
Introduction to Musical Corpus Studies

Release 0.0.1

Fabian C. Moss

Oct 22, 2020

CONTENT

1 Organization	3
1.1 Overview	4
1.2 Credits	5
1.3 Learning objectives	5
1.4 Deliverables	5
2 Introduction	7
2.1 About this course	7
2.1.1 About me	7
2.1.2 Focus of this course	8
2.2 What are Musical Corpus Studies?	8
2.3 Epistemological goals	8
2.4 Issues	8
2.5 MCS and traditional musicology	8
2.6 Basic representations	8
3 Folk Songs and the Melodic Arc	9
3.1 Data	9
3.2 Notes, Pitch Classes	9
3.3 Intervals	9
4 Solos in the Weimar Jazz Database	11
5 Renaissance Chansons and Masses	13
6 Harmony in Beethoven's String Quartets	15
6.1 Access the data	15
6.2 Harmonic Annotations	15
6.3 Chord Transitions	15
7 Harmonic Clusters in Pop Charts	17
7.1 The dataset	18
7.2 Research question	18
7.3 Operationalization	18
7.4 Clustering	18
7.4.1 General idea	18
7.4.2 Algorithm	18
8 Data-Driven Music History	19
8.1 Research Questions	19
8.2 A bit of theory	19

8.3	Data	20
8.3.1	A (kind of) large corpus: TP3C	20
8.4	Recovering the line of fifths from data	23
8.5	Historical development of tonality	26
8.6	If there is time: some more advanced stuff	28
8.7	Summary	31
8.8	Conclusion	32
8.9	The end	32
9	Harmony and Form in Brazilian Choro	33
9.1	Data	33
10	A Data-Driven History of Tonality	35
10.1	Musical Pieces as Tonal Pitch-Class Distributions	35
10.1.1	Distributions of TPCs	37
10.1.2	Principle Component Analysis	38
10.2	Historical Development	38
10.2.1	LOWESS	39
11	Malian Percussion Music	41
12	Electronic Music 1950–1990	43
13	Conclusion	45
14	Python Basics	47
14.1	Types	47
14.2	Lists	47
14.3	Reading and saving files	47
15	Bibliography	49
	Bibliography	51

Warning: This material is still (heavily) under construction and might change throughout the course!

You can help improving the course and [let me know](#) about any errors and inconsistencies that you find or suggest other ways of improving the course.

Welcome!

These pages present the content of the course “Introduction to Musical Corpus Studies” at the [Institute of Musicology](#), given at [University of Cologne](#) in Fall 2020.

In the last two decades *Musical Corpus Studies* evolved from a niche discipline into a veritable research area. The growing availability of digital and digitized musical data as well as the application and development of modern methodologies from computer science, machine learning, and data science cast new light on old musicological questions and generate entirely novel approaches to empirical music research.

Moreover, the general methodological and epistemological approach of Musical Corpus Studies allows to transcend traditional intra-musicological boundaries between its sub-disciplines (historical/systematic/ethnological/...) without sacrificing the respective specific viewpoints and perspectives.

This course offers a fundamental and practical introduction into these topics. It demonstrates, explores, and critically reflects central thematic areas and methods by means of a number of case studies. In the engagement with these topics the course also introduces elementary methods from natural language and music processing, as well as statistics, data analysis and visualization.

The course is aimed at students at the undergraduate level who have little or no empirical background and are curious about quantitative approaches to musicology.

CHAPTER
ONE

ORGANIZATION

1.1 Overview

No.	Date	Time	Room	Topics
1	Fr., 13.11.2020	16:00- 17:20 Uhr	Neuer Seminar- raum 1.315	Introduction / Background
2		17:40- 19:00 Uhr		Folk Songs, Melodies, Pitches and Intervals frequencies, mean, variance
3	Sa., 14.11.2020	09:00- 10:20 Uhr	Neuer Seminar- raum 1.315	Jazz Solos, Melodies, Regular Expressions
4		10:40- 12:00 Uhr		Beethoven's string quartets, harmony, <i>n</i> -grams, Markov models
		12:00- 13:00 Uhr		Lunch Break
5		13:00- 14:20 Uhr		Pop Charts Billboard 100, harmony, Clustering, <i>k</i> -means, [Hidden Markov Models]
6		14:40- 16:00 Uhr		Free group work
7	Fr., 11.12.2020	10:00- 11:20 Uhr	Alter Seminar- raum 1.408	Cadences in Renaissance Polyphony with guest researcher Richard Freedman
8		11:40- 13:00 Uhr		Brazilian Choro, harmony, form, context-Free Grammars
9	Sa., 12.12.2020	09:00- 10:20 Uhr	Neuer Seminar- raum 1.315	Malian Percussion Music, rhythm, meter
10		10:40- 12:00 Uhr		Electronic Music 1950-1990
		12:00- 13:00 Uhr		Lunch Break
11		13:00- 14:20 Uhr		Free group work
4		Uhr		Chapter 1. Organization
12		14:40- 16:00 Uhr		Recapitulation and conclusion

1.2 Credits

Note: Ich gehe in der Seminarplanung von 12 Semesterwochen à 2 SWS aus, für das gesamte Blockseminar also 24 SWS. Das Seminar wird mit 3 CP bewertet, was 90 Stunden aktiver Arbeit entspricht. Davon entfallen 24 SWS an die Präsenzzeit im Seminar plus 48 SWS an Vor- und Nachbereitung der Seminarsitzungen. Die verbleibenden 18 SWS sind für die Lektüre der Fachliteratur vorgesehen.

Course work consists of three parts: preparing the relevant literature (reading), completing the relevant exercises (group work), and critically engaging with the course materials in the form of a report written together with your group.

1.3 Learning objectives

1. content of the course units
2. work load management
3. organization
4. reading of scientific literature
5. writing academic reviews

1.4 Deliverables

Reading

For each session, the relevant literature is cited in the text. Careful preparation is required in order to be able to follow the content of the course. Because the course will mainly talk about methods and general points of musical corpus research, the content (and musical topic) will mainly be introduced by the literature.

I am aware that the reading workload is relatively high since the course will be taught as a block seminar and doesn't spread out over the entire semester. I hope that the fact that the course is finished before the end of the year compensates for this.

Group work

At the beginning of the course, you will be randomly assigned to a group. Together with your group, you will work on a number of exercises during the course, e.g. in Zoom breakout rooms.

Review

After the course has ended, your group will be randomly assigned a course topic. It is your task to write a review/report on this topic. What did you learn? Which concepts are not clear? Which methods did you (not) understand? What is missing? How can the textual descriptions be improved? Who in your group did what? Write about the organization of your group, challenges and benefits.

Important: Submit your report by **31 January 2021**.

