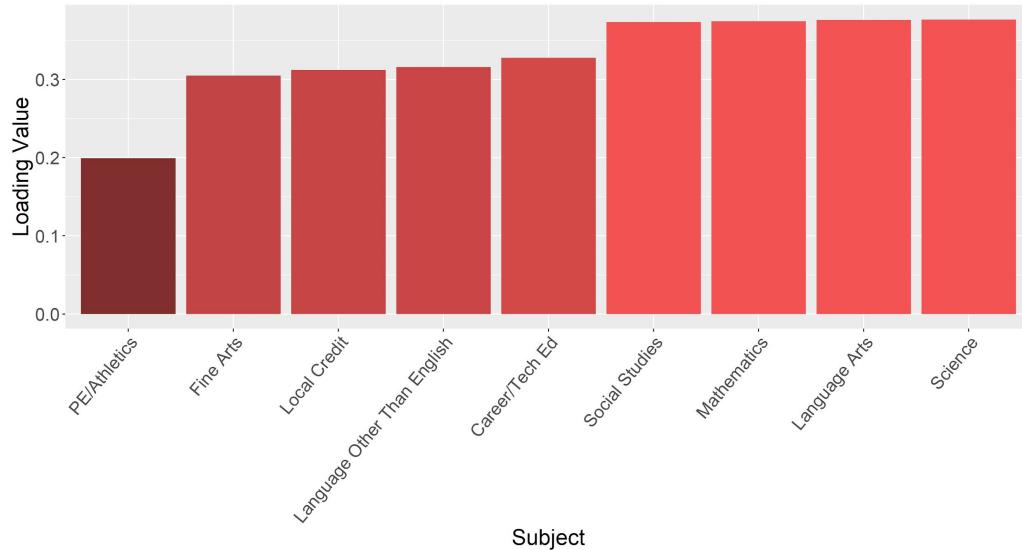


# Biplot





# PC1 Loadings

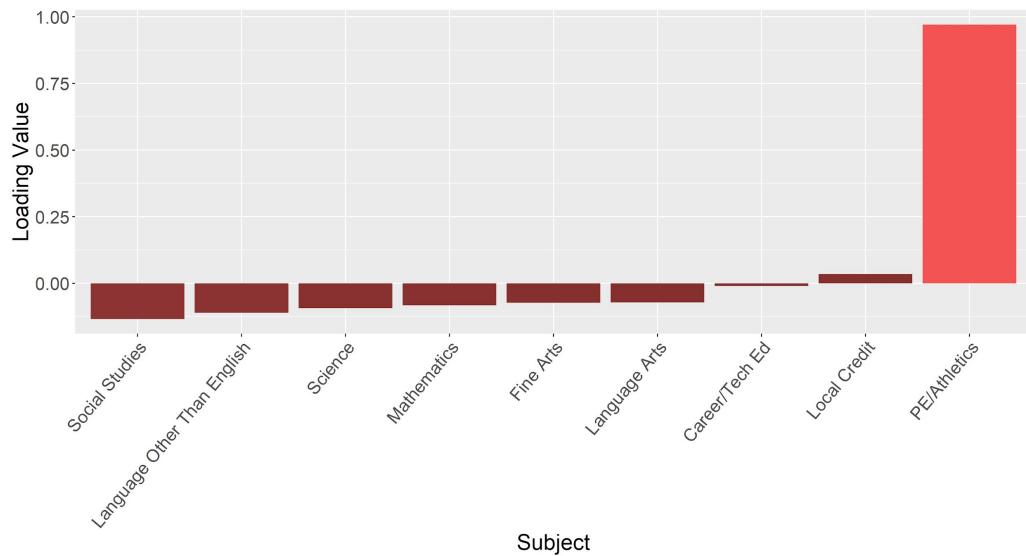


- High everything
- Separates students who get good grades from students who get bad grades





