# AFFECT PREDICTION IN MUSIC USING BOOSTED ENSEMBLE OF FILTERS

*Rahul Gupta, Naveen Kumar, Shrikanth Narayanan*

Signal Analysis and Interpretation Lab,
University of Southern California, Los Angeles, CA-90007, USA

## ABSTRACT

Music influences the affective states of its listeners. For this reason, music is extensively used in various media forms to enhance and induce emotional feeling. Automatic evaluation of affect from music can have impact on music design and can also aid further analysis of music. In this work, we present a novel scheme for affect prediction in music using a *Boosted Ensemble of Single feature Filters* (BESiF) model. Given a set of frame-wise features, the BESiF model predicts the affective rating as a weighted sum of filtered feature values. The BESiF model improves the Signal to Noise Ratio for arousal and valence prediction by a factor of 1.92 and 1.06, respectively, over the best baseline method. This performance is achieved using only 14 signal features for arousal (16 for valence). We further analyze the transformation of one of the features selected towards arousal prediction.

*Index Terms*— Affect, Arousal, Valence, Emotion in music, Boosting

## 1. INTRODUCTION

In recent years, considerable amount of research has gone into improving automatic understanding and indexing of music signals. This effort has been partly led by the data deluge in digital music and partly by the large number of new applications in multimedia such as information retrieval, automatic transcription and music fingerprinting. A majority of these applications require classifying songs into meaningful categories as the first step. Typically, grouping is done on the basis of genre or melody. However, in the context of a music recommendation system, it is desirable for these categories to be aligned with the listener's music preference or mood. Thus, studying the affective component in music signal is as important as studying its structural aspects.

Music has been shown to possess the ability to influence the emotional state of its listeners [1–3]. For example, consider the elaborate use of background soundtracks in movies to support the narrative being carried through speech and video. In movies, music plays a complementary role to the cinematography and dialog delivery [4]. In fact, previous studies have found that music in movies often plays a more important role in conveying emotion compared to other modalities [5]. Owing to its positive emotional influence, music has also been used for therapeutic purposes [6, 7].

The ability of music to affect human's emotion state has led to considerable interest in the study of affective features and prediction models for music emotion recognition. Previous studies have focused on predicting both static and dynamic emotion labels from music [8–10]. These emotion labels are typically measured along affective dimensions of arousal and valence, and collected from multiple human annotators [11, 12]. Predicting dynamic emotion ratings in music is considerably more difficult (as opposed to predicting a static overall rating) as it involves accounting for temporal evolution of emotion with music signal. The Emotion in Music Challenge at the 2014 Mediaeval Workshop [13] led to several investigations towards capturing the emotional content in music using low level frame-wise features. Some of the successful schemes involved using a recurrent neural network [14], multi-level regression systems [15] as well as state space models [16]. These methods perform well in predicting the emotional dimensions from low level features. However, they fail to explain the temporal evolution of emotion and its relation to the features. Moreover, all the above systems use a large number of features to predict the affective ratings, rendering feature analysis difficult. To overcome these shortcomings, we propose a gradient boosting-based [17] *Boosted Ensemble of Single feature Filter* (BESiF) method. Given a set of frame-wise features, the BESiF model sequentially learns filters on a selected set of features. The model later performs a weighted combination of the filtered feature values to provide the prediction for affect ratings. We obtain Signal to Noise Ratio (SNR) improvement by a factor of 1.92 and 1.06 for arousal and valence prediction when compared to the next best baseline algorithm. The BESiF model uses a small set of 14/16 features for arousal/valence prediction when compared to the available 6000+ features used by the baseline algorithms. We analyze the output from one of the filters for arousal prediction and interpret the transformation of feature values which contribute towards the final prediction.

In the next section we describe the database used in experiments, followed by the details of affect prediction in Section 3. Section 4 presents the results and conclusions are presented in Section 5.

