

Classical 846 LSTM

Aaron A. Gauthier, Avijeet Kartikay and Pedro Uria Rodriguez
Machine Learning II
GWU

April 24, 2019

Table of Contents

1 Introduction

- Objective

2 Data

- Data Format
- Dataset/s

3 Encoding

- General Thoughts
- Notes
- Rests and Holds
- Multivoice polyphony

Table of Contents

4 Network Model

- Structure
- Teaching
- Making music
- Issues with generation

5 Experimental Results

- Different training approaches
- Hyperparameters
- Hyperparameters Adjusted
- Trained on various authors, individually
- Future Work

6 Conclusion

- Some last words
- The End

Music Generation System

- Build a Neural Network that can generate new music

Music Generation System

- Build a Neural Network that can generate new music

Because of the temporal nature of music...

Music Generation System

- Build a Neural Network that can generate new music

Because of the temporal nature of music...

- We will use LSTM

Music Generation System

- Build a Neural Network that can generate new music

Because of the temporal nature of music...

- We will use LSTM
- We will train on many songs from an author and/or genre

Music Generation System

- Build a Neural Network that can generate new music

Because of the temporal nature of music...

- We will use LSTM
- We will train on many songs from an author and/or genre
- The ANN will learn the probability distribution of a sequence of notes, i.e, the music

Music Generation System

- Build a Neural Network that can generate new music

Because of the temporal nature of music...

- We will use LSTM
- We will train on many songs from an author and/or genre
- The ANN will learn the probability distribution of a sequence of notes, i.e, the music
- Using this distribution, we can predict the next note based on an arbitrary number of previous notes

MIDI

- Protocol that stores music data and metadata, and allows different instruments and software to communicate with each other.

MIDI

- Protocol that stores music data and metadata, and allows different instruments and software to communicate with each other.

It is made up by a series of events, with info regarding:

MIDI

- Protocol that stores music data and metadata, and allows different instruments and software to communicate with each other.

It is made up by a series of events, with info regarding:

- Location in time
- Duration in time
- Pitch, intensity and tempo
- Other metadata

Data Sources

- Classical Piano Midi Page: <http://piano-midi.de/>
- Folders Organized by Author and Genre
- We want the data to be well organized
- We also have other sources but...

Data Sources

- Classical Piano Midi Page: <http://piano-midi.de/>
- Folders Organized by Author and Genre
- We want the data to be well organized
- We also have other sources but...

“How do I go about finding a certain classical work in MIDI format in the Internet? Some MIDI archives in the Internet contain thousands of classical works. You can find the corresponding links to them on my Linkpage. **The quality of the pieces, however, may vary considerably**”

From the author of <http://piano-midi.de/>.

Overview

- MIDI files provide us with more info than we need

Overview

- MIDI files provide us with more info than we need
- We are going to use a many-hot encoding approach

Overview

- MIDI files provide us with more info than we need
- We are going to use a many-hot encoding approach

Thus...

Overview

- MIDI files provide us with more info than we need
 - We are going to use a many-hot encoding approach
- Thus...
- Tempo will not be encoded: Too many possible values

Overview

- MIDI files provide us with more info than we need
- We are going to use a many-hot encoding approach

Thus...

- Tempo will not be encoded: Too many possible values
- Intensity will not be encoded: Same reason

Overview

- MIDI files provide us with more info than we need
- We are going to use a many-hot encoding approach

Thus...

- Tempo will not be encoded: Too many possible values
- Intensity will not be encoded: Same reason

We are essentially losing expressiveness info in order to reduce the complexity of the network.

Many-One-Hot Encoding

- Python's `music21` library to read `.mid` files.

Many-One-Hot Encoding

- Python's `music21` library to read `.mid` files.
- Preprocess the `stream` objects to get the data for the time events.

