
COMS4060A/7056A: Assignment 3

Puseletso Nakana : 2341162
Ntando Ngobese : 2112256
Nthabiseng Thema : 2012016

Department of Computer Science
University of Witwatersrand

1 Data Cleaning

Outlier Analysis: We identified outliers in several categories, with the most significant numbers in:

3

- 4 • Salary (48 outliers)
- 5 • VORP (36 outliers)
- 6 • OWS (27 outliers)
- 7 • TOV% (21 outliers)

8 Dataset Reduction:

9

- 10 Original dataset: 433 rows
- 11 After initial outlier removal: 337 rows
- 12 After efficiency (WS/48) outlier removal: 333 rows

13 Odd Points:

- 14 a) Age: We identified 2 outliers. Given the range (19-37), these likely represent very young rookies
- 15 or veteran players.
- 16
- 17 b) Salary: The high number of outliers (48) suggests a wide disparity in player salaries,
- 18 which is common in the NBA.
- 19
- 20 c) ORP and OWSV: The high number of outliers in these advanced metrics might indicate
- 21 players with exceptional positive or negative impacts.

22 Decision on Outlier Removal:

- 23 a) Age outliers: Age extremes are valid in basketball and can provide valuable insights, so we did not
- 24 remove these outliers.
- 25
- 26 b) Salary outliers: Salary disparities are a reality in the NBA and often correlate with per-
- 27 formance or player status, hence we did not remove these outliers.
- 28
- 29 c) Performance metric outliers (VORP, OWS, BPM, etc.): We kept positive outliers as they
- 30 likely represent star players. We removed extreme negative outliers, as they might represent players
- 31 with very limited playing time.

32
33 d) Efficiency outlier (WS/48): The removal of 4 additional outliers seems reasonable to us,
34 as extreme efficiency values might skew analysis.

35 **Justification:**

36 We recognize that basketball often has legitimate statistical outliers due to varying roles, playing time,
37 and skill levels. We caution that removing too many outliers might eliminate important data points
38 representing star players or unique roles. We note that the current approach reduced the dataset by
39 about 23%, which is significant.

Table 1: Number of Outliers by Statistical Category

Category	Number of Outliers
Age	2
Salary	48
GP	0
MP	0
FG	17
TOV%	21
USG%	11
OVS	27
DWS	1
WS	11
WS/48	12
OBPM	16
DBPM	10
BPM	17
VORP	36

40

Table 2: Summary Statistics of Cleaned Dataset

Statistic	Salary	Age	GP	GS	MP	FG
Count	405.00	405.00	405.00	405.00	405.00	405.00
Mean	9.39M	26.04	53.19	25.32	21.31	3.64
Std	11.10M	4.37	21.55	27.60	8.81	2.41
Min	32.80K	19.00	1.00	0.00	2.80	0.20
25%	2.00M	23.00	39.00	2.00	14.00	1.90
50%	4.44M	25.00	59.00	11.00	20.30	3.00
75%	12.20M	29.00	70.00	53.00	29.00	4.60
Max	48.07M	42.00	83.00	83.00	37.40	11.20

Table 3: Advanced Statistics Summary

Statistic	TOV%	USG%	WS	WS/48	BPM	VORP
Count	405.00	405.00	405.00	405.00	405.00	405.00
Mean	12.28	18.45	2.60	0.09	-0.93	0.61
Std	3.83	5.53	2.53	0.06	2.83	1.22
Min	2.80	7.40	-0.50	-0.06	-12.10	-1.30
25%	9.70	14.40	0.60	0.05	-2.80	-0.10
50%	11.90	17.70	1.90	0.09	-1.10	0.20
75%	14.40	21.10	3.80	0.13	0.60	0.90
Max	35.50	38.80	12.60	0.28	9.20	6.40

Table 4: Age Outliers

Player Name	Team	Age
Andre Iguodala	GSW	39
Udonis Haslem	MIA	42

Table 5: Top 10 Salary Outliers

Player Name	Team	Salary
Stephen Curry	GSW	\$48,070,014
John Wall	LAC	\$47,345,760
Russell Westbrook	LAL/LAC	\$47,080,179
LeBron James	LAL	\$44,474,988
Kevin Durant	BRK/PHO	\$44,119,845
Bradley Beal	WAS	\$43,279,250
Kawhi Leonard	LAC	\$42,492,492
Paul George	LAC	\$42,492,492
Giannis Antetokounmpo	MIL	\$42,492,492
Damian Lillard	POR	\$42,492,492

2 Dimensionality Reduction

Dimensionality reduction in this NBA dataset context is crucial for several important reasons:

Data Complexity: The original dataset has 51 columns containing various player statistics (like Age, Salary, GP, MP, FG%, etc.). Working with such high-dimensional data can be:

- Computationally expensive
- Difficult to visualize
- Prone to the "curse of dimensionality" (data becomes sparse in high dimensions)

Feature Relationships: Many NBA statistics are likely correlated (e.g., points, field goal attempts, and usage rate). Reducing dimensions helps:

- Identify underlying patterns
- Remove redundant information
- Capture the most important variations in the data

Visualization: Reducing to two dimensions allows us to:

- Create meaningful visualizations
- Identify player clusters/similarities
- Spot trends and patterns that aren't visible in higher dimensions

