
Research of Efficient Transformers

Vladimir Baikalov¹ Kovaleva Maria¹ Konstantin Shlychkov¹ Vo Ngoc Bich Uyen¹

Abstract

The attention-based methods and transformers made a significant breakthrough in the deep learning area and greatly impacted NLP task solutions. Although recent works show that they could potentially improve results in different task domains, the application of transformers for financial data in particular transaction data is unexplored. While applying attention mechanisms, one can face the apparent restriction on input sequence length due to the method's quadratic complexity. Recent papers proposed different ways to overcome this problem, but we want to concentrate on two promising approaches: Informer and Performer. The Informer is the most current and prospective approach. Its main assumption is that the model should have an "infinite memory" and fit a sequence with arbitrary length. The Performer model shows good results in the NLP task but is not well-explored for other datatypes. Its main idea is to use some trigonometric approximation of the attention matrix to decrease memory consumption. To sum up, the project aims to compare several recent methods proposed to decrease the evaluation complexity in particular tasks predicting the user's gender based on transactions.

1. Introduction

The most popular sequence transduction models are built on complicated recurrent or convolutional neural networks with an encoder and a decoder. Transformer is proposed as a new fundamental network design based solely on attention techniques for drawing global dependencies between input and output. Transformers are strong neural network designs that have emerged as the standard in many fields of machine learning. These models surpass others in terms of quality while being more parallelizable and taking substantially less time to train. Transformer-based language models have

achieved impressive results by increasing the context size. However, whereas humans process information sequentially, continually updating their memories, and recurrent neural networks (RNNs) update a single memory vector across time, transformers do not - they comprehensively query every presentation connected with prior events. The regular Transformer scales quadratically with the number of tokens L in the input sequence, which makes huge datasets prohibitively expensive. Most approaches either limit the attention mechanism to local neighborhoods or add structural prior on attention. As a result, the amount of work they must perform rises in proportion to the length of the context. In this study, we present two possible solutions to this challenge: Informer and Performer.

First, we introduce the Informer or infinite-former - a transformer that extends the vanilla transformer with an unbounded long-term memory (LTM). By making use of a continuous-space attention framework that trades off the amount of information units that fit into memory (basis functions) with the granularity of their representations. This representation has two significant advantages: (i) the context can be represented using a smaller number of basis functions N than the number of tokens, lowering the attention computational cost; and (ii) N can be fixed, allowing unbounded context to be represented in memory without increasing its attention complexity. When expressing the input sequence as a continuous signal, however, choosing a lesser number of basis functions results in reduced accuracy. To address the issue of resolution loss, we offer the notion of "sticky memories," which is a process that ensures the preservation of critical information in the LTM. In "sticky memories," we assign a bigger space in the LTM signal to sections of memory that are regularly examined, allowing the model to collect long contexts without losing significant information.

Second, we present the Performer, the first Transformer designs capable of provably accurate and practical estimation of regular (softmax) full-rank attention. Other strategies that try to reduce the space complexity of Transformers include reversible residual layers that allow for one-time activation storage in training and shared attention weights. Performer is the first linear structures that are fully compatible with conventional Transformers

¹Skoltech University, Moscow, Russia. Correspondence to: Vladimir Baikalov <vladimir.baikalov@skoltech.ru>.

(with minor fine-tuning). Performer uses the Fast Attention Via positive Orthogonal Random features (FAVOR+) mechanism, leveraging new methods for approximating softmax and Gaussian kernels. FAVOR+ can also be used outside of the Transformer scope as a more scalable substitute for frequent attention.

Finally, in a specific task that predicts the gender of a customer using information about receipts and expenditures on a bank card, we execute three comparisons: baseline model, Performer model, and Informer model. The correctness of all models is then computed in terms of speed and memory usage.

2. Background

2.1. Vanilla Transformer

The model [1] follows the standard encoder-decoder architecture, but it omits the recurrent module and relies only on attention mechanism to compute representations of its input and output. The encoder is composed of several layers of multi-head self-attention mechanism followed by a feed-forward layer, along with residual connections and layer normalization. The decoder has similar structure but additionally uses the attention between the output from the previous decoder layer and the encoder output. To prevent the leftward information in decoder self-attention, the masking technique is used to avoid the illegal connections.

