

Chapters 3 + 4 (DRAFT - 3/1/16)

Hugh Zabriskie

16 February 2015

1 Introducing Neural Networks

I will now shift focus to supervised learning of the chorales using both non-recurrent and recurrent neural models, seeking to improve upon the baseline models presented in Chapter 2. Neural networks are able to engage in sophisticated decision-making by passing the input forward through the network and making more complex and abstract decisions in the hidden perceptron layers. As demonstrated in the literature review at the beginning of Chapter 2, neural networks are popular in computational musicology for their ability to perform highly complex pattern recognition over large datasets and generalize when given unexpected input patterns. In particular, recurrent neural networks (RNNs) show great promise in sequence labelling because of their ability to incorporate contextual information about previous computations. In the case of musical harmony and the task of next-step harmonic prediction, it is critical to understand the local harmonic progression in order to sense the direction of the music. All chords except the tonic generate a sense of tension that insist upon a goal of resolving to the tonic. Therefore, I hypothesize that is critical for the model to understand the tension generated by previous output harmonies when predicting the next harmony. This hypothesis will be evaluated by use of an "Oracle experiment."

2 Non-Recurrent Models

The "vanilla" neural network chosen as the non-recurrent model was implemented in Lua using torch7 library. The network has 3 layers, with one layer of varying sizes that was adjusted on the basis of evaluation on a validation set. The validation set consisted of 33 chorales randomly extracted from the training set. A lookup table was added as a layer of convolution that takes the input as a vector of indexed features and outputs a matrix where each column represents an embedding for a feature in the original input. After passing through the hidden layer, the intermediate output is transformed by a softmax to output a probability distribution over all possible output classes. The criterion used was negative log likelihood (NLL), and an NLL value closer to 0 indicates how strongly the model the correct class as its prediction. The objective is to minimize negative log likelihood during training.

For each subtask, this model was trained on the same input data X and output data Y used in Chapter 2 after GCT was introduced to replace Roman numeral classification. The number of output class for each subtask were 12, 130, 4, 12, and 12, respectively. These models all trained in under a minute, requiring no more than 10 epochs to reach convergence. NLL was reported as an average of over all inputs in the training and test sets.

Table 1: "Vanilla" neural network accuracy.

Task	Training NLL	Acc%	Test NLL	Acc%
Root	0.865	68.92%	1.106	60.86%
Base	1.402	62.10%	1.686	57.83%
Inversion	0.675	73.26%	0.785	69.67%
Alto	1.473	59.22%	1.710	56.27%
Tenor	1.320	52.91%	1.751	40.05%

Update these values. They might have been from the Oracle experiment and were mislabeled.

2.1 Full Harmonization

The baseline models and the neural network were also trained on the "full harmonization task", where the output represented a full set of decisions for all subtasks. Each full set of decisions in the training set is assigned a unique class index. For the test set, only inputs that map to an index in the training set are used. To understand why, imagine a model that took an image as input and predicted whether the image depicted a cat, dog, or bird. If during training the model, only images of cats and dogs are shown, then the model will gradually adjust its weights to reflect the fact that bird is never right answer. Therefore, if the model is tested on an image of a bird, it is extremely improbable that the model will guess correctly because of the data it was trained on. Similarly, the neural network will not guess the correct full harmonization class if it never appeared during training.

The network was trained until convergence around 20 epochs, and the code is in the script `eval-sm.lua`. Random forests and neural networks continue to outperform the other models, although Random forests consistently achieving a slightly higher accuracy. The results for multinomial naive Bayes, the worst performer on average, were therefore left out.

Table 2: "Vanilla" neural network accuracy, full harmonization task.

Classifier	Test NLL	Acc%
Multiclass Logistic	–	6.26%
Random Forest	–	31.35%
Neural Network	4.614	26.35%

There are two significant flaws to this approach. First, the accuracy metric here is harsh because by combining a series of decisions into a single class, a much more fine-grained decision must be made to choose the correct class. This could be improved by creating a more sophisticated accuracy metric that took into account its overall accuracy across all subtasks. Second, I hypothesized that the model would be improved significantly if after

63 the model output a decision for a subtask that that information would be provided as part
 64 of the input for the next subtask. This hypothesis reflects the musician’s method of harmo-
 65 nization, where, for example, knowing that the next chord should be in terms of its root and
 66 base will significantly influence the selection of the alto and tenor pitches.

