

Informer: Beyond Efficient Transformer for Long Sequence 时间序列预测

周浩一,¹ 张上航,² 彭洁琪,¹ 张帅,⁴ 李建新,¹
熊辉,³ 张万才,³
¹ 北京航空航天大学 ² 加州大学伯克利分校 ³ 罗格斯大学
⁴ 北京国网富达科技发展有限公司
{zhouhy, pengjq, zhangs, lijx}@act.buaa.edu.cn, shz@eecs.berkeley.edu, xionghui@gmail.com,
zhangwancai@sdee.sgcc.com.cn

抽象的

许多现实世界的应用程序需要长时间的预测
时序时间序列,如用电量计划等。长序列时间序列预测 (LSTF) 需求

模型的高预测能力,即能力
捕获精确的远程依赖耦合
有效地输出和输入。最近的研究表明,
Transformer 提高预测能力的潜力。
但是,Transformer 存在几个严重的问题
阻止它直接适用于 LSTF,例如
二次时间复杂度,高内存使用,以及编码器-解码器架构的限制。为了解决
这些问题,我们设计了一个高效的基于变压器的

LSTF 模型,名为 Informer,具有三个独特的特征:(i) ProbSparse Self-
attention 机制,它
在时间复杂度和内存使用上达到 $O(L \log L)$,
并且在序列依赖性方面具有相当的性能
结盟。(ii) 自注意力提炼通过级联层输入减半来突出主导注意力,并且
有效地
处理极长的输入序列。(iii) 生成
风格解码器,虽然概念上很简单,但可以预测长
一次前向操作的时间序列序列,而不是
循序渐进的方式,极大地改善了推理
长序列预测的速度。广泛的实验
在四个大规模数据集上表明 Informer 显著优于现有方法并提供了新的

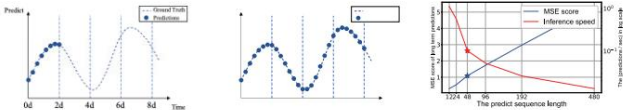
LSTF 问题的解决方案。

介绍

时间序列预测是许多领域的关键要素
域,例如传感器网络监控 (Papadimitriou
和 Yu 2006)、能源和智能电网管理、经济和金融 (Zhu 和 Shasha 2002) 和疾
病
传播分析 (Matsubara et al. 2014)。在这些场景中,我们可以利用大量的时
间序列
过去行为的数据,以进行长期预测,
即长序列时间序列预测 (LSTF)。然而,现有的方法是在有限的问题下设计的

设置,例如预测 48 个点或更少 (Hochreiter 和
施米德胡伯 1997;李等人。2018;于等人。2017;刘等人。
2019;秦等人。2017;温等人。2017)。越来越
长序列使模型的预测能力
有些人认为这种趋势正在保持

版权所有 c 2021,人工促进协会
情报 (www.aaai.org)。版权所有。



(a) 短序列
预测。
(b) 长序列
预测。
(c) 运行 LSTM
长序列。

图 1:(a) 短序列预测仅揭示
不久的将来。(b) 长序列时间序列预测可以
为更好的政策规划和
投资保护。(c) 现有方法的预测能力限制了长序列的性能,即

从长度 = 48 开始,MSE 上升到无法接受的高度,
推理速度迅速下降。

LSTF 研究。作为一个经验例子,图 (1)显示了
在真实数据集上预测结果,其中 LSTM 网络从短期内预测电力变压器站的每小
时温度 (12 分,0.5
天)到长期 (480点,20天)。这
当预测时,整体性能差距很大
长度大于 48 个点 (图 (1 (c)中的实心星)。
MSE 分数上升到不能令人满意的性能,推理速度急剧下降,LSTM 模型失败。

LSTF 的主要挑战是增强预测能力以满足日益增长的序列需求,这需要 (a)
非凡的远程对齐能力和 (b) 对长序列输入和输出的有效操作。最近,
Transformer 模型在捕获远程依赖方面表现出比

RNN 模型。自注意力机制可以减少
网络信号传播路径进入的最大长度
理论上最短的 $O(1)$ 并避免循环结构,
Transformer 显示出解决 LSTF 问题的巨大潜力。但另一方面,自注意力机制

违反要求 (b),因为它在 L 长度输入/输出上的 L 二次计算和内存消耗。

一些大型的 Transformer 模型倾注了资源和
在 NLP 任务上产生了令人印象深刻的结果 (Brown 等人,2020),
但是在数十个 GPU 上的训练和昂贵的部署
成本使这些模型在现实世界的 LSTF 上无法承受
问题。自注意力机制的效率和

Transformer框架成为应用瓶颈
他们解决了 LSTF 问题。因此,在本文中,我们试图回答一个问题:Transformer 模型是否可以改进为
计算、内存和架构的效率也很高
保持较高的预测能力?

