

Thesis

A Thesis
Presented to
The Division of Mathematics and Natural Sciences
Reed College

In Partial Fulfillment
of the Requirements for the Degree
Bachelor of Arts

Emily Palmer

May 2018

Approved for the Division
(Mathematics)

Andrew Bray

Acknowledgements

I want to thank a few people.

Table of Contents

| | |
|--|-----------|
| Chapter 1: | 1 |
| 1.1 Introduction: | 1 |
| 1.2 Brief history of text classification | 2 |
| 1.3 Methods for classification. | 2 |
| 1.4 Feature selection | 3 |
| 1.5 A Very Short Introduction to Music Theory | 3 |
| 1.6 Musical features | 4 |
| 1.6.1 Background on Variable and Feature Selection | 4 |
| 1.7 Previous research. | 5 |
| 1.7.1 Previous choices of features | 6 |
| 1.7.2 Previous applications | 7 |
| 1.8 Fanny and Felix Mendelssohn | 8 |
| Chapter 2: | 11 |
| 2.1 About the data and conversion process | 11 |
| 2.2 Pieces used | 11 |
| 2.3 Optical music recognition | 13 |
| 2.4 .krn to R | 17 |
| 2.5 About the functions - MuseR | 18 |
| Chapter 3: EDA | 21 |
| Chapter 4: | 23 |
| 4.1 Feature Selection | 23 |
| Chapter 5: | 25 |
| 5.1 Models | 25 |
| 5.2 Linear Methods for Classification | 25 |
| 5.2.1 Linear Regression: | 25 |
| 5.3 K- nearest - neaighbor | 27 |
| 5.3.1 Logistic Regression | 27 |
| Chapter 6: | 29 |
| 6.1 Model Fit | 29 |
| Chapter 7: | 31 |

| | |
|---|-----------|
| 7.1 Discussion | 31 |
| Conclusion | 33 |
| Appendix A: The First Appendix | 35 |
| References | 37 |

List of Tables

List of Figures

| | | |
|-----|--|----|
| 2.1 | A picture of a score | 12 |
| 2.2 | Shows how basic kern files correspond to sheet music | 16 |

Abstract

The preface pretty much says it all.

Second paragraph of abstract starts here.

Dedication

You can have a dedication here if you wish.

Chapter 1

1.1 Introduction:

In the digital age data sets are everywhere. Billions of new data are generated daily, in banking, social media, or in other scientific studies. These data can be in numerous forms. With the availability of the internet, text classification has become an interesting form of data. While text classification problems are very frequent, similar methods that instead classify music have not been explored as much.

Thinking of text or music as data has its own challenges. If we normally think of data as something with lots of numbers in a spreadsheet, and then running analysis on that spreadsheet. Analysis of text using machine learning can find patterns, authors, or categories of text. Similar methods can in fact be used to classify music.

The essential part of either text or music classification is feature selection. Unlike in a data set of numerical or categorical values, text and music must first go through a processing stage where hopefully features of interest can be extracted before any models can be fit.

For music, the interest is in building a model that uses the extracted features would be able to correctly classify a likely composer for that piece of music. For musically trained humans, this might be an easy task. Some are able to either by listening to a recording, or looking at sheet music, automatically distinguish a piece composed by Bach or Mozart. This becomes more difficult when composers are contemporaries. Mozart and Salieri might be distinguishable to a scholar or music fan, but it might be harder. Harder still is when a piece of music has a disputed composer history. These examples exist throughout music history, most notably in the Renaissance. Music classification has attempted to assign composers to pieces to be disputed to be Josquin Des Prez.

In both text and music classification, we must create features that can be calculated that would give some signal to indicate some unconscious tendency of the composer that would make them distinguishable.

There have been numerous other interesting applications of machine learning to music. Many studies have used neural nets etc trained on composers music to have the computer create a similar composition.

1.2 Brief history of text classification

One of the earliest instances of text classification was on the Federalist papers. (Mosteller & Wallace, 1964). The famous Federalist Papers were written under the pen name ‘Publius’. There are several disputed papers attributed to James Madison or Alexander Hamilton. The authors never admitted authorship, as some of the writings were contradictory to their later political platforms. (Adair, 1944) Historians have often examined the papers using styles of previously known writings of Madison and Hamilton. Their analysis is often partially based on the content of the letters, for example the existence of citing English history is a trait more common to Hamilton. (Ford & Bourne, 1897)

In contrast, using the frequency of words such as and ‘by’, ‘from’, and ‘upon’, Mosteller and Wallace trained the writings on a set of known writings by each author. These unconscious indicators were able to differentiate between the two writers, and when a model was trained, the model was able to identify the likely author of the disputed paper.

1.3 Methods for classification.

Initial approaches often use linear methods for classification. If we are trying to predict the author, $y \in \{\text{Fanny, Felix}\}$, or more generally, $y \in \{\text{list of composers}\}$, given features or predictors $\mathbf{X} = \{X_1 \dots X_p\}$, we can divide the input space* into a collection of regions labeled according to the classification.

We can create linear decision boundaries for K classes where the fitted linear model for the k th indicator response is $\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$. The decision boundary between classes i and j is the set of points for which $\hat{f}_i(x) = \hat{f}_j(x)$, or in other words, the set $\{x : (\hat{\beta}_{i0} - \hat{\beta}_{j0}) + (\hat{\beta}_i - \hat{\beta}_j)^T x = 0\}$ which defines a hyperplane.