CHAPTER
TWO

INTRODUCTION



2.1 About this course

2.1.1 About me

- Music and Mathematics education (Uni & HfMT Köln)
- MA Musicology (HfMT Köln)
- PhD Digital Humanities (EPFL)

2.1.2 Focus of this course

Programming introductions often boring. A lot of time lost in introducing basic concepts and techniques (important!) but quite remote from actual (!) applications. Examples are usually “toy examples” that work well, but the transition to real-world applications is difficult. Of course, the example studies discussed in this course work well, too. However, they are without exception taken from peer-reviewed, published, open access articles. They thus reflect actual, recent research questions that reflect current research.

This course takes thus the opposite approach to “toy examples”. We will not introduce many specific programming concepts. The course rather showcases what is possible with musical corpus studies. If this sparks your interest, it will be much easier to pick up the basics for yourself, knowing what they are *for* and being motivated intrinsically. If you are not particularly interested in doing this kind of work yourself, you will still see a broad range of applications that are much more useful to you than learning (or not learning) Python basics.

2.2 What are Musical Corpus Studies?

tbc... (text from diss?)

2.3 Epistemological goals

tbc...

2.4 Issues

tbc [[Coo06](#)][[Hon06](#)][[Hur13](#)][[Mar16](#)][[NR16](#)][[Pug15](#)][[Sch16](#)][[TV13](#)]

2.5 MCS and traditional musicology

tbc

2.6 Basic representations

- tones, notes
- (tonal/neutral) pitch classes
- meter (hierarchy)

FOLK SONGS AND THE MELODIC ARC

Tones are among the basic elements of music. Most musical styles combine tones in different ways to create songs, chants, instrumental pieces, or other elaborate compositions. In this chapter, we will analyze some basic aspects of songs by studying distributions of tones and intervals.

Huron... / MusThe Tutorial

3.1 Data

<http://kern.humdrum.org/data?f=zip&l=essen> or <http://kern.humdrum.org/help/data/>

Open and read *README.txt*

Essen Folksong Collection

3.2 Notes, Pitch Classes

https://github.com/DCMLab/DigitalMusicologyExercises/tree/master/tone_profiles

means, variance

also multidimensional (for later)

3.3 Intervals

https://github.com/DCMLab/DigitalMusicologyExercises/tree/master/interval_bigrams

maybe extend with Hansen and Pearce (2014) (but data not available?)

Note: In this chapter we covered the following musical terms:

- a
 - b
 - c
-

SOLOS IN THE WEIMAR JAZZ DATABASE



Fig. 4.1: Photo by Janine Robinson on Unsplash

The first project we will have a look at is the [Jazzomat](#) project. Transcriptions of Jazz solos [PFAbesser+17]. The *Weimar Jazz Database* (WJD) consists of 456 transcriptions of Jazz solos from diverse substyles. As all the corpora that we deal with here, it is freely available on the internet.¹

The WJD contains a number of tables:

¹ <https://jazzomat.hfm-weimar.de/dbformat/dboverview.html>

Table 4.1: Tables in the *Weimar Jazz Database*

Table name	Description
beats	Table for beat annotation of WJD melodies, referenced by melody (melid)
composition_info	Infos regarding the underlying composition of a WJD solo, referenced by melody (melid)
db_info	Information regarding the distributed database file like version information, license, etc
esac_info	EsAC infos for EsAC melodies, referenced by melody (melid)
melody	Main table for all melody events
melody_type	Indicated type of melody: WJD solos or EsAC (Folk songs using Essen Associative Code), referenced by melody (melid)
popsong_info	Pop song infos, referenced by melody (melid)
record_info	Infos regarding the specific audio recording of a WJD solo was taken from, referenced by melody (melid)
sections	All sections (phrase, chorus, form, chords, etc.), referenced by melody (melid)
solo_info	Solo infos for WJD solos, referenced by melody (melid)
track_info	Information specific to a track on a record (or CD)
transcription_info	Transcription infos for WJD solos, referenced by melody (melid)

Here, we focus on the main table `melody`. First, we download the entire database from <https://jazzomat.hfm-weimar.de/download/download.html> (under “Weimar Jazz Database”) and save it as the file `wjazz.db`.

```
import sqlite3 # for working with databases
import pandas as pd # for working with tabular data

# create connection to database
conn = sqlite3.connect("wjazzd.db")

# read all entries of the 'melody' table into a pandas DataFrame
df = pd.read_sql("SELECT * FROM melody", con=conn)

df.head()

>>> df.head()
output
```

The part of the code `SELECT * FROM melody` reads “Select all entries from the table ‘melody’”.

CHAPTER
FIVE

RENAISSANCE CHANSONS AND MASSES

[Fre14][FVC17]

HARMONY IN BEETHOVEN'S STRING QUARTETS

6.1 Access the data

The data lies on the GitHub repository [DCMLab/ABC](#). Either download the `.tsv` file directly and open it in pandas or load it from the URL as follows:

```
import pandas as pd
df = pd.read_csv("https://github.com/DCMLab/ABC/corpus.tsv", sep="\t")
```

The corpus is now stored in the variable `df`.

6.2 Harmonic Annotations

- regular expressions

[NHMR18][MNHM19]

6.3 Chord Transitions

- n-grams

CHAPTER
SEVEN

HARMONIC CLUSTERS IN POP CHARTS



Clustering analysis in [SVJ+20].

```
def zipf_mandelbrot(x, a, b, c):  
    """Zipf-Mandelbrot function of `x` given parameters `a`, `b`, `c`. """  
    z = a / ( (b + x)**c )  
    return z
```

7.1 The dataset

djfjfjf

7.2 Research question

Hypothesis: Corpus is not stylistically uniform but consists of several sub-styles, each with respective harmonic sub-grammars.

Discussion: what is harmonic grammar/syntax (Riemann/GTTM/Rohrmeier)

7.3 Operationalization

Harmonic grammar → bigrams (we know that already)

7.4 Clustering

7.4.1 General idea

Variance, minimizing *within-cluster sum of squares*, ...

7.4.2 Algorithm

For data $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$:

1. Choose $k \leq n$.
2. Randomly assign means μ_1, \dots, μ_k to points in \mathbf{X} .
3. Assign each point in \mathbf{X} to cluster C_i by determining closest mean μ_i .
4. Given a cluster C_i , calculate its *centroid* m_i :

$$m_i = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$$

5. Repeat steps 3 and 4 until the means do not change anymore.

Note that, while this is a relatively simple clustering algorithm, it is not guaranteed that it finds the global optimal solution. This means that different *initializations* (step 2 of the algorithm) might lead to different clustering results. More sophisticated clustering methods are, e.g. kNN...

DATA-DRIVEN MUSIC HISTORY

Traditionally, musicology has been divided into historical and systematic research agendas, encompassing qualitative-hermeneutic and quantitative-empirical methodologies, respectively. Innovations in the emerging and rapidly growing field of musical corpus studies question this fundamental divide and address, for instance, inherently historical questions with quantitative methods, fueled by the creation of ever larger and more appropriate datasets.