Vanilla Transformer (Vaswani et al. 2017) 在求解 LSTF 时存在三个显著限制:

- 1.self-attention的二次计算。这
自注意力机制的原子操作,即
规范点积,导致时间复杂度和
每层的内存使用量为 $O(L^2)$ 。
- 2.长时间堆叠层的内存瓶颈
输入。J编码器/解码器层的堆栈使
总内存使用量为 $O(J \cdot L^2)$,这限制了
接收长序列输入的模式可扩展性。
- 3.预测多头输出的速度骤降。这
vanilla Transformer 的动态解码使得
与基于 RNN 的模型一样慢的逐步推理,
建议在图 (1c)中。

有一些关于提高效率的前期工作
自我关注。稀疏变压器 (Child et al. 2019) ,
LogSparse Transformer (Li et al. 2019) 和 Longformer
(Beltagy,Peters 和 Cohan 2020)都使用启发式方法
解决限制 1 并将自我注意机制的复杂性降低到 $O(L \log L)$,其中它们的效率

收益是有限的 (Qiu et al. 2019) 。改革者 (基塔耶夫、凯撒、
和 Levskaya 2019)也通过局部敏感的哈希自注意力实现了 $O(L \log L)$,但它
仅适用于极长的序列。最近,Linformer (Wang

等。2020) 声称线性复杂度 $O(L)$,但该项目
对于实际的长序列,矩阵不能固定,这可能有退化到 $O(L^2)$
的风险。
Transformer-XL (Dai et al. 2019) 和 Compressive Transformer (Rae
et al. 2019) 使用辅助隐藏状态来捕获远程依赖,这可能会放大限制

- 1.不利打破效率瓶颈。一切
作品主要关注限制1,限制2&3
仍然存在于 LSTF 问题中。为了提高预测能力,我们将解决所有这些问题并实现改进
超出了建议的 Informer 的效率。
- 为此,我们的工作明确地深入研究了这三个问题。我们研究了自注意力机制的稀疏性,改进了网络组件,并进行了广泛的实验。本文的贡献

总结如下:

- 我们建议 Informer 成功增强 LSTF 问题的预测能力,这验证了
类似 Transformer 的模型在捕获长期之间的个体长期依赖方面的
潜在价值
时序时序输出和输入。
- 我们提出 ProbSparse 自注意力机制来
有效地取代了规范的自注意力和
它实现了 $O(L \log L)$ 时间复杂度和
 $O(L \log L)$ 内存使用量。
- 我们提出了在J-stacking 层中主导注意力分数的自注意力蒸馏操作
特权
并大幅降低总空间复杂度
 $O((2 -)L \log L)$ 。
- 我们提出生成式风格解码器来获取长

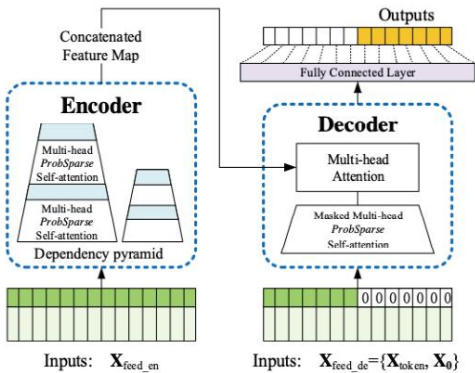


图 2:Informer 模型的整体图。这
左边是Encoder,它接收到海量的长序列
输入 (绿色系列)。我们已经替换了规范
使用建议的 ProbSparse 自注意力进行自注意力。
蓝色梯形是自注意力蒸馏操作
提取主导注意力,减少网络规模
尖锐地。层堆叠副本提高了鲁棒性。
对于右侧部分,解码器接收长序列输入,
将目标元素填充为零,测量加权
特征图的注意力组合,并立即以生成方式预测输出元素 (橙色系列) 。

只需要向前一步的序列输出,
同时避免累积误差扩散
在推理阶段。

初步的

我们首先提供问题定义。滚滚之下
具有固定大小窗口的预测设置,我们有
输入 $X \in \mathbb{R}^{L \times d} = \{x_t^1, \dots, x_t^d \mid x_t^i \in \mathbb{R}^d\}$ 在时间 t ,
输出是预测对应的序列 $y_t = \{y_t^1, \dots, y_t^d \mid y_t^i \in \mathbb{R}^d\}$ 。LSTF 问题鼓励一个
比以前的作品更长的输出长度 L_y (Cho et al.
2014; Sutskever,Vinyals 和 Le 2014)并且特征维度不限于单变量情况
($d_y \geq 1$) 。
编码器-解码器架构许多流行的模型是
设计将输入表示 X “编码”为隐藏状态表示 H_t 并从 $H_t = \{h_t^1, \dots, h_t^d\}$ “解码”
输出表示 Y_t 。推论

涉及一个名为“动态解码”的逐步过程,
其中解码器计算一个新的隐藏状态 h_{t+1} 从
前一个状态 h_t 以及来自第 k 个的其他必要输出
step 然后预测第 $(k+1)$ 个序列 y_{t+1} 。
输入表示统一的输入表示是
用于增强时间序列输入的全局位置上下文和局部时间上下文。为避免使描述
变得琐碎,我们将详细信息放在附录 B 中。

方法

现有的时间序列预测方法大致可以
分为两类1。经典时间序列模型
¹ 限于篇幅,完整的相关工作调查
附录 A 中提供。

els 是时间序列预测的可靠主力 (Box et al. 2015; Ray 1990; Seeger et al. 2017; Seeger, Salinas 和 Flunkert 2016)以及深度学习技术主要通过以下方式开发编码器-解码器预测范式使用 RNN 及其变体 (Hochreiter 和 Schmidhuber 1997;李等人。 2018;于等人。 2017)。我们提议的告密者保持编码器-解码器架构,同时针对 LSTF 问题。请参阅图 (2)的概述和以下部分了解详细信息。