Similarly quadratic decision boundaries can be used when we increase our predictor space to include squares and higher polynomials of X . We then fit linear decision boundaries, which then map down to quadratic functions in the original space.*.

Logistic regression is often used when the response is binary. It models the probability that Y belongs to either category. We use the logistic function $p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$. To calculate estimates of $\hat{\beta}_0, \hat{\beta}_1$, we use maximum likelihood. To do this, we choose $\hat{\beta}_0, \hat{\beta}_1$ to maximize the function ...

The Naive Bayes classifier is often used for musical classification as it is good when the dimension p of the features space is large, making density estimation unattractive*. It makes the (naive) assumption that all the features are independent for a given class i .

$$f_i(X) = \prod_{k=1}^p f_{ik}(X_k)$$

The idea of separating Hyperplanes is essential to Support Vector Machines (SVM)
 Linear discriminant analysis: Using
 GDA

1.4 Feature selection

Text analysis, such as in the federalist papers, is read one word after another. Information in piece of music, however, is read in a variety of ways. It can be read left to right note by note, but it can also be read vertically as the harmony, or the notes played together. Also in a piece with several instruments, the above happens at the same time for each instrument. There are also aspects that take place over large sections, such as phrasing, or cadencial patterns. There are rules of counterpoint that are followed throughout the entire piece. Thus we need to find features that can be measured for each piece, or perhaps each measure or instrument, that can describe a certain piece of music. Then we must decide which features are those of rules and practices of classical music, and where the creativeness and individuality of a composer happens. As there are so many different aspects of a piece, the melodic changes, harmonic ** etc, we can end up with very many (x) features to describe one piece. As in this instance when we only have ~ 150 total pieces, we have that $p > n$. We thus need to figure a clever way to decrease our predictor space. This can happen in choosing what features we want to use to model, or using a dimension reduction technique such as PCA.

Most of the musical stylometry papers have focused on composers in the Renaissance, Baroque, and Classical eras. The Mendelssohns were composing in the Romantic period. This choice might be because composers in earlier eras followed rules of counterpoint more exactly, or perhaps had less “expressive” allowances for their composing, thus making features easier, although this is just speculation. There are also more pieces with doubtful authorship in those eras. In addition Computer Assisted Research in the Humanities (CCARH) has a large corpus of encoded music from these times. * have a citation list of all the ones that use these*

1.5 A Very Short Introduction to Music Theory

Western sheet music is presented on a five line staff. The vertical distance between notes (also known as an interval) depends on how many half steps occur between two notes. There are 12 half steps in the western scale. Melodic intervals are defined as the number of half steps between two adjacent notes. Harmonic intervals are defined as the number of half steps between notes played at the same time. Cadences are a type of chord progression, usually occurring at the end of a phrase and especially at the end of a piece. Musical notes are in the set $\text{note} \in \{A, B, C, D, E, F, G\}$. The value of a note can be changed up or down a halfstep, by adding a sharp or flat. There are 12 unique note value per octave. Intervals and chords can be either dissonant or consonant. Consonant intervals sound nice to our ears, whereas dissonant intervals add a sense of tension or unease, which is used to shape the feel of the music. In addition chords and intervals can be minor or major. Minor intervals feel “sad”, whereas major intervals feel “happy”.

1.6 Musical features

Features calculated from music often closely follow music theory. In music, addition to deciding the features, there is also the decision of the scope at where those features take place. They can be features for a given instrument, the entire piece, or each measure. Also windowing techniques can be used where a “window” is created over a given number of bars or notes, and shifts across the whole piece. For each window, a feature is recorded. These can be overlapping windows by creating an “offset” of a number of beats or notes. This produces more data, as instead of one feature for each piece, there is a feature for each window.

Musical features are often thought about as either high level, or low level. While the exact definition of each often varies, low level is often understood to be features such as note frequency, etc. High level features are more about a broader sense of the piece, including chord progressions, etc. High level features are often what music experts use in their analysis, whereas low-level features are more easily done with computer analysis.

Especially in research regarding gene expression and text categorization, data sets have enormous numbers of variables. We use “feature” and “variable” interchangeably, with the exception when features are created from variables, and the distinction will be made in that case. (Guyon & Elisseeff, 2003). There are several variable selection algorithms that select the “important” variables. If we included every variable that we extract from the piece, our model would very likely be overfit.

The start of feature selection is domain knowledge. Thanks to John Cox in the music department for suggesting a list of valuable features. These will be described in Chapter 2.

1.6.1 Background on Variable and Feature Selection

Often, especially in musical research, before analysis is done, numerous features are extracted from the music without knowing a priori which ones will be helpful in identifying a composer’s style. Thus we have to choose which features we want to use in our model. There are many ways to do this. There are existing variable selection algorithms that can help with this process.

Several variable selection algorithms include variable ranking. Variable ranking uses a score function to assign a score to each possible variable. It is a computationally efficient method and is robust against overfitting as it introduces bias but may result in less variance. It is tempting to only include variables that have a high score. However, this possibly leads to redundancy. In addition, variables that are not important by themselves can have a significant performance improvement when considered with other variables. Popular variable ranking methods for classification are single variable classifiers and information theoretic ranking criteria.

Single variable classifiers rank the variable according to their individual predictive power. The predictive power can be measured in terms of error rate, or using the false positive or false negative rate (fpr, fnr). This classifier cannot distinguish variables that perfectly separate the data.