## 2. DATABASE

We use the music dataset provided in the *Emotion in Music Challenge* at the 2014 Mediaeval Workshop [13] to evaluate the BESiF algorithm. The data set consists of 1744 songs from different musical genres. 45 second clips were selected from each song in the dataset and assigned emotion labels by at least 10 annotators (at a rate of 2 frames/second). More
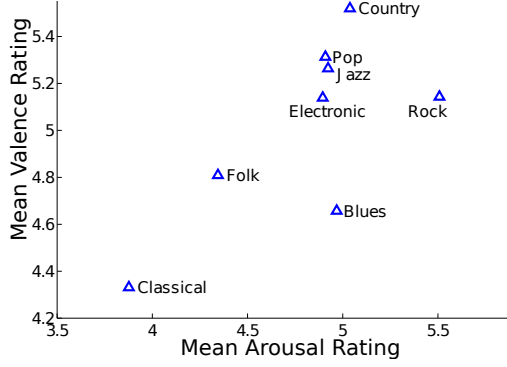
**Fig. 1**. Average arousal-valence ratings in each music genre

details about the data set and annotation process can be found in [18].

To further motivate our study of emotion in music, we analyze the relation between genre and emotional variability present in this database. We plot the average emotion ratings per genre in Figure 1. We observe that the annotated emotion ratings follow intuitive trends along high level music categories such as genres. As an example, notice how *rock* music displays high arousal, while *country* music is high valence. *Classical* music on the other hand has both low arousal and low valence. This suggests that a relation exists between the style of music and its emotional content. This further encourages automatic prediction of affect in music with potential to impact music design, recommendation and understanding music perception.

For the purpose of prediction, we use the mean dynamic annotations for arousal-valence as the gold truth in accordance with the challenge task [13]. Moreover, the first 15 seconds of annotations were excluded from consideration, to allow the dynamic annotations to stabilize. We use a set of 6000+ Opensmile [19] features supplied during the challenge [18]. These features are functionals of various spectral and frequency properties of music signals (Mel Frequency Bank, Fundamental frequency etc.), extracted at the rate of 2 frames/second (same as annotation frame rate.) Out of 1744 songs in the dataset, we use a split of 744, 300 and 700 songs as the train, development and test set respectively. The 744 files for training are as provided during the challenge. The development and testing set are randomly selected from a separate set of 1000 files. In the next section, we describe our training methodology to predict the affective ratings using the provided feature set.

## 3. AFFECT PREDICTION

Through our experiments, we not only aim to maximize the prediction quality, but also understand the relation between these low level signal features and the affective dimension. In this work, we focus on minimizing the mean squared error between predicted and true affect ratings.

We denote the true affect ratings (arousal/valence) for a file $f$ with $N$ frames (arousal/valence) as the row vector $\boldsymbol{t}(f) = [t_1^f, .., t_n^f, .., t_N^f]$ and the corresponding time series of feature vectors as $\boldsymbol{X}(f) = [\boldsymbol{x}_1^f, .., \boldsymbol{x}_n^f, .., \boldsymbol{x}_N^f]$, where $\boldsymbol{x}_n^f$ is

a $D$-dimensional feature vector. The $d^{\text{th}}$ row of $\boldsymbol{X}(f)$ represents the values over time for the $d^{\text{th}}$ feature and we represent that as $\boldsymbol{x}_d(f) = [x_{d,1}^f, ..., x_{d,N}^f]$. A function $M(\boldsymbol{x}_n^f)$ maps the feature vector $\boldsymbol{x}_n^f$ to the one dimensional affect rating. We represent the time series vector $[M(\boldsymbol{x}_1^f), ..., M(\boldsymbol{x}_N^f)]$ obtained after mapping as $\boldsymbol{M}(\boldsymbol{X}(f))$. The mean square error $L_f$ between the mapped and true ratings as obtained for the file $f$ is shown in equation 1 ($||\ ||_2$ represents the $L_2$ norm).

