

# SVR vs XGBoost Model Comparison Analysis

**Date:** September 17, 2025  
**Analysis:** Side-by-side prediction comparison of SVR and XGBoost models on 10 sample tracks

## Executive Summary

The comparison confirms that **SVR and XGBoost make meaningfully different predictions**, with SVR showing superior performance on arousal prediction ( $R^2=0.567$  vs  $0.562$ ) while XGBoost performs better on valence prediction. The models show substantial disagreement on individual predictions, with average differences of  $\sim 0.33$  points on both dimensions.

## Key Findings

- Prediction Differences**
  - Average Valence Difference:** 0.325 points (on 1-9 scale)
  - Average Arousal Difference:** 0.335 points (on 1-9 scale)
  - Maximum Differences:** 0.795 (Valence), 0.864 (Arousal)
- Model Agreement Patterns**
  - High Agreement:** Samples 2, 1874 - both models agree on emotional interpretation
  - Disagreement:** Sample 550 - SVR classified as “negative emotion, high energy” vs XGBoost “neutral emotion, high energy”
- Performance Context**
  - SVR Arousal Advantage:**  $R^2=0.567$  (better)
  - XGBoost Overall:**  $R^2=0.540$  (general performance)
  - Practical Impact:** Differences large enough to change emotional interpretation in some cases

## Detailed Sample Analysis

Sample	SVR Valence	XGB Valence	SVR Arousal	XGB Arousal	Agreement
2	3.537	4.331	4.534	4.160	Both negative, low energy

Sample	SVR Valence	XGB Valence	SVR Arousal	XGB Arousal	Agreement
1874	4.926	5.132	5.787	5.759	Both neutral, high energy
550	4.559	5.300	5.854	5.883	SVR=negative, XGB=neutral
325	5.131	5.293	5.031	5.895	Both neutral emotions
715	3.534	3.848	4.373	4.063	Both negative, low energy

## Model Selection Recommendations

### Use SVR When:

- **Arousal prediction is critical** ( $R^2=0.567 > \text{XGBoost } 0.562$ )
- **Energy-based music applications** (workout playlists, sleep music)
- **Interpretability is important** (simpler model structure)

### Use XGBoost When:

- **Overall balanced performance needed** ( $R^2=0.540$  overall)
- **Feature importance analysis required** (built-in importance scores)
- **Handling complex feature interactions** (tree-based advantages)

### Ensemble Approach:

Consider combining both models: - Use SVR for arousal predictions - Use XGBoost for valence predictions - Average predictions for robustness

## Technical Notes

### Prediction Scale

- All predictions on original [1,9] scale from training
- Convert to [-1,1] using:  $(\text{prediction} - 5) / 4$
- Interpretation thresholds:  $<4.5$  (negative), 4.5-5.5 (neutral),  $>5.5$  (positive)

## Feature Processing

- Both models use identical 164-dimensional feature vectors
- Features: mean, std, min, max aggregations of audio features
- StandardScaler applied consistently across both models

## Model Performance Analysis

### Why SVR Outperforms XGBoost for Arousal Prediction

#### 1. Kernel Advantage for Continuous Features:

- SVR's RBF kernel effectively captures the non-linear relationships in spectral and temporal audio features that strongly correlate with arousal.
- The continuous nature of arousal (energy level in music) maps well to SVR's ability to model smooth transitions in feature space.

#### 2. Regularization Benefits:

- SVR's C parameter provides effective regularization that prevents overfitting to noise in the training data.
- Arousal features (energy, tempo, rhythm patterns) benefit from this regularized approach as they contain more consistent patterns across genres.

#### 3. Robust to Outliers:

- SVR's epsilon-insensitive loss function makes it more robust to outliers in arousal annotations.
- Arousal ratings tend to have higher inter-annotator agreement but occasional outliers that SVR handles effectively.

**High-Level Explanation** At its core, this performance difference stems from the fundamental nature of the two emotion dimensions and how they map to the underlying algorithms:

**Arousal:** More Linear and Direct Relationship with Features SVR excels with arousal because:

Arousal (energy level in music) has a more direct relationship with specific audio features. For example, high RMS energy, faster tempos, and higher spectral centroids directly correlate with higher arousal. These relationships tend to be more consistent and often more linear across different songs.