The Information Theoretic Ranking Criteria is used in variable selection. They often rely on estimates of the mutual information between the predictor and response, as given by

$$I(i) = \int_{x_i} \int_y p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)} dx dy$$

where $p(x_i)$ and $p(y)$ are the probability densities of x_i the i^{th} predictor and y the response, and $p(x_i, y)$ is the joint density. $I(i)$ is a measure of dependency between the density of variable x_i and the density of the response y (reword)

After knowing the ranking of a variable we then select which variables will be useful for our model. This is known as variable subset selection. The three most common types of variable subset selection are wrappers, filters, and embedded methods. Filters do not involve any machine learning to create the criterion for variable subset selection. Wrappers on the other hand use the “performance of a learning machine trained using a given feature subset.” Embedded methods perform variable selection in the process of training and are usually specific to given learning machines

All possible subsets of variables is $2^p - 1$, which for large p is often computationally impossible. Strategies like best-first, branch and bound, simulated annealing, and genetic algorithms can help with the computational difficulties.

Wrappers are often thought of as brute force methods. This can be good, as it can reduce overfitting. Two types include forward selection and backward elimination. These both give nested subsets of variables

Often there is a need for dimensionality reduction. Is there a way to combine enough of the information given in the features in a smaller dimensional space? This results in feature creation; using the recorded variables to create new features to fit the model on. These include clustering, basic linear transformations of the input variables, such as PCA/SVD, and LDA. Also more sophisticated linear transformations like Fourier and Hadamard. Two basic goals of these feature creations are that we can achieve a good reconstruction of the data. The second is that we can be most efficient in making our predictors. The first is an unsupervised problem. The second is supervised.

Clustering is in fact a type of feature construction. The group of clustered points thereby becomes a feature. Examples of this include K-means and hierarchical clustering.

1.7 Previous research.

Research on musical stylometry focuses on two areas, data in the form of audio, and data in the form of sheet music. Our analysis uses data in the form of sheet music. To predict a composer, a training data set of pieces of known composer is needed. Then that model can be fit to a testing data set to predict composer. If that shows good predictions, that model can be used on the pieces of unknown authorship.

Musical stylometry can be used disputed authorship, but also to detect distinguishing musical styles of composer, even if there are no disputed pieces. One such

study did both. Backer et. al. first looked at differences in style between J.S. Bach, Telemann, and Handel, Haydn and Mozart. Next they looked at piece BWV 534 which is disputed to be composed by either J.S. Bach, J.L. Krebs, or W.F. Bach (J.S. Bach's son). (Backer & Kranenburg, 2005)

1.7.1 Previous choices of features

Deciding on and extracting features of music is the first step to analysis. Depending on the characteristics of the composer and time period, different features would be useful. Often, features are extracted en masse and then work is done later to determine which features are important or useful in identifying style.

In addition to what kinds of features, in music there is also the question of the scale at where those features take place. They can be features for a given instrument, the entire piece, or each measure. Also windowing techniques can be used where a "window" is created over a given number of bars or notes, and moves through the whole piece. For each window, a feature is recorded. These can be overlapping windows by creating an "offset" of a number of beats or notes. This produces more data, as instead of one feature for each piece, there is a feature for each window, and there can be tens of windows in each piece.

Common types of features used before in music analysis are: Frequencies or fractions of notes, chords, etc are a common low-level feature. These include the fraction of the score that consisted of dissonant sonorities, as well as the fraction of bars that begin with a dissonant sonority. Other features include the type of intervals or consonances present in a piece: perfect consonance, imperfect consonances, and dissonance. In polyphonic pieces, the four types of motion, (parallel, similar, oblique, and contrary) can also be used as features.

Features measuring "stability" are also popular. Stability is computed by dividing the standard deviation of the lengths of the fragment by the mean length of the fragments. It is normalized in this way to be comparable over differing time signatures. (Backer & Kranenburg, 2005)

Markov transition matrices for the rhythms of the pieces, and Markov transition matrices of the pitches in each piece

The above techniques were used to analyze the music of Bach, Handel, Telemann, Mozart and Haydn and compare J.S. Bach, W.F. Bach and J.L. Krebs in an attempt to classify BWV 534. (Backer & Kranenburg, 2005) They use overlapping windowing over each entire composition to produce more data, and avoid issues of dimensionality. They chose a window of 30 bars to create a high enough number of fragments per piece and a low enough variance of the feature values between fragments. They chose to extract 20 features including features of fractions and measuring stability, and entropy.

Additionally, a number of previous papers have focused on Josquin des Prez. This is likely due to the fact that there is a large training and testing data set available in easily analyzable format provided by the Josquin Research Project (citation). In addition there are a number of pieces of disputed authorship that have been attributed to him. Work by Brinkman et al. (Brinkman, Shanahan, & Sapp, n.d.) use machine

learning approaches to evaluate attribution of compositions by des Prez. They used both high level and low level features. The high level features were 9-8 suspensions, oblique motion, contrary motion, similar motion and parallel motion. The low level features were average melodic entropy, normalized pairwise variability index (?), and note-to-note transition probabilities.

Work by Speiser and Gupta (Speiser & Gupta, n.d.) analyzed Josquin and his contemporaries to attempt to classify unknown works. They extracted four categories of features, frequencies of individual notes, frequencies of pairwise interval combinations between each of the voices, Markov transition matrices for the rhythms of the pieces, and Markov transition matrices of the pitches in each piece. In total, this led to a total of 3000 features.