$$L_f = ||\boldsymbol{t}(f) - \boldsymbol{M}(\boldsymbol{X}(f))||_2^2 = \sum_{n=1}^{N} (t_n^f - M(\boldsymbol{x}_f^n))^2 \quad (1)$$

Given the set of files in the training set, we learn the function $M$ by optimizing the sum of squared error losses, $\mathcal{L}$, as defined below.

$$\mathcal{L} = \sum_{f \in \text{Training set}} \frac{1}{2} L_f = \sum_{f \in \text{Training set}} \frac{1}{2} ||\boldsymbol{t}(f) - M(\boldsymbol{X}(f))||_2^2 \quad (2)$$

One can assume any functional form for $M$ before optimizing the cost function $\mathcal{L}$. For the problem of interest, several schemes were proposed during the Mediaeval challenge [13]. However, these schemes are either too simple to capture the complex relationship between the acoustic features and the abstract affective space (e.g. linear regression) or are difficult to interpret (e.g. Recurrent Neural Networks). In this work, we present a new gradient boosting [17] based *Boosted Ensemble of Single-feature Filters*, which overcomes the shortcomings of both these categories of models. The BESiF model is an ensemble of filters trained sequentially on one feature at a time and the final prediction is given as the weighted sum of the filter outputs. We provide the BESiF model training algorithm below along with a brief description of the gradient boosting method.

### 3.0.1. Boosted Ensemble of Single feature Filters (BESiF)

The BESiF model consists of an ensemble of filters operating over the feature time series, learnt using gradient boosting. Gradient boosting is a general technique for learning an ensemble of weak learners applicable in the cases of arbitrarily differentiable loss functions (e.g. mean squared error loss). We represent an ensemble of $K$ weak learners $\{\tilde{h}_1, \tilde{h}_2, ..., \tilde{h}_K\}$ as $M_K$, where the prediction $\boldsymbol{M}_k(\boldsymbol{X}_f)$ is given as shown below.

$$\boldsymbol{M}_K(\boldsymbol{X}_f) = \sum_{k=0}^{K} \tilde{h}_k \quad (3)$$

The base learners $\{\tilde{h}_1, \tilde{h}_2, ..., \tilde{h}_K\}$ are learnt sequentially. The first base learner ($\tilde{h}_0$) is initialized to be a constant model obtained by solving the optimization problem shown in (4). $\bar{\boldsymbol{1}}$ represents a vector of ones of the size same as the target affect variable $\boldsymbol{t}(f)$.

$$\tilde{h}_0 = \gamma_0 = \arg\min_{\gamma_0} \sum_{f \in \text{training set}} ||\boldsymbol{t}(f) - \gamma_0 \times \bar{\boldsymbol{1}}||_2^2 \quad (4)$$

Subsequently, new regressors $\tilde{h}_k$ are added by solving the following optimization.

$$\tilde{h}_k = \arg\min_{h_k} \sum_{f \in \text{training set}} \big|\big| \boldsymbol{t}(f) - M_k(\boldsymbol{X}(f)) \big|\big|_2^2$$
$$= \arg\min_{h_k} \sum_{f \in \text{training set}} \big|\big| \boldsymbol{t}(f) - \Big( M_{k-1}(\boldsymbol{X}(f)) + h_k \Big) \big|\big|_2^2 \quad (5)$$

However the optimization problem in equation 5 is not easy to solve and, in practice optimization is performed iteratively using the steepest descent method [20]. This is equivalent to fitting base learners to a set of *pseudo residuals* (defined for the current problem in equation 6) and learning the weights of base learners using a one-dimensional optimization. For more details on gradient boosting please refer to [17]. In our case, we chose the set of base learners to be Finite Impulse Response (FIR) filters operating on a single feature. We chose this feature probabilistically, with the probability of selection proportional to its absolute correlation with the pseudo residuals. We summarize the training algorithm for the BESiF model below.