Many-One-Hot Encoding

- Python's `music21` library to read `.mid` files.
- Preprocess the `stream` objects to get the data for the time events.



$$\bar{p}_t = [0, 0, \dots, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, \dots, 0]$$

Many-One-Hot Encoding

- Python's `music21` library to read `.mid` files.
- Preprocess the `stream` objects to get the data for the time events.



$$\bar{p}_t = [0, 0, \dots, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, \dots, 0]$$

- Create a sequence where each input vector \bar{p}_t corresponds to the duration of the shortest note on the piece/s:
 $time\ step \equiv \Delta_t = t_1 - t_0 = t_2 - t_1 = \dots = cte$

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

time step = ♩ , event at $t_i = \text{♩}$

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

$$\text{time step} = \text{♩}, \text{ event at } t_i = \text{♩} \Rightarrow \bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots]$$

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

$time\ step = \text{♩}$, event at $t_i = \text{♩} \Rightarrow \bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots]$ and
 $\bar{p}_{t_{i+1}} = [0, 0, \dots, 1, 0, 0, \dots] = \bar{p}_{t_i}$

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

$time\ step = \text{♩}$, event at $t_i = \text{♩} \Rightarrow \bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots]$ and
 $\bar{p}_{t_{i+1}} = [0, 0, \dots, 1, 0, 0, \dots] = \bar{p}_{t_i} \Rightarrow \text{♩} \text{♩} \neq \text{♩}$

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

$time\ step = \text{♪}$, event at $t_i = \text{♪} \Rightarrow \bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots]$ and
 $\bar{p}_{t_{i+1}} = [0, 0, \dots, 1, 0, 0, \dots] = \bar{p}_{t_i} \Rightarrow \text{♪} \text{♪} \neq \text{♪}$

- Add *hold* component #89, which indicates that the notes played at a time event shall be held. We end up with:

Two extra dimensions

- Rests are essential to music. Add component #88 to encode rests.
- If a note is longer than the *time step*, it will be split into more than one vector. Suppose we have:

$time\ step = \text{♩}$, event at $t_i = \text{♩} \Rightarrow \bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots]$ and $\bar{p}_{t_{i+1}} = [0, 0, \dots, 1, 0, 0, \dots] = \bar{p}_{t_i} \Rightarrow \text{♩} \text{ ♩} \neq \text{♩}$

- Add *hold* component #89, which indicates that the notes played at a time event shall be held. We end up with:

$\bar{p}_{t_i} = [0, 0, \dots, 1, 0, 0, \dots, 1], \bar{p}_{t_{i+1}} = [0, 0, \dots, 1, 0, 0, \dots, 0] \Rightarrow \text{♩}$

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} + \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}} = [0, 1, \dots, 0, 1, 0, 1]$$

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} + \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}} = [0, 1, \dots, 0, 1, 0, 1]$$

- Problem: *hold* = 1 or *hold* = 0?

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} + \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}} = [0, 1, \dots, 0, 1, 0, 1]$$

- Problem: *hold* = 1 or *hold* = 0? \Rightarrow Duration of note for each hand is lost.

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} + \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}} = [0, 1, \dots, 0, 1, 0, 1]$$

- Problem: *hold* = 1 or *hold* = 0? \Rightarrow Duration of note for each hand is lost.
- Approach 2: Stack the left and right vectors together horizontally, effectively doubling the number of dimensions.

Combining Melody & Harmony

- Approach 1: Keep the same number of dimensions by just adding up the vectors from the left and right hands.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} + \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}} = [0, 1, \dots, 0, 1, 0, 1]$$

- Problem: *hold* = 1 or *hold* = 0? \Rightarrow Duration of note for each hand is lost.
- Approach 2: Stack the left and right vectors together horizontally, effectively doubling the number of dimensions.

$$\bar{p}_{t_i} = \overbrace{[0, 1, \dots, 0, 0, 1]}^{\text{left hand}} \overbrace{[0, \dots, 0, 1, 0, 0]}^{\text{right hand}}$$

Architecture

- Model:
 - LSTM
 - Sigmoid
 - Fully connected output layer
- Transfer Function: Sigmoid assumes there is no influence of one class on the other. Therefore preventing the emergence of one single dominant probability.
- Performance Index: Binary Cross Entropy is picked over categorical cross entropy as we wish to treat each dimension of the output to be separate.
- Optimizer: Adam Optimizer

Training

- Single Song:
 - Loss of the entire sequence: The loss of each item of the sequence is determined
 - Loss of the last sequence member: The loss of only the last element of the sequence is calculated. Once losses are computed, using BPTT algorithm the weights and biases are updated.
- Multiple Songs:
 - Sequence index (Approaches 1 & 2): Given a sequence index, all the losses across sequences are calculated and only then we go backwards to calculate to update weights.
 - Stacking all songs into a matrix (Approach 3): Irrespective of the songs, the encoded events are stacked one over another to form on large matrix and used to train.