This chapter first introduces some methodological and epistemological issues regarding empirical approaches to music history. It then presents a hands-on exercise on a case study. Finally, it invites critical discussion about the implications and relevance of the results for other subfields such as music psychology. In doing so, the workshop simulates (nearly) the entire life cycle of a research project, from an initial idea via selecting appropriate operationalisations and measures up to choosing suitable visualisations to communicate the results, e.g. in a research article or a blog post. At each point, participants will be invited to critically reflect the decisions taken. Along the way, more general methods for data analysis (e.g. data transformation, clustering, dimensionality reduction, and plotting) will be introduced. This is expected to benefit participants in a vast number of future projects.

```
import pandas as pd # for working with tabular data
pd.set_option('display.max_columns', 500)
import matplotlib.pyplot as plt # for plotting
plt.style.use("fivethirtyeight") # select specific plotting style
import seaborn as sns; sns.set_context("talk")
import numpy as np
```

8.1 Research Questions

- General: How can we study historical changes quantitatively?
- Specific: What can we say about the history of tonality based on a dataset of musical pieces?

8.2 A bit of theory

```
note_names = list("FCGDAEB") # diatonic note names in fifths ordering
note_names
```

```
['F', 'C', 'G', 'D', 'A', 'E', 'B']
```

```
accidentals = ["bb", "b", "", "#", "##"] # up to two accidentals is sufficient here
accidentals
```

```
['bb', 'b', '', '#', '##']
```

```
lof = [ n + a for a in accidentals for n in note_names ] # lof = "Line of Fifths"
print(lof)
```

```
['Fbb', 'Cbb', 'Gbb', 'Dbb', 'Abb', 'Ebb', 'Bbb', 'Fb', 'Cb', 'Gb', 'Db', 'Ab', 'Eb',
← 'Bb', 'F', 'C', 'D', 'A', 'E', 'B', 'F#', 'C#', 'G#', 'D#', 'A#', 'E#', 'B#',
← 'F##', 'C##', 'G##', 'D##', 'A##', 'E##', 'B##']
```

```
len(lof) # how long is this line-of-fifths segment?
```

```
35
```

We call the elements on the line of fifths **tonal pitch-classes**

8.3 Data

8.3.1 A (kind of) large corpus: TP3C

Here, we use a dataset that was specifically compiled for this kind of analysis, the Tonal pitch-class counts corpus (TP3C) [MNR20].

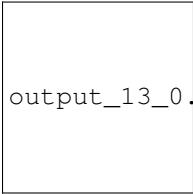
- 2,012 pieces
- 75 composers
- approx. spans 600 years of music history
- does not contain complete pieces but only counts of tonal pitch-classes

```
import pandas as pd # to work with tabular data

url = "https://raw.githubusercontent.com/DCMLab/TP3C/master/tp3c.tsv"
data = pd.read_table(url)

data.sample(10)
```

```
data.display_year.plot(kind="hist", bins=50, figsize=(15,6)); # historical overview
```



output_13_0.png

- it can be seen that there are large gaps and that some historical periods are underrepresented
- however, it is not so obvious how to fix that
- do we want a uniform distribution over time?
- do we want a “historically accurate” distribution?
- do we want to remove geographical/gender/class/instrument/etc. biases?

- on one hand, balanced datasets are likely not to reflect historical realities
- on the other hand, such datasets rather represent the “canon”, that is a contemporary selection of “valuable” compositions that may differ greatly from what was considered relevant at the time

-> There is no unique objective answer to these questions. It is important to be aware of these limitations and take them into account when interpreting the results

For this workshop we ignore all the metadata about the pieces (titles, composer names etc.) but only focus on their tonal material. Therefore, we don't need all the columns of the table.

```
tpc_counts = data.loc[:, :lof] # select all rows ":" and the lof columns
tpc_counts.sample(20)
```

```
piece = tpc_counts.iloc[10]

fig, axes = plt.subplots(2, 1, figsize=(20,10))

axes[0].bar(piece.sort_values(ascending=False).index, piece.sort_
    ↴values(ascending=False))
axes[0].set_title("'without theory'")

axes[1].bar(piece.index, piece)
axes[1].set_title("'with theory'")

plt.savefig("img/random_piece.png")
plt.show()
```

output_17_0.png

Let us have an overview of the note counts in these pieces!

If we would just look at the raw counts of the tonal pitch-classe, we could not learn much from it. Using a theoretical model (the line of fifths) shows that the notes in pieces are usually come from few adjacent keys (you don't say!).

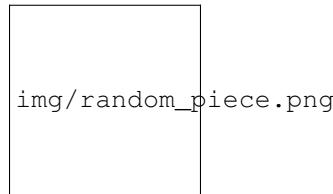


Fig. 8.1: Random piece

We probably have very long pieces (sonatas) and very short pieces (songs) in the dataset. Since we don't want length (or the absolute number of notes in a piece) to have an effect, we rather consider tonal pitch-class distributions instead counts, by normalizing all pieces to sum to one.

```
tpc_dists = tpc_counts.div(tpc_counts.sum(axis=1), axis=0)
tpc_dists.sample(20)
```

For further numerical analysis, we extract the data from this table and assign it to a variable X.

```
# extract values of table to matrix
X = tpc_dists.values

X.shape # shows (#rows, #columns) of X
```

```
(2012, 35)
```

Now, X is a 2012×35 matrix where the rows represent the pieces and the columns (also called “features” or “dimensions”) represent the relative frequency of tonal pitch-classes.

Thinking in 35 dimensions is quite difficult for most people. Without trying to imagine what this would look like, what can we already say about this data?

Since each piece is a point in this 35-D space and pieces are represented as vectors, pieces that have similar tonal pitch-class distributions must be close in this space (whatever this looks like).

What groups of pieces that cluster together? Maybe pieces of the same composer are similar to each other? Maybe pieces from a similar time? Maybe pieces for the same instruments?

If we find clusters, these would still be in 35-D and thus difficult to interpret. Luckily, there are a range of so-called *dimensionality reduction* methods that transform the data into lower-dimensional spaces so that we actually can look at them.

A very common dimensionality reduction method is **Principal Components Analysis (PCA)**.

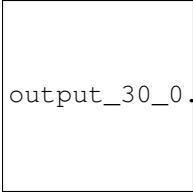
The basic idea of PCA is:

- find dimensions in the data that maximize the variance in this direction
- these dimensions have to be orthogonal to each other (mutually independent)
- these dimensions are called the *principal components*
- each principal component is associated with how much of the data variance it explains

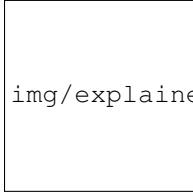
```
import numpy as np # for numerical computations
import sklearn
from sklearn.decomposition import PCA # for dimensionality reduction

pca = sklearn.decomposition.PCA(n_components=35) # initialize PCA with 35 dimensions
pca.fit(X) # apply it to the data
variance = pca.explained_variance_ratio_ # assign explained variance to variable
```

```
fig, ax = plt.subplots(figsize=(14,5))
x = np.arange(35)
ax.plot(x, variance, label="relative", marker="o")
ax.plot(x, variance.cumsum(), label="cumulative", marker="o")
ax.set_xlim(-0.5, 35)
ax.set_ylim(-0.1, 1.1)
ax.set_xlabel("Principal Components")
ax.set_ylabel("Explained variance")
plt.xticks(np.arange(len(variance)), np.arange(len(variance)) + 1) # because Python starts counting at 0
plt.legend(loc="center right")
plt.tight_layout()
plt.savefig("img/explained_variance.png")
plt.show()
```



output_30_0.png



img/explained_variance.png

Fig. 8.2: explained variance

```
variance[:5]
```

```
array([0.41144591, 0.23410347, 0.09063507, 0.07574242, 0.04436989])
```

The first principal component explains 41.1% of the variance of the data, the second explains 23.4% and the third 9%. Together, this amounts to 73.6%.

