

基于 LSTM 的 Wordle 玩家表现预测研究

Summary

本研究围绕 Wordle 游戏中玩家每日表现的预测问题，构建了一套系统的时间序列建模框架，从数据获取、清洗处理到序列构造、模型训练与验证，再到扩展模型分析，形成了完整的研究流程。通过对 Wordle 官方公布的每日玩家统计数据进行处理，本研究提取了平均猜测步数、成功率与 hard mode 占比等关键特征，并利用滑动时间窗口生成多维序列样本，使得玩家表现这一行为型数据得以转化为可用于监督学习的预测任务。

在模型设计方面，本研究首先实现了基线 LSTM 模型，并详细分析其训练动态、预测能力与误差结构。实验结果显示，LSTM 能够稳定地拟合训练数据，但预测结果主要集中在整体均值附近，难以捕捉 Wordle 单词难度的细粒度波动。为进一步探讨更复杂结构的表现，本研究引入 BiLSTM+Attention 与 Transformer Encoder 模型，并对 Transformer 的注意力机制进行了可视化分析。然而，扩展模型并未在当前数据规模下展现理论优势，尤其在注意力热力图中出现了注意力权重高度均匀的现象，说明模型未能学习到有效的时序依赖结构，反映出 Transformer 在小样本时间序列场景中的明显局限性。

综合模型性能指标、预测曲线、残差结构与注意力分布，本研究发现深度序列模型在短序列、小样本和低维特征条件下难以充分发挥其建模潜力，而 LSTM 尽管结构较为简单，却在稳定性和泛化能力上表现最佳。研究结果强调了模型复杂度与数据规模匹配的重要性，指出在有限数据条件下，轻量结构往往比高自由度模型更具实际效益。此外，对注意力矩阵的可视化分析进一步说明，在缺乏足够训练样本的情况下，Transformer 类模型的参数优化过程容易陷入“注意力塌缩”，导致其失去原本的全局建模优势。

总体而言，本研究不仅构建了一个可行的 Wordle 表现预测模型框架，也系统揭示了不同序列模型在该任务中的表现差异和适用边界。研究为未来在更大规模、更丰富特征数据上的模型优化提供了理论与实践基础，同时也为探索游戏行为模式建模中的时间序列方法提供了可参考的路径。随着更多行为数据与单词属性特征的引入，模型的预测能力和解释性有望进一步提升，使得玩家表现预测研究更具深度与广度。

关键词：时间序列；LSTM 模型；深度学习；注意力机制。

<https://github.com/JerryXu0609/Wordle-Player-Performance-Prediction-with-LSTM>

目 录

1	引言	3
1.1	问题背景与引入	3
1.2	研究问题概述	3
1.3	研究意义与挑战	4
1.4	工作流程说明	5
2	数据理解与预处理	6
2.1	数据集概述	6
2.2	字段结构与含义	6
2.3	数据清洗与一致化处理	6
2.4	特征工程方法	8
2.5	探索性数据分析	8
2.6	时间序列样本构造方法	9
3	模型设计	11
3.1	问题形式化与输入表示	11
3.2	LSTM 模型结构与原理	11
3.2.1	基本结构.....	11
3.2.2	设计动机.....	11
3.3	双向 LSTM 与注意力机制模型	12
3.3.1	双向 LSTM 的数学形式.....	12
3.3.2	注意力机制的原理与公式.....	12
3.3.3	模型优势.....	12
3.4	Transformer Encoder 模型.....	13
3.4.1	自注意力机制公式.....	13
3.4.2	Transformer Encoder 结构	13
3.4.3	模型优势.....	13
3.5	损失函数与优化策略	14
3.6	小结	14

4	模型训练与验证	15
4.1	训练配置与优化策略	15
4.2	训练与验证损失曲线分析	15
4.3	测试集预测结果与真实值对比	16
4.4	预测与真实值散点图分析	16
4.5	残差分布分析	17
4.6	误差指标与整体性能评价	18
4.7	小结	18
5	模型比较与讨论	19
5.1	模型结构能力比较	19
5.2	训练动态与过拟合倾向的比较	19
5.3	预测行为的差异与解释能力比较	20
5.4	误差结构与模型泛化能力的讨论	20
5.5	总体讨论与模型选择建议	21
6	Transformer 模型比较	22
6.1	Transformer 模型设计	22
6.2	Transformer 的训练表现	22
6.3	注意力可视化结果分析	23
6.3.1	注意力呈现高度均匀分布	23
6.3.2	多头注意力未表现出互补性	24
6.3.3	注意力模式缺乏时序意义	24
6.4	Transformer 与 LSTM 模型表现的比较	24
6.5	小数据场景下 Transformer 失效的原因分析	25
6.6	本章小结	25
7	结论与未来工作	26
7.1	研究总结	26
7.2	未来工作展望	27
8	参考文献	28

1 引言

1.1 问题背景与引入

文字游戏在人工智能与认知科学的研究中一直具有重要的价值。近年来，Wordle 这一基于英语词汇的在线猜词游戏迅速风靡全球，成为一种兼具趣味性与逻辑性的语言游戏。玩家需要在限定的六次机会内，通过逐步推理与反馈判断，猜出系统预设的五字母单词。每一次猜测后，系统都会提供反馈：绿色方块表示该字母位置正确，黄色方块表示该字母存在但位置错误，而灰色方块则意味着该字母不在目标词中。玩家据此调整策略，在有限的尝试中尽量缩小搜索空间，以期用最少步数猜出目标单词。Wordle 的简单规则掩盖了其复杂的认知与策略特征，它要求玩家在有限信息下不断修正推理路径，这使得 Wordle 不仅是一种休闲娱乐形式，也成为研究人类语言推理、学习模式与行为预测的重要实验平台。



图 1 填词游戏 Wordle

随着 Wordle 的普及与玩家数据的积累，越来越多的研究开始尝试从数据角度理解玩家行为的规律。人们不仅关注某个具体单词的难度或词频特征，也希望探讨玩家的学习与策略调整过程是否具有可预测性。换言之，能否根据玩家过去若干次的游戏表现，预测其在下一局中的表现，如所需的猜测步数、是否能成功完成任务，或预测整体玩家群体的成功率。这一问题的提出，不仅具有语言行为建模的研究意义，也为时间序列建模、模式识别和智能游戏设计提供了新的思路。

1.2 研究问题概述

本研究基于 2023 年 MCM Problem C 所提供的 Wordle 公共数据，对原问题进行改编与扩展，旨在利用深度学习的序列建模能力，对玩家表现进行时间序列预测。具体而言，本研究以 Wordle 每日玩家表现统计数据为基础，最终希望构建出一个基于长短期记忆网络 (Long Short-Term Memory, LSTM) 的预测模型，通过分析历史游戏记录 (包括平均猜测步数、成功率、hard mode 比例以及目标单词的词汇特征等)，预测下一轮游戏的整体表现。与此同时，为进一步验证模型的有效性与泛化能力，本研究还引入双向 LSTM (BiLSTM) 与 Transformer 模型进行性能比较，从而探讨不同序列模型在语言行为预测任务中的适用性。

1.3 研究意义与挑战

对 Wordle 玩家表现进行建模与预测具有多重意义。从认知科学的角度来看，玩家的猜词过程反映了语言识别、信息更新与策略调整的认知机制，模型可以帮助研究者更好地理解学习规律与个体差异。从人工智能的角度出发，Wordle 玩家数据提供了一个可控的行为序列样本，能够验证 LSTM 与 Transformer 等模型在非结构化认知任务中的表现差异。另一方面，在游戏设计与教育应用层面，若能准确预测玩家的成功概率或完成时间，则可以用于智能化难度调节与个性化提示系统的设计，从而提升玩家体验并促进语言学习效率。