Training algorithm for BESiF models:
- Initialize $M_0$ with a constant model $\tilde{h}_0$ (equation 4).
- For $k = 1$ to $K$
  - Computing the pseudo-residuals $\boldsymbol{r}_k(f) = [r_{k,1}^f, ..., r_{k,N}^f]$ for each file in the training set.

$$\boldsymbol{r}_k(f) = -\frac{\partial\left( \frac{1}{2}\big|\big| \boldsymbol{t}(f) - M\big(\boldsymbol{X}(f)\big) \big|\big|_2^2 \right)}{\partial M\big(\boldsymbol{X}(f)\big)} \Bigg|_{\substack{\text{at } M(\boldsymbol{X}(f))= \\ M_k(\boldsymbol{X}(f))}} \quad (6)$$
$$= \boldsymbol{t}(f) - M_k\big(\boldsymbol{X}(f)\big)$$

  - Randomly selecting a feature: In the next step, we randomly select one of the $D$ features. The probability of selecting a feature is proportional to the absolute correlation of feature with $\boldsymbol{r}_k(f)$. Let $d$ be index of the selected feature with values for the file $f$ represented as $\boldsymbol{x}_d(f) = [x_{d,1}^f, ..., x_{d,N}^f]$.
  - Learning a filter to predict the pseudo residuals using the selected feature: Given the filter length $L$, we learn the filter coefficients $\boldsymbol{w}_k = \{w_k^1, .., w_k^L\}$ to predict residuals. These filter coefficients are convolved with the selected feature to obtain with residuals as the target outputs. Coefficients $\boldsymbol{w}_k$ are obtained by solving the optimization problem mentioned in equation 8. $\langle \boldsymbol{w} * \boldsymbol{x}_d(f) \rangle$ represents the convolution of the selected feature with the filter coefficients and is denoted by $h_k(f)$.

$$h_k(f) = \langle \boldsymbol{w} * \boldsymbol{x}_d(f) \rangle \quad (7)$$

$$\boldsymbol{w}_k = \arg\min_{\boldsymbol{w}} \sum_{f \in \text{training set}} \big|\big| \boldsymbol{r}_k(f) - \langle \boldsymbol{w} * \boldsymbol{x}_d(f) \rangle \big|\big|_2^2 \quad (8)$$

  - Computing weights of the base learners: After obtaining the filter coefficients, we compute the scalar $\gamma_k$ to weigh the filter outputs $h_k(f)$. We solve the following one dimensional problem to obtain $\gamma_k$ using backtracking algorithm [21].

| Model | SNR ($\sigma_{\text{signal}}^2/\sigma_{\text{error}}^2$) | |
|---|---|---|
| | Arousal | Valence |
| Linear Regression + smoothing | 1.39 | 1.37 |
| Greedy Linear Regression + smoothing | 1.27 | 1.21 |
| Least squares boost + smoothing | 1.47 | 1.44 |
| BESiF | 2.83 | 1.53 |
| Energy in signal ($\sigma_{\text{signal}}^2$) | 0.11 | 0.06 |

**Table 1**. SNR values for affect rating prediction using the baseline and the proposed BESiF models.

$$\gamma_k = \arg\min_{\gamma} \sum_{\substack{f \in \\ \text{training set}}} \big|\big| \boldsymbol{t}(f) - \Big( M_{k-1}(\boldsymbol{X}(f)) + \gamma \times h_k(f) \Big) \big|\big|_2^2$$
$$(9)$$

  - Updating the model: After obtaining $\boldsymbol{w}_k$ and $\gamma_k$, the predicted outputs for a file $f$ are obtained as

$$M_k(\boldsymbol{X}_f) = M_{k-1}(\boldsymbol{X}_f) + \tilde{h}_k(f) =$$
$$M_{k-1}(\boldsymbol{X}_f) + \big( \gamma_k \times h_k(f) \big) \quad (10)$$

- End For

## 4. EXPERIMENTS AND DISCUSSION

We use the proposed BESiF model to predict the arousal and valence ratings from the low level signal features. As the function $M\big(\boldsymbol{X}(f)\big)$ can assume several functional forms, we chose three other methods as baseline models for comparison. The first baseline method performs linear regression followed by a smoothing operation as was proposed in our previous work [18]. The other two baselines involve techniques such as sequential selection of features and boosting, like the BESiF model. We describe these methods in detail below.