Almost three quarters of the variance in the dataset is retained by reducing the dimensionality from 35 to 3 dimensions (8.6%)! If we reduce the data to two dimensions, we still can explain $\approx 65\%$ of the variance.

This is great because it means that we can look at the data in 2 or 3 dimensions without loosing too much information.

8.4 Recovering the line of fifths from data

```
pca3d = PCA(n_components=3)
pca3d.fit(X)
```

```
X_ = pca3d.transform(X)
X_.shape
```

```
(2012, 3)
```

```
from mpl_toolkits.mplot3d import Axes3D

fig = plt.figure(figsize=(6, 6))

ax = fig.add_subplot(111, projection='3d')
ax.scatter(X_[:, 0], X_[:, 1], X_[:, 2], s=50, alpha=.25) # c=cs,
ax.set_xlabel("PC 1", labelpad=30)
ax.set_ylabel("PC 2", labelpad=30)
ax.set_zlabel("PC 3", labelpad=30)

plt.tight_layout()
plt.savefig("img/3d_scatter.png")
plt.show()
```

output_38_0.png

img/3d_scatter.png

Fig. 8.3: 3D Scatterplot

Each piece in this plot is represented by a point in 3-D space. But remember that this location represents ~75% of the information contained in the full tonal pitch-class distribution. In 35-D space each dimension corresponded to the relative frequency of a tonal pitch-class in a piece.

- What do these three dimensions signify?
- How can we interpret them?

Fortunately, we can inspect them individually and try to interpret what we see.

```
from itertools import combinations

fig, axes = plt.subplots(1, 3, sharey=True, figsize=(24, 8))

for k, (i, j) in enumerate(combinations(range(3), 2)):

    axes[k].scatter(X_[:, i], X_[:, j], s=50, alpha=.25, edgecolor=None)
    axes[k].set_xlabel(f"PC {i+1}")
    axes[k].set_ylabel(f"PC {j+1}")
    axes[k].set_aspect("equal")

plt.tight_layout()
plt.savefig("img/3d_dimension_pairs.png")
plt.show()
```

output_41_0.png

img/3d_dimension_pairs.png

Fig. 8.4: Principal Components

Clearly, looking at two principal components at a time shows that there is some latent structure in the data. How can we understand it better?

One way to see whether the pieces are clustered together systematically be coloring them according to some criterion.

As always, many different options are available. For the present purpose we will use the most simple summary of the piece: its most frequent note (which is the *mode* of its pitch-class distribution in statistical terms) and call this note its **tonal center**.

This will also allow to map the tonal pitch-classes on the line of fifths to colors.

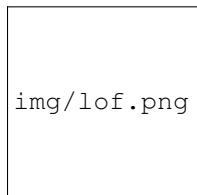


Fig. 8.5: Line of fifths coloring

```
tpc_dists["tonal_center"] = tpc_dists.apply(lambda piece: np.argmax(piece[lof].values) - 15, axis=1)
tpc_dists.sample(10)
```

```
from matplotlib import cm
from matplotlib.colors import Normalize

#normalize item number values to colormap
norm = Normalize(vmin=-15, vmax=20)

# cs = [ cm.seismic(norm(c)) for c in data["tonal_center"]]
cs = [ cm.seismic(norm(c)) for c in tpc_dists["tonal_center"]]
```

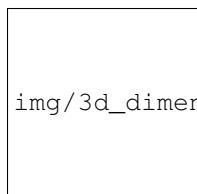


Fig. 8.6: Dimension pairs

```
from itertools import combinations

fig, axes = plt.subplots(1,3, sharey=True, figsize=(24,8))

for k, (i, j) in enumerate(combinations(range(3), 2)):

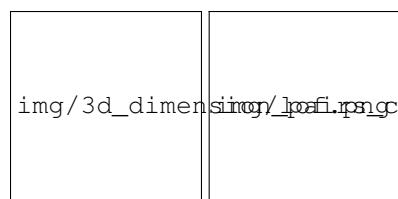
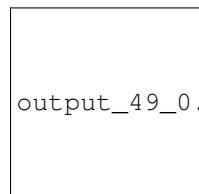
    axes[k].scatter(X_[:,i], X_[:,j], s=50, c=[ np.abs(c) for c in cs], edgecolor=None)
    axes[k].set_xlabel(f"PC {i}")
    axes[k].set_ylabel(f"PC {j}")
    axes[k].set_aspect("equal")

plt.tight_layout()
```

(continues on next page)

(continued from previous page)

```
plt.savefig("img/3d_dimension_pairs_colored.png")
plt.show()
```



8.5 Historical development of tonality

The line of fifths is an important underlying structure for pitch-class distributions in tonal compositions

But we have treated all pieces in our dataset as synchronic and have not yet taken their historical location into account.

Remember the tonal pitch-class distribution of an example piece above?

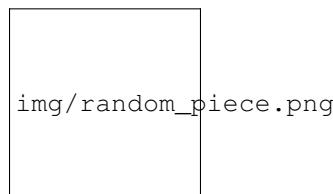


Fig. 8.7: Random piece

Let's assume the pitch-class content of a piece spreads on the line of fifths from F to A \sharp .

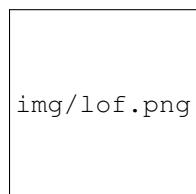


Fig. 8.8: Line of fifths

This means, its range on the line of fifths is $10 - (-1) = 11$. The piece covers eleven consecutive fifths on the lof.

We can generalize this calculation and write a function that calculates the range for each piece in the dataset.

```
def lof_range(piece):
    l = [i for i, v in enumerate(piece) if v!=0]
    return max(l) - min(l)
```

```
data["lof_range"] = data.loc[:, "lof"].apply(lof_range, axis=1) # create a new column
data.sample(20)
```

This allows us now to take the `display_year` (composition or publication) and `lof_range` (range on the line of fifths) features to observe historical changes.

```
fig, ax = plt.subplots(figsize=(18,9))
ax.scatter(data["display_year"].values, data["lof_range"].values, alpha=.5, s=50)
ax.set_xlim(0,35)
ax.set_xlabel("year")
ax.set_ylabel("line-of-fifths range")
plt.savefig("img/hist_scatter.png");
```

output_59_0.png

img/hist_scatter.png

Fig. 8.9: Historical scatterplot

We could try to fit a line to this data to see whether there is a trend (kinda obvious here).

```
g = sns.lmplot(
    data=data,
    x="display_year",
    y="lof_range",
    line_kws={"color":"k"},
    scatter_kws={"alpha":.5},
    #     lowess=True,
    height=8,
    aspect=2
)
g.savefig("img/hist_scatter_line.png");
```

output_62_0.png

But actually, this is not the best idea. Why should any historical process be linear? More complex models might make more sense.

