

Business Intelligence

Prediction the best shooters in the NBA

Instructor: Prof. Dr. Adem Karahoca
Student Ahmet Hacıoğlu
Student ID: 041801123



The main purpose of this project...



- *Predict the best shooters using key shooting statistics from basketball players.*



Why did I choose this topic?



- *Real-world application*
- *Interest-Relevance-Motivation*
- *A Huge Industry*
- *Valuable data analysis and modeling skills*



Attributes

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter

Choose Normalize -S 1.0 -T 0.0 Apply Stop

Current relation

Relation: NBA-Stats-202122-All-Player-Statistics-in-on... Attributes: 29
Instances: 217 Sum of weights: 217

Attributes

All None Invert Pattern

No.	Name	Label	Count	Weight
1	<input type="checkbox"/> FULL NAME	Precious Achiuwa	1	1
2	<input type="checkbox"/> TEAM	Steven Adams	1	1
3	<input type="checkbox"/> POS	Bam Adebayo	1	1
4	<input type="checkbox"/> AGE	Nickeil Alexander-W...	1	1
5	<input type="checkbox"/> GP	Grayson Allen	1	1
6	<input type="checkbox"/> MPG	Jose Alvarado	1	1
7	<input type="checkbox"/> MIN	Kyle Anderson	1	1
8	<input type="checkbox"/> USG	Giannis Antetokoun...	1	1
9	<input type="checkbox"/> TO	Thanasis Antetokoun...	1	1
10	<input type="checkbox"/> FTA	OG Anunoby	1	1
11	<input type="checkbox"/> FT%			
12	<input type="checkbox"/> 2PA			
13	<input type="checkbox"/> 2P%			
14	<input type="checkbox"/> 3PA			
15	<input type="checkbox"/> 3P%			
16	<input type="checkbox"/> eFG			
17	<input type="checkbox"/> TS			
18	<input type="checkbox"/> PPG			
19	<input type="checkbox"/> RPG			

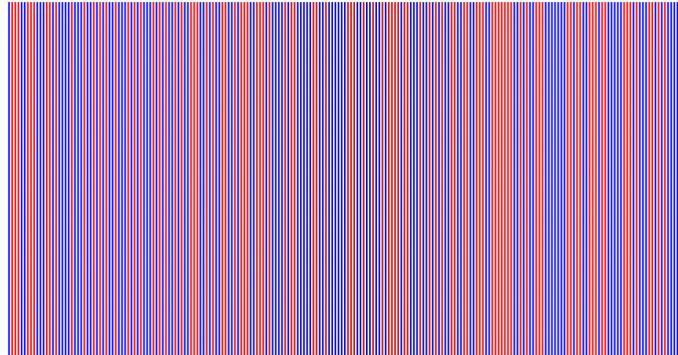
Remove

Status OK Log x 0

Selected attribute

Name: FULL NAME Type: Nominal
Missing: 0 (0%) Distinct: 217 Unique: 217 (100%)

Class: Results (Nom) Visualize All



The Important Attributes for Selecting Good Shooter



1. **Field Goal Percentage (FG%):** This measures the percentage of field goals a player makes compared to the number of attempts. This can give you a good sense of how accurate a player is when shooting.
2. **Three-Point Percentage (3P%):** This measures the percentage of three-point shots a player makes compared to the number of attempts. This can be particularly important if you're looking for players who are good at shooting from beyond the arc.
3. **Effective Field Goal Percentage (eFG%):** This is a more advanced metric that takes into account the fact that three-point shots are worth more than two-point shots. It adjusts a player's field goal percentage based on how many three-pointers they make.
4. **True Shooting Percentage (TS%):** This is another advanced metric that takes into account free throw shooting as well as two-point and three-point shooting. It gives you a sense of how efficient a player is at scoring.
5. **Usage Rate (USG%):** This measures how often a player is involved in their team's offensive plays. A player with a high usage rate may be more likely to take a lot of shots, but may also be more prone to turnovers.