Other work looking at renaissance and baroque composers looked specifically at differences in counterpoint. Since most composers in that era for the most part followed the rules of counterpoint, there is a question of if there are distinguishing differences. Using counterpoint movement types, dissonance distributions, parallel intervals of each kind, and vertical interval distributions, a classifier using a WEKA algorithm, as well as Naive Bayes and a Decision Tree was created that correctly predicted composer 2/3 of the time. (Mearns, Tidhar, & Dixon, 2010)

1.7.2 Previous applications

Most of the previous research has needed to do some kind of feature selection. A lot of features are extracted as a priori we don't know which features are distinguishing.

A modification of a forward selection (Floating Forward Selection(cite)) was used to extract features in order to identify distinguishing style between Bach, Handel, Telemann, Mozart, and Haydn, and then subsequently classify the authorship of BWV 535. (Backer & Kranenburg, 2005) Each composer was compared via creating comparisons of all possible class arrangements, ie (Bach)(Handel), (Bach)(Handel,Telemann), etc. The algorithm extracted features for each class arrangement that distinguished the groups the best. A decision boundary was used for Bach and not Bach, on the features Diss Part, Par thirds, and stab time slice. A k-nearest neighbors classifier was successful in comparing Bach and others as well as each individual composer. Decision trees to interpret the features used in decision making of the different class arrangements. To determine authorship of BWV 535, they train a quadratic Bayesian classifier to distinguish J.S. Bach, W.F. Bach and J.L. Krebs. They again compare every possible class arrangement as potential composers.

PCA was used to analyze the music of Bach, Handel, Telemann, Mozart and Haydn and compare J.S. Bach, W.F. Bach and J.L. Krebs in an attempt to classify BWV 534.(Backer & Kranenburg, 2005) Although only two PC's accounted for most of the variance, 5 PCs were used to account for more variance. Binary comparisons were used to compare composers. This resulted in a relatively clear separation between Bach and Josquin. For Josquin and his contemporaries, the PC's do not do as well a job of separation. The results of the principal component analysis run on all the composers, were used to train a classifier on all the composers. First a k-nearest neighbor classifier was used. To account for most of the variance, 27 PCs were used.

Next they trained a support vector machine classifier with a radial kernel. Finally they used a decision tree to determine which features were important in discerning the composers.

Speiser and Gupta (Speiser & Gupta, n.d.) scored each feature by the mutual information of each features. They then chose the top 50 features and ran GDA. They then ran PCA to attempt to remove some of the dependencies associated with musical features. They first fit a Naive Bayes for classification, but it had a large training error as the independence assumption does not work well with musical data. Next they used support vector machines with a Gaussian kernel and GCM learning algorithms.

** insert picture **

1.8 Fanny and Felix Mendelssohn

Most musical stylometry analysis focuses on music of the Renaissance and Baroque period, as there are more questions of authorship in that period. As the Romantic period is much more modern in comparison, there are many more surviving records of original manuscripts that include the composer.

Felix Mendelssohn, often considered a prodigy akin to Mozart, was a prolific composer. Before he was fourteen years old, he had already written over 100 compositions.

His lesser known sister Fanny Hensel was also a composer of incredible skill. The two were very close, for many years training and studying together. In their early education living in Berlin, Felix and Fanny received the same musical education, first piano lessons by Madam Bigot, a famous pianist esteemed by Haydn and Beethoven. Beginning in 1818, Carl Friedrich Zelter, a somewhat removed student of Bach and the most influential Berlin musician of the time, began to teach them both composition. In addition to music, the children were tutored by some of the finest scholars in Berlin in subjects such as languages, history, and drawing. Goethe himself claimed that Fanny was “as gifted as Felix”. (Tillard, 1996)

As Fanny grew up, her father started implying that she should focus her energy on the domestic sphere of her life. While the fact that she never became a world famous composer and performer is often attributed to the gender politics of her time, it is also likely due to her high class. (Reich, 1991) Especially considering the anti-semitic feelings of the time, and since the family had recently converted from Judaism to Christianity, the family did not want any other unusual characteristic such as a professional female composer to set them further apart from “polite” society.

Most of Fanny’s available work are *Lieder*, short pieces of voice accompanied by piano. They were accepted at the time as the more feminine, domestic compositions, acceptable for women to compose. Her brother moved on to more elaborate compositions such as operas and orchestral concertos. Her father pressured Fanny to remain composing *Lieder*. (Todd, 2003)

“Music will perhaps become his profession, while for you it can and must only be an ornament, never the root of your being and doing. We may therefore pardon him some ambition and desire to be acknowledged in a pursuit which appears very important to him, . . . while it does you credit

that you have always shown yourself good and sensible in these matters; ... Remain true to these sentiments and to this line of conduct; they are feminine, and only what is truly feminine is an ornament to your sex.”

Throughout their lives, Felix and Fanny maintained contact through letters until Fanny’s death in 1847 and Felix’s death shortly thereafter. These letters contain many instances of Felix asking for advice on his compositions (include quote)

Unlike Felix who conducted and performed piano and organ in some of Berlin’s most esteemed concert halls, most of Fanny’s performances were private, only performed in small circles of her friends and family at intimate parties. Similarly, although she was quite a prolific composer, under recommendation of her father Abraham Mendelssohn, and to a lesser extent Felix, Fanny did not publish her work until later in her life. In 1846 after her father’s death and though her brother’s disapproval, she published her first collection of *Lieder*. Many of Fanny’s unpublished notebooks are in private collections and are inaccessible

However, it is widely speculated (known?) that some of Fanny’s work was published under her brother’s name, Especially three pieces each in his Op 8 and 9 *Lieder*. Famously, when Felix met the Queen of England, she sang Felix’s *Lied* “*Italien*”, and Felix had to admit that in fact, it was his sister that had written it. In a letter to Felix, Fanny admits:

“I have just recently received a letter from Vienna, which contained basically nothing but the question of whether “*On Wings of Song*” was by me, and that I should really send a list of things that are running about in the world disguised, it seems that they aren’t clever enough themselves to separate the wheat from the chaff.” (Mace, 2013)

As she never made such a list, we are left to wonder if there are any other pieces of hers that have been published under her brother’s name and reputation.