67

68 A new approach was therefore taken for deciding a full harmonization at each time step.
 69 Figure 1 shows the sequence of decisions that comprises the complete harmonic decision
 70 for a given time step. This decision process acknowledges the *dependent* nature of the
 71 subtasks, and therefore the ordering of their completion influences the model performance.
 72 This ordering decides the harmony in the two subtasks, while the final three decide the
 73 specific notes to voice that harmony. The final two are hypothesized to be interchangeable,
 74 since the Alto and Tenor voices are traditionally inner voices are relatively equal weight.

Let’s assume for the sake of shorthand that there are only three subtasks. Then the
 objective is to select

$$\arg \max_{x,y,z} P(X = x, Y = y, Z = z | \mathbf{X})$$

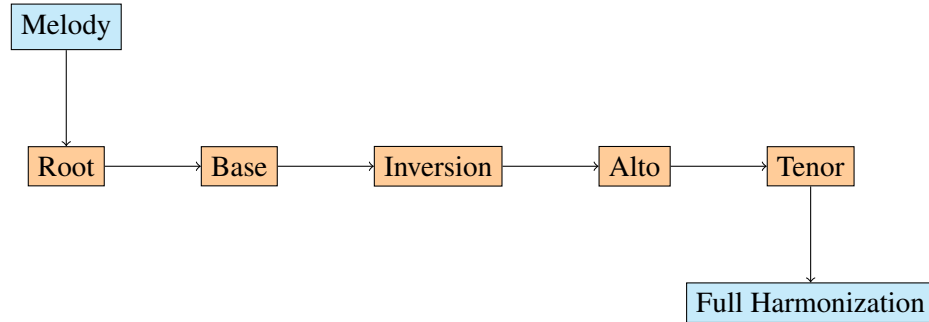
where X, Y, and Z are each random variables associated with a specific subtask and
 can only have the value of a class belonging to the variable’s subtask, and \mathbf{X} represents the
 initial input data. Let $\hat{x}, \hat{y}, \hat{z}$ be the predicted classes for the 3 subtasks. By the chain rule,
 we can show that Then by the chain rule,

$$\begin{aligned} \hat{x}, \hat{y}, \hat{z} &= \arg \max_{x,y,z} P(X = x, Y = y, Z = z | \mathbf{X}) \\ &= \arg \max_{x,y,z} P(X = x | \mathbf{X}, Y = y, Z = z) \cdot P(Y = y | \mathbf{X}, Z = z) \cdot P(Z = z | \mathbf{X}) \end{aligned}$$

75 More work needs to be done to evaluate this model and ensure it is working probably.
 76 During training, the each subtask is fed the original input data along with the correct classes
 77 for the previous subtasks in the sequence. During test time, its decision for each subtask is
 78 propagated through future decisions. The possible choices of x, y, z are selected from the

79 combinations of outputs seen in the training set. But as of now, this model achieves around
80 40% accuracy across subtasks.

Figure 1: **Sequence of harmonization subtasks with dependent decision making.**



81 **3 Recurrent Models**

82 **3.1 Oracle experiment**

83 Before, applying RNNs to this task, I wanted to know if having the context of previous
84 harmonies would in fact improve the model. An Oracle experiment is so-called termed
85 in the field of natural language processing because there is the addition of an "oracle"
86 in the data that always provides the correct answer. It is used to compare one system to
87 another system by including the "oracle" and evaluating how your model responds when
88 one component of it always provides the correct answer. Such an experiment is performed
89 here to determine whether the neural model presented in the previous section will perform
90 better with this knowledge. If so, it would strongly suggest that the context provided by a
91 recurrent model would improve results.

92 In `featurize.py`, the final three columns of X describe the correct harmony (root,
93 base, and inversion) at time step $t - 1$. This information simulates the additional informa-
94 tion provided in a recurrent model from the computation at the previous time step. The
95 baseline models were re-evaluated on this dataset as well as the neural model, using the
96 same training-test split as a GCT-updated dataset described in Chapter 2.