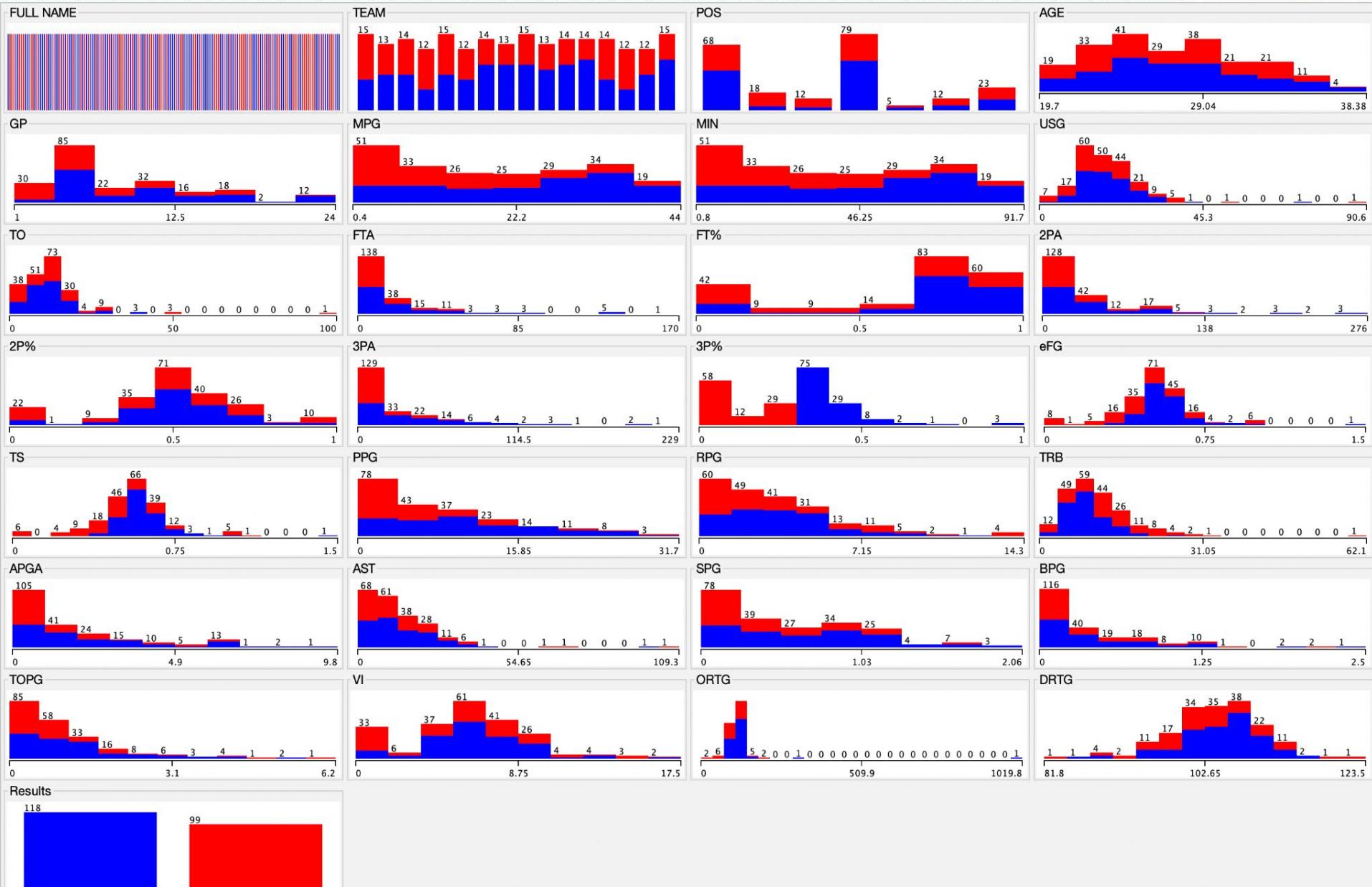
DataSet



- *The data set comes from 2021-2022 NBA play-offs.*
- *The data set have 29 different attributes and 217 instances.*