This project will use *Lieder* of Fanny and Felix Mendelssohn. Most of the available work by Fanny are *Lieder*, of which Felix also composed a great deal. We will see if there is a determinable difference in style of these siblings who grew up very close and received mostly the same musical education. We will then look at the (disputed?) Op 8 and 9. Additionally, using *Lieder* that have been decidedly written by Felix, we will see if any other of his earlier publications could have potentially been written by Fanny.

Chapter 2

2.1 About the data and conversion process

2.2 Pieces used

The majority of the pieces used in this paper were Lieder of Felix Mendelssohn and Fanny Hensel. While both composers composed different types of music, the majority of available scores from Fanny Hensel were Lieder. Felix Mendelssohn composed many different styles of music, orchestral, piano, etc. Fanny Hensel in contrast has an available existing corpus of mostly Lieder, although she did compose many works for solo piano and orchestra. Of Felix's music, Lieder of Op 8 (12), op 9 (12), op 19(6), op 34 (6), op 47 (6), op 57 (6), op 71 (6), and 6 pieces of Lied without opus numbers, also 56 lied without a collection. A total of 116 pieces.

Of Fanny's music, 23 lieder were used from her Lieder without Name collection, 10 from her *wo kommst du her* collection, and 10 from, a total of 43 pieces.

Data from JS Bach was also used. These data were available in Kern Score format. from the Center for Computer Assisted Research in the Humanities. (CCARH). The pieces used were the Well tempered clavier. These were written as training pieces and each book contains 24 pieces with one in every possible key. These were chosen as the data were more easily accessible (no scanning was required), and they were a similar format as the Mendelssohn songs, written as for solo piano(or harpsichord).

„Dieser Tage habe ich einen Brief aus Wien erhalten, der von
Anfrage ob „auf Flügeln des Gesanges“ von mir wäre, u. ist
verköpft
von W
blos i
m we
es ha
im D

23 Das Fenster
(Carl Klingemann)

63

[August 1826]

Andante

1. Es lausch-te das Laub... so...
2. Es lauscht aus dem Laub... so...

p *cresc.* *p*

dan - kel-grün dort in das Fen-ster hi - nein, die Son - ne da - rin... am -
dan - kel-grün, es strah-len dort Au-gen he - raus, es ran - ken die Re-ben im...

liebs - ten... schien, dort san - gen die Vö - ge - lein. Sie wähl-nen, es wer-de so
stül-len Be - mühen, um - gar - nen das ein - sa-me Haus. Es wähl-nen die Ar-men das

Beispiel EIB 8653

A
E
H
RO
B
nd
nd

Figure 2.1: A picture of a score

2.3 Optical music recognition

The vast majority of classical music is found solely in PDF or physical copies. Sheet music as a form of data requires a lengthy process of conversion before being able to be used in any analysis. Simply scanning the scores into, say, a PDF, gives no musical semantics and can only be viewed on screen or printed on paper. Thus, the two main steps in reading in data from sheet music are: using optical music recognition software to transform physical scores into digital formats, and to read the digital format in to R where subsequent analysis can be done.

The scores used in this paper were obtained from physical copies available in the Reed music library. These scores were then scanned using software designed for optical music recognition (OMR).

Optical music recognition requires learning from graphical and textual information. The main things the software must pick up are the locations of bar lines, notes, rests, slurs, dynamic markings, tempo markings, lyrics etc. Basic optical music recognition has been around since 1966.

Most commonly, the first step in optical music recognition is to remove the staff lines. The staff lines are critical, as they define the basis for the vertical definition distance of pitch, and the horizontal distance definition of rhythm. The staff gives a normalization that is helpful, essentially defining the size of what notes and rhythm look like. Staff removal methods include projections, histograms, run lengths, Candidates assemblage, contour tracking, and Graph path search. (Doermann, Tombre, & others, 2014)

The next step is music symbol extraction and classification. These methods include template matching, where the object in question is compared to existing known musical symbols, simple operators, such as analysis of bounding boxes and projections, joining graphical primitives, such as combining extracted objects such as notes, note heads, and note beams to connect them in a musically correct way to form chords etc. Other methods use statistical models for analyzing musical primitives (the objects its trying to classify) such as Neural Networks, Support Vector Machines, k-Nearest Neighbor, and Hidden Markov Models.

The next step OMR preforms is syntactical analysis and validation. This step essentially uses defined grammars describing the organization of music notation in terms of music symbols. This makes the classification problem simpler, as there are existing rules and relationships between musical symbols.

```
knitr::include_graphics("images/museScore.png")
```

The screenshot shows the MuseScore 2.10.0 interface. The title bar indicates 'MuseScore 2: scV71992'. The main window displays a musical score for '2. Die Nonne' by Mai [18J22 Berlin]. The score is in 3/8 time and features a vocal line with German lyrics. The lyrics are: 'Im stil- len Klo- ster- / O wohl mir, dai! ge- / Sic trat mit za- gem / Sic sank zu sei- nen'. The score is displayed on a staff with a treble clef and a key signature of one sharp (F#). The interface includes a toolbar at the top, a palette on the left, and a status bar at the bottom.