高效的自注意力机制

(Vaswani et al. 2017) 中的规范自注意力是在接收元组输入时定义 (查询、键、值) 并将缩放的点积执行为 $A(Q, K, V) =$

$\text{Softmax}(\frac{QK^T}{\sqrt{d}})V$, 其中 $Q \in \mathbb{R}^{L_Q \times d}$, $K \in \mathbb{R}^{L_K \times d}$, $V \in \mathbb{R}^{L_V \times d}$, d 为输入维度。进一步讨论自注意力机制,让 q_i, k_i, v_i 代表第 i 个分别排在 Q, K, V 中。按照中的公式 (Tsai et al. 2019), 第 i 个查询的注意力被定义为概率形式的核更平滑:

$$A(q_i, K, V) = \sum_j \frac{p(k_j | q_i)}{\sum_l p(k_l | q_i)} v_j = E_{p(k_j | q_i)} [v_j | k(q_i, k_l)], \quad (1)$$

其中 $p(k_j | q_i) = \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_l \exp(\frac{k(q_i, k_l)}{\sqrt{d}})}$ 和 $k(q_i, k_j)$ 选择非对称指数核 $\exp(\frac{k(q_i, k_j)}{\sqrt{d}})$ 。自己 d attention 结合价值并获得输出基于计算概率 $p(k_j | q_i)$ 。它需要二次方乘以点积计算和 $O(L_Q L_K)$ 内存使用时间,这是增强预测的主要缺点

容量。

之前的一些尝试表明,分布自注意概率具有潜在的稀疏性,并且它们设计了一些“选择性”计数策略 $p(k_j | q_i)$ 不会显着影响性能。这 Sparse Transformer (Child et al. 2019) 结合了两者的行输出和列输入,其中稀疏性产生于分离的空间相关性。对数稀疏 Transformer (Li et al. 2019) 注意到了自我注意并迫使每个细胞注意其先前的一个指数步长。Longformer (贝尔塔吉, Peters 和 Cohan 2020) 将前两部作品扩展到更多复杂的稀疏配置。但是,它们是有限的从以下启发式方法到理论分析和用相同的策略处理每个多头自注意力,这缩小了它的进一步改进。

为了激发我们的方法,我们首先执行定性对经典自我注意的学习注意模式的评估。“稀疏”自注意力得分形式

一个长尾分布 (详见附录 C),即很少有点积对引起主要关注,并且其他的可以忽略。那么,接下来的问题是如何区分它们?

从等式 (1) 查询稀疏度测量,第 i 个查询对所有键的注意力被定义为概率 $p(k_j | q_i)$, 输出是它与值 v 的组合。

占主导地位的点积对鼓励相应查询的注意力概率分布远离

均匀分布。如果 $p(k_j | q_i)$ 接近于均匀 $\text{dis } 1$ 贡献 $q(k_j | q_i) = \text{self-attention}$ 变成了 triv LK , 值 V 的总和,对于住宅投入来说是多余的。自然地,分布 p 和

q 可用于区分“重要”查询。我们通过 Kullback-Leibler 散度测量“相似度” $KL(q || p) = \ln LK \sum_{l=1}^{L_K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}{\sum_{j=1}^{L_K} \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l=1}^{L_K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}} - \ln LK$ 。删除常量,我们定义第 i 个查询的稀疏度测量为

$$M(q_i, K) = \ln \sum_{j=1}^{L_K} \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l=1}^{L_K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l=1}^{L_K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}}, \quad (2)$$

其中第一项是 q_i on 的 Log-Sum-Exp (LSE) 所有的键,第二项是算术平均值他们。如果第 i 个查询获得更大的 $M(q_i, K)$, 则其注意概率 p 更“多样化”,并且有很高的机会在标题字段中包含主要的点积对长尾自注意力分布。

ProbSparse Self-attention 基于建议的测量,我们通过允许每个键只关注 u 个主要查询来实现 ProbSparse Self-attention:

$$A(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d}})V, \quad (3)$$

其中 Q 是与 q 大小相同的稀疏矩阵,它仅包含稀疏度量 $M(q, K)$ 下的 Top- u 查询。由恒定的采样因子 c 控制,

我们设置 $u = c \cdot L_Q$, 这使得 ProbSparse self attention 只需要计算 $O(\ln L_Q)$ 点积每个查询键查找和层内存使用保持 $O(L_K \ln L_Q)$ 。

然而,遍历测量 $M(q_i, K)$ 的所有查询需要计算每个点积对,

即二次 $O(L_Q L_K)$, 并且 LSE 操作具有潜在的数值稳定性问题。受此启发,我们提出了查询稀疏度测量的近似值。

引理 1. 对于每个查询 $q_i \in \mathbb{R}^d$ 和 $k_j \in \mathbb{R}^d$ 键集 K , 我们的边界为 $\ln LK \leq M(q_i, K) \leq \frac{\max_{j \in K} \{ \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l \in K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}} \}}{L_K} + \ln LK$ 。当 $q_i \in K$ 时, 立。

根据引理 1 (证明在附录 D.1 中给出), 我们建议最大平均测量为

$$\overline{M}(q_i, K) = \max_j \{ \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l \in K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}} \} - \frac{1}{L_K} \sum_{j \in K} \frac{\exp(\frac{k(q_i, k_j)}{\sqrt{d}})}{\sum_{l \in K} \frac{\exp(\frac{k(q_i, k_l)}{\sqrt{d}})}}. \quad (4)$$