Generating

- Sequence Input: $[1, 2, 3] \rightarrow [2', 3', 4], [2, 3, 4] \rightarrow [3', 4', 5], \dots$
 - A sequence of inputs is sampled from the training data and the output sequence is predicted for the input sequence
 - The last note of the output sequence is added to the previous input sequence, while the first note of this sequence is removed, and then this new sequence is used as input to the model and the new outputs are generated
 - This is repeated until a satisfactory number of output notes is generated.
- Single Input: $[1] \rightarrow [2], [1, 2] \rightarrow [2', 3], [1, 2, 3] \rightarrow [2', 3', 4], \dots$
 - In this approach, a single input note is supplied to the model.
 - The generated output, together with the previous input, is then used to predict the next input for the model
 - We do this again and again, until we reach a sequence length equal to the one our network had been trained with. Then, we start doing the sequence input approach.

Problems

- Event Stuck:
 - Issue: Network sometimes got stuck repeating a particular note over and over again
 - Solution: We added an option (default to False) to change remove the last note from the generated sequence when the exact same note was generated 6 times in a row, and add a new note randomly sampled from one of the songs of the training set
- Sequence Stuck:
 - Issue: Network getting stuck repeating a particular sequence
 - Solution: Change the note when the number of generated notes is a multiple of the sequence length
- Rest Stuck:
 - Issue: Too many rests in between the notes
 - Solution: Remove rests of duration $4 * \text{time-step}$ or if the rests have been repeated during 4 time step.

- Approach #1 did not work. The network got confused and did not learn any patterns
- Approach #2 worked decently, but tweaks to the generator were needed.
- Approach #3 worked well for some data, although it also has its flaws.

Training Hyperparameters

- `time_step`: duration of each vector
- `seq_len`: number of time steps to input as a sequence
- `hidden_size`: number of neurons in LSTM layer/s
- `lr`: learning rate
- `n_epoch`: number of training iterations
- `use_all_seq`: uses all sequences of the pieces
- `use_n_seq`: used when all sequence is false

Best Values in General

- `time_step`: 0.25
- `seq_len`: 30-50
- `hidden_size`: 178 for Melody & Harmony, 89 for Melody.
- `lr`: 0.01
- `n_epoch`: 50 - 100
- `use_all_seq`: False
- `use_n_seq`: 100

Results

Albeniz:

- It was very hard for us to get a good result with this music - seems rather unpredictable
- Tuned multiple hyperparameters such as `hidden_size`, `time_step`, `lr`, `n_epochs`

Tschaikovsky:

- Model predicted well with this music - higher quality result
- Tuned multiple hyperparameters such as `hidden_size`, `time_step`, `lr`, `n_epochs`

Other authors:

- We tried many other authors, with various degrees of success
- Also needed to tune hyperparameters for most of them

So many things we could do

- Use a mix of different authors.
- We actually tried this but with a inherently flawed approach due to a lack of time
- We could also work on different genres of music
- Also different instruments, and potentially add more than two voices
- We could even generate only melodies via LSTM and then map these melodies with harmonies training a MLP/CNN on real melodies as inputs and their harmonies as outputs.
- Expand the network architecture by combining with other types

Final Thoughts

- Music generation is very complex!
- Mathematics and music are intricately intertwined.
- Multiple complexities and levels of difficulty associated with music: Time, Duration, Tempo
- Music Files: Quality of music files, Encoding, etc.
- Quality / Complexity of Model: Different models will yield different results
- Accuracy of Tools: PyTorch, Keras, etc.
- Infinite number of possible variables to tweak
- Improvements for the Future:
- Generate harmonies for LSTM-generated melodies using MLP/CNN
- New models for other genres of music
- Improving the quality of output of the model

Happy Music Generation!



Sources

- <https://arxiv.org/abs/1709.01620>
- <http://hagan.okstate.edu/NNDesign.pdf>
- MLII GWU class slides & notes.
- <https://web.mit.edu/music21/doc/>
- <https://pytorch.org/docs/stable/index.html>
- https://gombru.github.io/2018/05/23/cross_entropy_loss/
- <https://www.depends-on-the-definition.com/guide-to-multi-label-classification-with-neural-networks/>
- Classical Piano Midi Page

You can play with our system at:

https://github.com/PedroUria/DL-Music_Generation