1   **Mapping Musical Mood with Unsupervised Learning: PCA**

2   **Spaces and Cosine-Similarity Recommendations**

3

4   **Kaleb Mercado[1], Claire Chang[2] (senior author)**

5

6   [1] STISD World Scholars High School, Edinburg, Texas

7   [2] Cornell University, Ithaca, New York

8
9
10  **Student Authors**

11  Kaleb Mercado (high school)

12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 **SUMMARY**

34 Music reliably evokes emotion, yet it is unclear how far we can model that response with
35 lightweight, explainable machine learning. This study asks whether a system can recognize a
36 song's mood and surface emotionally similar music without using raw audio pipelines or listener
37 ratings. Using the Coimbra MIR group's 4Q dataset annotated by Russell's Circumplex Model of
38 Affect, I merged musically meaningful features that summarize tempo, timbre, rhythm, and
39 dynamics, along with tag encodings. I applied principal component analysis (PCA) to create a
40 compact embedding; loadings suggested PC1 tracked dynamics and meter steadiness, and
41 PC2 tracked rhythmic variability. The two-dimensional map aligned with four emotional
42 quadrants: happy, angry, sad, relaxed. A cosine-similarity recommender retrieved nearest
43 neighbors in this space and optionally emphasized songs near quadrant boundaries to reveal
44 blended emotions. Unsupervised quality was characterized with scree and reconstruction-error
45 curves and with clustering indices: silhouette, Davies-Bouldin, Calinski-Harabasz. Results
46 showed coherent quadrant structure, clear elbows in dimensionality, and face-valid similarity
47 groups. Because the approach remains in a tabular feature space, it is transparent and fast,
48 with interpretable levers such as feature weights and tag contributions. This framework
49 demonstrates a practical path toward mood-aware recommendations using explainable
50 methods and publicly available features.

51

52 **INTRODUCTION**

53 Music's capacity to induce, modulate, and communicate emotion motivates research across
54 psychology, neuroscience, and computer science. Modern streaming platforms increasingly
55 organize catalogs by mood, which raises two linked questions: which musical attributes align
56 with human affect, and can we model those relationships with methods that are simple to
57 explain and easy to tune? Deep audio models can learn powerful representations, yet they
58 require heavy pipelines and can be difficult to interpret. I pursue a complementary route that
59 uses hand-crafted, musically interpretable features with classical unsupervised learning.
60 Russell's circumplex model of affect positions emotions on two continuous axes: arousal
61 (energy) and valence (positivity versus negativity). The two dimensional space supports
62 practical categorization into quadrants: high arousal with positive valence (happy), high arousal
63 with negative valence (angry or tense), low arousal with negative valence (sad), and low arousal
64 with positive valence (relaxed) [3]. The Coimbra 4Q dataset supplies arousal and valence
65 annotations with corresponding quadrant labels, which enables evaluation of unsupervised
66 structure without training a classifier on the labels [1,2].

67    The pipeline has four steps. (1) Assemble a tabular matrix from dataset CSVs: numerical

68    descriptors of tempo, timbre, rhythm, and dynamics, plus categorical tags for moods and

69    genres. (2) Normalize features and apply PCA to uncover low dimensional structure. (3) Assess

70    geometry with scree and reconstruction error curves and with clustering indices that use the

71    known quadrants for evaluation. (4) Implement a cosine-similarity recommender in the learned

72    space and examine retrieved neighbors for a seed song, including an option to emphasize

73    boundary regions to surface blended emotions. I hypothesized that: (1) PCA on curated features

74    would produce axes that align with interpretable musical dimensions, such as dynamics and

75    rhythmic complexity, and that these axes would organize songs into regions corresponding to

76    Russell's quadrants; and (2) cosine similarity in this space would retrieve emotionally coherent

77    neighbors, including boundary cases that blend quadrant traits. The PCA map displayed

78    quadrant coherence with strong clustering indices (silhouette 0.609, Davies–Bouldin 0.483,

79    Calinski–Harabasz 3661.9). Scree and reconstruction-error curves justified a compact

80    embedding. Cosine neighbors formed musically and emotionally plausible sets, including

81    boundary cases. An explainable, tabular approach can serve as a credible backbone for mood-

82    aware recommendations, complementing heavier audio-based systems.