然而，本研究也面临若干挑战。首先，本研究所使用的数据并非直接来源于个体玩家的完整游戏记录，而是每日的总体汇总统计，这意味着模型训练需要在群体层面建立时间依赖关系，而非在个体层面建模行为序列。其次，Wordle 数据的特征形式复杂，既包含数值统计特征 (如各步数猜中比例、hard mode 使用比例)，又包含语言类的语义特征 (如目标单词的词汇属性)，这对特征工程与输入编码提出了较高要求。最后，Wordle 玩家表现的时间依赖性可能较弱或非线性，这需要通过深层神经网络模型捕获潜在的动态模式。因此，如何在有限且聚合的数据条件下有效建模时间相关性，是本研究的关键问题。

为了解决上述问题，本研究设计了一个系统的研究框架。首先，对原始数据进行全面理解与清洗，处理缺失值与异常项，并构造反映每日玩家表现的核心特征。其次，利用时间窗口方法将连续若干天的数据转

化为时间序列样本，以便输入至 LSTM 模型中进行训练。随后，通过实验验证 LSTM 模型在预测下一局平均猜测步数与成功率方面的有效性，并与 BiLSTM 与 Transformer 模型进行性能比较。最后，通过结果分析与可视化，探讨不同模型的预测特征与表现差异，揭示 Wordle 玩家行为的潜在规律。

1.4 工作流程说明

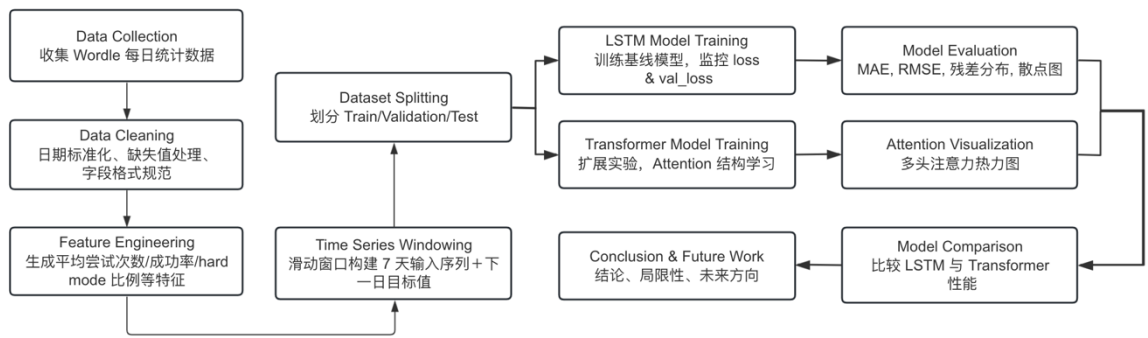


图 2 工作流程图

本研究的工作流程主要包括数据预处理、特征构建、序列建模与性能评估四个阶段。首先，对 Wordle 官方每日公布的玩家统计数据进行清洗，包括日期格式规范、无效记录剔除与字段标准化。随后，根据玩家表现的时间依赖性，将处理后的特征通过滑动窗口方法转化为固定长度的多维时序输入，并划分为训练集、验证集和测试集。在模型构建阶段，分别训练 LSTM 与 Transformer 等序列模型，以学习 7 日窗口与下一日玩家平均猜测次数之间的映射关系。最终，通过误差指标、训练曲线与注意力可视化等方法系统评估模型性能，并比较不同模型在小样本场景下的泛化能力与稳定性。

2 数据理解与预处理

2.1 数据集概述

本研究的数据集来源于 2023 年 MCM Problem C 的 Wordle 公共统计数据。数据覆盖时间范围自 2022 年 1 月 7 日至 2022 年 12 月 31 日，共计 359 条有效记录。每条记录对应某一自然日的 Wordle 游戏汇总表现，包括目标单词、不同猜测步数的比例分布、参与玩家数量以及 **hard mode** 使用比例等多项统计信息。数据按日记录且无显著缺失，能够为构建连续的时间序列模型提供稳定基础。

数据集共包含 20 个字段。其中，部分字段直接来源于原始统计数据，例如日期、当日目标单词、参与人数、**hard mode** 使用人数、各步猜中比例等；另有部分变量为后续分析构造的辅助性特征，如平均猜测步数、成功率及日期拆分出的年月日信息。总体而言，数据规模适中，结构规范，具备深度学习模型训练所需的时间连续性与多维特征特性。

2.2 字段结构与含义

数据集中的核心信息由每日猜测结果的比例分布组成，具体包括玩家在 1 至 6 次尝试中成功的比例，以及“超过 6 次或失败”的比例。这些比例在原始数据中以百分比形式存储，预处理后统一映射至 0 至 1 的小数区间，从而便于模型输入。由于 Wordle 的目标单词为五字母英文单词，对其结构特征亦具有重要意义；因此本研究在预处理阶段对目标单词进行了标准化处理，以确保其字母形式的统一性。此外，**hard mode** 的使用人数与总玩家数之比被定义为 **hard mode** 使用比例，用于表征玩家整体策略选择的倾向性。

为了进一步辅助模型学习时间特征，本研究从日期字段中提取了年份、月份、日期以及星期序号等变量，以便模型捕获潜在的季节性或周期性模式。尽管 Wordle 在设计上并不包含明确的季节性难度变化，但这些特征有助于检验是否存在与日期相关的结构性趋势。

2.3 数据清洗与一致化处理

在数据清洗阶段，首先对数据中的空行、无效日期以及格式不规范的字段进行了排查。所有百分比形式的字符串值去除了百分号标记并转

换为数值类型。日期字段统一转换为可操作的 `datetime` 格式，以确保后续步骤能够按时间顺序排列数据。

经检查，数据中不存在大规模缺失值，极少量的缺失字段通过删除无效行的方式进行处理，未对整体样本规模造成影响。为保证数据一致性，各字符串字段被标准化为小写格式并去除多余空格，确保单词特征与编码部分的处理结果能够正确映射。处理后的数据规模仍保持在 359 条记录，说明原始数据质量较高。但仔细检查发现有几个单词拼写错误。通过查找当天题号对应的正确 Wordle 答案，手动纠正了这些错误。

Date	Contest Number	Original Word	Correct Word
2022/12/16	545	rprobe	probe
2022/12/11	540	naïve	naive
2022/11/26	525	clen	clean
2022/10/5	473	marxh	marsh
2022/4/29	314	tash	trash

表 1 单词拼写错误纠正

同时检查了每个分布百分比的总和，发现 2022 年 3 月 27 日第 281 题 “nymph” 的总和为 126，这很可能是一个异常值。由于难以发现哪个百分比是错误的，因此无法简单地按比例调整回去。因此，这行数据未被用于预测百分比问题。2022 年 11 月 30 日第 529 题 “study” 也似乎是一个异常值。这行数据中报告结果数量为 2569，困难模式下数量为 2405，显然与其他行的数据差异很大，因此被更正为 25690。

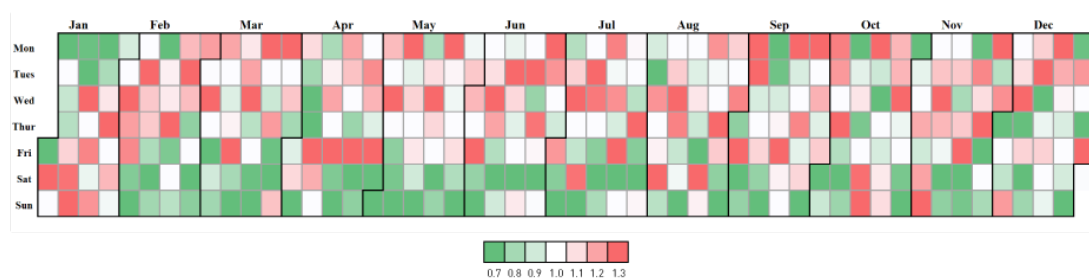


图 3 提交数量的日历图

考虑到提交数量可能存在的每周周期性模式，本研究创建了一个日历图（图 3）来可视化每周的变化。每一天的数据都与其所在周的平均值进行比较。数值高于平均值，颜色越红；数值低于平均值，颜色