Top- u 的顺序在边界松弛中成立与命题 1 (参见附录 D.2 中的证明)。在下面长尾分布,我们只需要随机抽样 $U = L_Q \ln LK$ 点积对计算 $M(q_i, K)$, 即用零填充其他对。我们从 $\overline{M}(q_i, K)$ 中的最大算子不太敏感为零值并且是数值稳定的。在实践中,查询和键的输入长度通常是等价的,即 $L_Q = L_K = L$ 这样总的 ProbSparse self-attention 时间复杂度和空间复杂度为 $O(L \ln L)$ 。

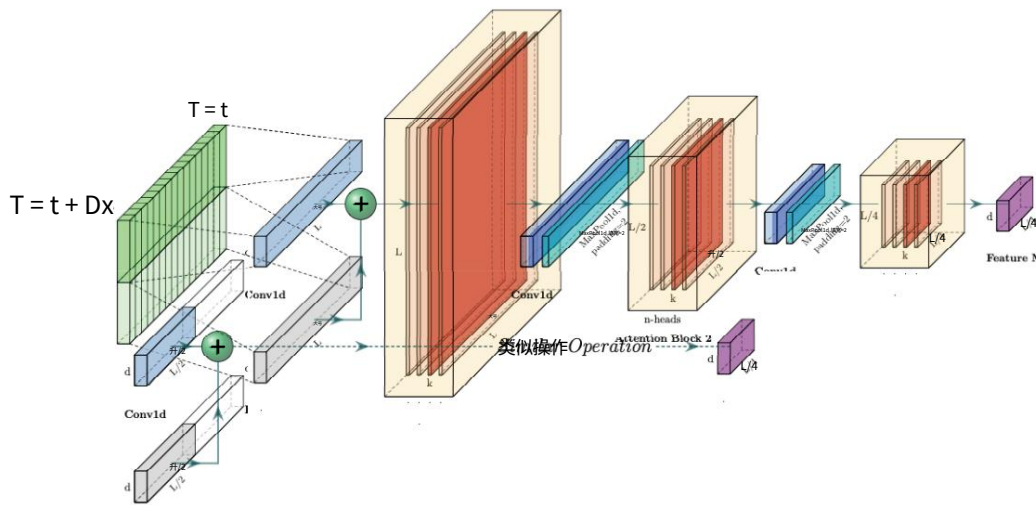


图 3:Informer 编码器的架构。(1) 每个水平堆栈代表一个单独的编码器副本在图 (2)中。(2) 上层栈是主栈,接收整个输入序列,而第二层栈占据一半输入的切片。(3) 红色层是self-attention机制的点积矩阵,级联减少在每一层上应用自注意力蒸馏。(4) 将两个堆栈的特征图连接起来作为编码器的输出。

命题 1.假设 $k_j \in N(\mu, \Sigma)$ 并且我们设 q_{kj} 表示集合 $\{(q_{ij} | j = 1, \dots, LK)\}$, 那么 $\forall M_m = \max_i M(q_i, K)$ 存在 $\kappa > 0$ 使得: 区间 $\forall q_1, q_2 \in \{q | M(q, K) \in [M_m, M_m - \kappa]\}$, 如果 $M(q_1, K) > M(q_2, K)$ 且 $\text{Var}(q_{k1}) > \text{Var}(q_{k2})$, 我们很有可能 $M(q_1, K) > M(q_2, K)$ 。到简化, 概率的估计在证明。

编码器:允许处理更长的顺序 内存使用限制下的输入

编码器旨在提取长序列输入的稳健远程依赖。在输入表示之后,第 t 个序列输入 X 已经被整形为

矩阵 $X_t \in \mathbb{R}^{L \times d_{\text{model}}}$ 。我们给出了 en 的草图,为清楚起见,图 (3)中的编码器。

Self-attention Distilling作为自然的结果在 ProbSparse Self-attention 机制中,编码器的特征图具有值 V 的冗余组合。我们使用蒸馏操作以赋予上级特权支配特征并在下一层制作集中的自注意力特征图。它修剪输入的时间维度急剧地,看到 n 头权重矩阵 (重叠红色图 (3)中的注意块的正方形)。受扩张卷积的启发 (Yu, Koltun 和 Funkhouser 2017; Gupta 和 Rush 2017), 我们的“蒸馏”程序从第 j 层变为第 $(j+1)$ 层为

$$X_{t,j+1} = \text{MaxPool ELU}(\text{Conv1d}([X_{t,j} \text{AB}])) \quad (5)$$

其中 $[\cdot] \text{AB}$ 包含 Multi-head ProbSparse 自注意力和注意力块中的基本操作, $\text{Conv1d}(\cdot)$ 使用 $\text{ELU}(\cdot)$ 激活函数 (Clevert, Unterthiner, and Hochreiter 2016) 在时间维度上执行一维卷积滤波器 (内核宽度 = 3)。

我们添加一个步幅为 2 和下采样的最大池化层 X_t 在堆叠一层后进入它的半片,这减少了整个内存使用量为 $O((2 - \epsilon) \log L)$, 其中 ϵ 是一个很小的数字。为了增强蒸馏操作的鲁棒性,我们构建了主堆栈的减半副本

并通过一次丢弃一层来逐渐减少自注意力蒸馏层的数量,就像金字塔一样

在图 (3)中,它们的输出维度是对齐的。因此,我们连接所有堆栈的输出并获得编码器的最终隐藏表示。

解码器:生成长序列输出 通过一个前向程序

我们使用标准解码器结构 (Vaswani et al. 2017) 图 (2), 它由两个相同的多头注意力层组成。然而,生成推理是

用于缓解长期预测中的速度骤降。

我们为解码器提供以下向量:

$$X_{t_0} = \text{连续}(X_t, X_{t_0}) \in \mathbb{R}^{(L_{\text{token}} + L_y) \times d_{\text{model}}}, \quad (6)$$