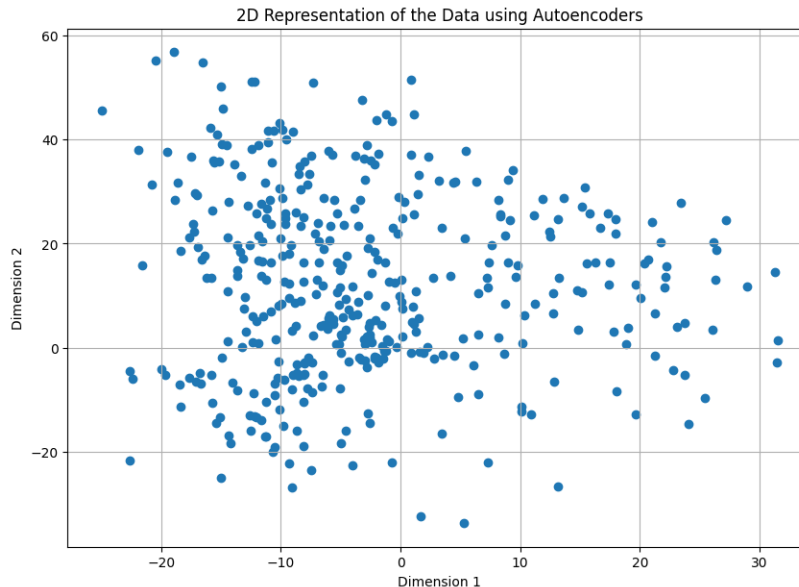
Analysis Efficiency: The different techniques mentioned (Autoencoders, t-SNE, UMAP, etc.) each offer unique advantages for:

- Understanding player similarities
- Identifying playing styles
- Finding patterns in performance metrics
- Creating more efficient player comparison systems

This dimensionality reduction will help us better understand player characteristics and relationships while maintaining the most important information from our original dataset.

Algorithm 1 NBA Player Data Dimensionality Reduction using Autoencoder

- 1: **Data Preparation:**
- 2: Select numerical columns from cleaned dataset
- 3: Normalize data using StandardScaler
- 4: Set input dimension = number of features
- 5: Set encoding dimension = 2
- 6: **Autoencoder Architecture:**
- 7: *Encoder:*
- 8: Input layer (shape = input_dim)
- 9: Dense layer (64 neurons, ReLU activation)
- 10: Dense layer (32 neurons, ReLU activation)
- 11: Dense layer (2 neurons, linear activation)
- 12: *Decoder:*
- 13: Dense layer (32 neurons, ReLU activation)
- 14: Dense layer (64 neurons, ReLU activation)
- 15: Dense layer (input_dim neurons, sigmoid activation)
- 16: **Model Configuration:**
- 17: Compile model with Adam optimizer
- 18: Use Mean Squared Error loss function
- 19: **Training:**
- 20: Train for 100 epochs
- 21: Batch size = 32
- 22: Validation split = 20%
- 23: Enable shuffle
- 24: **Dimensionality Reduction:**
- 25: Use encoder to transform data
- 26: Convert to DataFrame with 2 dimensions
- 27: Visualize reduced data using scatter plot



66 The scatter plot reveals the distribution of data points in two dimensions, allowing us to observe
 67 potential clusters or patterns that may exist within the dataset. The presence of distinct clusters
 68 suggests that certain groups share similar characteristics. The spread of points across the plot
 69 indicates the variance in the dataset, with tightly packed points representing similar observations.
 70

71

72 The scatter plot shows several distinct clusters of points, each concentrated in separate areas. This
73 indicates that there are groups of similar observations in the dataset. There are points which are densely
74 packed in certain areas, potentially overlapping. This suggests that there are many observations
75 with similar characteristics. The points form a linear or near-linear pattern. This indicates a strong
76 correlation between two features represented in the dimensions.

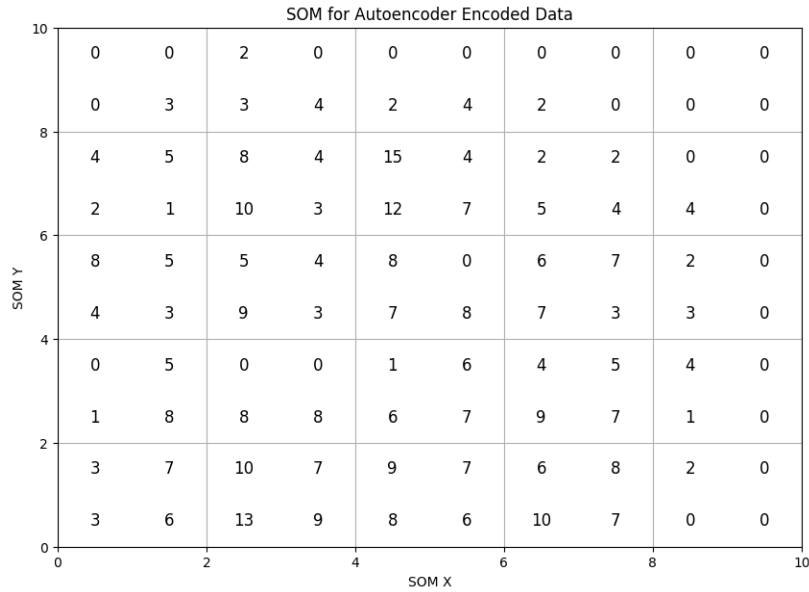
77 2.2 Autoencoders + self-organising maps (SOMs)