Visualize all attributes



Filter Choose

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter

Choose Normalize -S 1.0 -T 0.0 Apply Stop

Current relation

Relation: NBA-Stats-202122-All-Player-Statistics-in-on... Attributes: 29
Instances: 217 Sum of weights: 217

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> FULL NAME
2	<input type="checkbox"/> TEAM
3	<input type="checkbox"/> POS
4	<input type="checkbox"/> AGE
5	<input type="checkbox"/> GP
6	<input type="checkbox"/> MPG
7	<input type="checkbox"/> MIN
8	<input type="checkbox"/> USG
9	<input type="checkbox"/> TO
10	<input type="checkbox"/> FTA
11	<input type="checkbox"/> FT%
12	<input type="checkbox"/> 2PA
13	<input type="checkbox"/> 2P%
14	<input type="checkbox"/> 3PA
15	<input type="checkbox"/> 3P%
16	<input type="checkbox"/> eFG
17	<input type="checkbox"/> TS
18	<input type="checkbox"/> PPG
19	<input type="checkbox"/> RPG

Remove

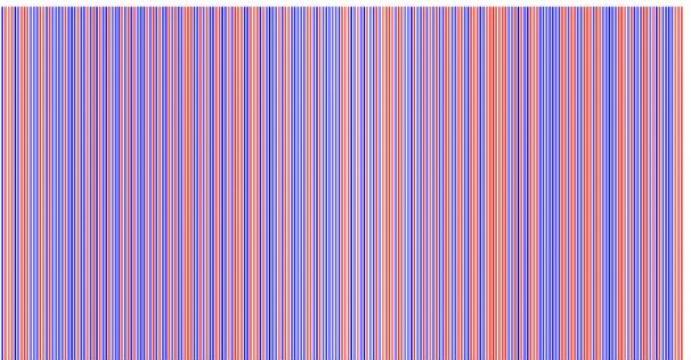
Status OK Log x 0

Selected attribute

Name: FULL NAME Type: Nominal
Missing: 0 (0%) Distinct: 217 Unique: 217 (100%)

No.	Label	Count	Weight
1	Precious Achiuwa	1	1
2	Steven Adams	1	1
3	Bam Adebayo	1	1
4	Nickeil Alexander-W...	1	1
5	Grayson Allen	1	1
6	Jose Alvarado	1	1
7	Kyle Anderson	1	1
8	Giannis Antetokoun...	1	1
9	Thanasis Antetokoun...	1	1
10	OG Anunoby	1	1

Class: Results (Nom) Visualize All



Classify - Which Algorithms?



1. *Naive Bayes*
2. *J48*
3. *Random Forest*
4. *SMO*
5. *K-nearest neighbor*
6. *Random Tree*



Naive Bayes



Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose NaiveBayes

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) Results

Start Stop

Result list (right-click for options)

05:14:34 - bayes.NaiveBayes

Classifier output

	0	1
precision	6.5372	6.5372

DRTG

	0	1
mean	104.7318	102.9194
std. dev.	5.3272	6.768
weight sum	108	72
precision	0.3446	0.3446

Time taken to build model: 0.01 seconds

==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances 174 80.1843 %
Incorrectly Classified Instances 43 19.8157 %
Kappa statistic 0.6009
Mean absolute error 0.2181
Root mean squared error 0.4266
Relative absolute error 43.9456 %
Root relative squared error 85.6421 %
Total Number of Instances 217

==== Detailed Accuracy By Class ====

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0	0.814	0.212	0.821	0.814	0.817	0.601	0.847	0.854	1
1	0.788	0.186	0.780	0.788	0.784	0.601	0.847	0.803	0
Weighted Avg.	0.802	0.200	0.802	0.802	0.802	0.601	0.847	0.831	

==== Confusion Matrix ====

	a	b	<-- classified as
a	96	22	a = 1
b	21	78	b = 0

Status OK

Log x 0

J48

x 0

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **J48 -C 0.25 -M 2**

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) Results

Start Stop

Result list (right-click for options)