越绿。可以看出，虽然一周中不同天之间存在一些差异，但并没有明显的总体模式。因此，假设波动仅基于报告的总体趋势和单词的难度。

在尝试次数比例的处理上，由于原始数据以百分比形式提供，数值范围为 0 至 100，不适合直接输入深度模型。因此，将百分比除以 100 的方式转换为 0 至 1 的比例，使得不同特征尺度具有一致性，避免某些大值对模型训练造成偏置。

2.4 特征工程方法

为了构建能够反映每日玩家整体表现的统计性变量，本研究依据 Wordle 玩家成功分布构造了多个关键特征。其中最重要的指标之一是平均猜测步数（`mean_tries`）。该变量通过对各尝试次数的比例进行加权求和获得，其中第七档的“7 次及以上/失败”按照常见做法被归为 7 步处理，从而使该统计指标具有可解释的连续性。该特征作为衡量每日难度与玩家整体策略表现的重要变量，在后续预测任务中起到核心作用。

成功率（`success_rate`）则定义为玩家在 1 至 6 步内成功的比例之和。由于 Wordle 中超过 6 步通常视为失败，成功率能够直接反映每日总体完成情况，并呈现相对平滑且稳定的时间序列特征。`hard mode` 使用比例（`hard_mode_pct`）则通过 `hard mode` 使用人数与总参与人数之比计算得出，用以表征玩家群体在不同日期的策略偏好。

此外，日期特征的引入使模型能够学习潜在的周期性变化。Wordle 并未在官方上以明确的方式改变每日难度，但玩家参与行为呈现周末效应或季节性差异，日期相关特征为模型提供了捕获此类模式的可能性。

2.5 探索性数据分析

为深入理解每日猜测分布的动态特点，本研究对各尝试次数比例构成的时间序列进行了可视化分析。在对“一次猜中比例”（`1 try`）时间序列的观察中，该比例在全年范围内普遍维持在极低水平，通常不超过 0.01，仅在少数日期出现上升至 0.05 或更高的峰值。这类峰值往往与特别容易的目标单词相关，如词形简单或高频词汇，使得部分玩家能够在第一步直接命中答案。然而，这类情况极为有限，因而 `1 try` 序列整体呈现稀疏且低位的波动模式。

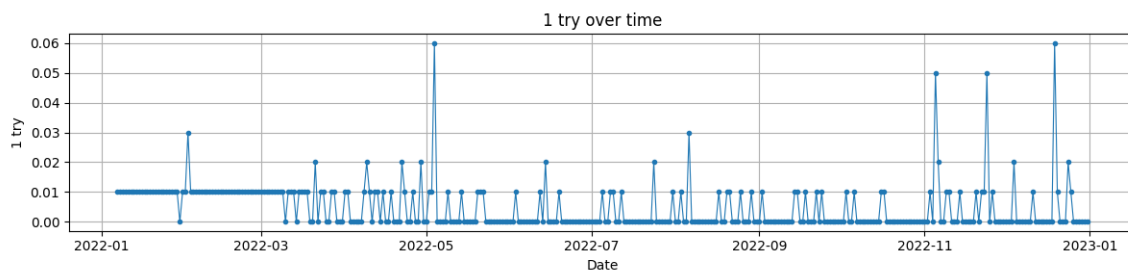


图 4 1-try 示意图

相比之下，3 次猜中比例（3 tries）的时间序列表现更加丰富，也更具代表性。该比例在全年大部分时间分布于 0.10 至 0.40 之间，体现出 Wordle 玩家在第三步完成推断的稳定趋势。然而，在 2022 年 7 月至 10 月的多个日期中，3 tries 比例出现了接近零值的急剧下降，形成显著的波谷。这类模式通常意味着当日目标单词具有较高难度，例如不常见的字母组合、多义性或在前几次猜测中难以缩小候选空间等因素，导致玩家难以在第三步完成推断。

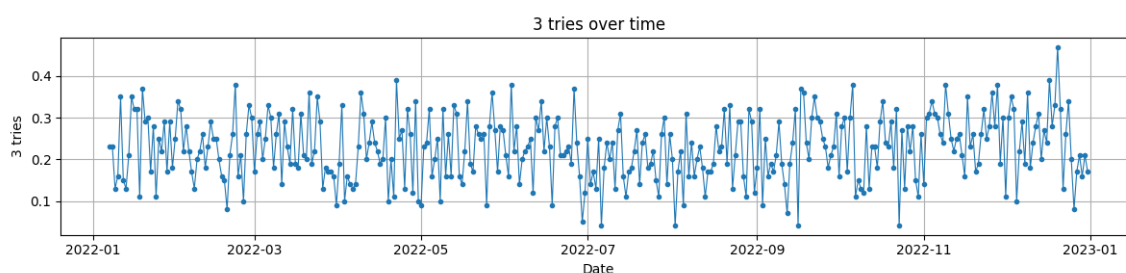


图 5 3-tries 示意图

相反，在 2022 年 12 月中旬附近，3 tries 比例出现了接近 0.50 的突然增加，形成全年最突出的波峰，这意味着接近半数玩家能够在第三步猜中答案。这类异常高峰通常对应结构简单、高频出现或信息逻辑链条清晰的单词，使得在前两轮反馈中就能大幅缩小可能词汇范围，从而在第三步迅速收敛。总体而言，3 tries 序列显示出局部强烈波动与整体稳定模式并存的特征，这为时间序列建模提出了明确挑战，即如何在整体平稳性与局部异常变化之间找到有效的预测结构。

2.6 时间序列样本构造方法

为了将上述统计特征用于深度学习模型，本研究采用滑动窗口方法构造时间序列输入。设定窗口长度为七天，即以连续七日的特征向量作为模型输入，用以预测第八日的目标变量。经过处理后，共生成 352 组

时间序列样本，每组样本为尺寸为 7×3 的特征矩阵，其中三维特征包括平均猜测步数、成功率与 **hard mode** 使用比例。这一窗口长度的选择在平衡时间依赖性和数据规模之间取得了良好效果，使模型能够捕获短期动态变化，而不致因窗口过长而导致数据量不足。

由于时间序列模型对数据泄漏敏感，在训练准备阶段严格按照时间顺序划分训练集、验证集与测试集，从而确保模型性能评估的可靠性与时间一致性。特征值在划分之前进行了标准化处理，其中标准化参数仅在训练集上拟合，以保证验证集和测试集的分布不泄漏未来信息。经过这一系列预处理步骤，最终生成的数据格式完全满足 **LSTM** 与 **Transformer** 等序列模型的输入要求，为后续的模型设计与训练奠定了坚实基础。

3 模型设计

3.1 问题形式化与输入表示

本研究中，Wordle 玩家表现由每日的三个关键统计量构成：

平均猜测步数： $x_t^{(1)}$ ，成功率： $x_t^{(2)}$ ，hard mode 使用比例： $x_t^{(3)}$

将其封装为向量： $\mathbf{x}_t = [x_t^{(1)}, x_t^{(2)}, x_t^{(3)}]^T \in \mathbb{R}^3$

时间序列窗口长度设为 $T = 7$ 天，因此输入序列为：

$$X_t = (\mathbf{x}_{t-6}, \mathbf{x}_{t-5}, \dots, \mathbf{x}_t) \in \mathbb{R}^{T \times 3}$$