The attention mechanism uses queries (Q), keys (K) and values (V) as its input and aims to estimate the contribution of each query for all values via measuring the similarity between keys and queries. The resemblance is measured with a standard scaled dot-product with further application of softmax function to obtain the weights on the values. We can describe this operation as follows:

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

where d is the dimension of hidden representations for queries and keys. Its usage may be illustrated with the following example: suppose ξ and η are independent random vectors, whose components have mean 0 and variance 1. Then their dot product (ξ, η) has mean 0 and variance d , which means the values may become large and hence result in computational instability. To avoid this we need to scale the dot product by \sqrt{d} .

To get more stable and generalized hidden representations the attention mechanism is separated in h groups. Queries, keys and values are linear projected into smaller subspaces with dimensions d_k , d_k and d_v respectively and then the attention mechanism is applied to each group.

Obtained values are concatenated and projected again to result in final result. This procedure may be represented as follows:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(h_1, h_2, \dots, h_z) W^l$$

where $h_i = \text{Attention}(QW_i^q, KW_i^k, VW_i^v)$ and $W_i^q, W_i^k \in \mathbb{R}^{d \times d_k}$, $W_i^v \in \mathbb{R}^{d \times d_v}$, $W^l \in \mathbb{R}^{d_v \times d}$. This allows not only jointly attend to information from different representation subspaces, but also to speed up the computation with parallelization.

However, there is enormous limitation in computational complexity as attention mechanism requires all pairwise dot products between two inputs and the model works extremely slow when handles with long sequences. Many approaches were proposed to overcome this issue, but we will investigate only two of them — performer and ∞ -former.

2.2. ∞ -former

Transformers are unable to model long-term memories effectively, since the amount of computation they need to perform grows with the context length. All efficient transformers have a finite memory capacity and are forced to drop old information. The ∞ -former [2] is a transformer, which extends the vanilla transformer with an unbounded long term memory. By making use of a continuous space attention mechanism to attend over the long-term memory, the ∞ -former's attention complexity becomes independent of the context length, trading of memory length with precision. In order to control where precision is more important, ∞ -former maintains "sticky memories," being able to model arbitrarily long contexts while keeping the computation budget fixed.

The main idea of this model is usage of continuous attention using N radial basis functions (RBF).

Lets consider details of using continuous attention. $X = [x_1, \dots, x_L]$ - input vector. Every x_i is associated with $t_i \in [0, 1]$, so the input vector can be represented as $\bar{X}(t) = B^T \psi(t)$, where $\psi(t)$ is a vector of N RBFs: $\psi_j(t) = N(t; \mu_j, \sigma_j)$ and B is coefficient matrix, which can be found by ridge regression: $B^T = X^T F^T (F F^T + \lambda I)^{-1} = X^T G$, where $F = [\psi(t_1), \dots, \psi(t_L)]$. As out input vector X is replaced by $\bar{X}(t) = B^T \psi(t)$ so the keys and values will be computed as $K_h = B_h W^{K_h}$ and $V_h = B_h W^{V_h}$ (where different h associated with different heads). And for every query $q_{h,i}$ we need to compute the parameters (which are given by neural network) of the normal distribution from which

we will sample our long-term memory (LTM) context vectors: $\mu_{h,i} = \text{sigmoid}\left(\text{affine}\left(\frac{K_h q_{h,i}}{\sqrt{d}}\right)\right)$, $\sigma_{h,i} = \text{softplus}\left(\text{affine}\left(\frac{K_h q_{h,i}}{\sqrt{d}}\right)\right)$, $p_{h,i} \sim N\left(t; \mu_{h,i}, \sigma_{h,i}^2\right)$. And having value function $\bar{V}_h(t) = V_h^T \psi(t)$ we will have $Z_{h,i} = \mathbb{E}_{p_{h,i}}[\bar{V}_h] = V_h \mathbb{E}_{p_{h,i}}[\bar{\psi}(t)]$ which forms the rows of matrix $Z_{\text{LTM},h}$. $Z_{\text{LTM}} = [Z_{\text{LTM},1}, \dots, Z_{\text{LTM},H}] W^o$ is LTM context representation. So the final context representations in this framework will be $Z = Z_T + Z_{\text{LTM}}$, where Z_T is e transformer context vector.

The key matrix K_h size depends only on the number of basis functions N , but not on the length of the context being attended to. Thus, the ∞ -former's attention complexity is also independent of the context's length. It corresponds to $O(L^2 + LN)$, while complexity of a vanilla transformer attending to the same context is $O(L(L + L_{\text{LTM}}))$.