1. Linear regressor + smoothing: In this method, we use linear regression on the entire feature set followed by a smoothing operation to predict the affective ratings. Our analysis in the past work [18] showed that the affective signals evolve rather smoothly. The linear regressor computes the affective rating using the provided features and smoothing is used to incorporate local temporal context. The smoothing operation is fundamentally a moving average operation where output at a frame is recomputed as an average of predictions over a local window. This operation also helps to remove any high frequency noise added during regression.

2. Greedy linear regressor + smoothing: This method is same as the previous baseline method, except for a greedy selection of a few features for regression. Similar to the BESiF training algorithm, features are added sequentially based on their correlation with the residual at each iteration. However, note that after addition of every new feature, the algorithm re-optimizes the regression coefficients for each selected feature. This may lead to the problems associated with *curse of dimensionality* and high computation cost. The total number of features added are tuned on the development set. BESiF model does not suffer from this problem as filter coefficients are determined only for a single feature at a time. The final outputs after regression are again smoothed using a moving average filter.

3. Least squares boost + smoothing: Least squares boost [22] is another class of boosting algorithm used to optimize
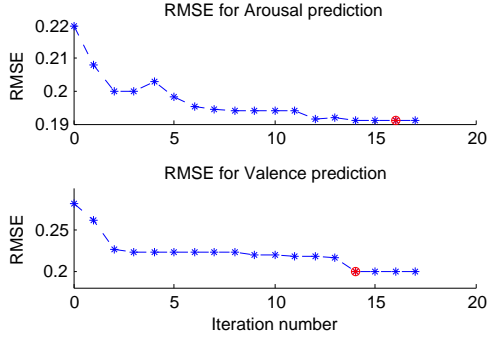
**Fig. 2**. RMSE ($\sigma_{\text{error}}$) of BESiF model on the test set against the number of base classifiers. The red point denotes the count of base filters chosen based on development set.



**Fig. 3**. (a) Target arousal values, (b) filtered feature values and (c) raw feature values for a selected file in the testing set.

squared error loss functions. However, this algorithm uses all the features at every iteration to predict the residuals. We use a regression tree [23] as our base learner in this case. Note that unlike BESiF, the regression trees can not account for the temporal relationship between the residuals and the feature time series. We again smooth the outputs using a moving average filter to account for smooth temporal evolution.

We present the Signal to Noise Ratio (SNR) for affective rating predictions using the baseline methods and the BESiF model in Table 1. SNR is computed as the ratio of energy in the true arousal/valence signal ($\sigma_{signal}^2$) and energy in the prediction error ($\sigma_{error}^2$). The length $L$ of the FIR filters for BESiF model and the length of moving average filters for the baseline methods are tuned on the development set. The number of base classifiers for BESiF and the least squares boost models are also tuned on the development set.