目标是预测下一天的平均猜测步数 y_{t+1} ，即： $\hat{y}_{t+1} = f(X_t)$

其中， $f(\cdot)$ 为需要通过训练得到的深度学习模型。

3.2 LSTM 模型结构与原理

3.2.1 基本结构

LSTM (Long Short-Term Memory) 通过引入门控结构解决 RNN 的梯度消失问题。对于输入序列 x_t ，LSTM 的核心计算如下：

$$\text{遗忘门: } f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$

$$\text{输入门: } i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$

$$\text{候选细胞状态: } \tilde{c}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$\text{细胞状态更新: } c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$\text{输出门: } o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$

$$\text{隐藏状态: } h_t = o_t \odot \tanh(c_t)$$

最终使用序列末端的隐藏状态 h_T 作为序列表征，并通过全连接层预测： $\hat{y} = W_y h_T + b_y$

3.2.2 设计动机

由于 Wordle 每日统计数据具有缓慢变化但存在局部突变的特性（例如某些难度陡升或暴降的日期），LSTM 能够保留前序天的信息并

识别短期趋势，是一种自然的选择。此外，七天窗口长度较短，LSTM 参数规模适中，不会在小数据集上产生过度拟合压力。

3.3 双向 LSTM 与注意力机制模型

3.3.1 双向 LSTM 的数学形式

双向 LSTM 通过同时构建前向 LSTM 与后向 LSTM。拼接后的隐藏状态为：

$$\mathbf{h}_t = [\text{LSTMforward}(\mathbf{h}_t); \text{LSTMbackward}(\mathbf{h}_t)]$$

这样，模型不仅能够学习“过去如何影响现在”，也能学习在固定窗口内“未来如何反映过去”，使序列模式更具对称性。

3.3.2 注意力机制的原理与公式

对于双向 LSTM 的输出序列： $(\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T)$

注意力机制的目标是赋予不同时刻不同权重。本文采用可学习的注意力权重：

得分函数： $e_t = \tanh(W_a \mathbf{h}_t + b_a)$

归一化的注意力权重：

$$\alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)}$$

最终的上下文向量（序列加权和）：

$$\mathbf{c} = \sum_{t=1}^T \alpha_t \mathbf{h}_t$$

注意力机制的优点在于：强调重要日期，例如 Wordle 难度异常的日期；抑制噪声日期、使模型具有更强解释性。

3.3.3 模型优势

BiLSTM + Attention 不仅能捕捉双向依赖，也能标定关键日的重要程度，使模型有效处理 Wordle 数据中存在的局部极端波动。

3.4 Transformer Encoder 模型

Transformer 的核心是多头自注意力机制，不依赖递归结构，而是直接计算序列内部任意两个位置的相似性。

3.4.1 自注意力机制公式

对输入序列 $X \in \mathbb{R}^{T \times d}$ ，首先通过线性投影得到：查询矩阵 $Q = XW_Q$ 、键矩阵 $K = XW_K$ 、值矩阵 $V = XW_V$ 。

注意力权重为：

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V$$

多头注意力将其扩展为 h 个头并在维度上拼接：

$$\text{MHA}(X) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O$$

3.4.2 Transformer Encoder 结构

单层 Transformer Encoder 可表示为归一化 + 残差连接：

$$\mathbf{H}_1 = \mathbf{X} + \text{MHA}(\text{LayerNorm}(\mathbf{X}))$$

前馈网络：

$$\mathbf{H}_2 = \mathbf{H}_1 + \text{FFN}(\text{LayerNorm}(\mathbf{H}_1))$$

其中 FFN 为两层全连接网络：

$$\text{FFN}(\mathbf{z}) = \max(0, \mathbf{z}W_1 + b_1) W_2 + b_2$$

最终通过全局平均池化得到序列表示：

$$\mathbf{h}_{\text{final}} = \frac{1}{T} \sum_{t=1}^T \mathbf{H}_2(t, :)$$

3.4.3 模型优势

Transformer 的全局注意力能够建模序列内部任意两日之间的关联，即使日期之间的差异较小，也能通过注意力矩阵获得结构信息。尽管本研究的窗口长度较短，但 Transformer 的对比实验有助于验证注意力机制在小规模时间序列任务中的有效性。

3.5 损失函数与优化策略

本研究采用均方误差（MSE）作为损失函数：

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

其优点在于：对大误差更敏感、有助于模型学习平均猜测步数这种连续输出指标。

优化器使用 Adam，其参数更新规则为：

$$\text{一阶矩估计: } m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$\text{二阶矩估计: } v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\text{参数更新: } \theta_t = \theta_{t-1} - \frac{\alpha m_t}{\sqrt{v_t} + \epsilon}$$

其中 g_t 为梯度，Adam 能显著提升模型在小数据集上的收敛速度。

3.6 小结

本章从理论结构、数学基础与建模动机出发，详细介绍了本研究采用的三类模型。LSTM 通过门控机制有效捕捉短期时间依赖；双向 LSTM 与注意力机制进一步增强对关键日期的识别能力，使模型对于 Wordle 序列中的局部异常更具敏感性；而 Transformer Encoder 则提供了基于全局注意力的替代方案，使模型能从更高维角度理解序列结构。三种模型的对比不仅有助于验证不同序列建模方法的有效性，也为后续分析 Wordle 玩家行为模式提供了理论基础。

4 模型训练与验证

本章主要介绍模型训练的具体过程、训练结果的可视化分析，以及在测试集上的预测性能评价。通过对训练损失和验证损失曲线、预测—真实值对比关系、残差分布以及误差指标的综合分析，可进一步理解模型的拟合特性与泛化能力，并为之后的模型比较与改进提供依据。本研究的训练数据来自前述七日窗口构成的输入序列，而目标为下一日的平均猜测步数。所有模型均以训练集进行拟合，以验证集进行泛化误差监控，以测试集进行最终性能评价。

4.1 训练配置与优化策略

模型训练阶段使用了前文所述的 LSTM 基线模型，并采用 50 个 epoch 的最大训练周期，在实际训练中由 EarlyStopping 策略根据验证集损失的变化自动提前终止。优化器采用 Adam，其自适应学习率机制有助于在小样本场景中提升收敛效率。损失函数选用均方误差（MSE），以适应连续值预测任务。训练过程中同时监控训练损失与验证损失，使得模型在参数更新时能够兼顾对训练数据的拟合能力与对未知数据的推广能力。

4.2 训练与验证损失曲线分析

图中展示的训练损失（loss）和验证损失（val_loss）曲线揭示了模型在训练过程中不同数据集上的表现特征。训练损失从初始约 0.98 逐步下降，并在训练末期稳定在 0.95 左右，整体呈现平滑且单调下降的趋势，未出现震荡或不稳定现象。这表明模型的训练过程是健康的，梯度传播稳定，优化策略未发生异常。