The OMR used in this paper was Photo Score. Photo Score works by scanning the score on a flatbed scanner at a high resolution. It then uses OMR techniques to output a musicXML file that can be read in by most music composing software, such as Sibelius, Finale or Muse Score. MusicXML is used commonly as a format for digital music, as it is conducive to representing sheet music and music notation, and it can be transferable to many different music software. Muse Score was chosen to be the music software for viewing digital scores, as it is a free software that can read MusicXML. After being read into Muse Score, each piece was proof-read and corrected, as there were often errors in the OMR, especially in recognizing triplets. Unfortunately, the scanning process is very lengthy and time consuming, as the scanning often gives a large number of mistakes. The score must be then scanned again. In addition the proof-reading process is lengthy. We must check each note and theme for errors against the original score, and change the afflicted notes using Muse Score. The corrected score must then be output as a musicXML file.

MusicXML on its own is not conducive to converting into a data frame as representing the single half note middle C looks like this:

```

<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<!DOCTYPE score-partwise PUBLIC
    "-//Recordare//DTD MusicXML 0.5 Partwise//EN"
    "http://www.musicxml.org/dtds/partwise.dtd">
<score-partwise>
  <part-list>
    <score-part id="P1">
      <part-name>Music</part-name>
    </score-part>
  </part-list>
  <part id="P1">
    <measure number="1">
      <attributes>
        <divisions>1</divisions>
        <key>
          <fifths>0</fifths>
        </key>
        <time>
          <beats>4</beats>
          <beat-type>4</beat-type>
        </time>
        <clef>
          <sign>G</sign>
          <line>2</line>
        </clef>
      </attributes>
      <note>
        <pitch>
          <step>C</step>
          <octave>4</octave>
        </pitch>
        <duration>4</duration>
        <type>whole</type>
      </note>
    </measure>
  </part>
</score-partwise>

```

We then need to convert into a format more easily readable into R. The Kern Score music format is much more easily readable. The below picture shows how a basic piece of music corresponds to a .krn file.

Start of kern data spine

G clef, **2**nd staff line of staff
key signature (**B-flat**)
key, **D minor**
meter, **2/2**

metric sign (**common time**)

first bar (styled invisible)
half note D4, stem up
half note A4, stem up

second bar
half note F4, stem up
half note D4, stem up


third bar
half note C#4, stem up
quarter note D4, stem up
quarter note E4, stem up


fourth bar
start tie half note F4, **su**
3th F4, **su**, **tie end**, **beam start**
eighth G4, stem up
eighth F4, stem up
8th E4, stem up, **beam end**


fifth bar
quarter D4, stem up


end of kern data spine


****kern**


***clefG2** ↔ 


***k[b-]** ↔ 


***d:** ↔ 


***M2/2** ↔ 


***met(c)** ↔ 


=1- ↔ 


2d/ ↔ 

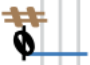
2a/ ↔ 


=2 ↔ 


2f/ ↔ 


2d/ ↔ 


=3 ↔ 


2c#/ ↔ 


4d/ ↔ 


4e/ ↔ 


=4 ↔ 


[2f/ ↔ 


8f]/L ↔ 

8g/ ↔ 

8f/ ↔ 

8e/J ↔ 

=5 ↔ 

4d/ ↔ 

*** -**

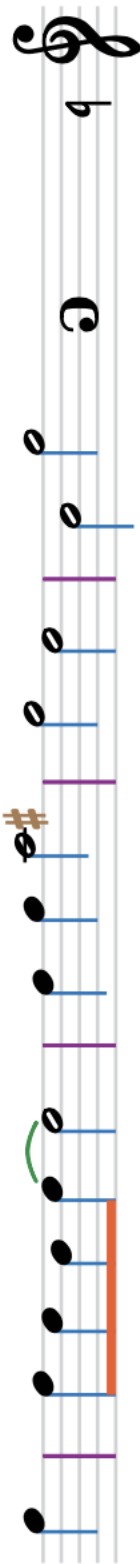


Figure 2.2: Shows how basic kern files correspond to sheet music

Each line of a krn file represents one value of a timebase. Kern files are based on the smallest (shortest) rhythm value of a note found in a piece. For example, if a piece was in 4/4 and there were sixteenth notes present there would be 16*4 rows for each measure. The “attack” of each note is the only note printed, the following time while the note is held is represented with dots in the remaining rows until a new note is sounded for that staff.

We do this by using Humdrum’s function `xml2hum` that converts a musicXML file into a .krn file. Humdrum is a computational music software. It is a command line tool that has many functions for music analysis. The kern file type can be read much more easily into R. Compared to above, the code for a single middle c whole note would be :

```
**kern
*clefG2
*k[c]
*M4/4
=1-
1c/
```

Each staff of a piece of music corresponds to one .krn spline. Each spline is represented in a column. The lowest bass staff is the first column and then progresses up to a soprano line.

The `import_xml_files.sh` file goes through the process of converting scores from musicXML to .krn. The CCARH has a large data base mainly baroque and Renaissance composers already in this format, which is where the Bach data came from. The files that were scanned need to be separated into having a separate file for each staff, which would mean a separate file for each instrument. In addition, since we are focused on musical style, the text of the pieces are removed in this stage. In our case for each piece there are always two or three files for each piece, which are voice, piano right hand, and piano left hand. This is necessary to avoid the bugs in `xml2hum` that have issues when staves don’t necessarily match up as a result of the conversion process, most often by rhythm.