A more versatile technique is *Locally Weighted Scatterplot Smoothing* (LOWESS) that locally fits a polynomial. Using this method, we see that a non-linear process is displayed.

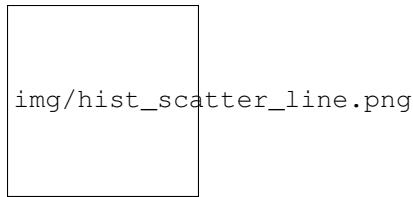


Fig. 8.10: Line Scatter

```
from statsmodels.nonparametric.smoothers_lowess import lowess

x = data.display_year
y = data.lof_range
l = lowess(y,x)

fig, ax = plt.subplots(figsize=(15,10))

ax.scatter(x,y, s=50)
ax.plot(l[:,0], l[:,1], c="k")
ax.set_ylabel("line-of-fifths range")
plt.savefig("img/hist_scatter_lowess.png")
plt.show()
```

output_66_0.png

A small black rectangular placeholder box with a thin white border, indicating the location where the file 'img/hist_scatter_lowess.png' would normally be displayed.

Fig. 8.11: Scatter Lowess

8.6 If there is time: some more advanced stuff

```
B = 200
delta = 1/10

fig, ax = plt.subplots(figsize=(16, 9))

x = data.display_year
y = data.lof_range
l = lowess(y,x, frac=delta)

ax.scatter(x,y, s=50, alpha=.25)
```

(continues on next page)

(continued from previous page)

```

for _ in range(B):
    resampled = data.sample(data.shape[0], replace=True)

    xx = resampled.display_year
    yy = resampled.lof_range
    ll = lowess(yy,xx, frac=delta)

    ax.plot(ll[:,0], ll[:,1], c="k", alpha=.05)

ax.plot(l[:,0], l[:,1], c="yellow")

## REGIONS
from matplotlib.patches import Rectangle

text_kws = {
    "rotation" : 90,
    "fontsize" : 16,
    "bbox" : dict(
        facecolor="white",
        boxstyle="round"
    ),
    "horizontalalignment" : "center",
    "verticalalignment" : "center"
}

rect_props = {
    "width" : 40,
    "zorder" : -1,
    "alpha" : 1.
}

stylecolors = plt.rcParams["axes.prop_cycle"].by_key()["color"]

ax.text(1980, 3, "diatonic", **text_kws)
ax.axhline(6.5, c="gray", linestyle="--", lw=2) # dia / chrom.
ax.add_patch(Rectangle((1960,0), height=6.5, facecolor=stylecolors[0], **rect_props))

ax.text(1980, 9.5, "chromatic", **text_kws)
ax.axhline(12.5, c="gray", linestyle="--", lw=2) # chr. / enh.
ax.add_patch(Rectangle((1960,6.5), height=6, facecolor=stylecolors[1], **rect_props))

ax.text(1980, 23.5, "enharmonic", **text_kws)
ax.add_patch(Rectangle((1960,12.5), height=28, facecolor=stylecolors[2], **rect_
    props))

ax.set_xlim(1300,2000)
ax.set_ylim(0,35)

ax.set_ylabel("line-of-fifths range")
plt.savefig("img/final.png", dpi=300)
plt.show()

```

output_69_0.png

Using bootstrap sampling we achieve an estimation of the local variance of the data and thus of the diversity in the note usage of the musical pieces.

img/final.png

Fig. 8.12: Final Result

We also can distinguish three regions in terms of line-of-fifth range: diatonic, chromatic, and enharmonic.

Grouping the data together in these three regions, we see a clear change from diatonic and chromatic to chromatic and enharmonic pieces over the course of history.

```

epochs = {
    "Renaissance" : [1300, 1549],
    "Baroque" : [1550, 1649],
    "Classical" : [1650, 1749],
    "Early\nRomantic" : [1750, 1819],
    "Late Romantic/\nModern" : [1820, 2000]
}

strata = [
    "diatonic",
    "chromatic",
    "enharmonic"
]

widths = data[["display_year", "lof_range"]].sort_values(by="display_year").reset_index(drop=True)

df = pd.concat(
    [
        widths[
            (widths.display_year >= epochs[e][0]) & (widths.display_year <= epochs[e][1])
        ][["lof_range"]].value_counts(normalize=True).sort_index().groupby(
            lambda x: strata[0] if x <= 6 else strata[1] if x <= 12 else strata[2]
        ).sum() for e in epochs
    ], axis=1, sort=True
)

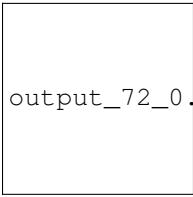
df.columns = epochs.keys()
df = df.reindex(strata)
df.T.plot(kind="bar", stacked=True, figsize=(12,5))
# plt.title("Epochs")
plt.legend(bbox_to_anchor=(1.3, 0.75))

```

(continues on next page)

(continued from previous page)

```
plt.gca().set_xticklabels(epochs.keys(), rotation="horizontal")
plt.tight_layout()
plt.savefig("img/epochs_regions.png")
plt.show()
```



output_72_0.png



img/epochs_regions.png

Fig. 8.13: Epochs

- Renaissance: largest diatonic proportion overall but mostly chromatic
- Baroque: almost completely chromatic
- Classical: enharmonic proportion increases -> more distant modulations
- This trend continues through the Romantic eras

8.7 Summary

1. We have analyzed a very specific aspect of Western classical music.
2. We have used a large(-ish) corpus to answer our research question.
 3. We have operationalized musical pieces as vectors that represent distributions of tonal pitch-classes.
 4. We have used the dimensionality-reduction technique Principal Component Analysis (PCA) in order to visually inspect the distribution of the data in 2 and 3 dimensions.
 5. We have used music-theoretical domain knowledge to find meaningful structure in this space.
 6. We have seen that pieces are largely distributed along the line of fifths.
 7. We have used Locally Weighted Scatterplot Smoothing (LOWESS) to estimate the variance in this historical process.
 8. We have seen that, historically, composers explore ever larger regions on this line and that the variance also increases.

8.8 Conclusion

1. Data-driven approaches to music analysis offer new ways of studying music history.
2. One of the largest obstacles is the lack of appropriate data (maybe you could help improve the situation?)
3. It is difficult to operationalize/formalize musical concepts.
4. Good news: there is a lot to be done for Master/PhD students!