05:14:34 - bayes.NaiveBayes
05:15:08 - trees.J48

Classifier output

J48 pruned tree

3P% <= 0.3: 0 (99.0)
3P% > 0.3: 1 (118.0)

Number of Leaves : 2

Size of the tree : 3

Time taken to build model: 0.01 seconds

==== Stratified cross-validation ====
==== Summary ====

Correctly Classified Instances	217	100	%
Incorrectly Classified Instances	0	0	%
Kappa statistic	1		
Mean absolute error	0		
Root mean squared error	0		
Relative absolute error	0	%	
Root relative squared error	0	%	
Total Number of Instances	217		

==== Detailed Accuracy By Class ====

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
1	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	1
0	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	0
Weighted Avg.	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	

==== Confusion Matrix ====

a	b	<-- classified as
118	0	a = 1
0	99	b = 0

Status OK Log

RandomForest

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1**

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) Results

Start Stop

Result list (right-click to copy): Starts the classification
05:14:34 - bayes.NaiveBayes
05:15:08 - trees.J48
05:15:33 - trees.RandomForest

Classifier output

Test mode: 10-fold cross-validation

==== Classifier model (full training set) ====
RandomForest
Bagging with 100 iterations and base learner
weka.classifiers.trees.RandomTree -K 0 -M 1.0 -V 0.001 -S 1 -do-not-check-capabilities

Time taken to build model: 0.05 seconds

==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances 216 99.5392 %
Incorrectly Classified Instances 1 0.4608 %
Kappa statistic 0.9907
Mean absolute error 0.2598
Root mean squared error 0.2741
Relative absolute error 52.3527 %
Root relative squared error 55.0236 %
Total Number of Instances 217

==== Detailed Accuracy By Class ====

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0	0.992	0.000	1.000	0.992	0.996	0.991	1.000	1.000	1
1	1.000	0.008	0.990	1.000	0.995	0.991	1.000	1.000	0
Weighted Avg.	0.995	0.004	0.995	0.995	0.995	0.991	1.000	1.000	

==== Confusion Matrix ====
a b <-- classified as
117 1 | a = 1
0 99 | b = 0

Status OK Log x 0

SMO

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **SMO** -C 1.0 -L 0.001 -P 1.0E-12 -N 0 -V -1 -W 1 -K "weka.classifiers.functions.supportVector.PolyKernel -E 1.0 -C 250007" -calibrator "weka.classifiers.function

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

Classifier output

```
+ 0.463 * (normalized) TOPG
+ 0.4043 * (normalized) VI
+ -0.1876 * (normalized) ORTG
+ 0.0757 * (normalized) DRTG
+ 1.69
```

Number of kernel evaluations: 19065 (94.207% cached)

Time taken to build model: 0.08 seconds

==== Stratified cross-validation ====
==== Summary ====

Correctly Classified Instances	191	88.0184 %
Incorrectly Classified Instances	26	11.9816 %
Kappa statistic	0.7561	
Mean absolute error	0.1198	
Root mean squared error	0.3461	
Relative absolute error	24.1448 %	
Root relative squared error	69.4909 %	
Total Number of Instances	217	

==== Detailed Accuracy By Class ====

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.941	0.192	0.854	0.941	0.895	0.761	0.874	0.835	1	
0.808	0.059	0.920	0.808	0.860	0.761	0.874	0.831	0	
Weighted Avg.	0.880	0.131	0.884	0.880	0.879	0.761	0.874	0.833	

==== Confusion Matrix ====

a	b	<-- classified as
111	7	a = 1
19	80	b = 0

Status OK Log x 0

K-Nearest Neighbors

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""**

Test options

Use training set
 Supplied test set
 Cross-validation Folds
 Percentage split %

Classifier output

DRTG
Results
Test mode: 10-fold cross-validation

==== Classifier model (full training set) ====
IB1 instance-based classifier
using 1 nearest neighbour(s) for classification

Time taken to build model: 0 seconds

==== Stratified cross-validation ====
==== Summary ===

	Correctly Classified Instances	139	64.0553 %
Incorrectly Classified Instances	78	35.9447 %	
Kappa statistic	0.2648		
Mean absolute error	0.3609		
Root mean squared error	0.5965		
Relative absolute error	72.7214 %		
Root relative squared error	119.7544 %		
Total Number of Instances	217		

==== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.746	0.485	0.647	0.746	0.693	0.269	0.628	0.622	1
	0.515	0.254	0.630	0.515	0.567	0.269	0.628	0.554	0
Weighted Avg.	0.641	0.380	0.639	0.641	0.635	0.269	0.628	0.591	

==== Confusion Matrix ===

		a b <-- classified as	
		a = 1	
88	30		a = 1
48	51		b = 0

Status OK Log x 0

Random Tree

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **RandomTree -K 0 -M 1.0 -V 0.001 -S 1**

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) Results

Start Stop

Result list (right-click for options)

- 05:14:34 - bayes.NaiveBayes
- 05:15:08 - trees.J48
- 05:15:33 - trees.RandomForest
- 05:16:26 - functions.SMO
- 05:16:55 - lazy.IBk
- 05:17:15 - trees.RandomTree**

Classifier output

```
FULL NAME = Robert Williams III : 0 (1/0)
FULL NAME = Grant Williams : 1 (1/0)
FULL NAME = Patrick Williams : 1 (1/0)
FULL NAME = Zaire Williams : 1 (1/0)
FULL NAME = Delon Wright : 1 (1/0)
FULL NAME = Thaddeus Young : 0 (1/0)
FULL NAME = Trae Young : 0 (1/0)
FULL NAME = Omer Yurtseven : 0 (1/0)
```

Size of the tree : 218

Time taken to build model: 0 seconds

== Stratified cross-validation ==

== Summary ==

	Correctly Classified Instances	176	81.106 %
Incorrectly Classified Instances	41	18.894 %	
Kappa statistic	0.6113		
Mean absolute error	0.2166		
Root mean squared error	0.3856		
Relative absolute error	43.6539 %		
Root relative squared error	77.4035 %		
Total Number of Instances	217		

== Detailed Accuracy By Class ==

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.932	0.333	0.769	0.932	0.843	0.629	0.865	0.854	1
	0.667	0.068	0.892	0.667	0.763	0.629	0.864	0.813	0
Weighted Avg.	0.811	0.212	0.825	0.811	0.806	0.629	0.865	0.835	