Algorithm 2 NBA Player Data Dimensionality Reduction using Autoencoder + SOM

```

1: procedure AUTOENCODERSOMREDUCTION(data_scaled)
2:   Autoencoder Configuration:
3:   encoding_dim  $\leftarrow$  10                                 $\triangleright$  Intermediate dimension
4:   Define Encoder:
5:   input_layer  $\leftarrow$  Input(shape=input_dim)
6:   encoder  $\leftarrow$  Dense(64, activation='relu')(input_layer)
7:   encoder  $\leftarrow$  Dense(32, activation='relu')(encoder)
8:   encoder_output  $\leftarrow$  Dense(encoding_dim, activation='relu')(encoder)
9:   Define Decoder:
10:  decoder  $\leftarrow$  Dense(32, activation='relu')(encoder_output)
11:  decoder  $\leftarrow$  Dense(64, activation='relu')(decoder)
12:  decoder_output  $\leftarrow$  Dense(input_dim, activation='sigmoid')(decoder)
13:  Create Models:
14:  autoencoder  $\leftarrow$  Model(input_layer, decoder_output)
15:  encoder_model  $\leftarrow$  Model(input_layer, encoder_output)
16:  Compile and Train Autoencoder:
17:  autoencoder.compile(optimizer='adam', loss='mse')
18:  autoencoder.fit(
19:    data_scaled,
20:    data_scaled,
21:    epochs=100,
22:    batch_size=32,
23:    validation_split=0.2
24:  )
25:  Encode Data:
26:  encoded_data  $\leftarrow$  encoder_model.predict(data_scaled)
27:  Initialize and Train SOM:
28:  som_size  $\leftarrow$  10                                 $\triangleright$  Size of the SOM grid
29:  som  $\leftarrow$  MiniSom(
30:    som_size,
31:    som_size,
32:    encoding_dim,
33:    sigma=1.0,
34:    learning_rate=0.5
35:  )
36:  som.train(encoded_data, num_iterations=1000)
37:  Get Winning Nodes:
38:  win_map  $\leftarrow$  som.win_map(encoded_data)
39:  Visualize Results:
40:  for x in range(som_size) do
41:    for y in range(som_size) do
42:      PlotNode(x + 0.5, y + 0.5, win_map[x,y])
43:    end for
44:  end for
45:  return som, win_map
46: end procedure

```



78

79 The grid shows an integration a Self-Organizing Map (SOM) following the autoencoder, allowing
 80 us to visualize the encoded data more effectively. The grid is visualized by counting the number of
 81 encoded data points associated with each node, providing insight into the density and distribution of
 82 the encoded representations.

83 The grid shows several nodes with high counts, indicating clusters of similar encoded data points. This
 84 suggests that certain groups in the dataset share similar characteristics. The zeros in the top right
 85 corner suggest that there are combinations of features that are either rare or non-existent. The rest of
 86 the grid having high counts while the top right is empty implies that the data points cluster around
 87 specific characteristics/combinations that are more common. This pattern suggests that there are
 88 certain features that, when combined, do not appear together frequently.

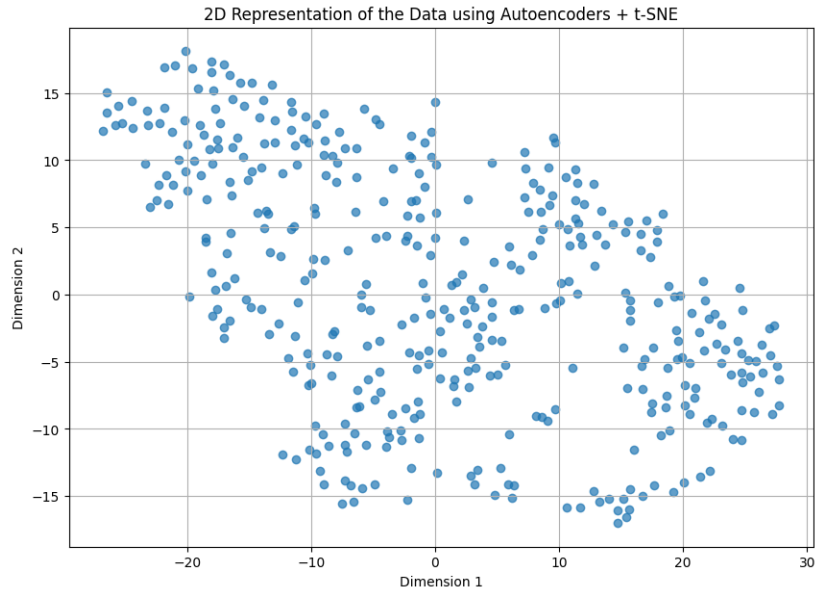
89

Algorithm 3 NBA Player Data Dimensionality Reduction using Autoencoder + t-SNE

```

1: procedure AUTOENCODERTSNEREDUCTION(data_scaled)
2:   Autoencoder Configuration:
3:   encoding_dim  $\leftarrow$  10 ▷ Intermediate dimension
4:   Define Encoder:
5:   input_layer  $\leftarrow$  Input(shape=input_dim)
6:   encoder  $\leftarrow$  Dense(64, activation='relu')(input_layer)
7:   encoder  $\leftarrow$  Dense(32, activation='relu')(encoder)
8:   encoder_output  $\leftarrow$  Dense(encoding_dim, activation='relu')(encoder)
9:   Define Decoder:
10:  decoder  $\leftarrow$  Dense(32, activation='relu')(encoder_output)
11:  decoder  $\leftarrow$  Dense(64, activation='relu')(decoder)
12:  decoder_output  $\leftarrow$  Dense(input_dim, activation='sigmoid')(decoder)
13:  Create Models:
14:  autoencoder  $\leftarrow$  Model(input_layer, decoder_output)
15:  encoder_model  $\leftarrow$  Model(input_layer, encoder_output)
16:  Compile and Train:
17:  autoencoder.compile(optimizer='adam', loss='mse')
18:  autoencoder.fit(
19:    data_scaled,
20:    data_scaled,
21:    epochs=100,
22:    batch_size=32,
23:    validation_split=0.2
24:  )
25:  Encode Data:
26:  encoded_data  $\leftarrow$  encoder_model.predict(data_scaled)
27:  Apply t-SNE:
28:  tsne  $\leftarrow$  TSNE(
29:    n_components=2,
30:    random_state=42,
31:    perplexity=30
32:  )
33:  tsne_results  $\leftarrow$  tsne.fit_transform(encoded_data)
34:  Create Visualization DataFrame:
35:  tsne_df  $\leftarrow$  DataFrame(tsne_results, columns=['Dimension 1', 'Dimension 2'])
36:  Visualize Results:
37:  plt.figure(figsize=(10, 7))
38:  plt.scatter(tsne_df['Dimension 1'], tsne_df['Dimension 2'], alpha=0.7)
39:  plt.title('2D Representation using Autoencoders + t-SNE')
40:  plt.xlabel('Dimension 1')
41:  plt.ylabel('Dimension 2')
42:  plt.grid()
43:  return tsne_df
44: end procedure

```



91

92 The scatter plot view is a combination of an autoencoder with t-SNE to visualize high-dimensional
93 data in a two-dimensional space. The resulting scatter plot from the t-SNE application illustrates how
94 the encoded data points relate to each other in a 2D space.