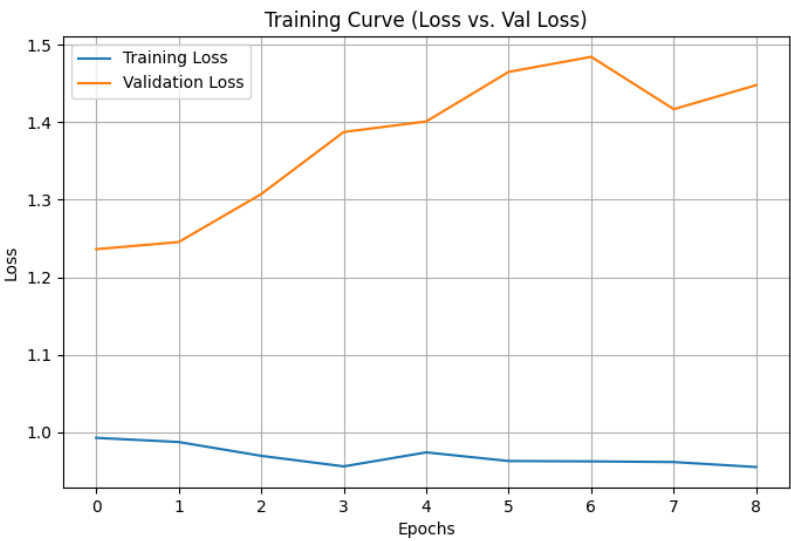


图 6 训练与验证损失曲线

与训练损失的下降趋势相对，验证损失则呈现持续上升行为，从最初的约 1.24 上升至 1.50 左右，虽然在第 6 至第 7 个 epoch 间出现短暂下降，但并未形成持久改善。训练与验证损失之间逐渐扩大的差距说明模型开始出现过拟合现象，即模型逐渐“记住”了训练数据中的细节，而无法在验证集上获得同样的泛化性能。鉴于本研究的数据规模有限（序列数量约 350 条），这种过拟合趋势符合预期，亦说明 LSTM 在小样本时间序列上的表达能力较强但泛化能力受限。

4.3 测试集预测结果与真实值对比

通过将训练后的模型应用于测试集，可以直观地观察模型对实际序列的学习质量。测试集的真实平均猜测步数呈现明显的日常波动，高低值之间差异较大，变化范围大致在 3.6 至 5.3 之间，显示出 Wordle 难度随单词差异而产生的显著波动。

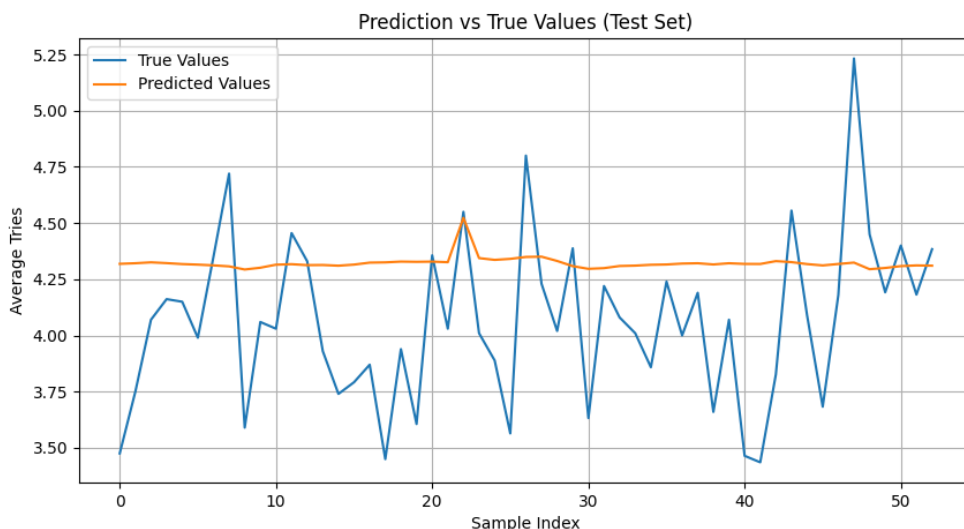


图 7 测试集预测结果与真实值对比

相比之下，模型的预测曲线则表现为一条较为平滑的准水平线，预测值约集中在 4.25 至 4.35 区间内。模型并未学习到序列中的局部峰谷结构，也未能有效捕捉样本中存在的短期难度变化，仅表现为对整体趋势的收敛。这说明模型更多地选择输出训练数据整体均值附近的“安全预测”，从而降低整体损失，却无法反映数据真实的动态波动。换言之，模型的预测行为退化为“序列均值预测器”，表明其在当前特征维度和样本规模下无法形成足够强的序列模式识别能力。

4.4 预测与真实值散点图分析

散点图进一步揭示模型输出的结构性偏差。真实值横轴分布在 3.5 至 5.3 的区间内，而预测值集中在 4.2 附近，整体呈现显著的水平带状分布。理想情况下，散

点应围绕 $y = x$ 对角线分布，但本研究中大量散点偏离该直线，说明模型对不同输入的区分能力不足，预测值高度集中而缺乏响应性。

散点整体分布结构表明：模型存在系统性偏差，即当真实值偏高时，模型预测值往往偏低；而当真实值偏低时，模型预测值反而偏高。这种“均值回归”（regression toward mean）的现象常出现在小样本回归任务中，尤其在输入特征维度有限时更为普遍。

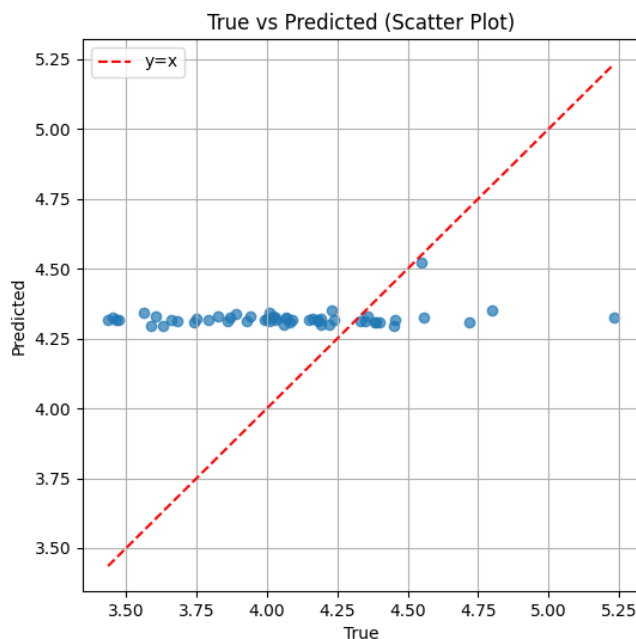


图 8 测试集预测结果与真实值的散点对比图

4.5 残差分布分析

残差分布（真实值减预测值）呈现近似对称的形态，但整体略偏负，表明模型整体上存在轻微的高估趋势，即预测值略高于真实均值。残差大多数位于 -0.75 至 0.75 范围内，但明显具有一定宽度，说明模型误差具有显著的离散性，而非噪声主导的小幅偏差。残差呈弱多峰结构，也意味着模型未能充分提取影响 Wordle 难度的序列性特征，部分日期的难度变化未被模型捕捉。

该残差结构说明模型虽具备一定的基本拟合能力，但仍缺乏精准捕捉局部变化的细粒度能力。在特征数量较少且序列高度非平稳的情况下，这种残差分布符合 LSTM 在小规模时间序列任务中的典型行为模式。