其中 $X_t \in \mathbb{R}^{L_{\text{token}} \times d_{\text{model}}}$ 为起始令牌, $X_{t_0} \in \mathbb{R}^{L_y \times d_{\text{model}}}$ 是目标序列的占位符 (设置标量为 0)。Masked multi-head attention 应用于通过将掩码点积设置为 $-\infty$ 来进行 ProbSparse 自注意计算。它阻止每个职位参加到即将到来的位置,这避免了自回归。一个完全连接层获取最终输出,其超大 d_y 取决于我们执行的是单变量预测还是多变量预测。

Generative Inference Start token 是 NLP 的“动态解码” (Devlin et al. 2018) 中的一种有效技术, 我们将其扩展为一种生成方式。我们不是选择一个特定的标志作为令牌,而是在输入序列中采样一个 L_{token} 长序列,它是之前的一个更早的切片

输出序列。以预测 168 个点为例（图 2(b)）中的 7 天温度预测,我们将

将目标序列前 5 天的已知时间作为“开始标记”,并为生成式推理解码器提供
Xfeed_{de} = {X5d, X0}。X0 包含目标序列的时间戳,即目标周的上下文。请注意,我们的建议的解码器通过一个前向预测所有输出程序并且免于耗时的“动态解码”事务在普通的编码器-解码器架构中。详细的性能比较在

计算效率部分。
损失函数我们选择 MSE 损失函数对目标序列进行预测,损失被传播

从整个模型的解码器输出返回。

实验

数据集

我们凭经验对四个数据集进行了实验,包括 2 个收集的 LSTF 真实世界数据集和 2 个公共基准数据集。

ETT（电力变压器温度）²: ETT 为电力长期部署的重要指标。我们从两个独立的县收集了 2 年的数据

在中国。为了探索 LSTF 问题的粒度,我们为 1 小时级别创建单独的数据集,如 {ETTh1, ETTh2}和15 分钟级别的 ETTm1。每个数据点由目标值“油温”和6个功率负载组成

特征。train/val/test 是 12/4/4 个月。
ECL（Electricity Consuming Load）³:收集321个客户的用电量（Kwh）。由于失踪数据（Li et al. 2019）,我们将数据集转换为 2 年的每小时消耗量,并将“MT 320”设置为目标值。train/val/test 是 15/3/4 个月。

天气⁴:该数据集包含当地气候数据从 2010 年到 2013 年的 4 年,近 1,600 个美国地点,每 1 小时收集一次数据点。每个数据点由目标值“湿球”和 11 个气候特征组成。训练/验证/测试是 28/10/10 个月。

实验细节

我们简要总结了基础知识,有关网络组件和设置的更多信息在附录 E 中给出。

基线:给出了网络组件的详细信息在附录 E.1 中。我们选择了 5 种时间序列预测方法作为比较,包括 ARIMA (Ariyo, Adewumi 和 Ayo 2014),先知 (Taylor 和 Letham 2018)、LSTMa (Bahdanau,Cho 和 Bengio 2015)和 LSTnet (Lai et al. 2018) 和 DeepAR (Flunkert, Salinas, 和 Gasthaus 2017)。为了更好地探索 ProbSparse 自注意力在我们提出的 Informer 中的性能,我们采用了规范的自注意力变体 (Informer†),

高效变体改革者 (Kitaev,Kaiser 和 Levskaya

²我们收集了 ETT 数据集并将其发布在 <https://github.com/zhouhaoyi/ETDataset>。
³ECL 数据集在 <https://archive.ics.uci.edu/ml/> 获得 datasets/ElectricityLoadDiagrams20112014 天
⁴气数据集在 <https://www.ncdc.noaa.gov/> 获得 订单/qclcd/

2019)和最相关的工作LogSparse self-attention (Li et al. 2019) 在实验中。
超参数调优:我们进行网络搜索超参数和详细范围在附录 E.3 中给出。Informer 包含一个 3 层堆栈和一个 2 层编码器中的堆栈（1/4 输入）,2 层解码器。我们提出的方法使用 Adam 优化器进行优化,其学习率从1e-4 开始,每下降 10 倍

2 epochs,总 epochs 为 10。我们按照推荐设置比较方法,批量大小为 32。每个数据集的输入都是零均值归一化的。在下面 LSTF 设置,我们延长预测窗口大小L_y 渐进式,即 {ETHh, 中的 {1d, 2d, 7d, 14d, 30d, 40d} ECL, Weather}, {6h, 12h, 24h, 72h, 168h} 在 ETTm 中。指标:我们使用了两个评估指标,包括 MSE = $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ 和 MAE = $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$ 在每一个上预测窗口（多元预测的平均）,并以 stride = 1 滚动整个集合。平台:全部模型在单个 Nvidia V100 32GB 上进行训练/测试显卡。源代码位于 <https://github.com/周浩毅/Informer2020>。

结果与分析

表 1 和表 2 总结了单变量/多变量所有方法在 4 个数据集上的评估结果。我们逐渐延长预测范围作为更高的要求的预测能力。为了公平比较,我们有精确控制问题设置以使 LSTF 为每种方法都可以在一个 GPU 上处理。最佳结果以粗体突出显示。

