
A Supervised Approach To Musical Chord Recognition

Pranav Rajpurkar
Brad Girardeau
Takatoki Migimatsu

Stanford University, Stanford, CA 94305 USA

PRANAVSR@STANFORD.EDU
BGIRARDE@STANFORD.EDU
TAKATOKI@STANFORD.EDU

Abstract

In this paper, we present a prototype of an online tool for real-time chord recognition, leveraging the capabilities of new web technologies such as the Web Audio API, and WebSockets. We use a Hidden Markov Model in conjunction with Gaussian Discriminant Analysis for the classification task. Unlike approaches to collect data through web-scraping or training on hand-labeled song data, we generate symbolic chord data programmatically. We improve the performance of system by substituting the usually tried Chroma features with a novel set of Chroma DCT-Reduced log Pitch features to push test accuracy on clean data to 99.96%. We finally propose a set of modifications to have the system predict with speed and accuracy in realtime.

1. Introduction

There is significant value in an automated tool to determine chords from audio. Knowing the progressions of chords underlying the melodies is an essential part of understanding, playing, and building on the music. To a curious learner of music, such a tool creates the opportunity to play a new pop song without meticulously hand-labelled chord tags. Equally useful to a learner is being able to receive feedback concerning the accuracy with which a chord was played, making such a system a good automated feedback tool, capable of being plugged into an MOOC music class. To a song writer, the system is useful for exploring chords supporting the melodic content of the song.

Furthermore, the use of such a system extends into

other machine learning tasks. The tasks of identifying a song from its waveform data, and of classifying the genre of song can be linked to finding the chord progressions underlying the harmonic content of the song. Hand-labelling chord names and marking chord changes in a song takes a lot of manual time and effort. An automated tool for this process saves time, and allows the development of new musical tools and research.

There has been progress in chord recognition research. A few have built real-time systems that have shown to achieve promising results (??). However, these have not leveraged the web to make a chord-recognition system accessible online. We build a real-time online chord recognition system that makes use of the novel HTML5 capabilities such as the WebAudio api and WebSockets, and detail the offline training strategies and online challenges posed by the novel adaptation.

2. Data Generation

The chord prediction task is a multiclass classification task. In music, a chord is a harmonic set of notes sounding simultaneously. We choose the minor and major chords, the two most common sets of chords in popular music to classify on. Using the traditional twelve pitch scale (C, C#, D, D#, E, F, F#, G, G#, A, A#, B), we have 24 such distinct chords.

There are different ways of playing the same chord. The Cmaj chord, for instance, is the set of three notes C, E, and G played simultaneously. On a piano, these notes can be played on different octaves. Played on the fourth octave, the Cmaj would have the notes C4, E4, G4. It is also possible to play Cmaj in with E as the lowest note - E4, G4, C5 (first inversion form) and with G as the lowest note G4, C5, E5 (second inversion form).

To train the system, we generate training data programmatically. This has found to have advantages over hand-labelling song data in its ability to gener-

ate lots of training data (?). We generate midi files for each of the 24 chords, taking into account 8 octaves, and 3 inversion forms, to generate a total of 576 midi files. We then use an audio synthesizer Timidity++ to convert the 576 generated midi files, in conjunction with 3 soundfonts for piano, guitar, and violin, to generate a total of 1728 audio files in wav format.

In musical chord recognition, feature extraction operates over frames. The generated wav files, which are an average of 4 seconds in length, are first split into frames of window size 100ms each, and a n-dimensional feature vector is extracted for each frame. We label each frame with the label of the chord on its corresponding sound file, to generate 69,120 examples. We use 80

3. Paralleling Speech Recognition

The pipeline of a chord recognition system is very similar to that of a speech recognition one and relies on techniques that were originally applied to speech recognition tasks. The use of the Hidden Markov Models (?), and the division of a sound file into frames on which the prediction task is performed, are two such techniques which have been reproduced in identifying chords. However, the task of finding chords is also different from speech tasks in a few ways, and these differences can be exploited to specialize a system in the task of chord recognition.

3.1. Feature Extraction

One important difference between the two surfaces in the choice of features for the tasks. Mel-frequency cepstrum coefficients (MFCC) have been the dominant features used in speech recognition systems. These represent the short-term power spectrum of a sound. It has been found that MFCCs are closely related to timbre - a quality that is considered able to distinguish between a voice, string instrument, wind instrument, and percussion instrument. These have traditionally been seen as poor features for chord recognition, since they discard the pitch content of the sound, but are useful in setting a baseline benchmark for such a system.