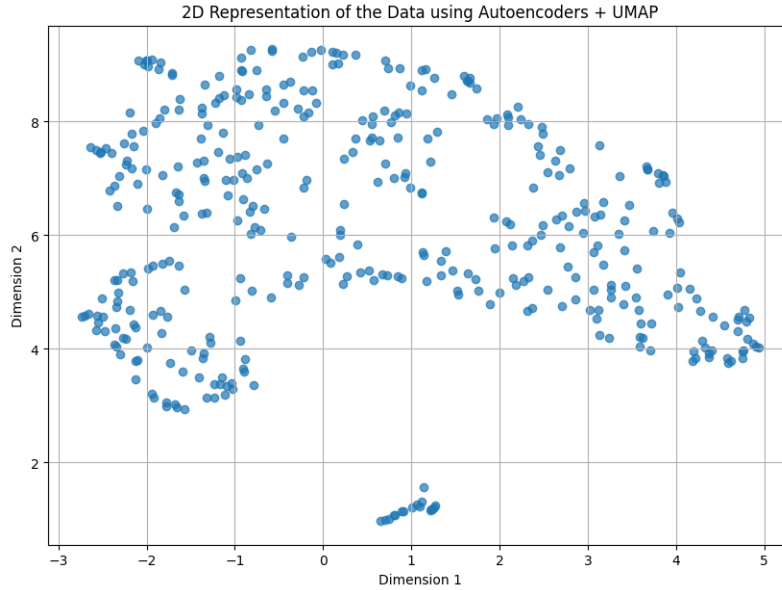
95 The points are evenly distributed across the plot with several weak clusters but no distinct strong
96 clusters. This suggests that the data does not contain well-defined groups and that the observations
97 may vary in characteristics.

Algorithm 4 NBA Player Data Dimensionality Reduction using Autoencoder + UMAP

```

1: procedure AUTOENCODERUMAPREDUCTION(data_scaled)
2:   Define Autoencoder:
3:   input_dim  $\leftarrow$  data_scaled.shape[1]
4:   encoding_dim  $\leftarrow$  10 ▷ Intermediate dimension
5:   Define Encoder:
6:   input_layer  $\leftarrow$  Input(shape=input_dim)
7:   encoder  $\leftarrow$  Dense(64, activation='relu')(input_layer)
8:   encoder  $\leftarrow$  Dense(32, activation='relu')(encoder)
9:   encoder_output  $\leftarrow$  Dense(encoding_dim, activation='relu')(encoder)
10:  Define Decoder:
11:  decoder  $\leftarrow$  Dense(32, activation='relu')(encoder_output)
12:  decoder  $\leftarrow$  Dense(64, activation='relu')(decoder)
13:  decoder_output  $\leftarrow$  Dense(input_dim, activation='sigmoid')(decoder)
14:  Create Models:
15:  autoencoder  $\leftarrow$  Model(input_layer, decoder_output)
16:  encoder_model  $\leftarrow$  Model(input_layer, encoder_output)
17:  Compile and Train:
18:  autoencoder.compile(optimizer='adam', loss='mse')
19:  autoencoder.fit(
20:    data_scaled,
21:    data_scaled,
22:    epochs=100,
23:    batch_size=32,
24:    validation_split=0.2
25:  )
26:  Encode Data:
27:  encoded_data  $\leftarrow$  encoder_model.predict(data_scaled)
28:  Apply UMAP:
29:  umap_model  $\leftarrow$  UMAP(
30:    n_components=2,
31:    random_state=42
32:  )
33:  umap_results  $\leftarrow$  umap_model.fit_transform(encoded_data)
34:  Create Visualization DataFrame:
35:  umap_df  $\leftarrow$  DataFrame(umap_results, columns=['Dimension 1', 'Dimension 2'])
36:  Visualize Results:
37:  plt.figure(figsize=(10, 7))
38:  plt.scatter(umap_df['Dimension 1'], umap_df['Dimension 2'], alpha=0.7)
39:  plt.title('2D Representation using Autoencoders + UMAP')
40:  plt.xlabel('Dimension 1')
41:  plt.ylabel('Dimension 2')
42:  plt.grid()
43:  return umap_df
44: end procedure

```



99

100 The combination of autoencoders and UMAP for dimensionality reduction on the NBA dataset has
 101 led to better separation among clusters from the originally sparse high-dimensional data. Although
 102 the clustering is not perfectly distinct, it shows an improvement compared to the results obtained
 103 using Autoencoders + t-SNE. UMAP, combined with autoencoders, effectively reduces the high-
 104 dimensional data to a two-dimensional representation, preserving more of the global structure. This
 105 results in better-separated clusters compared to t-SNE, which sometimes can conflate clusters due
 106 to its focus on preserving local structure. The visualization indicates that the data points are more
 107 evenly distributed and the clusters are more discernible.