单变量时间序列预测在此设置下,每种方法都在单个变量中获得预测时间。从表 1 中,我们观察到: (1)提出的 model Informer 大大提高了推理性能（最后一列中的获胜次数）在所有数据集中,以及他们的预测误差在不断扩大预测范围,这证明了成功 Informer 在提高 LSTF 中的预测能力问题。 (2) Informer 击败了它的规范降级 Informer†主要在获胜计数中,即 28>14,它支持查询稀疏假设,以提供可比较的注意力特征图。我们提出的方法也执行了最相关的工作 LogTrans 和 Reformer。我们

请注意,Reformer 保持动态解码并且在 LSTF 中表现不佳,而其他方法则受益于生成风格解码器作为非自回归预测器。 (3) Informer 模型显示出明显优于循环神经网络 LSTMa 的结果。我们的方法

MSE 下降了 41.5%（在 168）,60.7%（在 336）和 60.7%（在 720 处）。这揭示了一个较短的网络路径自注意力机制比基于 RNN 的模型获得更好的预测能力。 (4)我们提出的方法取得比 DeepAR,ARIMA 和 Prophet 更好的结果在 MSE 上降低 20.9%（在 168）,61.2%（在 336）和平均为 51.3%（在 720 处）。在 ECL 数据集上,DeepAR 在较短的视野（≤ 336）上表现更好,我们的方法在更长的视野中超越。我们将此归因于特定的例如,其中预测能力的有效性是体现为可扩展性问题。

自注意力机制。我们设置样本因子 $c = 5$ (红线)在实践中。堆叠层数: Layers的replica对self-attention是互补的蒸馏,我们研究图 (4 (c))中每个堆栈 $\{L, L/2, L/4\}$ 的行为。越长的堆栈越敏感投入,部分原因是接收到更多的长期信息。我们方法的选择 (红线),即结合 L 和 $L/4$,是最稳健的策略。

表 3:ProbSparse 机制的消融

预测长度	336			720		
编码器的输入	336	720	1440	720	1440	2880
告密者	MSE 0.243	MAE 0.225	0.212	0.258	0.238	0.224
	0.487	MSE 0.214	0.404	0.503	0.399	0.387
告密者†	MAE 0.369	MSE 0.205	-	0.235	-	-
	0.256	MAE 0.496	0.364	0.401	-	-
日志传输	MSE 1.848	MAE 0.233	-	0.264	-	-
	1.054	0.412	-	0.523	-	-
改革者	-	1.832	1.817	2.094	2.055	2.032
	-	1.027	1.010	1.363	1.306	1.334

1 Informer †使用规范的自注意力机制。 2 “-”表示内存不足失败。

表 4:自注意力蒸馏的消融

预测长度	336			480		
编码器输入	336	480	720	336	480	720
告密者†	MSE 0.101	0.175	0.215	0.185	0.172	0.136
	MAE 0.360	0.335	0.366	0.355	0.321	0.282
告密者‡	MSE 0.187	0.182	0.177	0.208	0.182	0.168
	0.304	-	-	-	-	-

1 Informer ‡消除了 Informer † 中的自注意力提炼。 2 “-”表示内存不足失败。

表 5:生成式解码器的消融

预测长度	336			480		
预测偏移	+0	+12	+24	+48	+0	+48
告密者‡	MSE 0.101	0.102	0.103	0.103	0.155	0.158
	MAE 0.215	0.218	0.223	0.227	0.317	0.397
告密者§	MSE 0.152	MAE 0.294	-	-	-	-
	-	-	-	-	-	-

1 Informer §将我们的解码器替换为 Informer ‡ 中的动态解码器。 2 “-”表示不可接受的度量结果失败。

消融研究:Informer 如何工作?

我们还对 ETHh1 进行了额外的实验
消融考虑。

ProbSparse 自注意力机制的性能在整体结果表 1 和表 2 中,我们限制了

问题设置以使内存使用对
规范的自注意力。在本研究中,我们将我们的方法与 LogTrans 和 Reformer 进行比较,并深入探索
他们的极端表现。高效隔离内存
问题,我们首先将设置减少为{batch size=8,heads=8, dim=64},并在单变量情况下保持其他设置。
在表 3 中,ProbSparse self-attention 表现出比同类产品更好的性能。 LogTrans 进入 OOM
公开实施的极端案例是
全注意力,仍然有 O(L2)内存使用。我们的
提出的 ProbSparse self-attention 避免了这种情况,因为 Eq.(4) 中的
查询稀疏假设带来了简单性,
参考附录 E.2 中的伪代码,并达到
较小的内存使用量。

self-attention distilling 的表现在此
研究中,我们使用 Informer †作为基准来消除
ProbSparse 自注意力的附加效果。另一个
实验设置与单变量时间序列的设置一致。从表 4 中可以看出,Informer †完成了
所有实验,并在利用序列输入后取得了更好的性能。比较方法 Informer ‡
去除了蒸馏操作和