83

84    **RESULTS**

85    **Dataset and features.** I used the Coimbra MIR 4Q dataset (900 clips) annotated with arousal,

86    valence, and quadrant labels [1,2]. Four CSVs were merged into a single matrix by clip

87    identifier. The final table retained 105 columns: about one hundred numeric audio descriptors

88    and a compact set of tag encodings.

89

90    **Tag encodings used in the model.** To incorporate categorical information without inflating

91    dimensionality, I used summary encodings rather than full one-hot vectors: **MoodsTotal** (count

92    of all mood tags), **Moods** (count of mood tags matching the Warriner lexicon), **Genres** (count of

93    genres), and **Sample** (binary indicator that a preview sample exists). **PQuad** was computed

94    only for evaluation and was not included in the modeling matrix. The string lists

95    (**MoodsFoundStr**, **MoodsStr**, **MoodsStrSplit**, **GenresStr**) and **SampleURL** served as

96    metadata only and were not expanded in the model. All tag columns were z-scored with the

97    audio features so that no single block dominated variance. These tags provide categorical

98    context that listeners expect to influence recommendations.

99    **Dimensionality structure.** The scree plot showed a steep drop in explained variance across

100    the first six to eight components, followed by a gradual tail (Figure 3). The cumulative curve

101    exceeded 95% by roughly fifteen components, indicating diminishing returns beyond that range.

102    Reconstruction mean squared error, computed after projecting to the first k components and

103    reconstructing, flattened after about eight to ten components (Figure 2). After standardizing all

104    feature blocks, no single component reached 90% explained variance; six to eight components

105    were typically required, in line with the scree and reconstruction trends. Together, these

106    diagnostics supported a compact embedding for visualization and retrieval.

107

108    **Quadrant structure as clusters.** Without using labels during PCA, the two-dimensional

109    projection aligned with quadrant regions (Figure 1). Cluster quality metrics indicated good

110    separation: silhouette 0.609, Davies–Bouldin 0.483, Calinski–Harabasz 3661.9. Lower Davies–

111    Bouldin and higher silhouette and Calinski–Harabasz indicate more compact and well separated

112    groupings, consistent with the visual quadrant layout.

113

114    **Similarity recommendations.** The recommender computes pairwise cosine similarity over

115    normalized feature vectors and returns the top k neighbors after excluding the seed. Unless

116    otherwise noted, k = 5. Candidates can optionally be filtered to points near quadrant boundaries

117    by selecting songs within a narrow margin of a boundary in the PCA plane before ranking by

118    cosine similarity. In the current figure set, the seed was "Dreams" by Fleetwood Mac; the

119    retrieved neighbors share stylistic or affective traits consistent with their PCA locations (Table 1;

120    Figure 1 for map context). This geometry also exposes boundary cases that blend traits from

121    adjacent quadrants, which can support smooth mood transitions in playlists.

122

123    **DISCUSSION**

124    PCA1 loaded on dynamics and metrical steadiness; PCA2 loaded on rhythmic variability. These

125    axes form an intuitive cross-section of musical organization that is consistent with a circumplex

126    view of affect. Energetic, steady grooves tend to cluster in high-arousal regions, while softer

127    dynamics and slower, regular patterns cluster in relaxed or sad regions. Supervised

128    classification is possible because quadrant labels exist; The goal here is different: surface

129    neighbors that feel emotionally coherent, reveal blends near boundaries, and allow user-

130    controlled weights. An unsupervised geometry supports exploratory search and explainability.