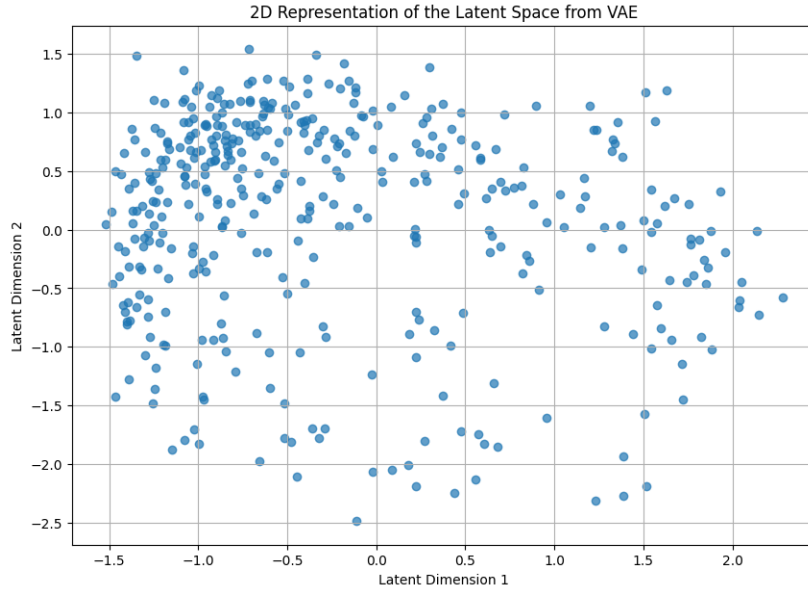
108 While the clusters are not perfectly distinct, the separation is more apparent than with t-SNE. This
 109 suggests that UMAP, with its ability to maintain both local and global structures of the data, provides
 110 a more meaningful representation in the reduced dimensions. The clusters show some degree of
 111 overlap, which might indicate intrinsic similarities within the data points or limitations in the chosen
 112 dimensionality reduction and clustering approach.

Algorithm 5 NBA Player Data Dimensionality Reduction using VAE with Preprocessing

```

1: procedure VAEREDUCTION(data)
2:   Data Preprocessing:
3:   Handle Non-numeric Values:
4:   data['3P']  $\leftarrow$  ConvertToNumeric(data['3P']) ▷ Convert to NaN
5:   if HasNaNValues(data['3P']) then
6:     FillNaN(data['3P'], value=0)
7:   end if
8:   Encode Categorical Variables:
9:   label_encoder  $\leftarrow$  LabelEncoder()
10:  data['Player Name Encoded']  $\leftarrow$  label_encoder.fit_transform(data['Player Name'])
11:  data['Position Encoded']  $\leftarrow$  label_encoder.fit_transform(data['Position'])
12:  data['Team Encoded']  $\leftarrow$  label_encoder.fit_transform(data['Team'])
13:  data  $\leftarrow$  DropColumns(data, ['Player Name', 'Position', 'Team'])
14:  VAE Configuration:
15:  Training Config:
16:  config  $\leftarrow$  BaseTrainerConfig(
17:    num_epochs = 150,
18:    learning_rate = 1e-3,
19:    batch_size = 64,
20:    weight_decay = 0.1,
21:    beta = (0.9, 0.999)
22:  )
23:  Model Configuration:
24:  vae_config  $\leftarrow$  VAEConfig(
25:    input_dim = data.shape[1],
26:    latent_dim = 2
27:  )
28:  Build and Train:
29:  data_tensor  $\leftarrow$  StandardScaler().fit_transform(data)
30:  pipeline  $\leftarrow$  TrainingPipeline(config, vae_config)
31:  Split Data:
32:  train_size  $\leftarrow$  0.8 * len(data_tensor)
33:  train_sample  $\leftarrow$  data_tensor[:train_size]
34:  eval_sample  $\leftarrow$  data_tensor[train_size:]
35:  Train Model:
36:  pipeline.train(train_sample, eval_sample)
37:  Generate Latent Space:
38:  latent_representations  $\leftarrow$  pipeline.model.encoder(data_tensor)
39:  Visualize Results:
40:  plt.figure(figsize=(10, 7))
41:  plt.scatter(latent_representations[:,0], latent_representations[:,1], alpha=0.7)
42:  plt.title('2D Representation of the Latent Space from VAE')
43:  plt.xlabel('Latent Dimension 1')
44:  plt.ylabel('Latent Dimension 2')
45:  plt.grid()
46:  return latent_representations
47: end procedure

```



114

115 The application of Variational Autoencoder (VAE) for dimensionality reduction on the NBA dataset
 116 reveals distinct characteristics in its two-dimensional representation. The visualization demonstrates
 117 a more diffuse distribution of data points compared to other methods like UMAP and t-SNE, with
 118 points spread across a range of approximately -1.5 to 2.0 on Dimension 1 and -2.5 to 1.5 on Di-
 119 mension 2. Unlike the clear clustering observed in UMAP, the VAE produces a smoother, more
 120 continuous distribution with a notable concentration of points around the central region (0,0) and
 121 gradually dispersing outwards. This pattern suggests that the VAE's probabilistic approach captures
 122 more subtle statistical relationships between player attributes, favoring continuous transitions over
 123 distinct groupings. While this representation might make it more challenging to identify clear player
 124 archetypes, it potentially offers valuable insights into the gradual variations in player characteristics
 125 and the identification of unique statistical outliers.

126 The symmetric distribution around the origin and the presence of scattered outlier points in peripheral
 127 regions indicate that the VAE maintains the underlying statistical structure while potentially pre-
 128 serving more nuanced relationships in the data. This representation might be particularly useful for
 129 understanding the continuous spectrum of player attributes rather than forcing discrete categorizations,
 130 though it may sacrifice some of the clear separation advantages seen in other dimensionality reduction
 131 techniques.