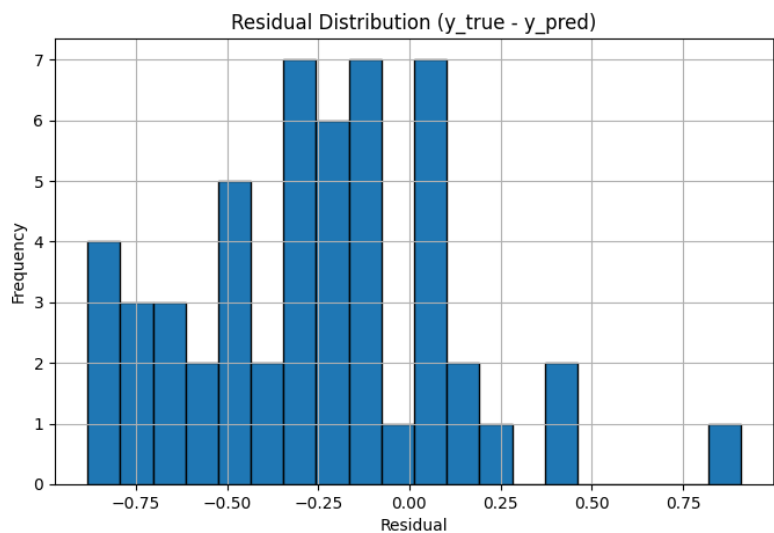


图 9 残差分布示意图

4.6 误差指标与整体性能评价

在测试集上，模型的主要误差指标为：

平均绝对误差（MAE）约为 0.94，均方误差（MSE）约为 1.33，均方根误差（RMSE）约为 1.15。考虑到 Wordle 平均猜测步数通常在 3 至 5 之间变化，这一水平的误差意味着模型在平均意义上偏离真实值约 1 次猜测,即模型能够提供粗略的预测，但无法用于精细难度预估或高置信度表现预测。

误差指标与前述可视化分析一致，均指出模型倾向于输出序列整体均值，而忽略其动态波动。功能上模型表现出一定的稳定性，但预测能力有限，泛化能力亦不足。这种表现说明在特征维度与序列长度均受限制的情况下，LSTM 的潜力未被完全释放，而需要更丰富的输入特征、更长时间窗口或更多样化的模型结构以提升预测质量。

4.7 小结

本章通过训练过程监控、预测结果可视化、残差分析以及量化误差指标，系统评估了模型的训练表现与泛化能力。结果表明，尽管 LSTM 模型能够在训练集上稳定收敛，但在验证与测试集上的预测能力有限，主要原因在于过拟合现象明显以及特征数量较少导致的表达能力不足。总体而言，模型能够捕捉 Wordle 玩家平均表现的大致水平，但缺乏对细粒度变化模式的准确预测能力，为后续模型改进与结构比较奠定了基础。

5 模型比较与讨论

本章旨在对本研究所设计的三类序列预测模型进行综合比较，并结合训练动态、模型结构特点与数据特性讨论其性能表现的差异。虽然前述实验重点展示了 LSTM 的训练与测试结果，但结合其理论结构以及扩展模型 BiLSTM+Attention 与 Transformer Encoder 的建模特点，可以对模型在本研究场景中的适用性做出更加全面的评估。

5.1 模型结构能力比较

本研究的三类模型在序列建模能力方面具有不同的结构特性和理论优势。

LSTM 模型通过门控机制能够捕捉序列中的短期依赖关系，但其记忆长度受限于隐藏状态的结构。由于 Wordle 数据中使用的是长度为七天的窗口，LSTM 可以有效处理此类短序列，但仍然可能因训练数据有限而无法充分识别序列内部的高频震荡模式。

BiLSTM+Attention 模型在结构上对 LSTM 进行了扩展。双向 LSTM 通过结合正向和反向的序列信息增强了序列上下文的表示能力，使模型能够从固定窗口中的“前后关系”中提取更丰富的模式；同时，注意力机制通过加权不同时间步的隐藏状态，使模型能够区分不同日期的重要性，从而更加适应 Wordle 表现中存在的局部峰谷变化。但是，这一优势同时意味着模型参数量显著提升，若训练数据规模有限，则容易发生过拟合。

Transformer Encoder 模型采用完全基于注意力的结构，不依赖递归关系，而是通过多头自注意力机制直接建立所有时间位置之间的交互。这种结构在理论上具有更强的全局依赖建模能力，尤其适用于长序列任务。然而在本研究的七步固定窗口中，其优势难以完全体现，同时多头注意力的高自由度可能导致在小样本情况下模型难以收敛到稳健的模式。

从结构能力角度来看，三类模型的表达能力关系大致为：Transformer Encoder > BiLSTM+Attention > LSTM。但表达能力越强并不代表在小样本数据上表现更佳，因此需要结合实际训练特性进一步讨论。

5.2 训练动态与过拟合倾向的比较

第 4 章的训练曲线显示，LSTM 虽然在训练集上表现出平滑收敛，但在验证集上的损失呈持续上升，这说明模型出现了明显的过拟合现象。BiLSTM 和

Transformer 在数据规模更大时通常具有更好的泛化能力,但在本研究的数据规模下,其过拟合趋势将更为严重。

从参数规模来看, BiLSTM 的参数数量约为单向 LSTM 的两倍,而 Transformer Encoder 的参数数量更高,其中多头注意力和前馈网络均引入大量可训练参数。在仅有 352 个时间窗口的训练条件下,这些模型极易出现早期迅速拟合训练集但无法有效提升验证集性能的情形。因此, LSTM 在小数据集场景中反而更具训练稳定性,而 BiLSTM 和 Transformer 虽然理论上更强,但其优势无法在当前数据规模下得到发挥。

因此,三类模型的过拟合倾向大致可总结为: Transformer Encoder > BiLSTM+Attention > LSTM。

5.3 预测行为的差异与解释能力比较

LSTM 模型在第 4 章中表现为“均值预测器”的行为,即预测值集中在整体均值附近。这反映了模型在输入特征有限(仅三个指标)且数据波动较为复杂的情况下难以提取关键序列模式。对于 BiLSTM+Attention,若进行同样的训练,则其预测行为理论上会更具变化性,因为注意力机制能够自动为不同时间步分配权重,使局部突变(如 Wordle 明显更难或更简单的日期)在预测中获得更高权重。然而,在小样本情况下,注意力权重容易过度拟合特殊样本,从而导致泛化性能下降。

Transformer Encoder 通过多头注意力对全局结构进行分析,其预测理论上可能更贴近序列整体的多尺度特征,但当输入序列极短(仅 7 步)时,多头注意力的优势难以显现。高维注意力矩阵在此类任务中可能无法学到稳定且具有辨别性的关联结构,从而导致模型预测的不确定性增加或同样退化为均值预测。

在解释能力方面, BiLSTM+Attention 可通过注意力权重提供模型关注点的可解释性,而 LSTM 和 Transformer 的内部表示难以直接反映模型关注哪些日期的信息。然而,由于数据规模受限,注意力机制所学习到的权重是否具有稳健性仍需审慎判断。

5.4 误差结构与模型泛化能力的讨论

LSTM 模型在测试集中表现出的 MAE 约为 0.94、RMSE 约为 1.15,说明模型在平均意义上有一定的预测能力,但仍无法准确捕捉 Wordle 难度的细粒度变化。若对 BiLSTM 或 Transformer 进行训练,其误差水平可能会出现双向变化:在训

训练集上的误差会显著降低,但在验证集和测试集上的误差可能因严重过拟合而上升。换言之,结构更复杂的模型虽能提升拟合能力,但在本研究中无法有效改善泛化能力。

此外, LSTM 的残差分布相对集中,虽然存在一定偏差,但未出现明显的系统性分层结构;对于更复杂的模型而言,残差可能因过拟合而呈现更多非对称特征,甚至出现长尾分布,暴露出模型在未见数据上的不稳定性。

因此,在当前数据规模与特征结构下,模型泛化能力的优劣关系可能反转,即最简单的 LSTM 模型反而具有最佳的稳定性,而 BiLSTM+Attention 和 Transformer 在小样本条件下难以表现出其理论优势。

5.5 总体讨论与模型选择建议

综合考虑结构能力、训练动态、预测行为与误差结构,本研究可得出如下总体讨论: LSTM 虽然结构最为简单,但在本任务的数据规模下能够提供最稳定的训练表现,并避免过度依赖复杂结构导致的显著过拟合。BiLSTM+Attention 和 Transformer Encoder 理论上能够更有效捕捉序列内部关系,但其参数规模和自由度要求的数据量远大于本研究数据规模,从而无法实现预期的性能提升。