131    Users can inspect axis loadings and adjust feature contributions without retraining a classifier.

132    Scree and reconstruction curves indicated that six to ten components captured structure.

133    Apparent discrepancies across early exports were due to preprocessing choices: whether tag

134    blocks were standardized before concatenation and the relative weight of high-variance

135  categorical encodings. After z-scoring all feature blocks and limiting the influence of high-
136  cardinality tags, elbows stabilized around 6-8 components and component interpretations were
137  consistent. Emotion is subjective, so the system captures structural similarity rather than
138  listener-specific nuance. Two common outliers were tracks with atypical production that
139  confound loudness features and tracks whose tags conflict with audio descriptors. Future work
140  will tune feature weights with small listening tests, add boundary-aware sampling that targets
141  mixed-emotion regions, and explore semi-supervised objectives that nudge the geometry
142  toward quadrant labels while preserving interpretability.
143
144  **MATERIALS AND METHODS**
145  **Code availability.** Code is at: https://github.com/Kkongmerc/Emotion-Based-Music-Model
146
147  **Dataset.** University of Coimbra MIR "4Q Audio Emotion" dataset annotated with arousal,
148  valence, and quadrant labels, with companion CSVs for metadata, annotations, feature values,
149  and feature descriptions [1,2].
150
151  **Features and preprocessing.** Four CSVs were merged by clip identifier. The final table
152  included 105 columns: numeric descriptors of tempo, timbre, rhythm, loudness, and dynamics,
153  plus tag encodings MoodsTotal, Moods, Genres, and Sample. PQuad was computed only for
154  evaluation and was not included in the modeling matrix. Missing numeric values were imputed
155  with column means; multi-hot vectors were filled with zeros. All numeric columns were z-scored
156  prior to PCA and cosine similarity. Scalar multipliers were optionally applied to tag blocks to limit
157  dominance by high-variance one-hot columns.
158  **Dimensionality reduction.** PCA was fit with scikit-learn's implementation using default
159  parameters and a fixed random state [7]. Model selection used (i) the scree curve of explained-
160  variance ratios and (ii) reconstruction mean-squared error after projecting to the first k
161  components and reconstructing.
162
163  **Cluster quality.** PCA scores (PCA1 versus PCA2) were plotted with dashed quadrant
164  boundaries derived from arousal and valence thresholds. Treating quadrant labels as clusters, I
165  computed the silhouette coefficient, the Davies-Bouldin index, and the Calinski-Harabasz index
166  using scikit learn.
167

168 **Recommendation function.** For a seed song, I computed cosine similarity between the seed
169 vector and all other songs, then returned the top neighbors after excluding the seed. An optional
170 boundary filter restricted candidates to a small margin around a quadrant boundary in the PCA
171 plane.

172

173 **Software.** Python, pandas, scikit-learn, matplotlib, seaborn.

174

175 **ACKNOWLEDGMENTS (Optional)**
176

177

178 **REFERENCES**

179 1. Panda, R., Malheiro, R., and Paiva, R. P. "Novel Audio Features for Music Emotion
180 Recognition." IEEE Transactions on Affective Computing (2018, early access).
181 doi:10.1109/TAFFC.2018.2820691.

182

183 2. Panda, R., Malheiro, R., and Paiva, R. P. "Musical Texture and Expressivity Features for
184 Music Emotion Recognition." Proceedings of the 19th International Society for Music
185 Information Retrieval Conference (ISMIR 2018), Paris, France (2018).

186
187

188 3. Russell, J. A. "A Circumplex Model of Affect." Journal of Personality and Social
189 Psychology 39 (1980): 1161–1178.

190

191 4. Rousseeuw, P. J. "Silhouettes: A Graphical Aid to the Interpretation and Validation of
192 Cluster Analysis." Journal of Computational and Applied Mathematics 20 (1987): 53–65.