2.4 .krn to R

Once we have .krn files to represent each piece we use regular expressions to extract key information. For scanned music (Felix and Fanny music), there are as many files as there are staves, usually three. MuseR’s `kern2df()` and `piece_df()` functions read in .krn files and output a data frame in R for each piece. First the data in .krn format are read in line by line using R’s `readLines()` function. This takes every line of the .krn file and converts it into a vector. Each entry contains the rythem value and note value. If there are multiple notes played at the same time, they are all in one line. Then each entry is seperated out into the rythem and note value for each note. Each line contains the following columns: measure, rhythm value [4], rhythm name [quarter note], note value [5], note name (octave inclusive)[cc], note name (octave exclusive)[C

sharp]. In addition, for the whole piece the key signature and meter are recorded as columns.

A lot of data included in the .krn files are not necessary. For example, we assume that whether or not a note has a step up or stem down offers no help in classifying composer style, so this information is removed when converting to an R data frame.

Inspired by the .krn file type, each row of the data frame contains one time base value. For a given piece, the time base represents the shortest note duration value. For example, if the shortest note a piece contained was a sixteenth note, the time base would be 16. Each measure then would contain 16 rows. This results in many rows of NA for certain instruments, when a note is still being voiced, but it is not the instance of the note being attacked.

2.5 About the functions - MuseR

Once the scores are converted into R raw data, feature creation begins. These features are mostly features suggested by John Cox.

To the best of my knowledge, there is currently no package of R that has been built to analyze sheet music. There are existing packages (such as `tuneR`) that examine audio formats of music. The intention of this thesis was to create a package, `museR`, that takes sheet music in the proper form (musicXML or .krn) and does all of the analysis using R.

Melodic intervals

Melodic intervals, or the interval between two successive notes, are found using the `mel_ints()` function. This first calculates the top line of any instrument.

c-dur

This outputs the proportion of each melodic interval happening over the whole piece.

Similarly the `consonance()` function outputs the proportion of consonant (perfect, imperfect, dissonant) intervals over the whole piece by calling `mel_ints()`.

Major__minor

For most musical analysis, the key of the piece is important in determining chords, etc. The key is based on the key signature, which is always given in a .krn file. Kern files from CCARH have the key of the piece given, but

Chords

Suspended chords are currently not supported by `MuseR`. Chords that begin, or are “attacked” at the same time count.

First, the key of the piece is found, as different chords depend on the key. (really?)

Next, the times notes are played at the same time are extracted into a list. Then the number of unique notes played at once is found. If there are two notes played at

once, `harm_int()` calculates the harmonic interval. This is done by calculating

$$note_1 - note_2 \mod (12)$$

This gives the number of half steps between each note. That number is then matched with the index of the interval. Work is being done to have this include augmented and diminished interval, but unfortunately that has not been completed at this time.

The possible triad chords are all defined by the intervals between each note. For example, a Major triad is given by the base note, a major third above the base, and a perfect fifth above the base. This corresponds to 4 half steps then 4 half steps. Alternatively a minor triad is given by the base, a minor third above, and a perfect fifth above the base, which is 3 half steps, then 5 half steps.

A similar process is done for seventh chords.

Resolutions

Density

Par thirds, par fourths, par sixths

Prop scale degree

Chapter 3

EDA

Chapter 4

4.1 Feature Selection

Next they computed the entropy of the probability of occurrence of ways of thinking about chords; chords are the same no matter what scale degree they are on, and distinguishing chords differently. Next the entropy was calculated given the probability of each pitch in the score. Entropy was calculated by $-\sum_{i=1}^N p_i \log p_i$ where p_i is the probability of occurrence, and N is the total number. Next they the average number of active voices at one time. This represents the voice density of the piece. Then for every interval, the duration of that interval was divided by the total duration of all intervals. Next the total duration of parallel thirds, fourths, and sixths divided by the total duration of all intervals was measured. Finally a measure of suspensions was found. (Backer)

Chapter 5

5.1 Models

For classification, we are interested in using a set of features to predict the response, or composer. We denote the features X , where X_i is a certain feature. Each song has a composer or response, known or unknown, denoted by Y . The i^{th} song has features x_i and composer Y_i , where x_i is a vector of p features. The composer takes values in a discrete set, which is the possible composers. Thus we can always divide the input space into a collection of regions labeled according to the classification

5.2 Linear Methods for Classification

5.2.1 Linear Regression:

A naive model for classification is linear regression. If our predictor space G has K classes, we code the response K different indicator responses Y_k where $Y_k = 1$ if $G = k$ and 0 otherwise. We can use the resulting k hyperplanes as a decision boundary. We find the coefficients for the line by finding coefficients β to minimize the residual sum of squares.

$$RSS(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)$$

where \mathbf{X} is an $N \times p$ matrix with each row an input vector, and \mathbf{y} is an N -vector of the outputs of the training set. This gives the unique solution:

$$\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

This gives us;

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

If there are K classes, and we have that the fitted linear model for the k^{th} class is $\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$, the decision boundary between class k and class l is the set of points for which $\hat{f}_k(x) = \hat{f}_l(x)$ which is equivalent to the set $\{x : (\hat{\beta}_{k0} - \hat{l}_0) + (\hat{\beta}_k - \hat{\beta}_l)^T x = 0\}$ which is the hyperplane.