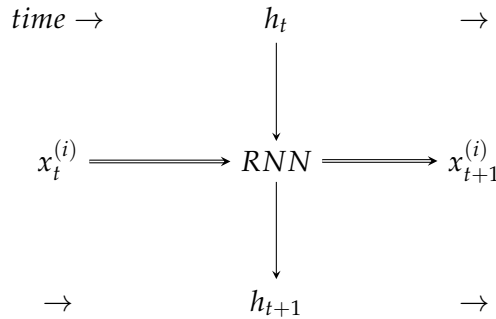
Table 3: **Baseline and neural model test accuracy, Oracle experiment**

Classifier	Root	Base	Inversion	Alto	Tenor
Multi-Class Logistic	59.75%	56.49%	64.43%	43.36%	38.93%
Multinomial Naive Bayes	52.62%	50.08%	63.05%	38.58%	35.93%
Random Forests	68.85%	61.37%	68.45%	51.96%	49.01%
Neural Network	60.86%	57.83%	69.67%	56.27%	40.05%

97 **TODO: re-run neural network values for Oracle experiment**

98 3.2 RNN Implementation

99 Based on the promising results of the Oracle experiment, I constructed a recurrent neural
100 network in Torch to determine if RNNs would improve accuracy over a *sequence* of inputs.
101 As discussed in Chapter 1, RNNs include recurrent connections that allow for feedback
102 from the computation that occurred at the previous time step. During the forward pass, at
103 each time step in the data sequence, the RNN input layer receives an external input from
104 the i th sequence $x_t^{(i)}$ as well as internal feedback from the hidden layer h_t , and then outputs
105 a distribution $x_{t+1}^{(i)}$ for the next time step. Here, each element of the sequence is a feature
106 vector describing a single melody note, drawn from the same dataset used for previous
107 models.



108

109 Each chorale represents an independent sequence that is fed into RNN and outputs
110 a synced sequence of distributions that correspond to the probability of each possibility
111 harmonization at each time step. The most likely harmonization for each time step is chosen
112 by taking the $\arg \max$ of that time step's output distribution, $x_{t+1}^{(i)}$.

$$\arg \max x_{t+1}^{(i)} = \arg \max_k \begin{bmatrix} x_{t+1_1}^{(i)} \\ x_{t+1_2}^{(i)} \\ \vdots \\ x_{t+1_k}^{(i)} \end{bmatrix}$$

Internally, h_t is stored as a "memory" of the melody features at time $t + 1$, which will be used to predict the harmonization x_{t+2} .

As a preprocessing step, each chorale was "padded" to the length of the longest chorale (in this dataset, this was 208 segments). The padding was added to all chorale sequences (both input and output vectors) using a unique "padding" feature in order not to confuse the padding with the actual musical section of the sequence. The output sequence was the same as the full harmonization task described in section 2.1 of Chapter 3, where each feature denotes a full set of harmonic decisions, including the notes for the harmonizing notes as well as a complete harmonic labelling using the GCT notation. Table 4 contains the results.

Table 4: **Neural and baseline model accuracy, full harmonization task**

Model	Train NLL	Train	Test
Multinomial Logistic	—	42.74%	25.15%
Random Forests	—	89.83%	33.41%
LSTM	0.807	38.01%	29.35%

Given that the MFFC frequency is 5%, these results show a significant increase in accuracy over previous results. The random forest model continues to outperform neural models based on accuracy, potentially for reasons related to imbalanced data, as discussed in Chapter 2.

128 **4 Chapter 4 - Looking forward**

129 **4.1 Further Exploration: The Bach Inventions**

130 The reality is a model trained to harmonize chorales has limited utility in modern times
131 (although it might have made Bach's work as *Kapellmeister* a bit easier). Musicians only
132 look to Bach's chorales as early model compositions for 4 voices and for theory-related
133 exercises. Chorales also have both a limited harmonic complexity as well as rhythmic
134 complexity, since they remain harmonically close to the tonic key and are easily segmented
135 into beat-long sections. This simplicity makes the chorales advantageous for statistical
136 learning and allowed my models to output a uniform set of harmonizations over all inputs.
137 But as chorale harmonization is increasingly well-understood as a computational task, it is
138 worth recognizing the ways in which these models can be extended to more general and
139 common musical processes. As one potential application of these models, I will explore
140 the task of composing a counterpoint line given an input melody, and how this might be
141 approached from a computational standpoint. This task is intentionally general and could be
142 applied to a wide variety of well-known musical corpuses. But here, I will focus on Bach's
143 15 Inventions because they are a body of exemplary two-voice counterpoint compositions
144 by the same composer as the Chorales. Despite being written explicitly as musical exercises,
145 they are also more melodically and rhythmically intriguing than the Chorales.