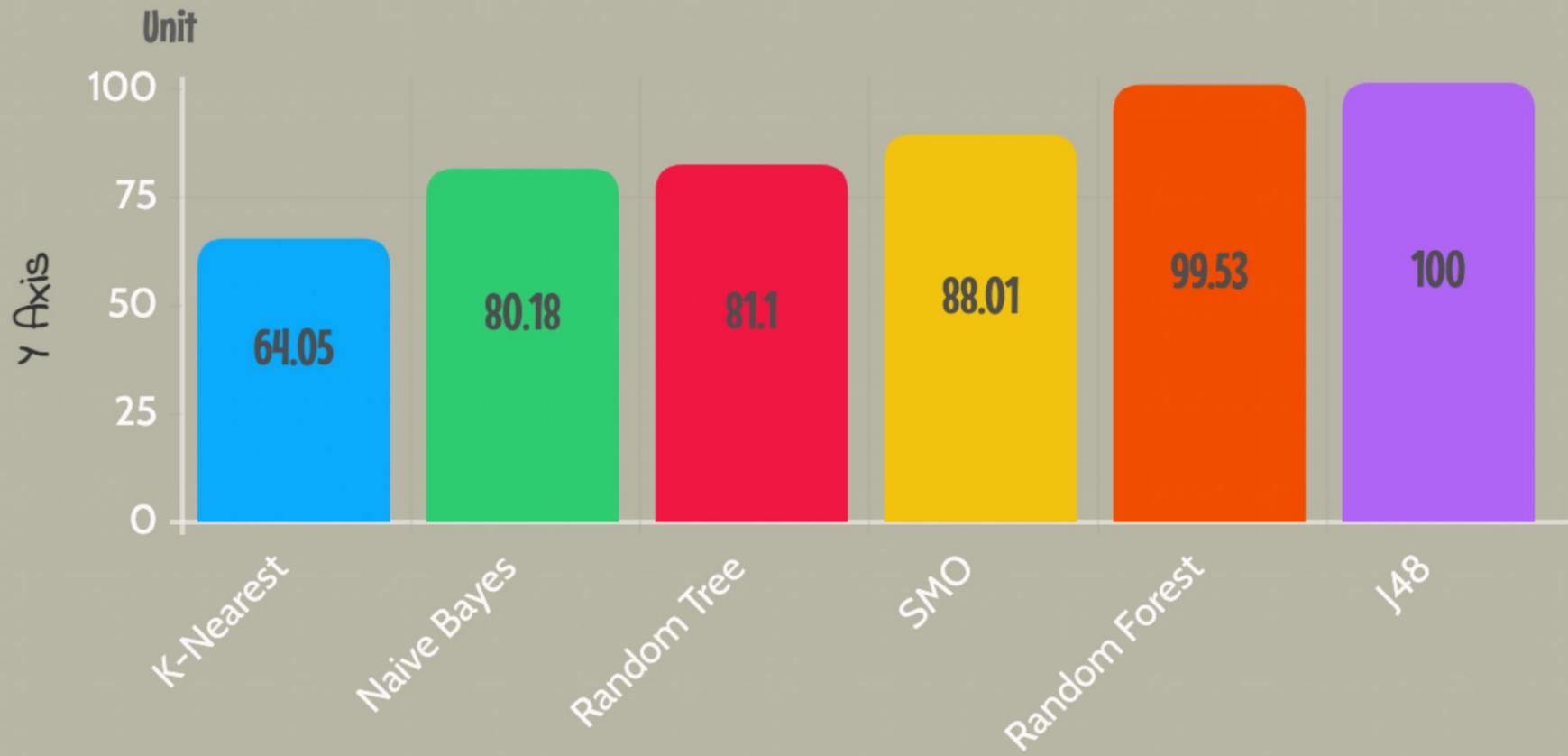
== Confusion Matrix ==

```
a    b    <-- classified as
110   8 |   a = 1
 33   66 |   b = 0
```

Status OK

Log x 0

Compare Classify Algorithms



- *J48 and Random Forest algorithms accuracy rates are perfect.*

Cluster - Which Algorithms?



1. *Simple k-means*
2. *Make Density Based Clusterer*
3. *Farthest First*
4. *Filtered Clusterer*



Simple K-means

Weka Explorer

Preprocess Classify Cluster **Cluster** Associate Select attributes Visualize

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -l 5

Cluster mode

Use training set

Supplied test set Set...

Percentage split % 66

Classes to clusters evaluation

(Nom) Results

Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

05:29:15 - SimpleKMeans

Clusterer output

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster#
	(217.0)	0
	(124.0)	1
FULL NAME	Precious Achiuwa	Steven Adams
TEAM	Tor G	Atl G
POS	27.4306	27.8165
AGE	8.7143	10.379
GP	19.4304	25.1508
MPG	40.4733	52.3839
MIN	18.935	19.2194
USG	12.7868	12.2323
T0	18.1014	25.9516
FTA	0.6232	0.7213
FT%	39.4147	56.1371
2PA	0.495	0.5059
2P%	27.7281	41.7339
3PA	0.2733	0.4007
3P%	0.5212	0.5665
eFG	0.5538	0.6001
TS	8.4571	11.7339
PPG	3.4037	3.9726
RPG	10.7829	9.2194
TRB	1.8286	2.529
APGA	14.0111	15.3694
AST	0.5788	0.7297
SPG	0.3585	0.4381
BPG	1.0796	1.4102
TOPG	6.0396	6.8081
VI	117.5078	116.5578
ORTG	103.9956	104.6311
DRTG	1	1
Results		0

Status OK Log x 0

Make Density Based Clusterer

Weka Explorer

Preprocess Classify Cluster **Associate** Select attributes Visualize

Clusterer

Choose **MakeDensityBasedClusterer** -M 1.0E-6 -W weka.clusterers.SimpleKMeans -- -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25

