

音乐生成（Music Generation with RNNs）

实验报告

一、实验目的

- 理解序列建模（Sequence Modeling）与循环神经网络（RNN/LSTM）的基本原理。
- 掌握音乐数据的数值化（MIDI → 序列 → 模型输入）的过程。
- 熟悉 Notebook 中的模型结构与训练流程。
- 尝试通过调整模型结构或参数，提升音乐生成的质量。
- 输出一段可播放的旋律文件（.mid）。

二、实验方法

（一）基础环境与工具

实验基于 Python 编程语言，依托 Google Colab 与 Jupyter Notebook 开发环境，主要依赖库包括：TensorFlow/PyTorch（模型构建与训练）、midiutil（MIDI 文件生成）、numpy（数据处理）、Comet.ml（实验日志与 metrics 记录）等。实验过程与结果通过 Comet.ml 实时上传记录，可通过实验 URL 回溯查看。

（二）核心流程

- 数据预处理：**首先对 MIDI 音乐文件进行解析，提取音符、节奏和时长等关键特征并转换为字符序列。通过构建字符到索引的映射关系，将 817 首有效训练歌曲合并为统一的数值序列，并根据设定的序列长度构造输入序列与目标序列的训练样本对，为模型训练做好数据准备。
- 模型构建：**采用基于 LSTM 的神经网络架构，包含嵌入层、LSTM 层和全连接输出层三个核心组件。嵌入层将字符索引转换为 256 维的密集向量，LSTM 层处理序列数据并可根据实验调整隐藏层维度，最终通过全连接层输出词汇表大小的预测结果。模型使用交叉熵损失函数和 Adam 优化器进行参数优化。
- 模型训练：**固定训练迭代次数为 3000 次，批次大小为 8，通过调整学习率和隐藏层维度等超参数进行多组对比实验。训练期间利用 Comet.ml 平台实时记录损失变

化，每 100 次迭代保存模型检查点，训练完成后同步上传训练指标、参数配置和生成的音频文件。

- 4. **音乐生成：**以特定起始字符串作为输入，基于训练完成的模型通过多轮采样生成 1000 个字符的 ABC 格式音乐文本。从生成文本中提取有效歌曲片段并合成为音频波形，最终通过 Comet.ml 平台记录生成歌曲的数量和质量评估结果。
- 5. **改进与对比：**分析通过调整隐藏层维度、学习率和序列长度这些关键参数，深入探究不同参数组合对训练时长、损失收敛趋势以及生成音乐质量的影响。通过对比实验结果，优化模型配置以提升音乐生成的多样性和连贯性，为模型性能改进提供依据。

（三）改进方案设计

本次实验选取两项核心参数调整，同时记录序列长度的变化（原始参数与改进参数实验中 seq_length 存在差异），具体方案如下：

改进方向	原始参数配置（实验名： shallow_tower_1358）	调整后参数配置（实验名： azure_cloudberry_6273）
训练参数优化 （学习率）	learning_rate = 0.005, seq_length = 100	learning_rate = 0.002, seq_length = 50
模型结构调整 （隐藏层维度）	hidden_size = 1024	hidden_size = 2048

三、调参过程

（一）基础模型运行（原始参数）

- 1. **参数配置：**hidden_size = 1024, learning_rate = 0.005, seq_length = 100, batch_size = 8, num_training_iterations = 3000, embedding_dim = 256;
- 2. **训练过程：**训练总时长 1 分 34 秒（31.78 it/s），损失范围为（0.694060206413269, 5.051288604736328），共记录 3300 个损失数据点；
- 3. **生成效果：**仅生成 1 首歌曲（Generated song 0），上传 1 个 MIDI 资产文件（40.37 MB）；生成的音乐存在音符衔接生硬、节奏连贯性不足的问题，偶尔出现突兀音符。

（二）改进模型运行（调整学习率 + 隐藏层维度）

- 1. 参数配置：hidden_size = 2048, learning_rate = 0.002, seq_length = 50, batch_size = 8, num_training_iterations = 3000, embedding_dim = 256;
- 2. 训练过程：因隐藏层维度提升，训练总时长延长至 3 分 01 秒（16.51 it/s），损失范围缩小至 (0.4434005618095398, 5.569828987121582)，同样记录 3300 个损失数据点；训练后期损失下降更明显，最低损失较原始参数降低约 36%；
- 3. 生成效果：成功生成 3 首歌曲（Generated song 0/1/2），上传 3 个 MIDI 资产文件（总大小 40.12 MB）；生成的音乐旋律连贯性显著提升，突兀音符减少，节奏规律更清晰。

四、结果分析

（一）关键指标对比

实验指标	原始参数实验 (hidden_size=1024, lr=0.005)	改进参数实验 (hidden_size=2048, lr=0.002)	变化幅度
训练时长	1 分 34 秒	3 分 01 秒	+97.9%
训练速度 (it/s)	31.78	16.51	-48.1%
损失最小值	0.694	0.443	-36.2%
损失最大值	5.051	5.570	+10.3%