When representing the memory as a discrete sequence, to extend it, we need to store the new hidden states in memory. In a vanilla transformer, this is not feasible for long contests due to the high memory requirements. However, the ∞ -former can attend to unbounded context without increasing memory requirements by using continuous attention. To be able to build an unbounded representation, we first sample M locations in $[0, \tau]$, where $\tau < 1$ and evaluate $\bar{X}(t)$ at those locations. Then, we concatenate the corresponding vectors with the new vectors coming from the short-term memory $X = [X_{\text{past}}^T, X_{\text{new}}^T]^T$. Finally, we simply need to perform multivariate ridge regression to compute the new coefficient matrix B , via $B^T = X^T G$, where $G \in \mathbb{R}^{(M+L)N}$.

Also sticky memories concept was proposed by authors of ∞ -former. The main idea of this concept is not to use linearly spaced M locations, but sample them according to histogram based on the attention given to each interval of the signal on the previous step.

2.3. Performer

Computation of attention matrix in vanilla transformer architecture requires quadratic space and time complexity with number of tokens / sequence length and it does not rely on any priors such as data sparsity or specific data structure. Paper *Rethinking attention with performers* [3] proposes a proven way for accurate estimation of regular (softmax) full-rank attention, but with linear space and time complexity. Performers uses above mentioned mechanism called FAVOR+. More precisely, authors of this paper show that it is possible to approximate attention matrix A having $O(Ld^2 \log(d))$ time complexity.

The main idea with is the basis of this paper states that it is possible to reconstruct in practice most kernels ϕ in the following way:

$$\phi(x) = \frac{h(x)}{\sqrt{m}} (f_1(\omega_1^\top x), \dots, f_l(\omega_l^\top x), \dots, f_l(\omega_m^\top x))$$

where functions $f_1, \dots, f_l : \mathbb{R} \rightarrow \mathbb{R}$, function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ and every $\omega_i \stackrel{\text{iid}}{\sim} \mathcal{P}(\mathbb{R}^d)$ is a deterministic vector. For example with $h(x) = 1$, $l = 2$, $f_1 = \sin$, $f_2 = \cos$ corresponds to shift-invariant kernels, and with $\mathcal{P}(\mathbb{R}^d) = \mathcal{N}(0, \mathbf{I}_d)$ we can get Gaussian kernel.

The key *softmax-kernel* itself may be written in the following form:

$$SM(x, y) \stackrel{\text{def}}{=} \exp(x^\top y) = \exp\left(\frac{\|x\|^2}{2}\right) K_{\text{gauss}}(x, y) \exp\left(\frac{\|y\|^2}{2}\right)$$

This equation can be estimated by using $h(x) = \exp\left(\frac{\|x\|^2}{2}\right)$, $l = 2$, $f_1 = \sin$, $f_2 = \cos$ and it can be called $\hat{SM}_m^{\text{trig}}(x, y)$. But since \sin and \cos may get negative values it may lead to high variance and abnormal behaviors.

To tackle this issue author propose two new estimators:

$$h(x) = \exp\left(-\frac{\|x\|^2}{2}\right)$$

with

$$l = 1, f_1(u) = \exp(u)$$

$$h(x) = \frac{1}{\sqrt{2}} \exp\left(-\frac{\|x\|^2}{2}\right)$$

with

$$l = 2, f_1(u) = \exp(u), f_2(u) = \exp(-u)$$

which is called \hat{SM}_m^+ and $\hat{SM}_m^{\text{hyp+}}$ respectively. For both estimators $\mathcal{D} = \mathcal{N}(0, \mathbf{I}_d)$.

Also, authors show that orthogonal choice of all random features ω will lead to better result.

To sum up, authors of *Rethinking attention with performers* [3] presented models called *Performer*, a new type of Transformers, relying on our Fast Attention Via positive Orthogonal Random features (FAVOR+) mechanism to significantly improve space and time complexity of regular Transformers. Their mechanism provides the first effective unbiased estimation of the original softmax-based Transformer with linear space and time complexity.

References

1. Attention is all you need / A. Vaswani [et al.] // Advances in neural information processing systems. — 2017. — Vol. 30.
2. *Martins P. H., Marinho Z., Martins A. F.* ∞ -former: Infinite Memory Transformer // arXiv preprint arXiv:2109.00301. — 2021.
3. Rethinking attention with performers / K. Choromanski [et al.] // arXiv preprint arXiv:2009.14794. — 2020.