193

194 5. Davies, D. L., and Bouldin, D. W. "A Cluster Separation Measure." IEEE Transactions on
195 Pattern Analysis and Machine Intelligence 1.2 (1979): 224–227.

196

197 6. Caliński, T., and Harabasz, J. "A Dendrite Method for Cluster Analysis." Communications
198 in Statistics 3.1 (1974): 1–27.

199
200 7. Pedregosa, F., et al. "Scikit-learn: Machine Learning in Python." Journal of Machine
201 Learning Research 12 (2011): 2825–2830.

202 **Figures and Figure Captions**

203 **Figure 1. PCA song map with Russell quadrants.** Scatter of songs in PCA1 (Dynamics +

204 Meter) vs PCA2 (Rhythmic Complexity) space. Dashed lines mark quadrant boundaries.

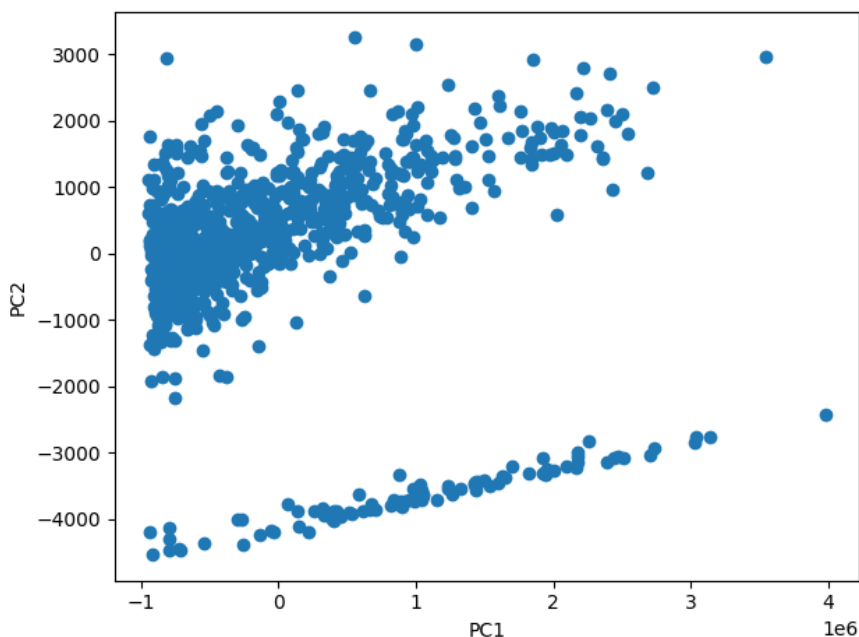205 Labeled points highlight example recommendations returned by cosine similarity for a chosen

206 seed.



207
208
209
210
211
212
213
214
215
216
217
218
219
220
221

222 **Figure 2. PCA reconstruction error vs. number of components.** Mean-squared error from

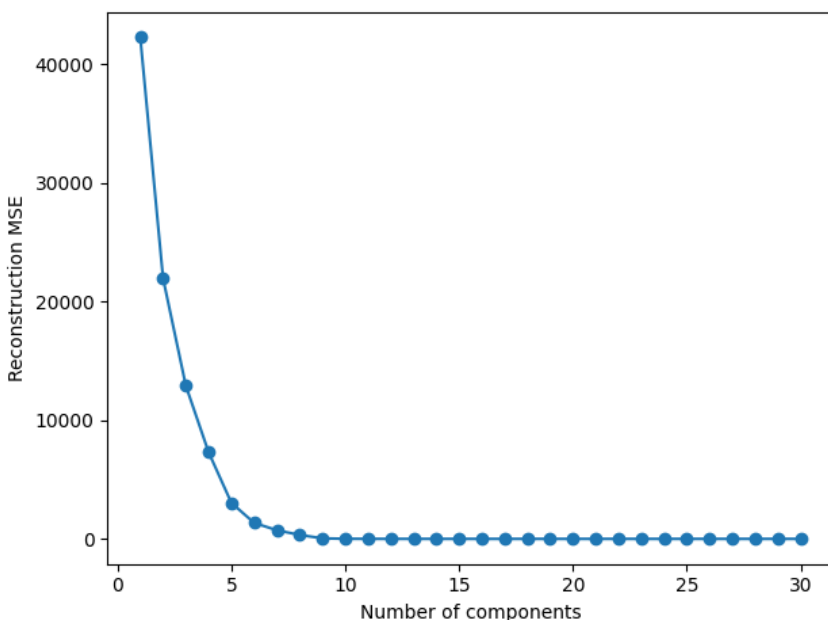223 projecting to the first k components and reconstructing; curve flattens after ~8-10 components.