### 4.1. Discussion

From the results in Table 1, we observe that a substantial gain in arousal using the BESiF model is achieved over all other baselines. In our previous work [18], we showed that the smoothing operation added context from neighboring windows, thus improving the prediction. However the regression design was decoupled from smoothing and the choice of filter during smoothing was ad-hoc. In the BESiF model, we overcome this drawback by incorporating filter design within the regression framework. Our base learners, i.e., the single-feature filters not only learn the mapping from the features to the affective dimension, but also incorporate the temporal context into account during prediction. We observe that the performance is particularly poor for greedy linear regression. We performed further investigation into this system and observed that even a backward feature selection (sequential removal of features starting from all features [24]) leads to degradation in performance. This suggests that removal of any feature leads to a degradation in performance of linear regression. The least squares boost algorithm in closest to BESiF in terms of performance. In general, boosting algorithms lead to strong regressors, therefore the better performance. However, decoupling of regressor design and smoothing again leads to poorer performance for least squares boost when compared to the BESiF model.
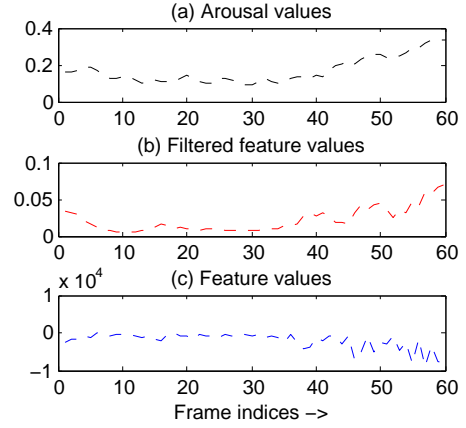
We plot the Root Mean Square Error (RMSE, $\sigma_{\text{error}}$) of the BESiF model on the test set against the number of base filters in Figure 2. We observe that the performance of BESiF model saturates approximately within 15 single feature filters for both valence and arousal. This observation is particularly interesting from the point of view of understanding the relation between the low level features and affective dimensions. We observed that 15 (out of 16) and 13 (out of 14) features selected for arousal and valence respectively are spectral features (Mel Frequency Cepstral Coefficients, Mel Filter Bank energies [19]) with a fundamental frequency (F0) statistic feature appearing once in both the cases. This reflects that most of the emotional information in music is associated with the evolution of spectral characteristics of the music. Since the final prediction is a weighted sum of the filtered feature values, one can also analyze a feature of interest and its contribution to the final prediction. For instance, we plot the target arousal rating, a selected feature (mean of absolute values for an MFB coefficient), and filtered feature values in Figure 3. We observe that the filtered feature values follow the same trend as the target values (despite the scales being different). Moreover, the filtered values are smoother than the feature itself, indicating removal of high frequency noise after the filtering operation. This is consistent with the premise of smooth evolution of affective ratings.

## 5. CONCLUSION

Music signals have been shown to carry emotional information. In this work, we present a novel BESiF scheme to predict the affective dimension of arousal and valence from low level audio signal features in music. This scheme is designed not only to better predict the affective ratings, but also to add insights and interpretability to the prediction process. We show that the BESiF system beats several comparable baseline methods, using only a handful of features. We interpret patterns as observed in a feature time series after filtering and compare it to the target value.

In the future, we aim to enhance our system by jointly predicting the affective dimensions instead to allow understanding of joint dynamics between valence and arousal. With availability of more data, one could also analyze the relation of system parameters (e.g. filters, base learner weights) with music categories such as genres. Further improvement in pre-

diction is also possible by using better feature selection techniques (e.g. using dynamic programming [25]) and filter design methods.

## REFERENCES

[1] Lars-Olov Lundqvist, Fredrik Carlsson, Per Hilmersson, and Patrik Juslin, "Emotional responses to music: experience, expression, and physiology," *Psychology of Music*, 2008.

[2] Patrik N Juslin and John A Sloboda, *Music and emotion: Theory and research.*, Oxford University Press, 2001.

[3] Jaak Panksepp, "The emotional sources of" chills" induced by music," *Music perception*, pp. 171–207, 1995.

[4] Tuomas Eerola and Jonna K Vuoskoski, "A comparison of the discrete and dimensional models of emotion in music," *Psychology of Music*, 2010.

[5] Nikos Malandrakis, Alexandros Potamianos, Georgios Evangelopoulos, and Athanasia Zlatintsi, "A supervised approach to movie emotion tracking," in *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*. IEEE, 2011, pp. 2376–2379.

[6] Susan Hallam and John Price, "Research section: can the use of background music improve the behaviour and academic performance of children with emotional and behavioural difficulties?," *British Journal of Special Education*, vol. 25, no. 2, pp. 88–91, 1998.