Chroma features are commonly used for chord recognition tasks (?). It is a representation in which the entire spectrum of sound frequencies is distributed into 12 bins representing the traditional twelve pitch scale. An advantage of Chroma features is that they are invariant to octaves and inversions. We use the Matlab Chroma Toolbox to extract Chroma features for the frames INSERT CITATION. Figure 1 shows the extracted Chroma features for the C major chord. The

Table 1. Accuracies for MFCC and Chroma on the binary classification task of distinguishing between major and minor chords

FEATURES	TEST ACCURACY
MFCC	51.0%
CHROMA	97.7%

energy spikes at C, E, and G, the notes constituting the C major, supporting the idea that Chroma encodes the harmonic content of a chord.

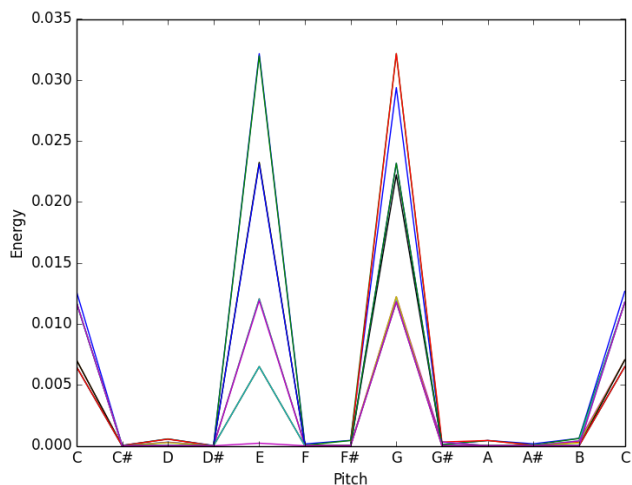


Figure 1. Chroma features for C major in various octaves and inversions

To test the performance of Chroma features against MFCC features, we start with a binary classification problem of distinguishing major chords from minor chords. An SVM with RBF kernel

$$K(x, z) = \exp(-\gamma \|x - y\|^2)$$

is trained, with regularization and kernel parameters ($\gamma = 1$ and $C = 100$). Table 1 summarizes the results, and confirms that chroma features are much better suited to the task of chord recognition than MFCCs.

4. Initial Models

4.1. Frame Model

With Chroma established as good features for the chord recognition model, we can now extend to the multiclass classification problem of determining the exact chord. We first use multinomial logistic regression, also called softmax regression, as our initial frame

Table 2. Softmax Regression frame model train and test accuracies

SET	ACCURACY
TRAINING	51.2%
TESTING	32.1%

Table 3. Comparisons of accuracies of mixer models with softmax frame model

MODEL	TEST ACCURACY
MIDDLE FRAME	6.7%
MAX COUNT	33.3%
INDEPENDENCE MIXER	48.3%

model. The frame model is responsible for making predictions on individual frames. Table 2 shows the accuracies achieved by the softmax classifier on the training and test set.

4.2. Mixer Model

Our frame model outputs a prediction for each frame. Our final classification task, however, is on an audio file, which consists of a sequence of f frames. Let us first make the simplifying assumption that a test sound file consists of a single chord being played.

We now define a mixer model, which is a model for collecting and using the results on individuals frames outputted by the frame model. A simple mixer model, the Middle Frame model, outputs the result for the entire file based on the frame model’s output for middle frame in the file. Another simple model, the Max Count model, counts the most frequent prediction made across all of the frames.

Consider another such model, we call the Independence Mixer model, which assumes that the prediction on each frame is independent of the prediction on other frames. The probability that chord y is the single chord played in the file is calculated by considering the probability that y is the chord played at each frame. For a test example, our predicted output is $y_p = \underset{y}{\operatorname{argmax}} p(y|X) = p(y|x_1, x_2, \dots, x_f) = \prod_{i=1}^f p(y|x_i)$.

Note that each $p(y|x_i)$ is given by our softmax frame model. The accuracies achieved with different Mixer Models are summarized in Table 3.