以更长的输入 (> 720) 达到 OOM。关于
LSTF 问题中长序列输入的好处,我们

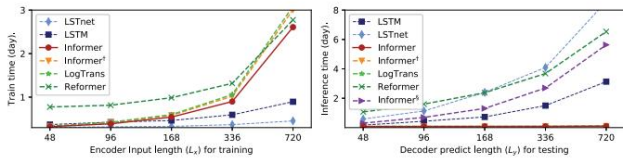


图 5:训练/测试阶段的总运行时间。

表 6:每层的 L 相关计算静态

方法	训练		测试
	时间复杂度	内存使用情况	脚步
告密者	$O(L \log L)$	$O(L \log L)$	1
变压器	$O(L^2)$	$O(L^2)$	大
日志传输	$O(L \log L)$	$O(L^2)$	1
改革者	$O(L \log L)$	$O(L \log L)$	大
长短期记忆	$O(L)$	$O(L)$	大

1 LSTnet 很难有一个封闭的形式。

得出结论,self-attention distilling 值得采用,
特别是当需要更长的预测时。

生成式解码器的性能
研究中,我们证明了我们的解码器在获得 “生成”结果方面的潜
在价值。与现有方法不同的是,
标签和输出在训练中被迫对齐,并且
推断,我们提出的解码器的预测完全依赖于
时间戳,可以用偏移量进行预测。从表 5 可以看出,一般预测性能

Informer ‡的抵制随着偏移量的增加而增加,而
对应的动态解码失败。它证明了解码器捕获个体远程依赖的能力

任意输出之间,避免错误累积
推论。

计算效率

使用多变量设置和每种方法的当前
最好的实现,我们在图 (5)中进行了严格的运行时比较。在训练阶段,Informer

(红线)达到最好的训练效率
基于变压器的方法。在测试阶段,我们的
方法比其他具有生成风格的方法快得多
解码。理论时间复杂度的比较
和内存使用情况总结在表 6 中,Informer 的性能与运行时实验一致。

请注意,LogTrans 专注于自注意力机制,我们在 LogTrans 中应用我们提出的
解码器
公平比较 (表 6 中的)。

结论

在本文中,我们研究了长序列时间序列预测问题,并提出了 Informer 来预测长
序列。具体来说,我们设计了 ProbSparse 自注意力机制和蒸馏操作来处理

vanilla Transformer 中二次时间复杂度和二次内存使用的挑战。此外,精心设
计的生成解码器缓解了传统编码器-解码器架构的局限性。上的实验

真实世界的的数据证明了 Informer 的有效性
用于提高 LSTF 问题的预测能力。

致谢

这项工作得到了自然科学基金的支持
中国基金会 (U20B2053, 61872022 和 61421003)
软件开发环境国家重点实验室 (SKLSDE-2020ZX-12) 。感谢北京大数据与
脑计算先进创新中心提供的计算基础设施。相应的

作者是李建新。

道德声明

提议的 Informer 可以处理长输入并生成
有效的长序列推理,可应用于
解决具有挑战性的长序列时间序列预测 (LSTF) 问题。一些重要的实际应用是
传感器网络监控 (Papadimitriou 和 Yu

2006)、能源和智能电网管理、疾病传播分析 (Matsubara et al. 2014)、经济和
金融预测 (Zhu and Shasha 2002)、农业生态系统的演变、气候变化预测和空气
气变化

污染。作为一个具体的例子,在线卖家可以预测
每月产品供应,有助于优化长期库存管理。我们的贡献不仅限于 LSTF 问题。除
了获取长序列外,我们的方法还可以为其他人带来实质性的好处

领域,例如文本的长序列生成,音乐,
图片和视频。
在伦理考虑下,任何从历史数据中学习的时间序列预测应用程序都会运行

产生有偏见的预测的风险。它可能给财产/资产的真正所有者造成无法弥补的
损失。领域
专家应该指导我们方法的使用,他们
从长序列预测中具有潜在的好处。
以应用我们的电力变压器温度预测方法为例,经理将检查

结果并决定未来的电力部署。如果长
有足够的预测,将有助于管理者在早期防止不可逆转的失败。除了识别偏差数
据之外,一种有前途的方法是

采用迁移学习。我们捐赠了收集到的
数据 (ETT 数据集) ,用于进一步研究相关主题
(如供水管理、5G网络部署) 。
另一个缺点是我们的方法需要高性能
GPU,这限制了它在欠发达地区的应用。

参考

阿里约,AA;阿德乌米,AO;和 Ayo, CK 2014。股票
使用 ARIMA 模型进行价格预测。在第 16 届计算机建模与仿真国际会议上,
106-112。 IEEE。

巴赫达瑙,D,赵,K,和 Bengio, Y. 2015。联合学习对齐和翻译的神经机器翻译。
在 ICLR 2015 中。

; 贝尔塔吉,我,彼得斯,我,和 Cohan, A. 2020。Longformer:
长文档转换器。 CoRR abs/2004.05150。
通用电气公司;詹金斯,总经理; Reinsel,GC;和 Ljung,总经理

2015。时间序列分析:预测与控制。约翰
威利父子。

布朗,结核病;曼恩,B。 ;莱德,N,苏比亚,M,卡普兰,J。
达里瓦尔,P,尼拉坎坦,A;希亚姆,P,萨斯特里,G,阿斯科尔,
一种。 ;阿加瓦尔,S,赫伯特-沃斯,A;克鲁格,G,黑根汉,
T;孩子,R,拉梅什,A,齐格勒,DM;吴,J,冬天,
C。 ;黑森,C,陈,M,西格勒,E,利特温,M,格雷,S。
国际象棋,B,克拉克,J,伯纳,C,麦坎德利什,S,拉德福,
一种。 ;苏茨克韦尔,我,和 Amodei, D. 2020。语言模型
是少数人的学习者。 CoRR abs/2005.14165。