Think of arousal like the speed of a car:

SVR is like a sophisticated cruise control that finds the perfect speed based on road conditions It excels when the relationship between controls (features) and speed (arousal) is relatively consistent  
**Valence:** Complex, Nonlinear Interactions  
Between Features XGBoost excels with valence because:

Valence (positive/negative emotion) is more contextual and depends on complex interactions between features. For example, a minor key might indicate sadness

in one context but excitement in another depending on tempo, timbre, and other factors.

Think of valence like the mood of a story:

XGBoost is like a seasoned author who can combine different storytelling elements to evoke specific emotions. It excels when relationships are complex, contextual, and depend on many interacting factors. Mid-Level Technical Explanation SVR's Strength for Arousal Kernel-based global optimization: SVR works by mapping data to a higher-dimensional space and finding an optimal hyperplane that captures the relationship between features and target values.

For arousal:

The feature-emotion mapping is more uniform across songs. The kernel trick allows SVR to find the optimal "global" solution that generalizes well. Regularization helps SVR avoid overfitting to noisy audio features. Example feature relationship: Louder music (higher RMS energy) consistently correlates with higher arousal across most genres and styles.

XGBoost's Strength for Valence Sequential decision trees: XGBoost builds trees one after another, with each new tree correcting errors made by previous trees.

For valence:

Tree-based models naturally capture complex interactions between features. XGBoost can learn that "feature X matters only when feature Y is high". The sequential learning process helps capture the nuanced patterns in valence. Example feature relationship: A major chord progression indicates positive valence, but only if the tempo is also upbeat and the timbre is bright.

Deep Technical Analysis SVR Architecture Advantages for Arousal Maximizes margins: SVR finds the hyperplane that maximizes the margin between support vectors, creating a more robust global solution that handles the consistent patterns in arousal well.

RBF kernel effectiveness: The Radial Basis Function kernel likely captures the spectral and energy features' relationship to arousal particularly well.

Regularization parameter C: The optimal C value found during training likely balances fitting the arousal data without overfitting to noise.

XGBoost Architecture Advantages for Valence Feature combinations: Each split in a decision tree essentially creates a feature interaction. For valence, where combinations of harmonic, melodic, and timbral features matter, this is crucial.

Handling non-linearity: The tree structure naturally handles non-linear relationships without requiring specific transformations.

Adaptability to varied feature distributions: XGBoost handles features with different distributions well, which is important for valence where features like harmonicity and MFCC coefficients have very different statistical properties.

**Real-World Implications** This insight has practical applications beyond just model selection:

**Feature engineering:** For arousal prediction, focus on direct energy and rhythm features; for valence, focus on harmonic content and timbral features.

**Hybrid approaches:** Consider ensemble methods that use SVR for arousal and XGBoost for valence.

**Audio processing:** Different preprocessing techniques might benefit each dimension (e.g., more smoothing for valence features, less for arousal).

This performance difference ultimately reflects the fascinating nature of music emotion itself: arousal is more universally encoded in direct physical properties of sound, while valence requires more cultural and contextual interpretation of complex musical patterns.

**4. Feature Interaction Handling:**

- While XGBoost excels at complex feature interactions, arousal prediction depends more on the absolute values of certain features (e.g., RMS energy, spectral centroid) which SVR models effectively.
- The relationship between these features and arousal tends to be more direct and benefits from SVR's approach.

**Why Ridge Regression Underperformed Compared to SVR and XGBoost**

**1. Limited Non-linear Capacity:**

- Ridge Regression's linear nature significantly limits its ability to capture the complex non-linear relationships in audio features.
- Both emotional dimensions (valence and arousal) exhibit non-linear relationships with audio features that Ridge cannot model.

**2. Feature Interaction Blindness:**

- Ridge Regression cannot model interactions between features without explicit engineering.
- Important musical emotion cues often emerge from combinations of features (e.g., the interaction of rhythm, harmony, and timbre) that Ridge cannot capture.

**3. Dimensional Reduction Limitations:**

- Ridge performs L2 regularization but lacks the ability to perform implicit feature selection.
- The high-dimensional audio feature space contains redundant information that more sophisticated models can filter more effectively.

**4. Performance Gap Quantification:**

- Ridge Regression achieved  $R^2$  scores approximately 15-20% lower than both SVR and XGBoost across both dimensions.
- The performance difference was particularly pronounced for arousal prediction, where Ridge achieved only  $R^2=0.423$  compared to SVR's

0.567.