146 **4.1.1 Added complexity**

147 Generating harmonizations for inventions introduces several new layers of musical com-
148 plexity that were not considered when harmonizing the chorales. One is rhythmic com-
149 plexity, since Invention melodies cannot be easily abstracted into units of a single rhythmic
150 length. One method for describing rhythm in the generated melody is by generating predic-
151 tions at the level of the smallest rhythmic unit (say a choice of pitch or rest for every 16th
152 beat), and then determining the duration of the pitch based on the number of iterations that
153 the same pitch is consecutively output by the model. Eck and Schmidhuber (2002) found

154 this approach effective for composing blues melodies. An similar but alternative model was
155 proposed by Franklin (2006), where separate LSTMs are trained to decide pitch and dura-
156 tion. An additional layer complexity is the large-scale structure of the Invention. Inventions
157 typically consist of a series of thematic statements interleaved with episodes - modulatory
158 sections that are melodically derived from the subject and create transitions between the-
159 matic statements in different keys. The remaining sections might be termed "elaborations",
160 which include cadences at the ends of episodes and codettas that occasionally come at the
161 end of an Invention. Each type of section has a specific set of conventions that govern its
162 structure as well as the relationship between the two voices. A final layer of complexity
163 to note is the need for subject identification. In contrast to the Chorales, a single melodic
164 idea (the "invention", or subject) governs the development of the entire work, and it is
165 consistently restated in various transformations. Identification of the subject is essential
166 for determining the melodic, harmonic, and rhythmic material used in the Invention, and
167 features from the subject should be extracted to create contrapuntal melodies.

168 **4.1.2 Subject identification**

169 Bach's Inventions are two-voice compositions, where the lower voice can be characterized
170 as a *function* of the upper voice. In order for the model to learn the harmonic and rhythmic
171 material that Bach employs in the work, the model should first identify and extract features
172 from the Invention's primary *subject*. Identifying the subject generally straightforward for
173 the listener, since the primary subject is typically the first statement in the upper voice.
174 However, inventions can contain multiple subjects with varying levels of importance. And
175 in Invention No. 6 in E major, Bach introduces two subjects of equal importance and devel-
176 ops them both consistently throughout. Therefore, a better approach might be to identify
177 subjects by selecting melodic ideas based on the frequency with which they are restated.
178 Dreyfus (1996) notes that the frequency of a subject is directly correlated with its impor-
179 tance in the Invention - a principle of repetition that is generally true of melodies in Western
180 music. One candidate algorithm for subject identification based on repetition is hierarchal

181 agglomerative clustering (HAC), which can learn to "cluster" melodic segments based on
 182 a supplied distance metric. Nagler (2014) applied HAC to the task of motivic analysis in
 183 Schubert's song cycle *Die Winterreise* and found promising results in extracting the main
 184 motive from each song. However, a significant difficulty in subject identification would be
 185 the harmonic and rhythmic transformations that are often applied to the subject when it is
 186 restated. For this task, it could be useful to encode the melody as a series of intervals rather
 187 than pitches in order to identify subjects in their original form as well as their transformed
 188 state.

189

Figure 2: The subject and transformations in Invention No. 4 in D Minor, BWV 775.

The figure shows a musical score for Invention No. 4 in D Minor, BWV 775, in 3/8 time. The score is divided into six measures. Measures 1 and 2 are labeled 'm. 1 X (original statement)' and contain the interval sequences [0, [0, 3, 7]] and [7, [0, 4, 7, 10]]. Measures 3 and 4 are labeled 'm. 18' and contain [4, [0, 4, 7]] and [10, [0, 4, 7, 10]]. Measures 5 and 6 are labeled 'm. 22 SHIFT(INV(X), 10)' and contain [4, [0, 4, 7]] and [0, [0, 4, 7, 10]]. A bracket under measures 3 and 4 is labeled 'SHIFT(X, 4)'.

190 4.1.3 Form identification

191 Once rhythmic and melodic features are extracted from the most prominent subjects, the
 192 model can use these to analyze the large-scale structure of the Invention and divide the work
 193 into a series of sections. One classifier could be trained to create the general boundaries of
 194 each section, while a separate classifier could label each section as a thematic statement (*T*),
 195 episodes (*S*), and elaborations (*E*) by analyzing if and how the subject is presented in the
 196 upper voice. Complete, sequential statements of the subject suggest a thematic statement,
 197 while partial and transformed subject statements will suggest an episode. Elaborations will
 198 most likely differentiate themselves by the unique and less repetitive melodic material.