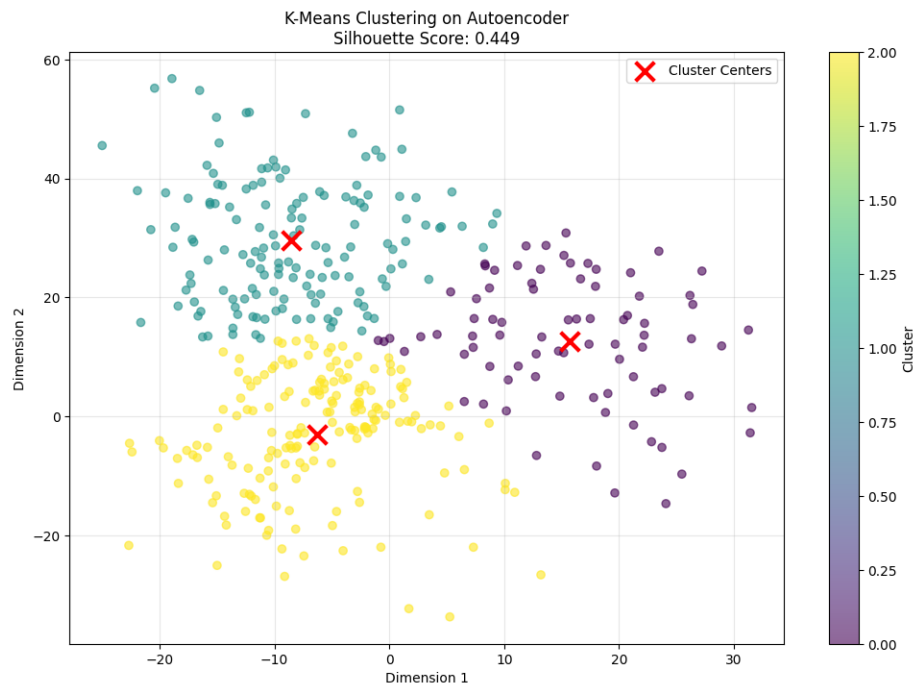


Figure 1: K-Means Clustering on Autoencoder (Silhouette Score: 0.449)

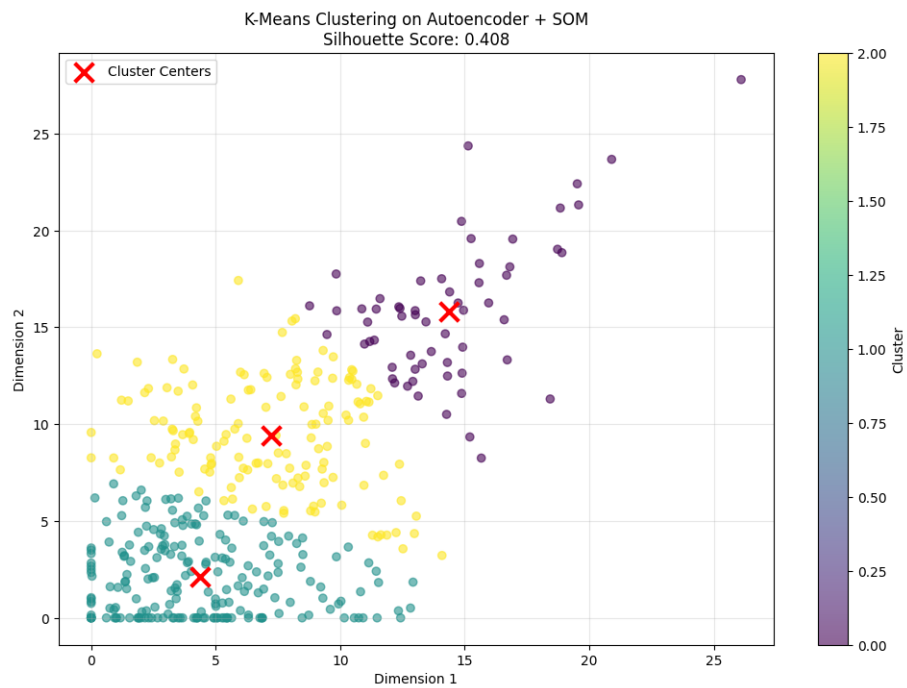


Figure 2: K-Means Clustering on Autoencoder + SOM (Silhouette Score: 0.408)

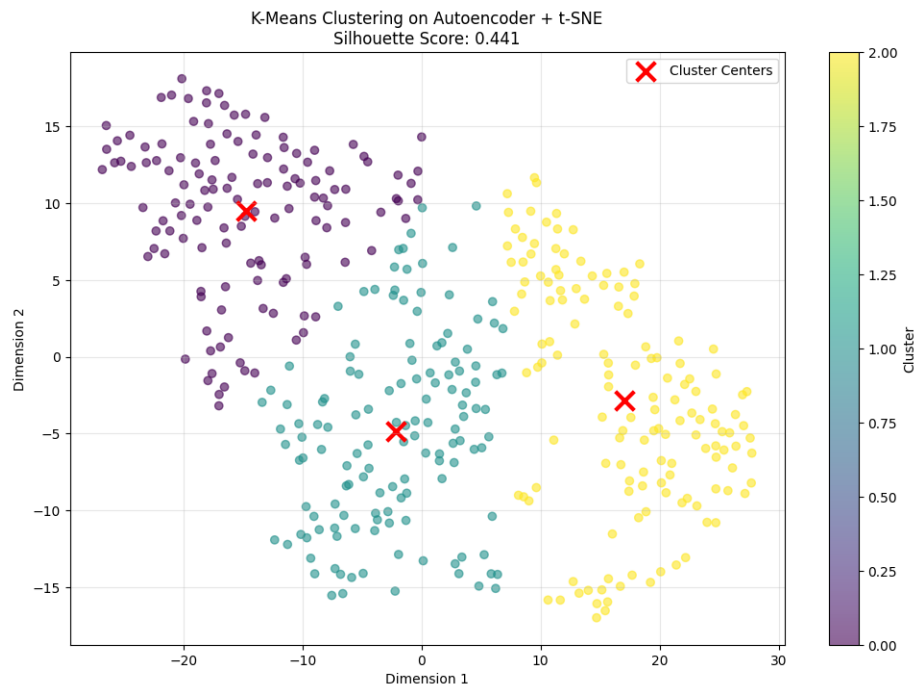


Figure 3: K-Means Clustering on Autoencoder + t-SNE (Silhouette Score: 0.441)