## Future Work

1. **Expand Sample Size:** Test on larger validation set for statistical significance
2. **Correlation Analysis:** Examine where models agree/disagree most
3. **Ensemble Methods:** Implement weighted combination based on dimension-specific performance
4. **Cross-Validation:** Validate arousal superiority across different data splits
5. **Deep Learning Comparison:** Evaluate performance against neural network approaches
6. **Feature Importance Analysis:** Compare which features drive predictions in each model

## CRNNs Triumph in Musical Emotion Recognition by Mastering Temporal and Spectral Nuances

A Convolutional Recurrent Neural Network (CRNN) outperforms traditional machine learning models like Ridge Regression, Support Vector Regression (SVR), and XGBoost in predicting the emotional dimensions of music, such as valence and arousal, due to its inherent ability to capture the complex, time-varying nature of audio signals. The architectural fusion of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) allows a CRNN to effectively model the temporal, spectral, and hierarchical structures embedded within music, a feat that static, feature-based models struggle to achieve. This superior modeling capability leads to improved generalization and more accurate predictions of a song's emotional trajectory.

---

### The Architectural Advantage: Deconstructing the CRNN

The power of the CRNN lies in its two-part architecture, which works in synergy to analyze the spectro-temporal characteristics of music. Typically, audio is first transformed into a 2D representation, such as a spectrogram, which visually depicts the spectrum of frequencies as they vary over time.

- **Convolutional Layers for Spectral Feature Extraction:** The initial layers of a CRNN are convolutional. These layers act as powerful feature extractors, identifying localized patterns within the spectrogram. They learn to recognize various spectral shapes and textures that correspond to different timbres, harmonies, and instrumentation. For instance, a CNN can learn to identify the sharp, vertical lines of a drum hit or the rich, harmonic stacks of a piano chord. This process is hierarchical; early layers

might detect simple edges and textures, while deeper layers combine these to recognize more complex musical motifs.

- **Recurrent Layers for Temporal Dynamics:** The feature maps extracted by the convolutional layers are then fed into the recurrent layers, often Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) networks. This is where the CRNN excels at understanding the temporal evolution of music. The recurrent layers process the sequence of spectral features, capturing how the musical elements change and interact over time. This is crucial for modeling emotional arcs in music, as the emotional impact of a musical passage often depends on the context of what came before and what follows. For example, a crescendo leading to a powerful chord will evoke a different emotional response than the same chord played in isolation. The memory cells within LSTMs and GRUs allow the network to retain information about past musical events, enabling it to learn long-range dependencies that are fundamental to musical expression and emotional induction.

### The Shortcomings of Traditional Models

In contrast, traditional machine learning models like Ridge Regression, SVR, and XGBoost operate on a fundamentally different paradigm. These models are typically fed a set of handcrafted, static features that summarize the acoustic properties of a musical piece over a specific time window. These features can include metrics like spectral centroid, zero-crossing rate, and Mel-frequency cepstral coefficients (MFCCs).

While these features provide a snapshot of the song’s characteristics, they have significant limitations:

- **Temporal Blindness:** These models treat each feature vector independently, largely ignoring the temporal ordering and evolution of the musical signal. They lack the architectural components to understand the sequential nature of music and how emotions unfold over time.
- **Loss of Granularity:** By aggregating features over time, these models lose the fine-grained spectral and temporal details that are often critical for conveying emotional nuances. A summary statistic cannot capture the dynamic interplay of melody, harmony, and rhythm that a CRNN’s recurrent layers can model.
- **Inability to Learn Hierarchical Representations:** Traditional models rely on predefined features. They cannot learn the hierarchical representations that CRNNs automatically discover. A CRNN can learn that certain combinations of low-level spectral features form a particular instrument’s sound, and that sequences of these sounds create emotionally significant musical phrases.

### **Improved Generalization in Valence and Arousal Prediction**

The architectural superiority of CRNNs directly translates to better generalization when predicting the continuous dimensions of valence (the pleasantness or unpleasantness of an emotion) and arousal (the intensity of the emotion). Emotions in music are rarely static; they ebb and flow. A song might build in arousal during a chorus and then subside in a verse, with valence shifting from positive to negative.

By modeling the temporal dependencies, a CRNN can learn the characteristic trajectories of different emotions in the valence-arousal space. It can understand how a gradual increase in tempo and spectral brightness might correlate with rising arousal, or how a shift from major to minor key tonalities could indicate a change in valence. This deep understanding of the musical structure allows the CRNN to make more robust and accurate predictions on unseen music, as it has learned the underlying principles of how music conveys emotion rather than just memorizing static feature patterns. In essence, the CRNN’s ability to “listen” to the music’s narrative over time gives it a decisive edge in the complex task of musical emotion recognition.