Cluster mode

Use training set
 Supplied test set Set...
 Percentage split % 66
 Classes to clusters evaluation
(Nom) Results ▾
 Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

05:29:15 - SimpleKMeans
05:30:15 - MakeDensityBasedClusterer

Clusterer output

```
Normal Distribution. Mean = 4.0882 StdDev = 4.0167
Attribute: RPG
Normal Distribution. Mean = 2.6452 StdDev = 2.7641
Attribute: TRB
Normal Distribution. Mean = 12.8677 StdDev = 8.6035
Attribute: APGA
Normal Distribution. Mean = 0.8946 StdDev = 1.1594
Attribute: AST
Normal Distribution. Mean = 12.2 StdDev = 15.608
Attribute: SPG
Normal Distribution. Mean = 0.3777 StdDev = 0.4117
Attribute: BPG
Normal Distribution. Mean = 0.2524 StdDev = 0.4323
Attribute: TOPG
Normal Distribution. Mean = 0.6388 StdDev = 0.8072
Attribute: VI
Normal Distribution. Mean = 5.0151 StdDev = 3.8952
Attribute: ORTG
Normal Distribution. Mean = 118.7745 StdDev = 98.1812
Attribute: DRTG
Normal Distribution. Mean = 103.1481 StdDev = 5.9386
Attribute: Results
Discrete Estimator. Counts = 4 91 (Total = 95)

Time taken to build model (full training data) : 0.01 seconds
== Model and evaluation on training set ==
Clustered Instances
0 122 ( 56%)
1 95 ( 44%)

Log likelihood: -71.88017
```

Status OK Log x 0

Farthest First

Weka Explorer

Preprocess Classify Cluster **Cluster** Associate Select attributes Visualize

Clusterer
Choose **FarthestFirst -N 2 -S 1**

Cluster mode
 Use training set
 Supplied test set Set...
 Percentage split % 66
 Classes to clusters evaluation
(Nom) Results
 Store clusters for visualization

Ignore attributes

Start **Stop**

Result list (right-click for options)
05:29:15 - SimpleKMeans
05:30:15 - MakeDensityBasedClusterer
05:31:15 - FarthestFirst

Clusterer output

```
TRB
APGA
AST
SPG
BPG
TOPG
VI
ORTG
DRTG
Results
Test mode: evaluate on training data
```

==== Clustering model (full training set) ====

FarthestFirst
=====

Cluster centroids:

Cluster 0 Khris Middleton Mil F 30.85 2.0 35.8 74.5 22.0 30.3 3.0 1.0 10.0 0.4 14.0 0.429 0.5
Cluster 1 Willy Hernangomez Nor C-F 28.06 1.0 1.9 4.0 90.6 0.0 0.0 0.0 3.0 0.333 1.0 0.0 0.25

Time taken to build model (full training data) : 0.01 seconds

==== Model and evaluation on training set ====

Clustered Instances

0	139	(64%)
1	78	(36%)

Status OK Log x 0

Filtered Clusterer

Weka Explorer

Preprocess Classify Cluster **Cluster** Select attributes Visualize

Clusterer

Choose **FilteredClusterer** -F "weka.filters.AllFilter" -W weka.clusterers.SimpleKMeans -- -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.

Cluster mode

Use training set

Supplied test set % 66

Percentage split % 66

Classes to clusters evaluation

(Nom) Results

Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

05:29:15 - SimpleKMeans
05:30:15 - MakeDensityBasedClusterer
05:31:15 - FarthestFirst
05:31:45 - FilteredClusterer