8.9 The end

- Thank you very much for participating in this workshop
- I would appreciate it if you would send me some feedback (mail: fabian.moss@epfl.ch; Twitter: [@fabian-moss](<https://twitter.com/fabianmoss>))
- Please get in touch if you are interested in working on a small project
- Special thanks to Diana Kayser for organization and making everything possible!!!
- My funding: École Polytechnique Fédérale de Lausanne (EPFL) and Swiss National Science Foundation (SNSF)

HARMONY AND FORM IN BRAZILIAN CHORO



Fig. 9.1: A 2019 musical tribute to renowned choro composer Jacob do Bandolim, on the 50th anniversary of his death in Sao Paulo, Brazil. (c) Governo do Estado de Sao Paulo / CC BY 2.0

9.1 Data

The *Choro Songbook Corpus* contains transcriptions of the chord symbols of 295 Choro pieces [MSFR20][Mos20]. It is available on [Zenodo](#) or [Github](#). See this link for a performance of the piece “Lamentos”

As before, we can import our data from a public resource.

```
url = "https://raw.githubusercontent.com/DCMLab/choro/1.1/data/choro.tsv"
df = pd.read_csv(url, sep='\t")
```

And display the first rows of the data frame:

```
df.head()
```

CHAPTER
TEN

A DATA-DRIVEN HISTORY OF TONALITY



10.1 Musical Pieces as Tonal Pitch-Class Distributions

Musical pieces are composed of notes which have a certain pitch, a certain duration, and are located in a specific position in a piece. Here, we disregard both the location and the duration of notes and consider only the pitch dimension of musical notes when counting them. It is moreover common to consider notes to be equivalent if their respective pitches are related by one or multiple octaves, and thus to speak of *pitch classes*. Pitch classes come in two varieties. The first, most commonly used representation in computational musicology and music information retrieval, distinguishes twelve different pitch classes. This representation system also assumes the enharmonic equivalence of certain notes, e.g. F♯ and G♭, C♯ and B, etc. The assumption of enharmonic equivalence is usually a consequence of music encoding formats that assume twelve-tone equal temperament, e.g. MIDI, in which enharmonically equivalent notes are indistinguishable. The assumptions of octave and enharmonic equivalence allow to represent pitch classes

as residuals in \mathbb{Z}_{12} and to arrange them on a circle. This arrangement of pitch classes is shown in Fig. 10.1, called the *circle of fifths*.

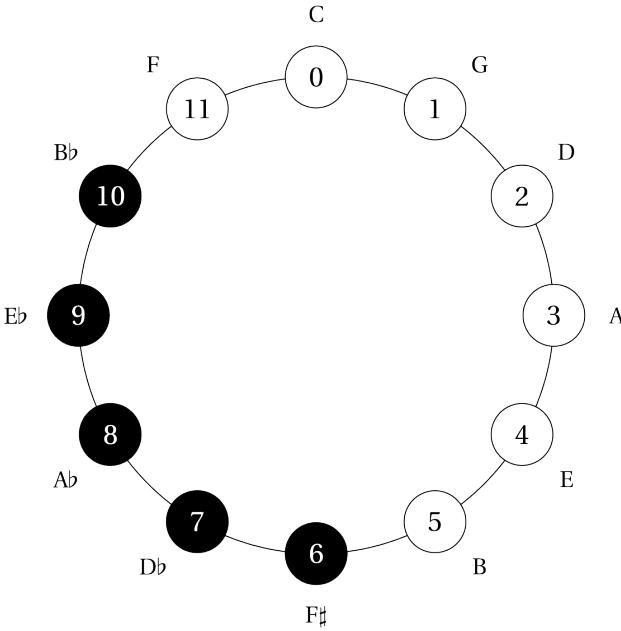


Fig. 10.1: Schematic depiction of the twelve neutral pitch-classes in \mathbb{Z}_{12} on the circle of fifths. One representative of each neutral pitch-class is shown as a tonal pitch-class label next to the node. The coloring of the nodes corresponds to the colors of the keys on the piano.

The numbers correspond to pitch classes in the order of fifths that can be transformed to chromatic ordering, the chromatic circle, by the mapping $t \mapsto 7t \bmod 12$. Pitch classes that correspond to white keys on the piano are shown in white and pitch classes that correspond to black keys on the piano are shown in black. The second variant of pitch classes does not assume enharmonic equivalence but only that octave-related notes are equivalent, and is hence more general. This representation allows to arrange pitch classes on the *line of fifths* [Tem00]. This linear ordering of tonal pitch-classes has been used by a number of music theorists, e.g. [Web51][Rie00][Han48][Bar56]. It is shown in Fig. 10.2.

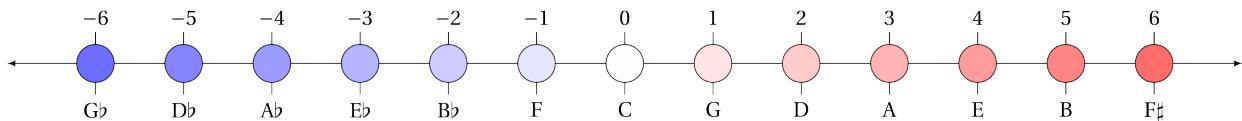


Fig. 10.2: Schematic depiction of the tonal pitch-classes on the line of fifths mapped to integers in \mathbb{Z} .

Tip: Think of other possible equivalence relations for pitch classes. E.g. which relation do we need to distinguish only two classes, high and low with respect to middle C?

Following [Tem00], we call the first representation *neutral pitch-classes* and the second one *tonal pitch-classes*. The linear structure of the line of fifths allows to associate each tonal pitch-class with an integer $k \in \mathbb{Z}$ such that this integer represents the number of perfect fifths that lie between this pitch-class and C (Gárdonyi and Nordhoff, 2002). In other words, the integer $k \in \mathbb{Z}$ corresponds to the number of flats (negative integers) or sharps (positive integers) that the major key has in which this tonal pitch-class is the root. For example, D is mapped to +2 because D major has a key signature with two sharps, Ab is mapped to -4 because the key signature for Ab major has four flats, and C is mapped to 0 because its key signature does not have any accidentals. Yet another benefit of this representation

can be seen in the way in which we associate each tonal pitch-class with a color. Positive integers ('sharpened' tonal pitch-classes) are associated with increasingly darker shades of red, negative integers are associated with increasingly darker shades of blue ('flattened' tonal pitch-classes), and C is associated with white as the neutral origin of the line of fifths. This color mapping is used throughout this part. The line of fifths does not only contain all tonal pitch-classes but also a number of central musical scales. For example, pentatonic scales are segments of length 4 (containing five pitch classes) on the line of fifths, e.g. from G \flat to B \flat , diatonic scales are defined by segments of length 6, e.g. from F to B, the early extensions of the natural diatonic scale by B \flat and F \sharp correspond to the segment spanning the eight fifths on the line between these two tonal pitch-classes, and the two whole-tone scales correspond to the odd and even numbers, respectively. Theoretically, the line of fifths extends to infinity in both directions but in actual compositions only a small segment of it is used. In the corpus that is used here, we consider only the segment from F $\flat\flat$ to B $\sharp\sharp$ because no piece in the corpus contains tonal pitch-classes outside this range. The vocabulary size of the corpus is thus $V = 35$, consisting of the seven natural pitch-classes F, C, G, D, A, E, B with two, one, or zero sharps or flats, respectively. While the transformation of tonal into neutral pitch-classes is achieved by deterministically mapping a tonal pitch-class $t \in \mathbb{Z}$ to a neutral pitch-class $t \mapsto t \bmod 12 \in \mathbb{Z}_{12}$ (in fifths ordering), the reverse direction involves some kind of inference and is called the problem of pitch spelling ([Tem01][Cam03][SRU04][CC05][Mer06]).