- 1. 训练效率与性能平衡：改进后模型因 hidden_size 翻倍，参数总量增加，导致训练时长延长近 1 倍、训练速度下降约 48%，但换来更优的拟合效果——损失最小值降低 36.2%，说明模型对音乐规律的学习更充分；
- 2. 生成能力提升：改进后模型生成歌曲数量从 1 首增至 3 首，且单首 MIDI 文件体积略有降低，说明模型在生成有效性与数据压缩效率上均有提升；
- 3. 损失稳定性分析：改进后损失最大值略有上升（+10.3%），推测因 seq_length 从 100 降至 50，模型对长时节奏规律的捕捉能力暂时减弱，但通过降低学习率，有效缓解了梯度震荡，确保整体损失趋势更优。

（二）生成音乐质量对比

评价维度	原始参数实验	改进参数实验
旋律连贯性	较差，音符跳跃突兀，无完整段落感	良好，段落衔接自然，突兀音符减少
节奏规律性	混乱，节拍间隔不稳定，易出现无规律停顿	清晰，基本保持稳定节拍，停顿位置符合音乐逻辑
风格统一性	单一且不明确，难以识别固定风格	每首歌曲风格相对统一，且 3 首间存在轻微风格差异，多样性提升

五、心得体会

这次实验里，把学习率从 0.005 降到 0.002，同时把隐藏层维度从 1024 提高到 2048，这两个改变一起起作用了。学习率降了，高维度的模型就不会出现梯度震荡了；维度提高了，模型的表达能力也更强了。结果就是损失减少了，生成的数量也变多了。还有，虽然把 `seq_length` 从 100 降到 50 可能会影响捕捉长时的节奏，但通过调整其他参数，总体上还是没受太大影响。

Comet.ml 这个工具记录了实验的所有参数、损失曲线、生成的素材这些重要信息，尤其是实验的回溯 URL 和详细的指标，这样对比参数的时候就更客观了，也不会出错。以后再做实验，可以用它的可视化功能，更直观地看损失的变化。

虽然提高了隐藏层维度，生成了更好的东西，但训练时间也翻倍了。这提示我们，以后优化的时候要考虑实际情况。如果想要快速迭代，可以适当降低 `hidden_size`，然后调整其他参数（比如增加 `epochs`）；如果想要高质量输出，就得平衡计算资源和训练时间。

六、实验反思问题回答

1. 模型为什么能学会“旋律规律”？

音乐的音符序列有时间顺序上的依赖关系，模型通过“嵌入层+LSTM层”的结构来学习规律：首先，嵌入层把离散的音符符号转换成低维向量，捕捉音符之间的语义关系

（比如相邻音阶的向量距离更近）；其次，LSTM 层通过门控机制（输入门、遗忘门、输出门）记住前面的音符信息，比如记住“C 大调中 C 音符后面常接 G 或 F 音符”的规律；最后，在 3000 次迭代训练中，模型通过最小化交叉熵损失，不断调整参数，优化“前序序列→下一个音符”的映射关系，最终学会音乐的旋律和节奏规律。

2. 为什么温度参数（temperature）会影响生成多样性？

温度参数决定了模型输出结果的随机性：温度等于 1 时，按照原始概率来采样，结果会符合模型学习的规律；温度大于 1 时，概率分布会被“拉长”，高概率的音符被选中的概率变大，低概率的音符几乎被忽略，生成结果更确定但多样性减少；温度小于 1 时，概率分布会被“压缩”，低概率的音符被选中的概率提升，生成结果更随机但可能偏离规律。这次实验虽然没有直接调整温度参数，但生成的歌曲数量从 1 首增加到 3 首，间接说明改进后的模型在基础概率分布学习上更优秀，为后续的温度参数调整打下了更好的基础。

3. 您的改进在哪些方面提升了音乐的自然度或节奏感？

提升音乐质量主要有两个方面的改进：一是把学习率从 0.005 降到 0.002，减少训练时的波动，让模型参数更稳定，避免出现突兀的音符（比如突然的高音跳跃），让旋律更连贯；二是把隐藏层的维度从 1024 提高到 2048，让模型能更好地提取特征，捕捉更复杂的节奏模式（比如 4/4 拍的强弱规律、八分音符和十六分音符的搭配），同时把 seq_length 从 100 降到 50，让模型在短时间内的节奏学习上更专注，最终让节拍更稳定，节奏更清晰。

4. 如何判断“音乐质量”的好坏？是否存在客观指标？

评价音乐好坏，要看“客观指标”和“主观感觉”两方面，这两方面得互相配合：

客观指标：① 损失指标（比如实验中损失值从 0.694 降到 0.443，说明模型对音乐规律的匹配有多好）；② 节奏稳定性（算一下 MIDI 文件里节拍的间隔，改进后这个间隔变化更小，节奏就更稳了）；③ 音符分布相似度（比一比生成的音乐和训练数据的音符出现频率，改进后更接近，风格就更统一了）。

主观感觉：包括旋律的连贯性（有没有突兀的地方）、节奏的规律性（是否符合常见的节拍）、风格的统一性（有没有明确的音乐风格）、听起来是否舒服（会不会觉得刺耳或杂乱）。