Clusterer output

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster# 0	1
	(217.0)	(124.0)	(93.0)
FULL NAME	Precious Achiuwa	Precious Achiuwa	Steven Adams
TEAM	Tor G	Gol G	Atl G
POS	27.4306	27.8165	26.9161
AGE	8.7143	10.379	6.4946
GP	19.4304	25.1508	11.8032
MPG	40.4733	52.3839	24.5925
MIN	18.935	19.2194	18.5559
USG	12.7868	12.2323	13.5262
T0	18.1014	25.9516	7.6344
FTA	0.6232	0.7213	0.4925
FT%	39.4147	56.1371	17.1183
2PA	0.495	0.5059	0.4805
2P%	27.7281	41.7339	9.0538
3PA	0.2733	0.4007	0.1035
3P%	0.5212	0.5665	0.4609
eFG	0.5538	0.6001	0.4921
TS	8.4571	11.7339	4.0882
PPG	3.4037	3.9726	2.6452
RPG	10.7829	9.2194	12.8677
TRB	1.8286	2.529	0.8946
APGA	14.0111	15.3694	12.2
AST	0.5788	0.7297	0.3777
SPG	0.3585	0.4381	0.2524
BPG	1.0796	1.4102	0.6388
TOPG	6.0396	6.8081	5.0151
VI	117.5078	116.5578	118.7745
ORTG	103.9956	104.6311	103.1481
DRTG			0
Results	1	1	

Status OK Log x 0

Compare Clusters

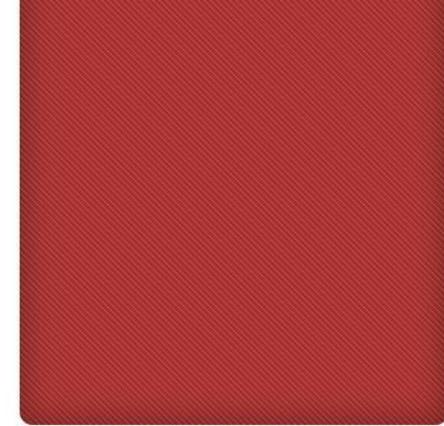
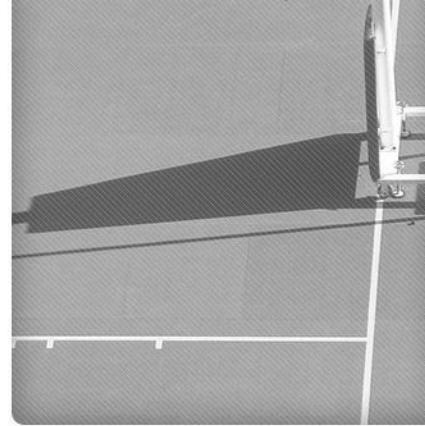
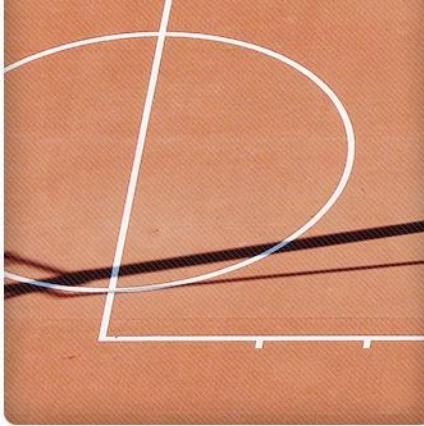


- *The algorithms could not cluster correctly.*
- *Simple k-means clustering and Filtered Clusterer results were same.*
- *The most worst clustering algorithm is Farthest First.*
- *The best clustering algorithm was make density based clustering but this algorithm good with poor results.*



References

1. *Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2017). The WEKA Machine Learning Workbench. Online appendix for "Data Mining: Practical Machine Learning Tools and Techniques". Fifth Edition, Morgan Kaufmann, 2017.* Available at <http://www.cs.waikato.ac.nz/ml/weka/book.html>.
2. *Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009). The WEKA data mining software: an update. ACM SIGKDD Explorations Newsletter, 11(1), 10-18.*
3. *[BilgisayarKavramları]. (2015, May 10). WEKA Eğitim Serisi [Video].Youtube.com.* https://www.youtube.com/watch?v=5s8lgMfH698&list=PLh9ECzBB8tJP8vpIM91W1k_qoGI10ikh1
4. *NBA (n.d.). NBA Advanced Stats.* <https://www.nba.com/stats>



THANK YOU