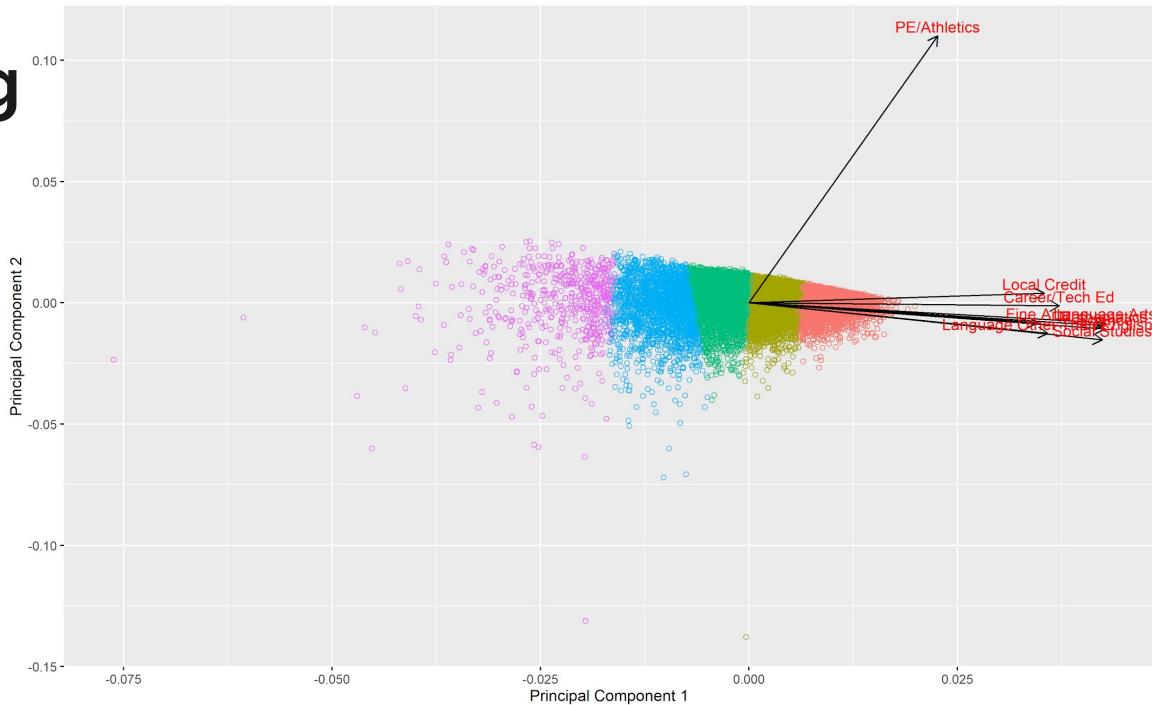
## PC2 Loadings



- High PE/Athletics
- Low everything else
- Separates students that do well in PE/Athletics from those that do not do well in PE/Athletics