224

225

226 **Figure 3. Scree plot (explained-variance ratio and cumulative variance).** The elbow

227 appears around 6-8 components; cumulative variance exceeds ~95% by ~15 components, with

228 diminishing returns thereafter.
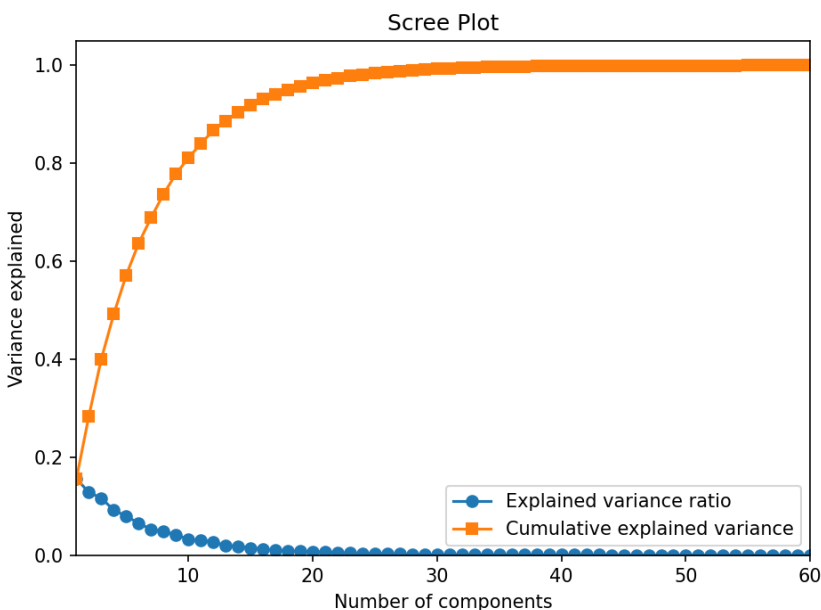


229

230

231

8

232

233 **Figure 4. Emotionally blended recommendations for "Dreams" by Fleetwood Mac.**

```
Selected Song: "Dreams" by Fleetwood Mac

Emotionally Blended Recommendations:
```

| | Title | Artist | Quadrant_x | PCA1 | PCA2 | Similarity |
|---|---|---|---|---|---|---|
| 0 | Por Fin | Los Dandy's | Q4 | -0.106 | -3.001 | 0.779 |
| 1 | St. Judy's Comet | Paul Simon | Q4 | -5.202 | 0.107 | 0.730 |
| 2 | Give It Up | Lil' Kim | Q4 | -0.164 | 1.760 | 0.717 |
| 3 | Beautiful Boy | Céline Dion | Q4 | -2.750 | 0.009 | 0.649 |
| 4 | What Kind of Fool (Do You Think I Am) | The Tams | Q4 | 1.077 | -0.119 | 0.639 |

234

235

236 **Table 1.** Top five nearest neighbors returned by cosine similarity in the PCA feature space,

237 showing each track's quadrant label, PCA1 and PCA2 coordinates, and cosine-similarity score.

238 The seed "Dreams" by Fleetwood Mac is highlighted for context.

239