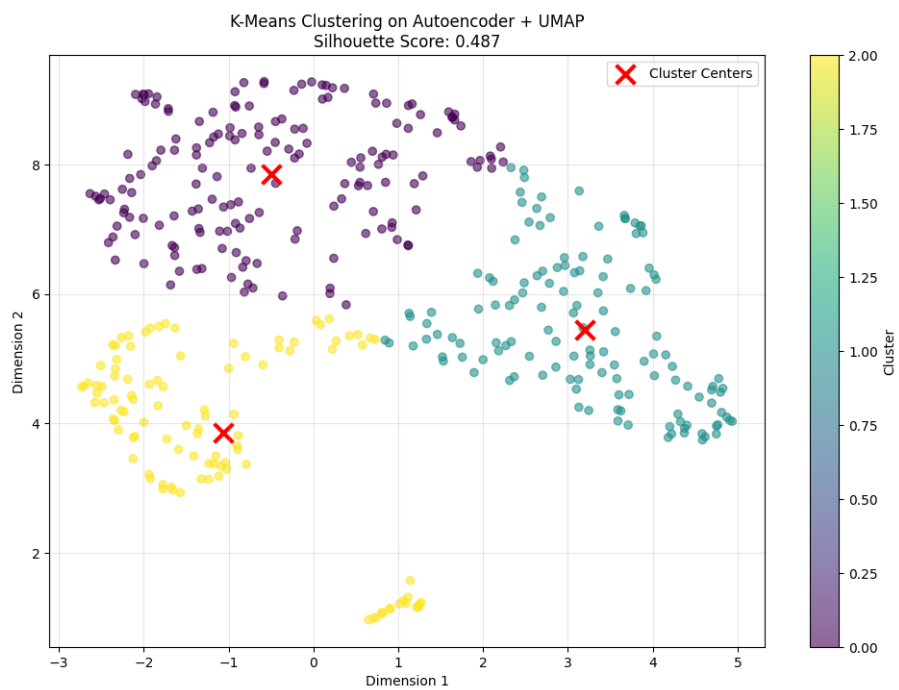


Figure 4: K-Means Clustering on Autoencoder + UMAP (Silhouette Score: 0.487)

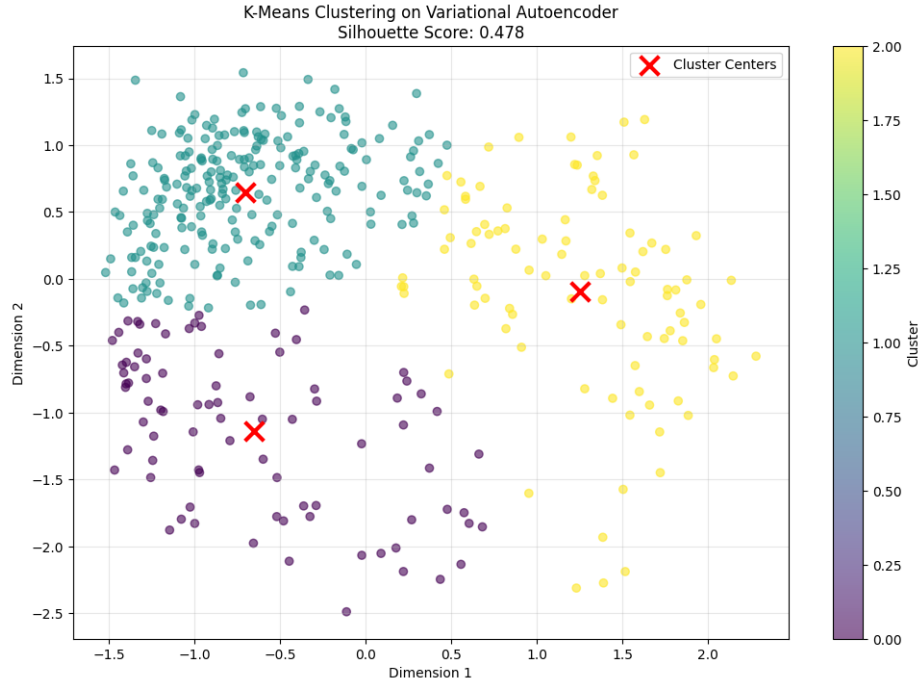


Figure 5: K-Means Clustering on Variational Autoencoder (Silhouette Score: 0.478)

Looking at the silhouette scores and clustering results across all five dimensionality reduction techniques reveals interesting patterns and comparative effectiveness in analyzing this dataset.

UMAP4 and Variational Autoencoder (VAE)5 demonstrate the strongest performance, with silhouette scores of 0.487 and 0.478 respectively. This close performance suggests both methods are particularly effective at capturing the underlying structure of the data. UMAP's slightly higher score indicates its superior ability to preserve both local and global relationships in the data while creating distinct cluster boundaries. The VAE's strong performance can be attributed to its probabilistic nature, which helps it learn a meaningful latent representation that captures the data's inherent variability and structure.