# Clustering



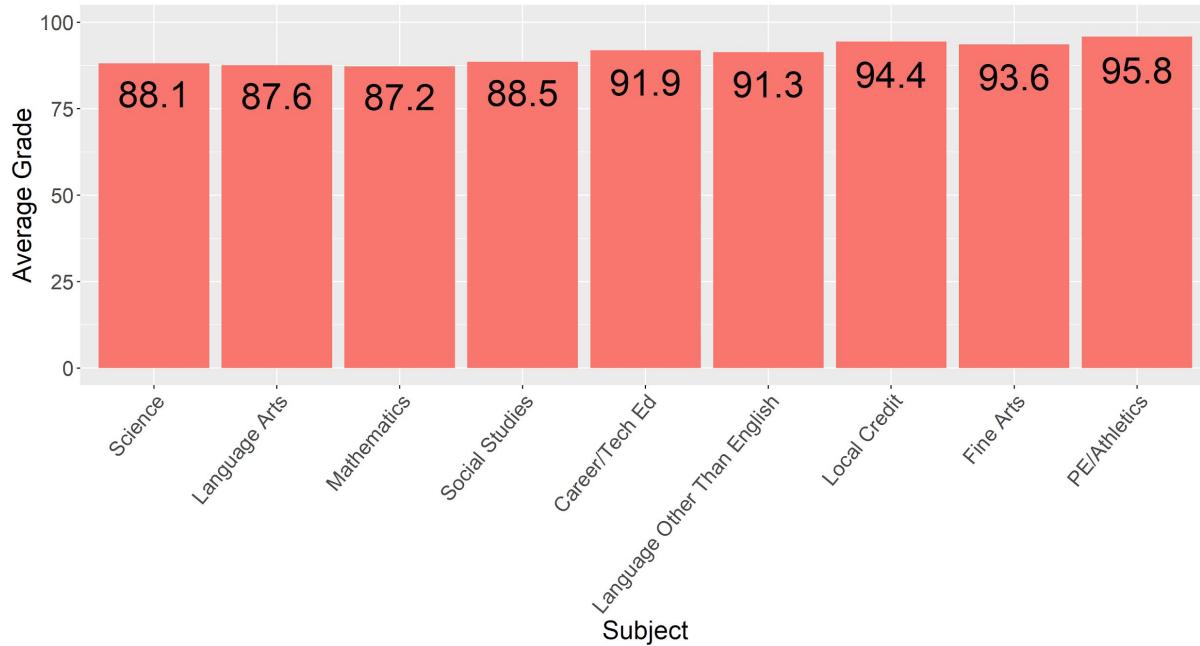
Cluster

- A/High B Student
- Low B Student
- High C Student
- Low C Student
- Low D/F Students