helix

Due to the encoding of the corpus, we adopt the representation of pieces as bags of notes (see Section~ref{sec:bagofnotes}) and represent each one as a distribution over the $V = 35$ tonal pitch-classes. That is, in the corpus used here with $D = 2012$ pieces, the tonal pitch-class distribution of a musical piece x_d is given by the relative frequencies of the tonal pitch class in that piece, for $d \in \{1, \dots, D\}$. Each piece is thus described by a V -dimensional vector, $x_d \in \mathbb{R}^V, \forall d \in \{1, \dots, D\}$, containing positive real numbers that sum to 1. In this chapter, we do not make any further assumptions about the process that generated this distribution and will postpone these considerations until Chapter~ref{chap:tonal_profiles}.

In this view, pieces simply correspond to points in the $V - 1$ -simplex

$$\Delta^{V-1} = \left\{ x_d \in \mathbb{R}^V \mid \sum_{i=1}^V x_{d,i} = 1; x_{d,i} \geq 0 \right\}.$$

In this space, those pieces with very different tonal pitch-class distributions will be very distant, whereas pieces that have similar tonal pitch-class distributions will be closer to one another and form clusters in the $V - 1$ -simplex. It is important to note that the bag-of-notes representation relies on the assumption that the V dimensions are independent, meaning that this model does not *a priori* assume any particular order between the tonal pitch-classes.

The average tonal pitch-class distribution of all pieces in the corpus is shown in Figure~ref{fig:pca_mean}. Figure~ref{fig:tpc_mean_dist_periods} in the appendix shows this distribution separately for each century from 1361 to 1943. Although we know that one can in principle order all tonal pitch-classes on the line of fifths (Fig. 10.2), we do not incorporate this assumption into the model but will show instead that it can be inferred from the data.

10.1.1 Distributions of TPCs

What are (discrete) distributions?

entropy

Hypothesis: similar pieces have similar distributions

operationalization: distances in vector space

Assumption (for coloring): tonal center is approximated by most frequent note

Taking tonal spaces into account: the line of fifths (dimensionality reduction, clustering, t-SNE, PCA)

10.1.2 Principle Component Analysis

Since the dimensionality of Δ^{V-1} is $|\{C, D, E, F, G, A, B\}| \cdot |\{\flat\flat, \flat, \natural, \sharp, \sharp\sharp\}| = 35$, it is impossible to visualize this space and pieces in it to see whether their arrangement contains any meaningful information. We address this problem by using a method for *dimensionality reduction* called *Principal Component Analysis* (PCA) [Bis06] that projects the data into a lower-dimensional space while at the same time maintaining as many characteristic properties of the original space as possible. PCA thus can aid to achieve a better understanding of the global structure of the space.¹

PCA considers the data to be represented as a matrix

$$X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_D^\top \end{bmatrix} \in [0, 1]^{D \times V}, \quad x_d \in \Delta^{V-1},$$

where the rows are given by the D data points, the pieces in the corpus, and the columns are given by the V features, the number of distinct tonal pitch-classes in the vocabulary. All entries in X range between 0 and 1 and all rows of X sum to 1 because the pieces are represented as distributions. PCA determines the $M \leq V$ largest directions and magnitudes of the variance in the data in X by first calculating the *covariance matrix*

$$K_X = \text{cov}[X, X] = E[(X - \bar{x})(X - \bar{x})^\top] \in \mathbb{R}^{D \times D},$$

where E denotes the expected value and $\bar{x} \in [0, 1]^V$ is the mean of the columns of X . The main directions of the variance in the data and their magnitude is given by the eigenvectors w_i and eigenvalues λ_i of K_X which can be calculated by solving

$$K \cdot w_i = \lambda \cdot w_i.$$

The projection into the lower-dimensional space is then achieved by selecting the M largest eigenvalues and their corresponding eigenvectors and transform the data to X' , the dimensionality reduction of X by

$$X' = X \cdot [w_1, \dots, w_M] \in \mathbb{R}^{D \times M}$$

The sum of all eigenvalues λ_i is the total amount of variance in the data and the variance explained by each principal component is given by $\lambda_i / \sum_j \lambda_j$.

In the present context, X was transformed to have zero mean before applying PCA but that the variance was not standardized to 1. This was done because the features all are on the same scale and because the differences in the variance between the respective tonal pitch-classes is of particular interest here.

10.2 Historical Development

historical expansion on line of fifths,

¹ While PCA is one of the most commonly used methods for dimensionality reduction, there are many others (sometimes relying on particular assumptions about the distribution of the data). In a qualitative comparison, *Locally Linear Embedding* [RS00] achieved a similar picture, whereas *t-distributed Stochastic Neighbor Embedding* (t-SNE) [VDMH08] that does emphasize the local over the global structure of the data did not. We opted here for PCA because it preserves most of the global structure and the interpretation of the results is straight-forward.

10.2.1 LOWESS

CHAPTER
ELEVEN

MALIAN PERCUSSION MUSIC



Fig. 11.1: Studio session in Bamako, May 2007. Left to right: Madu Jakite (Dundun), Sedu Balo (first Djembe) und Drissa Kone (second Djembe). Photo taken from <https://www.gmth.de/zeitschrift/artikel/908.aspx>.

Reproduce results from [LPJ17]. (Essentially same study in [P JL16])

CHAPTER
TWELVE

ELECTRONIC MUSIC 1950–1990



Fig. 12.1: Photo by Antoine Julien on Unsplash.

CHAPTER
THIRTEEN

CONCLUSION

Final thoughts, critical discussion...

[Some image]

CHAPTER
FOURTEEN

PYTHON BASICS

14.1 Types

- int
- float
- str - escape characters, using multiple quotes, encoding (utf-8)
- bool

Other types, e.g. complex numbers, are not covered here.