[7] Michael K Hul, Laurette Dube, and Jean-Charles Chebat, "The impact of music on consumers' reactions to waiting for services," *Journal of Retailing*, vol. 73, no. 1, pp. 87–104, 1997.

[8] Erik M Schmidt and Youngmoo E Kim, "Modeling musical emotion dynamics with conditional random fields.," in *ISMIR*, 2011, pp. 777–782.

[9] Yi-Hsuan Yang and Homer H Chen, "Machine recognition of music emotion: A review," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 3, no. 3, pp. 40, 2012.

[10] Youngmoo E Kim, Erik M Schmidt, Raymond Migneco, Brandon G Morton, Patrick Richardson, Jeffrey Scott, Jacquelin A Speck, and Douglas Turnbull, "Music emotion recognition: A state of the art review," in *Proc. ISMIR*. Citeseer, 2010, pp. 255–266.

[11] Jacquelin Speck, Erik M Schmidt, Brandon G Morton, and Youngmoo E Kim, "A comparative study of collaborative vs. traditional annotation methods," *ISMIR, Miami, Florida*, 2011.

[12] Mohammad Soleymani, Micheal N Caro, Erik M Schmidt, Cheng-Ya Sha, and Yi-Hsuan Yang, "1000 songs for emotional analysis of music," in *Proceedings of the 2nd ACM international workshop on Crowdsourcing for multimedia*. ACM, 2013, pp. 1–6.

[13] Anna Aljanaki, Y.H. Yang, and M. Soleymani, "Emotion in music task at mediaeval 2014," in *Mediaeval 2014 Workshop, Barcelona, Spain, October 16-17*, 2014.

[14] Eduardo Coutinho, Felix Weninger, Björn Schuller, and Klaus R Scherer, "The munich LSTM-RNN approach to the mediaeval 2014 emotion in music task," in *Mediaeval 2014 Workshop, Barcelona, Spain, October 16-17*, 2014.

[15] Yuchao Fan and Mingxing Xu, "Mediaeval 2014: Thu-hcsil approach to emotion in music task using multi-level regression," in *Mediaeval 2014 Workshop, Barcelona, Spain, October 16-17*, 2014.

[16] Konstantin Markov and Tomoko Matsui, "Dynamic music emotion recognition using state-space models," 2014.

[17] Jerome H Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of statistics*, pp. 1189–1232, 2001.

[18] Naveen Kumar, Rahul Gupta, Tanaya Guha, Colin Vaz, Maarten Van Segbroeck, Jangwon Kim, and Shrikanth Narayanan, "Affective feature design and predicting continuous affective dimensions from music," in *Mediaeval Workshop, Barcelona, Spain*, 2014.

[19] Florian Eyben, Martin Wöllmer, and Björn Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proceedings of the international conference on Multimedia*. ACM, 2010, pp. 1459–1462.

[20] Jan Snyman, *Practical mathematical optimization: an introduction to basic optimization theory and classical and new gradient-based algorithms*, vol. 97, Springer Science & Business Media, 2005.

[21] Larry Armijo et al., "Minimization of functions having lipschitz continuous first partial derivatives," *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.

[22] Jerome H Friedman, "Stochastic gradient boosting," *Computational Statistics & Data Analysis*, vol. 38, no. 4, pp. 367–378, 2002.

[23] Roger J Lewis, "An introduction to classification and regression tree (cart) analysis," in *Annual Meeting of the Society for Academic Emergency Medicine in San Francisco, California*, 2000, pp. 1–14.

[24] Isabelle Guyon and André Elisseeff, "An introduction to variable and feature selection," *The Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, 2003.

[25] Dimitri P Bertsekas, Dimitri P Bertsekas, Dimitri P Bertsekas, and Dimitri P Bertsekas, *Dynamic programming and optimal control*, vol. 1, Athena Scientific Belmont, MA, 1995.