Linear Discriminant Analysis

Linear regression on a categorical variable that has multiple variables has issues when there isn't a natural ordering with the categories. For large K and small p , groups can be masked. When there is a binary response, we can calculate $P(Y|X)$, but linear regression can give predictions that arent valid probabilities, namely negative probabilities or probabilities greater than 1.

Knowing the class posteriors $P(Y = k|X)$ gives us an optimal classification. If we assume $f_k(x)$ is the class-conditional density of X in class $G = k$ and that π_k is the prior probability of class k with $\sum_{k=1}^K \pi_k = 1$. We can then model $P(Y = k|X)$ by modeling the distribution of the features X separately in each response class, and then use Bayes' theorem to calculate $P(Y = k|X)$ which gives us the following:

$$P(Y = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l}$$

We thus must have a model to find $f_k(x)$. Linear and quadratic discriminant analysis assume a multivariate Gaussian density, given by:

$$f_k(x) = \frac{1}{(2\pi)^{p/2}|\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)}$$

Linear discriminant analysis (LDA) assumes that the covariance matrix is equal for every k : $\Sigma_k = \Sigma \forall k$. Quadratic discriminant analysis does not have this assumption. In addition we assume $\hat{\pi}_k = N_k/N$ where N_k is the number of class - k observations, $\hat{\mu}_k = \sum_{g_i=k} x_i/N_k$, and $\hat{\Sigma} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T / (N - K)$

For LDA we can look at the log ratio comparing two classes k and l and can show:

$$\log \frac{P(Y = k|x = x)}{P(Y = l|X = x)} = \log \frac{f_k(x)}{f_l(x)} + \log \frac{\pi_k}{\pi_l} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) + x^T \Sigma^{-1}(\mu_k - \mu_l)$$

This is a linear equation, so the classes will be separated by hyperplanes. From the above, we can find that the predicted class for any x is :

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

These functions are known as *linear discriminant functions*. We predict the class by finding the maximum value of the discriminant functions of all k .

For QDA we get the following discriminant functions:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k) + \log \pi_k$$

Linear discriminant analysis is helpful when the classes are well separated, when n is small and the distribution of the predictors X is approximately normal in each of the classes, and when there are more than two response classes.

Logistic Regression

Logistic regression differs from linear discriminant analysis by directly modeling $P(Y = k|X)$ by using the logistic function.

5.3 K- nearest - neighbor

Another method for classification is k-nearest neighbor methods. It uses observations in the training set closest to x to form \hat{Y} , the outputs. We often use Euclidean distance as a metric for closeness, although other methods exist. It is defined as

$$\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

where $N_k(x)$ is the k closest points, or neighbors, x_i in the training set. This is equivalent to taking the average of k observations with x_i closest to x . This gives a predicted class of taking the mode of the k nearest neighbors.

5.3.1 Logistic Regression

For logistic regression, we are interested in measuring the probability a certain occurrence is in a group. In our case, this corresponds to if a certain piece is composed by a certain composer. Linear regression is not a good application in this case, as it gives probabilities that are not in 0 to 1. Logistic regression uses

Chapter 6

6.1 Model Fit

Chapter 7

7.1 Discussion

Conclusion

If we don't want Conclusion to have a chapter number next to it, we can add the `{-}` attribute.

More info

And here's some other random info: the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.

Appendix A

The First Appendix

This first appendix includes all of the R chunks of code that were hidden throughout the document (using the `include = FALSE` chunk tag) to help with readability and/or setup.

In the main Rmd file

```
# This chunk ensures that the thesisdown package is  
# installed and loaded. This thesisdown package includes  
# the template files for the thesis.  
if(!require(devtools))  
  install.packages("devtools", repos = "http://cran.rstudio.com")  
if(!require(thesisdown))  
  devtools::install_github("ismayc/thesisdown")  
library(thesisdown)
```

In Chapter ??:

References

- Adair, D. (1944). The authorship of the disputed federalist papers. *The William and Mary Quarterly: A Magazine of Early American History*, 98–122.
- Backer, E., & Kranenburg, P. van. (2005). On musical stylometry—a pattern recognition approach. *Pattern Recognition Letters*, 26(3), 299–309.
- Brinkman, A., Shanahan, D., & Sapp, C. (n.d.). Musical stylometry, machine learning, and attribution studies: A semi-supervised approach to the works of josquin.
- Doermann, D., Tombre, K., & others. (2014). *Handbook of document image processing and recognition*. Springer.
- Ford, P. L., & Bourne, E. G. (1897). The authorship of the federalist. *The American Historical Review*, 2(4), 675–687. Retrieved from <http://www.jstor.org/stable/1833983>
- Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3(Mar), 1157–1182.
- Mace, A. R. (2013). *Fanny hensel, felix mendelssohn bartholdy, and the formation of the “mendelssohnian” style* (PhD thesis).
- Mearns, L., Tidhar, D., & Dixon, S. (2010). Characterisation of composer style using high-level musical features. In *Proceedings of 3rd international workshop on machine learning and music* (pp. 37–40). ACM.
- Mosteller, F., & Wallace, D. (1964). *Inference and disputed authorship: The federalist*. Addison-Wesley.
- Reich, N. B. (1991). The power of class: Fanny hensel. *Mendelssohn and His World*, 86–99.
- Speiser, J., & Gupta, V. (n.d.). Composer style attribution. *Project Report for CS*, 229.
- Tillard, F. (1996). *Fanny mendelssohn*. Hal Leonard Corporation.
- Todd, R. L. (2003). *Mendelssohn: A life in music*. Oxford University Press.