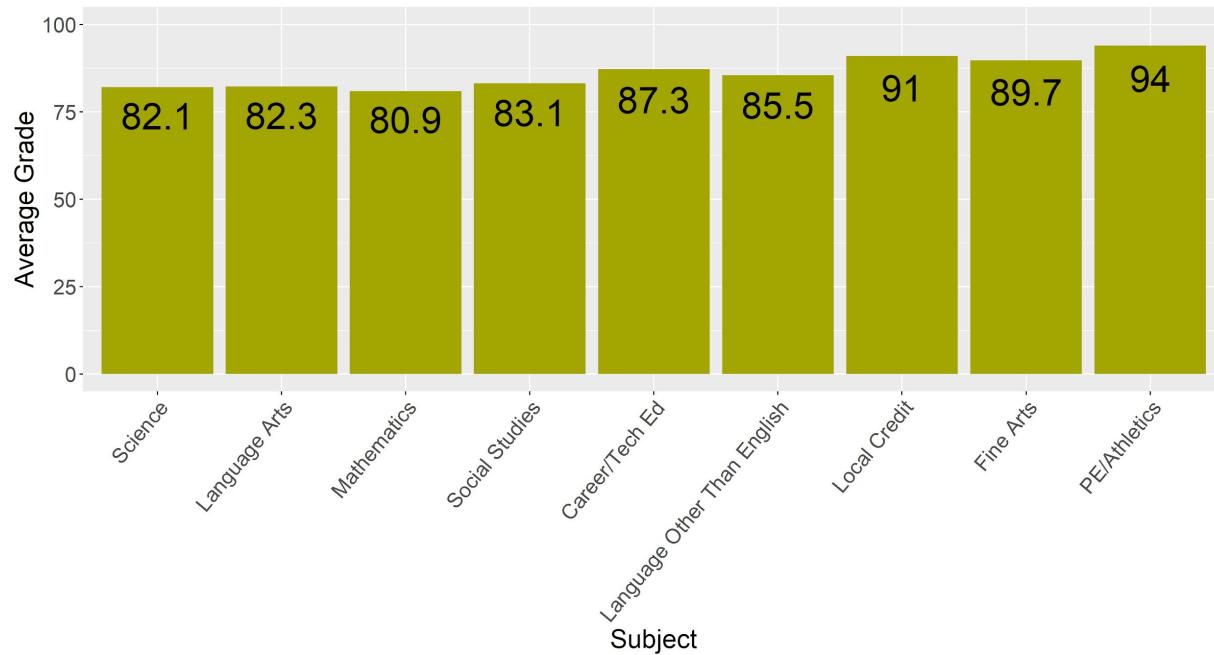


# A/High B Students



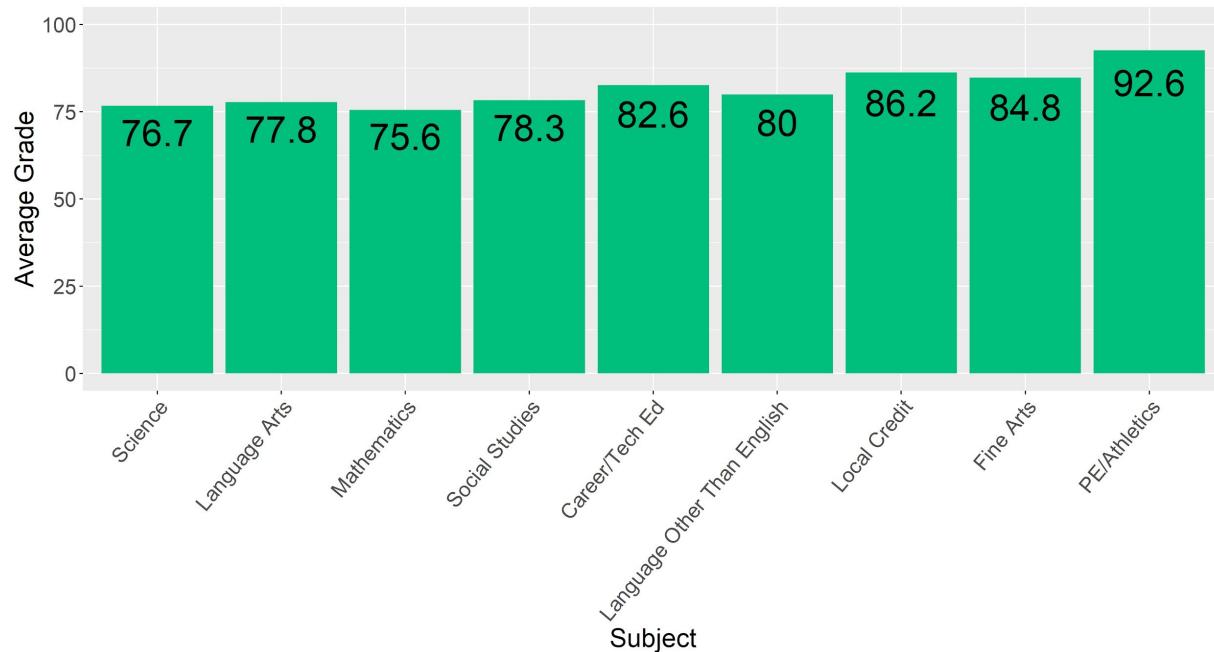


## Low B Students



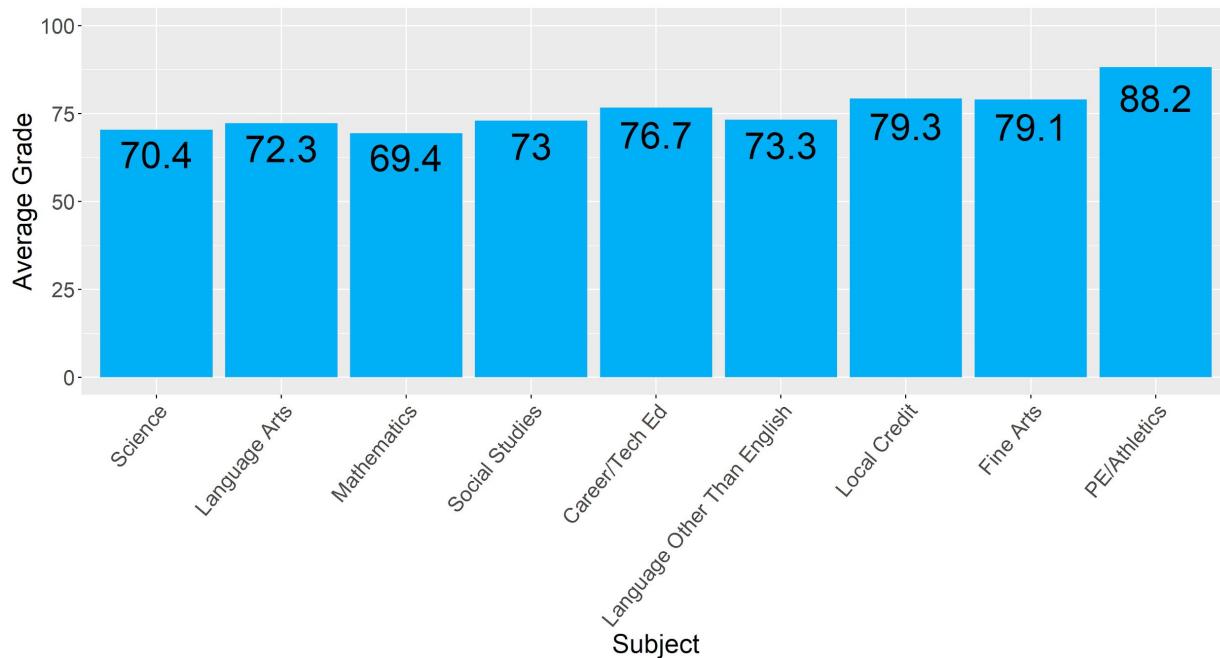


# High C Students



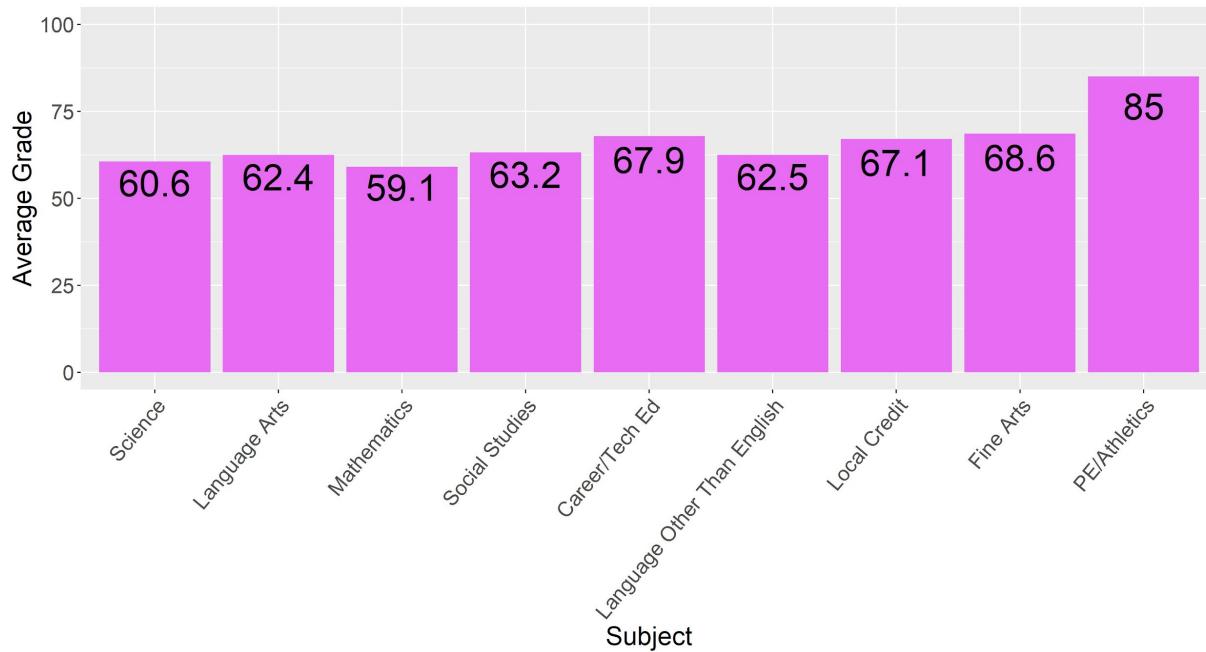


## Low C Students





# Low D/F Students



---

# Main Takeaways

- The best way to differentiate Aldine students is by separating students that get higher grades from students that get lower grades
- The second-best way to differentiate Aldine students is by separating students that do well in PE/Athletics and students that do not
- Student grades form a spectrum, rather than groups

Low Grades



High Grades



---

# Conclusion

- Clustering is a great way to understand and visualize the structure of a big set of data
- One way to cluster data points is by using Principal Components Analysis combined with K-means clustering
  - PCA emphasizes variation in the data, which makes natural groupings more apparent
  - K-means clustering decides what the groups are
- With a better understanding of our data set, our future analysis and decision-making is smarter and more informed