14.2 Lists

incl. list comprehension

14.3 Reading and saving files

CHAPTER
FIFTEEN

BIBLIOGRAPHY

BIBLIOGRAPHY

- [Bis06] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [Bar56] Layos Bárdos. Natürliche Tonsysteme: Methode ihrer Messung. In B. Rajeczky and L. Vargyas, editors, *Studia Memoriae Belae Bartók Sacra*, pages 209—248. Akadémiai Kiadó, 1956.
- [Cam03] E. Cambouropoulos. Pitch spelling: a computational model. *Music Perception: An Interdisciplinary Journal*, 20(4):411–429, 2003. doi:[10.1525/mp.2003.20.4.411](https://doi.org/10.1525/mp.2003.20.4.411).
- [CC05] E. Chew and Y.-C. Chen. Real-time pitch spelling using the spiral array. *Computer Music Journal*, 2005. doi:[10.1162/0148926054094378](https://doi.org/10.1162/0148926054094378).
- [Coo06] Nicholas Cook. Border Crossings: A Commentary on Henkjan Honing’s “On the Growing Role of Observation, Formalization and Experimental Method in Musicology”. *Empirical Musicology Review*, 1(1):7–11, 2006.
- [Fre14] Richard Freedman. The Renaissance chanson goes digital: digitalduchemin.org. *Early Music*, 42(4):567–578, 2014. doi:[10.1093/em/cau108](https://doi.org/10.1093/em/cau108).
- [FVC17] Richard Freedman, Raffaele Viglianti, and Adam Crandell. The Collaborative Musical Text. *Music Reference Services Quarterly*, 20(3–4):151–167, 2017. doi:[10.1080/10588167.2017.1404306](https://doi.org/10.1080/10588167.2017.1404306).
- [Han48] Jacques Handschin. *Der Toncharakter: Eine Einführung in die Tonpsychologie*. Atlantis Verlag, 1948.
- [Hon06] Henkjan Honing. On the Growing Role of Observation, Formalization and Experimental Method in Musicology. *Empirical Musicology Review*, 1(1):2–6, 2006.
- [Hur13] David Huron. On the Virtuous and the Vexatious in an Age of Big Data. *Music Perception: An Interdisciplinary Journal*, 31(1):4–9, 2013. doi:[10.1525/mp.2013.31.1.4](https://doi.org/10.1525/mp.2013.31.1.4).
- [LPJ17] Justin London, Rainer Polak, and Nori Jacoby. Rhythm histograms and musical meter: a corpus study of Malian percussion music. *Psychonomic Bulletin & Review*, 24:474–480, 2017. doi:[10.3758/s13423-016-1093-7](https://doi.org/10.3758/s13423-016-1093-7).
- [Mar16] Alan Marsden. *Music Analysis by Computer: Ontology and Epistemology*. Springer, 2016.
- [Mer06] David Meredith. The ps13 pitch spelling algorithm. *Journal of New Music Research*, 35(2):121–159, 2006. doi:[10.1080/09298210600834961](https://doi.org/10.1080/09298210600834961).
- [Mos20] Fabian C. Moss. Choro songbook corpus. 2020. doi:[10.5281/zenodo.3881347](https://doi.org/10.5281/zenodo.3881347).
- [MNHM19] Fabian C. Moss, Markus Neuwirth, Daniel Harasim, and Rohrmeier Martin. Statistical characteristics of tonal harmony: a corpus study of Beethoven’s string quartets. *PLoS ONE*, 14(6):e0217242, 2019. doi:[10.1371/journal.pone.0217242](https://doi.org/10.1371/journal.pone.0217242).
- [MNR20] Fabian C. Moss, Markus Neuwirth, and Martin Rohrmeier. Tonal Pitch-Class Counts Corpus (TP3C). 2020. doi:[10.5281/zenodo.4015177](https://doi.org/10.5281/zenodo.4015177).

- [MSFR20] Fabian C. Moss, Willian Souza Fernandes, and Martin Rohrmeier. Harmony and form in Brazilian Choro: A corpus-driven approach to musical style analysis. *Journal of New Music Research*, 2020. doi:10.1080/09298215.2020.1797109.
- [NHMR18] Markus Neuwirth, Daniel Harasim, Fabian C. Moss, and Martin Rohrmeier. The Annotated Beethoven Corpus (ABC): A Dataset of Harmonic Analyses of All Beethoven String Quartets. *Frontiers in Digital Humanities*, 5(July):1–5, 2018. doi:10.3389/fdigh.2018.00016.
- [NR16] Markus Neuwirth and Martin Rohrmeier. Wie wissenschaftlich muss Musiktheorie sein? Chancen und Herausforderungen musikalischer Korpusforschung. *Zeitschrift der Gesellschaft für Musiktheorie [Journal of the German-speaking Society of Music Theory]*, 13(2):171–193, 2016. doi:10.31751/915.
- [PFAbesser+17] Martin Pfleiderer, Klaus Frieler, Jakob Abeßer, Wolf-Georg Zaddach, and Benjamin Burkhart, editors. *Inside the Jazzomat - New Perspectives for Jazz Research*. Schott Campus, 2017.
- [PJL16] Rainer Polak, Nori Jacoby, and Justin London. Kulturelle Diversität in der empirischen Rhythmusforschung: Drei Analysen eines Audio-Korpus von Percussion-Ensemblemusik aus Mali. *Zeitschrift der Gesellschaft für Musiktheorie [Journal of the German-speaking Society of Music Theory]*, 13(2):195–235, 2016. doi:10.31751/908.
- [Pug15] Laurent Pugin. The challenge of data in digital musicology. *Frontiers in Digital Humanities*, 2(4):1–3, 2015. doi:10.3389/fdigh.2015.00004.
- [Rie00] Hugo Riemann. *Musik-Lexikon*. Max Hesse Verlag, 5 edition, 1900.
- [RS00] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. doi:10.1126/science.290.5500.2323.
- [Sch16] Kris Schaffer. What is computational musicology? online, 2016. URL: <https://medium.com/@krishaffer/what-is-computational-musicology-f25ee0a65102>.
- [SVJ+20] Kris Schaffer, Esther Vasiete, Brandon Jacquez, Aaron Davis, Diego Escalante, Calvin Hicks, Joshua McCann, Camille Noufi, and Paul Salminen. A cluster analysis of harmony in the McGill Billboard dataset. *Empirical Musicology Review*, 14(3–4):146–162, 2020. doi:10.18061/emr.v14i3-4.5576.
- [SRU04] J. Stoddard, C. Raphael, and P. E. Utgoff. Well-tempered spelling: a key-invariant pitch spelling algorithm. In *Proceedings of the 5th International Conference on Music Information Retrieval (ISMIR 2004)*. 2004. doi:10.1.1.100.9663.
- [Tem00] David Temperley. The line of fifths. *Music Analysis*, 19(3):289–319, 2000. doi:10.2307/854457.
- [Tem01] David Temperley. *The Cognition of Basic Musical Structures*. MIT Press, 2001.
- [TV13] David Temperley and Leigh VanHandel. Introduction to the Special Issue on Corpus Methods. *Music Perception: An Interdisciplinary Journal*, 31(1):1–3, 2013.
- [VDMH08] L. Van Der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(November):2579–2605, 2008.
- [Web51] Gottfried Weber. *The Theory of Musical Composition, Treated with a View to a Naturally Consecutive Arrangement of Topics*. Messrs. Robert Cocks and Co., 1851.