孩子,R,格雷,S,拉德福德,A,和 Sutskever, I. 2019。
使用稀疏变换器生成成长序列。
arXiv:1904.10509。

赵,K,范梅林博尔,B;巴赫达瑙,D,和本吉奥,
Y. 2014。关于神经机器翻译的属性:编码器-解码器方法。在诉讼中

SSST@EMNLP 2014,103-111。

聪明的,D。 ;Unterthiner,T,和 Hochreiter, S. 2016。快速
通过指数线性单元 (ELU) 进行准确的深度网络学习。在 ICLR 2016 中。

戴,Z,杨,Z,杨,Y,卡内尔,J,乐,QV;和
Salakhutdinov, R. 2019。Transformer-xl:专注的语言
超出固定长度上下文的模型。 arXiv:1901.02860。
德夫林,J,张,M,-W,韭葱。 ;和 Toutanova, K. 2018。
Bert:用于语言理解的深度双向转换器的预训练。 arXiv:1810.04805。

弗朗克特,五,萨利纳斯,D,和 Gasthaus, J. 2017。DeepAR:
使用自回归循环网络进行概率预测。 arXiv:1704.04110。

古普塔,A,和 Rush,AM 2017。用于模拟长距离基因组依赖关系的扩张卷积。
arXiv:1710.01278。

霍克赖特,S,和 Schmidhuber, J. 1997。长期短期
记忆。神经计算 9 (8) :1735-1780。
基塔耶夫,N,凯撒,L,和 Levskaya, A. 2019。改革者:
高效变压器。在 ICLR 中。

赖,G,常,W,-C,杨,Y,和 Liu, H. 2018。使用深度神经网络建模长期和短期时间
模式
网络。在 ACM SIGIR 2018 年,95-104。 ACM。

李,S,金,X。 ; 轩,Y,周,X,陈,W,王,Y,-X。
and Yan, X. 2019。增强地方性,打破传统
时间序列预测中 Transformer 的内存瓶颈。 arXiv:1907.00235。

李,Y,于,R,沙哈比,C,和 Liu, Y. 2018。扩散卷积递归神经网络:数据驱动流量
预测。在 ICLR 2018 中。

刘,Y,龚,C,杨,L,和 Chen, Y. 2019。DSTP-RNN:
基于双阶段两阶段注意力的循环神经网络适用于长期和多变量时间序列预测。
CoRR abs/1904.07464。

松原,Y,樱井,Y。 van Panhuis, WG;和 Falout sos, C. 2014 年。漏斗:自动挖
掘空间 coe 流行病。在 ACM SIGKDD 2014,105-114。

帕帕迪米特里乌,S.;和 Yu, P. 2006。时间序列流中的最佳多尺度模式。在 ACM SIGMOD 2006, 647–658. ACM。

秦,Y。宋,D。陈,H。程,W。江,G。和 Cot trell, GW 2017。用于时间序列预测的双阶段基于注意力的递归神经网络。在 IJCAI 2017,2627–2633。

邱,J。马,H。;利维,O。易,SW-t。王,S。和 Tang, J. 2019。用于长文档理解的 Blockwise Self-Attention。 arXiv:1911.02972。

雷,JW;波塔彭科,A.;贾古玛,SM;和 Lillicrap,TP 2019。用于远程序列建模的压缩变压器。 arXiv:1911.05507 。

Ray, W. 1990。时间序列 :理论和方法。皇家统计学会杂志 :A 系列 (社会统计学)153 (3) :400-400。

西格,M。兰加普拉姆,S.;王,Y。萨利纳斯,D。盖斯豪斯,J。亚努肖夫斯基,T。和 Flunkert, V. 2017。
用于大规模间歇性需求预测的线性状态空间模型中的近似贝叶斯推理。
arXiv:1709.07638 。

西格,兆瓦;萨利纳斯,D。和 Flunkert, V. 2016 年。大库存的贝叶斯间歇性需求预测。在 NIPS 中,4646–4654。

苏茨克韦尔,我。;乙烯基,O。和 Le,QV 2014。使用神经网络进行序列到序列学习。在 NIPS,3104-3112。

泰勒,SJ;和 Letham, B. 2018。大规模预测。美国统计学家 72 (1) :37-45。

蔡,Y.-HH;白,S。山田,M。莫伦西,L.-P.;和 Salakhutdinov, R. 2019。变压器解剖 :通过内核镜头对变压器注意力的统一理解。在 ACL 2019 中,4335–4344。

瓦斯瓦尼,A。沙泽尔,N。帕尔马,N。乌兹科雷特,J。琼斯,L。戈麦斯,AN;凯撒,。和 Polosukhin, I. 2017。注意力就是你所需要的。在 NIPS,5998-6008。

王,S。李,B。哈布萨,M。方,H。和马,H。 2020. Linformer:具有线性复杂性的自我注意。 arXiv:2006.04768。

温,R。托尔科拉,K。纳拉亚纳斯瓦米,B。和 Madeka, D. 2017 年。多水平分位数循环预报器。 arXiv:1711.11053。

于,F。科尔通,V。和 Funkhouser, T. 2017。扩张残差网络。在 CVPR 中,472–480。

于,R。郑,S。阿南德库马尔,A。和 Yue, Y. 2017。使用张量训练 rnns 进行长期预测。 arXiv:1711.00073
.

朱,Y。和 Shasha, DE 2002。StatStream:实时统计监控数以千计的数据流。在 VLDB 2002 中,358–369。