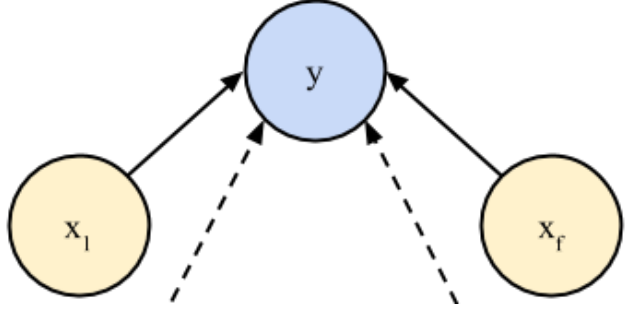


Figure 2. Bayesian network of Independence Mixer model

Table 4. GDA frame model accuracies on an 80/20 train/test split

SET	ACCURACY
TRAINING	67.8%
TESTING	53.9%

5. Improved models

5.1. Improved frame model

Softmax regression is a learning algorithm that models $p(y|x)$, the conditional distribution of the chord given the extracted frame features. We now look at a model that tries to model $p(x|y)$ and $p(y)$: Gaussian Discriminant Analysis (GDA). (?). We model $p(x|y)$ using a multivariate normal distribution. Although we model each gaussian with different means, we assume they all share the same covariance matrix: $x|y = i \sim \mathcal{N}(\mu_i, \Sigma)$. Since our aim is to make the system independent of any specific genre, we model $p(y) = 1/24$, a model in which all chords are equally likely.

With this improved frame model, we achieve a test accuracy of INSERT%. Figure INSERT shows the confusion matrix. The most common misclassifications are ones between a major chords and its corresponding minor version. This is explained by the fact that a major and a minor chord for any note have two out of three notes in common.

Softmax regression is a learning algorithm that models $p(y|x)$, the conditional distribution of the chord given the extracted frame features. We now look at a model that tries to model $p(x|y)$ and $p(y)$: Gaussian Discriminant Analysis (GDA). (?). We model $p(x|y)$ using a multivariate normal distribution. Although we model each gaussian with different means, we assume they all share the same covariance matrix: $x|y = i \sim \mathcal{N}(\mu_i, \Sigma)$. Since our aim is to make the system independent of any

specific genre, we model $p(y) = 1/24$, a model in which all chords are equally likely. With this improved frame model, we achieve a test accuracy of INSERT%. Figure INSERT shows the confusion matrix. The most common misclassifications are ones between a major chords and its corresponding minor version. This is explained by the fact that a major and a minor chord for any note have two out of three notes in common.

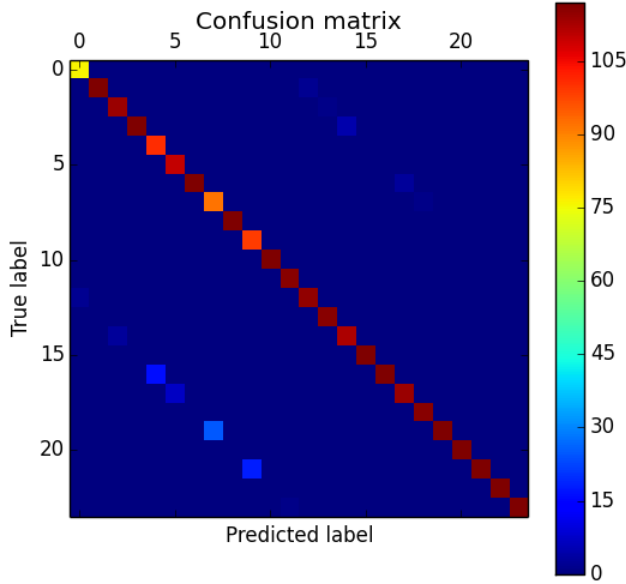


Figure 3. Confusion matrix of improved frame model

5.2. Improved mixer model

Earlier, we had imposed the constraint that chords could not change in a wav file. Our next model allows us to drop that constraint. We now use a Hidden Markov Model to predict the chord sequence in sound files, allowing us to determine chord changes in a file (?). Firstly, the emission probabilities $p(x|y)$ are modelled by our GDA frame model. While state transitions are usually learned in chord recognition tasks, (?), since each genre of music has a different distribution of transitions, assuming uniform state transitions allows us to remain flexible to any genre of music. We determine the most likely state sequence by the Viterbi decoding. Table 4 summarizes the accuracies acheived by the improved mixer model trained on different sets of instruments.