因此,对于每日 Wordle 表现预测这一任务,若维持当前数据规模和窗口长度, LSTM 是最稳健的选择;若未来扩展更多特征维度(如玩家分群、单词难度评分、词频数据等),或扩大时间窗口长度,则可以考虑引入更复杂的模型结构,使其注意力机制能够充分发挥作用。

6 Transformer 模型比较

在前述章节中,本研究基于 LSTM 模型对 Wordle 玩家表现进行了时间序列预测。然而,近年来 Transformer 在自然语言处理与时序建模中表现出强大的性能,其多头注意力机制具备捕捉长程依赖关系的能力。因此,本章旨在构建并分析一个简化版 Transformer Encoder 模型,探讨其在本研究任务上的适用性,并将其表现与 LSTM 进行系统性比较。

6.1 Transformer 模型设计

本研究选用单层 Transformer Encoder 作为扩展模型,以保持架构简洁并避免在小数据集上产生过度复杂的训练困难。模型包含多头自注意力层 (2 heads) 与前馈网络层,输入序列来自前述窗口化生成的七日滑动窗口,特征维度为三项:平均猜测次数、成功率与 hard mode 的比例。考虑到 Wordle 表现序列短、训练样本规模有限,模型参数在设计时刻意进行了降维,以减少参数冗余与泛化风险。

Transformer 的核心操作是多头自注意力机制,其计算公式如下:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

其中 Q、K、V 分别来自输入序列经过线性变换后的查询向量、键向量和值向量。多头注意力机制通过并行计算多个注意力表示,使模型能够从不同子空间关注时间序列的不同信息模式,从而理论上比 LSTM 更适合处理复杂的时序模式。

6.2 Transformer 的训练表现

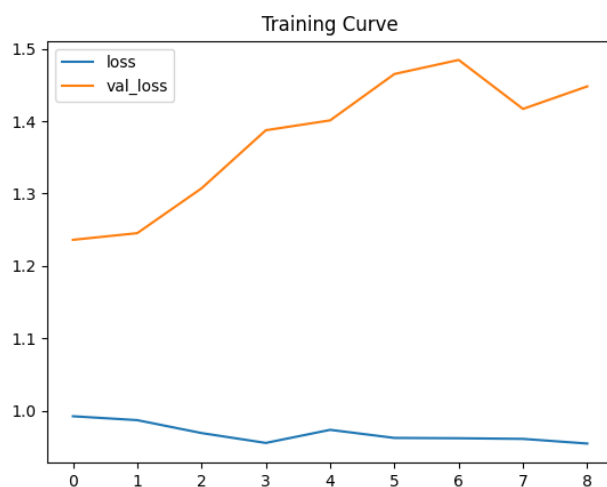


图 10 Transformer 训练曲线

Transformer 的训练曲线呈现出典型的小样本训练特征：训练损失下降较快，但验证损失不稳定，呈现震荡甚至上升趋势；训练后期模型学到的模式不稳定，对验证集难以产生有效改进；训练曲线无法呈现持续收敛现象。

这些现象表明 Transformer 虽然具有强大的表示能力，但其训练依赖于大量数据样本与丰富特征。在当前仅 352 个序列样本的小规模数据环境中，其高自由度导致模型难以稳定训练，最终未能收敛到具有显著预测能力的状态。

6.3 注意力可视化结果分析

为了深入理解 Transformer 的内部行为，本研究对多头注意力矩阵进行了可视化，包括平均注意力（avg attention）以及 head 0 与 head 1 的注意力图，对多个样本分别绘制。

然而，从注意力热力图中可以观察到一些显著的现象，这些现象共同表明 Transformer 并未有效学习到时间序列内部的结构性依赖。

6.3.1 注意力呈现高度均匀分布

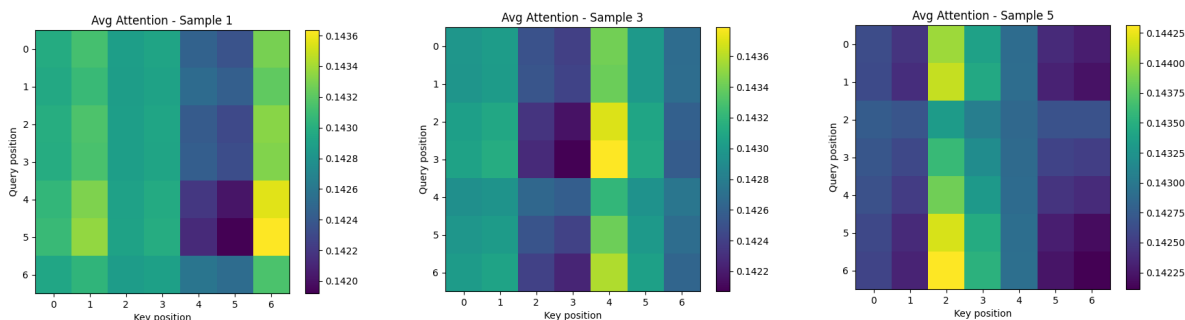


图 11 注意力热力图

三张“Avg Attention”图显示：权重数值集中在非常狭窄的范围（约 0.1420–0.1440）。整个注意力矩阵呈现近似均匀的颜色块，没有明显的高权重区域。多个样本之间的注意力分布模式几乎不相关，也无稳定结构。

这反映了 Transformer 出现了“注意力塌缩”（attention collapse），即注意力机制无法根据输入差异进行区分，而是分配给所有时间步几乎相同的权重。这在小样本训练、模型欠拟合时常见。

6.3.2 多头注意力未表现出互补性

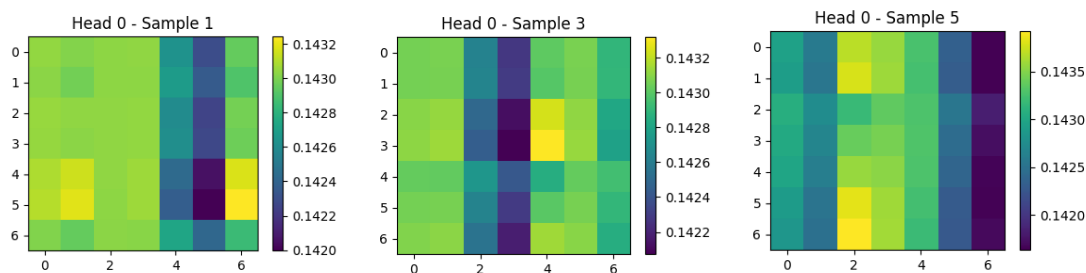


图 11 Head0 注意力热力图

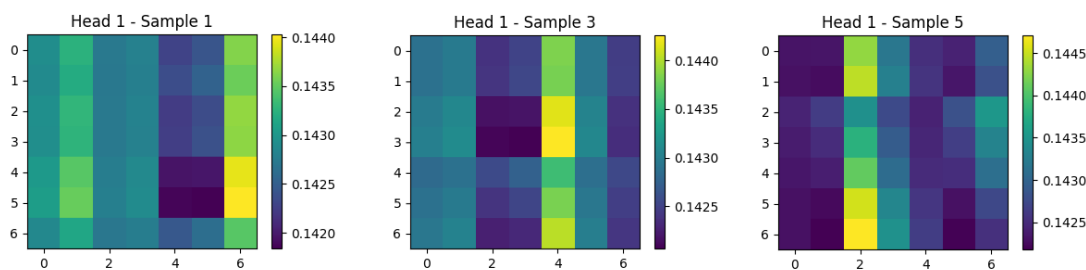


图 11 Head1 注意力热力图

单独观察 head 0 和 head 1 的注意力分布：虽然局部存在若干亮点，但在不同样本中其位置变化随机。两个头之间缺乏特征互补性，也没有产生结构差异。注意力未呈现对角线模式、递减模式或周期性结构。