199 4.1.4 Harmonic identification

200 Before generating the melody, a useful preprocessing step would be to perform harmonic
201 analysis over the upper voice. The classifier would be trained to accept the upper voice as
202 input and output a sequence of corresponding harmonic encodings. This task is well-suited
203 to the model presented in Chapter 2 and 3, where information about a segment of the melody
204 would be classified as a GCT-encoded harmony. While many pleasant melodies can be gen-
205 erated for the counterpoint melody, the most crucial constraint is harmonic agreement. A
206 crucial purpose of the contrapuntal lower voice is to support the upper voice in harmony.
207 The upper voice typically guides the progression, while the lower voice may echo and rein-
208 force harmonic transitions, as Figure 3a illustrates.

209

210 The primary technique for implying harmonies is the use of scalar and triadic passages.
211 Not coincidentally, many of Bach’s subjects are constructed to facilitate various transfor-
212 mations (*SHIFT*, *INV*, *AUGMENT*, etc.) and be harmonically indicative. However, given
213 a scalar or triadic passage, it can be difficult to determine which of many potential chords
214 is most strongly implied. There are several important indicators, including the harmony in
215 the previous melodic unit and the beat strength of the notes in the current unit. A recurrent
216 model would likely perform well on this task because of the sequential format of the data
217 (representing the upper voice as a *sequence* of melodic units) and the feedback loop that
218 would provide information about preceding harmonies. The size of each melodic unit, how-
219 ever, is an additional parameter that would require adjustment since the implied harmonies
220 have varying durations (Figure 3b).

221

222 4.1.5 Contrapuntal melodic generation

223 Using the structural information gathered during preprocessing, a final classifier is trained
224 to generate contrapuntal melodies. The classifier would accept input features describing the
225 current relevant subject, the current section of the form, and the current implied harmony,

Figure 3: Excerpts from Invention No. 13 in A Minor, BWV 784.

A) **B)**

The image displays two musical excerpts, labeled A) and B), from Invention No. 13 in A Minor, BWV 784. Both excerpts are in 4/4 time. Excerpt A) features a sequence of chords: Am, F#ø7/A, F#ø7, and B7/D#. Excerpt B) features a sequence of chords: Bø7, G#ø, Am, and E7(b9)?.

226 as well as melodic features related to the upper voice near time t . The output feature would
227 describe the most likely pitch for time t , with an implied duration of the length of the time
228 step. Consecutive time steps of the same pitch could be merged to represented longer notes.
229 The generalization of this model allows for a wide range of corpuses to be incorporated
230 into the training data. Stylistically similar candidates include Bach's Sinfonias (the 3-
231 voice equivalent of the Inventions) and the Fugues. Stylistically separate candidates include
232 Palestrina's masses and motets, which feature significant contrapuntal composition and de-
233 velopment of melodic motives in multiple voices. The larger dataset available for training
234 will very likely improve model performance, in comparison to the computationally small
235 dataset of chorales used in this paper. Further work should be done to explore this extension
236 of the harmonization task to various musical corpuses.

237 **4.2 Conclusion**

238 In this paper, the objective was to apply a variety of neural and non-neural models to dif-
239 ferent musical processes involved in the task of chorale harmonization. These models are
240 meant to supplement existing work on computational approaches to generative and com-
241 pletive musical tasks, not to replace human musical analysis or composition. The Random
242 forest model and recurrent neural models proved most effective in learning the various sub-
243 tasks associated with harmonization, and in the future they should be evaluated on similar
244 musical tasks, such as the generation of contrapuntal voices. The accuracy results from
245 experiments with these computational models demonstrate a promising ability to learn mu-
246 sical structure and connect temporally distant events.

247 **4.3 Code and related files**

248 All code and data related to this paper is freely available on GitHub at
249 <https://github.com/glasperfan/thesis/>.

250 **References**

- 251 Dreyfus, Laurence (1996). *Bach and the Patterns of Invention*. Harvard University Press.
- 252 Eck, Douglas and Jurgen Schmidhuber (2002). “Finding temporal structure in music: Blues
253 improvisation with LSTM recurrent networks”. In: *Neural Networks for Signal Pro-
254 cessing, 2002. Proceedings of the 2002 12th IEEE Workshop on*. IEEE, pp. 747–756.
- 255 Franklin, Judy A (2006). “Jazz melody generation using recurrent networks and reinforce-
256 ment learning”. In: *International Journal on Artificial Intelligence Tools* 15.04, pp. 623–
257 650.
- 258 Nagler, Dylan Jeremy (2014). “SCHUBOT: Machine Learning Tools for the Automated
259 Analysis of Schubert’s Lieder”. Honors thesis. Harvard University.