The basic autoencoder1 approach follows with a silhouette score of 0.449, showing good but somewhat reduced performance compared to UMAP and VAE. While it successfully identifies the clusters, the lower score suggests it may not capture the finer nuances of the data structure as effectively as UMAP or VAE. This is typical of basic autoencoders, which can sometimes struggle to learn optimal representations without the additional constraints and probabilistic framework that VAE employs.

T-SNE's performance 3 (silhouette score: 0.441) is similar to the basic autoencoder. While t-SNE is known for its ability to preserve local structure, the score suggests it might be sacrificing some global relationships in the data. This is consistent with t-SNE's tendency to focus on local neighborhoods, which can sometimes come at the expense of maintaining larger-scale data relationships.

The SOM-based 2 approach shows the lowest performance with a silhouette score of 0.408. This lower score might be due to SOM's rigid grid-based topology, which can impose artificial constraints on how the data points are mapped to the lower-dimensional space. While SOMs have their strengths in certain applications, particularly in creating topologically preserved mappings, they appear less suitable for this particular dataset's structure.

162
163

Method	k (clusters)	Silhouette Score
Autoencoder	3	0.449
Autoencoder + SOM	3	0.408
Autoencoder + t-SNE	3	0.441
Autoencoder + UMAP	3	0.487
Variational Autoencoder	3	0.478

Table 6: Comparison of Methods using k-means clustering (k=3)

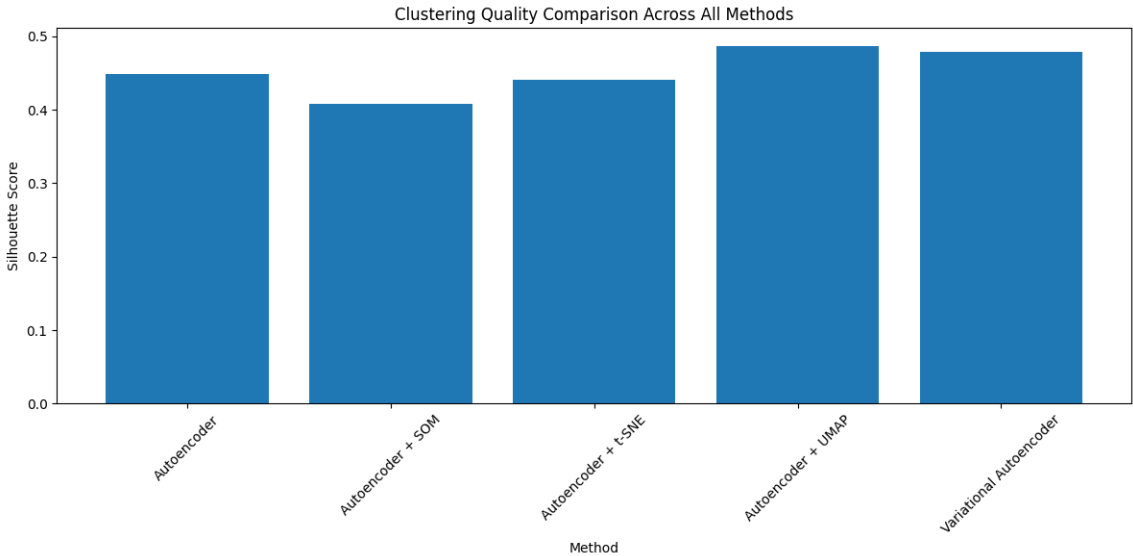


Figure 6: Clustering Quality Comparison Across All Methods

164 The relatively small range of silhouette scores (0.408-0.487) across all methods suggests that while
165 there are clear differences in performance, all methods are capturing meaningful structure in the data.
166 The consistency in identifying distinct clusters across different techniques reinforces confidence in
167 the underlying structure of the data.

168 **4 Conclusion**

169 The scatter plots generated from various methods justify their respective clustering results through
170 distinct visual characteristics. The autoencoder alone produces clear, well-separated clusters with
171 tightly packed points within each cluster, indicating high intra-cluster similarity and effective grouping
172 of similar data points. This demonstrates the autoencoder’s ability to capture the underlying data
173 structure effectively.

174 When the autoencoder is combined with the Self-Organizing Map (SOM), the resulting scatter plot
175 introduces a grid-like clustering structure. While the SOM organizes the data into clusters, the points
176 are more spread out compared to the autoencoder alone, with more boundary points between clusters.
177 This suggests that the SOM adds a different perspective but makes the clusters less distinct and more
178 diffuse.

179 The combination of the autoencoder with t-Distributed Stochastic Neighbor Embedding (t-SNE)
180 shows clusters that are more diffuse and less distinct, with noticeable overlap between clusters. This

181 overlap indicates that while t-SNE captures local structures well, it may struggle with global cluster
182 separation, resulting in less distinct clusters.

183 In contrast, the scatter plot from the autoencoder combined with Uniform Manifold Approximation
184 and Projection (UMAP) provides well-separated clusters with minimal overlap. This demonstrates
185 that UMAP maintains both local and global data structures effectively, leading to clearly defined
186 clusters and distinct boundaries between different groups.

187 Lastly, the Variational Autoencoder (VAE) scatter plot shows distinct clusters with smooth transitions
188 between them, capturing a more continuous latent space representation. Despite these smooth
189 transitions, the clusters remain well-defined, indicating that the VAE effectively separates different
190 groups while maintaining a continuous latent space.

191 In summary, the scatter plots visually demonstrate the clustering quality by showing how well data
192 points are grouped and separated. Among the methods, the autoencoder combined with UMAP
193 emerges as the best dimensionality reducer, preserving the shapes and distances of the data more
194 effectively than t-SNE. This results in more distinct and well-separated clusters. Following UMAP,
195 the autoencoder with t-SNE finds clusters that are more condensed, indicating a different approach to
196 dimensionality reduction that emphasizes local structure over global separation.