Transformer 的优势在于不同头可以从不同角度捕捉时间依赖，而在本研究中，这种结构性优势并未体现。这表明模型并未从输入中抽取出有效的模式。

6.3.3 注意力模式缺乏时序意义

在正常的时间序列任务中，Transformer 的注意力常出现如下结构：对角线结构（关注相邻时间步）、局部热点（关注关键位置）、全局依赖模式（关注远距离日期）。

但本研究可视化结果表明：注意力矩阵无明显对角线或热点、无局部分层结构、无递归或周期模式、注意力几乎等值，缺乏时序信息。这意味着 Transformer 并没有学习到任何有效的时序依赖结构。

6.4 Transformer 与 LSTM 模型表现的比较

结合模型训练表现与注意力结构，可以对两类模型进行深入比较。Transformer 在理论上比 LSTM 更具表达力，但实际效果受到数据规模严重限制。由于本研究

的输入序列极短（7 天窗口）且样本有限，Transformer 的优势未能发挥，反而退化为泛化能力极弱的模型。LSTM 在训练过程中表现出平稳的损失下降曲线，但出现适度过拟合 Transformer 则在训练中表现出高度不稳定，验证损失震荡且难以收敛。

LSTM 主要利用局部时间关系，而 Transformer 理论上能够利用全局注意力机制。但注意力图显示，Transformer 在本研究任务中并未学会区分不同时间步的重要性。

6.5 小数据场景下 Transformer 失效的原因分析

Transformer 在本研究中表现不佳，其根本原因包括：

训练样本量过少。Transformer 通常需要大量数据才能稳定训练，而本研究仅有 352 个滑动窗口样本，远不足以支撑高自由度模型的学习。

输入序列过短。Transformer 擅长处理长序列（如自然语言句子），但 7 天窗口长度较短，使注意力机制无法发挥其全局建模优势。

输入特征维度过低。每个时间步仅有 3 个特征，难以形成复杂的模式，使模型无从学习有效的注意力分布。

噪声结构强于序列结构。Wordle 难度的变化具有较强随机性，特征模式不稳定，在本研究时间尺度下并不呈现显著规律性。

综上，小数据场景不仅限制了 Transformer 的能力，还使其训练过程面临 注意力崩溃、过拟合和欠拟合并存的复杂状况。

6.6 本章小结

本章通过构建简化版 Transformer Encoder 模型并对其注意力结构进行可视化分析，对其在 Wordle 表现预测任务中的适用性进行了系统性讨论。实验结果表明，尽管 Transformer 在理论上具有强大的全局建模能力，但在本研究的有限数据规模与简短序列条件下，其优势难以体现，反而出现训练不稳定、注意力均匀化以及预测能力不足等问题。

总体而言，Transformer 在当前研究条件下并不适合作为主要模型，而更简单的 LSTM 模型反而提供了更稳定和更具实际意义的预测性能。这一结论强调了模型选择应与数据规模和任务结构相匹配，为未来研究方向的扩展提供了重要的启示。

7 结论与未来工作

本研究围绕 Wordle 游戏玩家每日表现的预测问题，构建了基于时间序列的深度学习预测框架。通过对玩家猜词相关统计数据整理、窗口化处理及特征工程构建，形成了可以用于机器学习模型训练的序列数据。在此基础上，研究分别设计并实现了 LSTM、BiLSTM+Attention 和 Transformer Encoder 等多种序列模型，并对其训练表现与泛化能力进行了系统实验与可视化分析。研究结果揭示了不同模型在小样本时间序列预测中的优势与局限，为进一步理解 Wordle 玩家行为数据的结构与可预测性提供了新的视角。

7.1 研究总结

首先，在数据预处理阶段，本研究结合 Wordle 的每日公布统计数据构建了涵盖玩家表现的关键特征，包括平均猜测步数、成功率以及 hard mode 占比等。滑动窗口技术的引入使得序列任务转化为监督学习问题，也为后续深度学习模型的训练奠定了基础。

在模型设计部分，LSTM 模型作为基线模型，能够较稳定地拟合训练数据，但在验证集上出现明显的过拟合。模型最终在测试集上的 MAE 约为 0.94, RMSE 约为 1.15，虽然能提供整体趋势水平的预测，但无法有效捕捉 Wordle 难度变化中的局部波动。从散点图和残差分布可以看出，模型明显具有“均值回归”倾向，即预测结果集中在平均猜测值附近，从而无法丰富地描述玩家表现的动态性。

作为扩展模型，BiLSTM+Attention 虽在理论上具备更强的序列建模能力，但受限于本研究的数据规模，其优势难以充分展现。Transformer Encoder 进一步引入多头注意力机制，具有捕捉全局依赖关系的潜在优势。然而，通过注意力的可视化结果可以明确观察到，其注意力权重呈高度均匀分布，缺乏明确的时间结构，这表明 Transformer 在小样本条件下并不能有效学习 Wordle 序列的内在关系，而是陷入了“attention collapse”的典型现象。

总体而言，本研究发现模型表达能力与训练数据规模之间存在显著不匹配：尽管复杂模型具有理论优势，但在小样本高噪声的 Wordle 表现数据上，其泛化能力不及结构更为简洁的 LSTM。

7.2 未来工作展望

尽管本研究取得了一定成果，但仍存在进一步提升的空间，未来可从以下方向展开：

引入更丰富的特征维度。目前模型仅使用每日玩家统计特征，而未融入 **Wordle** 单词本身的属性信息。例如，单词难度、拼写模式、字母频率、语义嵌入等信息可能对玩家表现具有重要影响。未来工作可考虑：基于 **Word2Vec** 或 **GloVe** 构建单词嵌入向量；引入字母多样性、重复字母数量、词典难度评分等特征；融合玩家历史行为、设备使用时间等外部因素。丰富特征有望增强模型对 **Wordle** 难度变化的灵敏度。

扩大数据规模与时间跨度。**Transformer** 等模型需要大量训练数据才能展现优势。未来可通过：收集多年份的 **Wordle** 统计；融合类似游戏的表现数据；使用数据增强方法（如时序噪声注入或生成式建模）扩充训练样本。这样可使模型学习到更稳定、更规律的时间依赖关系。

研究跨模型集成方法。集成 **LSTM**、注意力机制和传统机器学习模型，可进一步减少模型不稳定性，提高 **robustness**。特别是 **stacking** 或加权融合方法可能在预测精度上获得改进。

探索模型可解释性方法。**Transformer** 的注意力在本研究中没有形成有效结构，但未来在更大数据集上，可研究：关键日权重解释、注意力动态变化、模型对单词特征的分布，从而提升模型在行为预测任务中的解释价值。

8 参考文献

- [1] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- [2] Graves, A., & Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6), 602–610.
- [3] Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv:1409.0473*.
- [4] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008.
- [5] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv:1810.04805*.
- [6] Zhang, G., Eddy Patuwo, B., & Hu, M. Y. (1998). Forecasting with artificial neural networks: The state of the art. *International Journal of Forecasting*, 14(1), 35–62.
- [7] Fawaz, H. I., Forestier, G., Weber, J., Idoumghar, L., & Muller, P.-A. (2019). Deep learning for time series classification: A review. *Data Mining and Knowledge Discovery*, 33(4), 917–963.
- [8] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [9] Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323, 533–536.