Earlier, we had imposed the constraint that chords could not change in a wav file. Our next model allows us to drop that constraint. We now use a Hidden Markov Model to predict the chord sequence in sound files, allowing us to determine chord changes in

Table 5. HMM accuracies training and testing on different data sets

TRAINING DATA	TESTING DATA	ACCURACY
PIANO GUITAR	PIANO + ... GUITAR + ...	INSERT% INSERT%

a file (?). Firstly, the emission probabilities $p(x|y)$ are modelled by our GDA frame model. While state transitions are usually learned in chord recognition tasks, (?), since each genre of music has a different distribution of transitions, assuming uniform state transitions allows us to remain flexible to any genre of music. We determine the most likely state sequence by the Viterbi decoding. Table 4 summarizes the accuracies acheived by the improved mixer model trained on different sets of instruments.

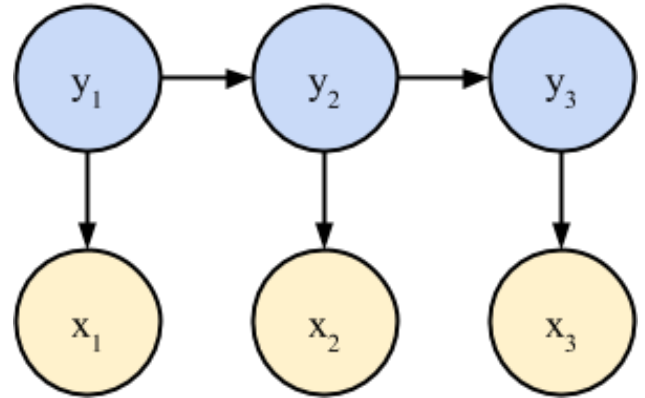


Figure 4. Hidden Markov Model to determine most likely chord state sequences

6. Improving features

Chroma features, in their invariance to octave and inversions, make great features for the chord recognition task. To boost the accuracy further would require features were invariant of instruments (?). In CRP (Chroma DCT-Reduced log Pitch) is a chroma-based audio feature that boosts the degree of timbre invariance. The general idea is to extract Chroma features, and then discard timbre-related information similar to that expressed by MFCCs, in effect leaving the information related to pitch. Table 5 summarizes the accuracies achieved by the new CRP features in relation to the Chroma features.

INSERT GRAPH SHOWING CHROMA INVARIANCE TO INSTRUMENTS VS CRP INVARIANCE

Table 6. Chroma vs CRP

TRAINING DATA	TESTING DATA	ACCURACY
PIANO GUITAR	PIANO GUITAR	INSERT% INSERT%

ing on hand-labeled song data, we generate symbolic chord data programmatically. We improve the performance of system by substituting the usually tried Chroma features with a novel set of Chroma DCT-Reduced log Pitch features to push test accuracy on clean data to 99.96%. We finally propose a set of modifications to have the system predict with speed and accuracy in realtime.

7. Live system consideration

A live system presents new challenges for chord recognition, while pushing the extents of possible uses of the system. One of the most important considerations for a live system to take into account is noise. It is important for a system to not predict any chord in this timespan.

Our realtime system makes use of the Web Audio API to capture live audio. Every 800ms, we encode the audio in wav format, and transfer it to a server using WebSockets. We then extract features for every 100ms frame in the wav file, and predict the most likely chord sequence in the 800ms using the HMM.

7.1. Handling Noise

At this stage, we are able to determine whether the 800ms segment is a noise or a chord. Using the knowledge that it usually takes a few seconds before chords change, we can then post-process the output by looking at the number of times chord changes were predicted by the HMM in the 800ms segment. If any chord lasts for more than 400ms in the prediction, then we output the chord as our prediction for the 800ms segment. Otherwise, we understand that segment of sound as consisting of noise.

7.2. Window Size

800ms was seen to be the optimal time for recording. Increasing the window time for collecting the recording from 800ms was found to create a highly noticeable delay in the live system. On the other hand, reducing the window time decreased accuracy. It was also harder to distinguish between noise and chords with a smaller window.

8. Conclusion

In this paper, we present an implementation of a real-time chord recognition system. Using novel web technologies, we are able to apply/ We use a Hidden Markov Model in conjunction with Gaussian Discriminant Analysis for the classification task. Unlike approaches to collect data through web-scraping or train-