Sequence Modeling

S.Calderara, R. Di Carlo

UNIMORE

December 3, 2020

S.Calderara, R. Di Carlo (UNIMORE)







2 CNN



4 Temporal Fusion Transformer

5 N-Beats

6 Unsupervised Anomaly Detection on Time Series

7 Credits



A sequence model is a model that takes in input a sequence of items (words, letters, time series, audio signals, etc) and outputs one item or another sequence of items.

$$\langle x_1 \rangle \langle x_2 \rangle \langle x_3 \rangle \longrightarrow \langle y_1 \rangle \langle y_2 \rangle \langle y_3 \rangle \langle y_4 \rangle$$

Je suis étudiant I am a student

Some of the tasks using seq-to-seq models are: Machine translation, Time-series forecasting, Video Captioning, Audio generation, Text Summarization and more.

Even if sequence modeling is very often associated with recurrent neural networks, yet recent results indicate that other neural network architectures can outperform RNN on many tasks.

CNN



The standard Convolution Neural Network was first introduced in Lenet-5 architecture. Conv2D is generally used on Image data. It is called 2 dimensional CNN because the kernel slides along 2 dimensions on the data.



In **Conv1D** the kernel slide only in one dimension. Which kind of data requires kernel sliding in only one dimension and have spatial properties?

1D Convolution





We can apply it to time series, as well as to text, since we can represent each word as a vector of fixed length.

The discrete 1D convolution is defined as:

$$(\mathbf{x} * \omega)(n) = \sum_{i=0}^{k-1} \mathbf{x}(n-i)\omega(i)$$

Where \boldsymbol{x} is our input vector of length \boldsymbol{n} , and $\boldsymbol{\omega}$ our kernel of length \boldsymbol{k} .

1D Convolution - example



Asimov	0.2	0.1	-0.3	0.4
scrisse	0.5	0.2	-0.3	-0.1
molti	-0.1	-0.3	-0.2	0.4
racconti	0.3	-0.3	0.1	0.1
sui	0.2	-0.3	0.4	0.2
robot	0.1	0.2	-0.1	-0.1
positronici	-0.4	-0.4	0.2	0.3

Apply a filter(or kernel) of size 3

Asimov scrisse molti	-1
scrisse molti racconti	-0.5
molti racconti sui	-3.6
racconti sui robot	-0.2
sui robot positronici	0.3

3	1	2	-3
-1	2	1	-3
1	1	-1	1

S.Calderara, R. Di Carlo (UNIMORE)

Sequence Modeling



A major disadvantage of using CNN for sequences is that in order to achieve a long effective history size, we need an extremely deep network or very large filter.



To mitigate this problem we can increase the receptive field using dilated convolution.

$$(\mathbf{x} * \omega)(n) = \sum_{i=0}^{k-1} \mathbf{x}(n-d \cdot i)\omega(i)$$

Where d is the dilation factor

Temporal Convolutional Network



1D-Conv layers are used in many models as a valid alternative to RNN in different tasks.

^(b) The **TCN** [1] is a model composed of residual blocks of dilated causal convolutions ^(a) Each layer use a different dilation factor d = 1,2,4 and filter size k = 3



Transformer



Sequence-to-sequence modelling using RNNs? No thanks, attention is all I need. [5]

Challenges with RNNs

- Long range dependencies
- Gradient vanishing and explosion
- Large # of training steps
- Recurrence prevents parallel computation

Transformer networks

- Facilitate long range dependencies
- No gradient vanishing and explosion
- Fewer training steps
- No recurrence, facilitate parallel computation

Transformer architecture





Figure 1: The Transformer - model architecture.

Transformer networks are composed by two main stacks:

- Encoder
- Decoder

Both are composed of a sub-stacks of *N* modules consisting mainly of **Multi-Head Attention** and **Feed Forward** layers.

Attention



Let's start with the description of the attention-mechanism. It's can be described by the following equation:

Scaled Dot-Product Attention



Multi-Head Attention



$$Attention(oldsymbol{Q},oldsymbol{K},oldsymbol{V})=Softmax\left(rac{oldsymbol{Q}oldsymbol{K}^{ op}}{\sqrt{oldsymbol{d}_k}}
ight)oldsymbol{V}$$

Where **Q**, **K**, **V** are vectors representing query, keys and values. d_k is just a scaling factor to prevent dot product growing large in magnitude, pushing the softmax into regions where it has small gradients.

Query, keys and values



Query, keys and values are generated using self-attention from each of the encoder's input vectors.



 $oldsymbol{Q} = oldsymbol{X}oldsymbol{W}^Q$ $oldsymbol{K} = oldsymbol{X}oldsymbol{W}^K$ $oldsymbol{V} = oldsymbol{X}oldsymbol{W}^V$

Query, keys, and values have dimension d_k . All sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{model} = 512$.

$$d_k = d_{model}/N_{heads}$$

Positional encoding



Since the Transformer contains no recurrence and no convolution, in order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens.

$$\overrightarrow{p_t}^{(i)} = f(t)^{(i)} := \left\{ egin{array}{c} \sin\left(\omega_k \cdot t
ight), & ext{if } i = 2k \ \cos\left(\omega_k \cdot t
ight), & ext{if } i = 2k+1 \end{array}
ight.$$
 where $\omega_k = rac{1}{10000^{2k/d}}$



Multi-head attention





A **multi-head attention** layer is composed of many attention layers, each of which has its own set of queries, keys and values obtained by randomly initializing the weights matrices.

Then, after training, each set is used to project the input embeddings into a different representation subspace.

 $MultiHead(\boldsymbol{Q},\boldsymbol{K},\boldsymbol{V}) = Concat(head_1,...,head_h)\boldsymbol{W}^O$

where:

*head*_{*i*} = Attention
$$\left(\boldsymbol{Q}\boldsymbol{W}_{i}^{Q}, \boldsymbol{K}\boldsymbol{W}_{i}^{K}, \boldsymbol{V}\boldsymbol{W}_{i}^{V}\right)$$

Decoding



0.11072	0.08901	0.10447	0.07113	0.14626
0.11533				
0.12546				
0.21091				
0.15397				

Attention matrix

0.11072		
0.11533		
0.12546		
0.21091		
0.15397		

Attention masked

1.00000		
0.48492		
0.32996		
0.26071		
0.19957		

Attention Softmax

During the training phase, the input of the decoder stack is the GT sequence shifted right to add the start token.

 $\langle y_1 \rangle \langle y_2 \rangle \langle y_3 \rangle \langle y_4 \rangle \langle y_5 \rangle$ SOS I am a student

Also we need to prevent leftward information flow in the decoder to preserve the auto-regressive property. We implement this inside of scaled dot-product attention by masking out (setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections.

Decoding





In "encoder-decoder attention" layers, the queries Q come from the previous decoder layer, and the memory keys K and values V come from the encoder output.



At inference time we do not have the GT sequence to provide as input to the decoder. So we use the transformer in auto-regressive mode.

Encoder input (x)	Decoder input (\hat{y}_{t-1})	Decoder Output (\hat{y}_t)
Je suis étudiant	<sos></sos>	1
Je suis étudiant	<sos>I</sos>	l am
Je suis étudiant	<sos>I am</sos>	l am a
Je suis étudiant	<sos>I am a</sos>	l am a student
Je suis étudiant	<sos>I am a student</sos>	I am a student <eos></eos>

Greedy Search





When we predict a sentence, the last Softmax layer gives us a probability distribution on the words in our dictionary. If we always choose the one with the highest probability, we are applying a greedy search.

A greedy search is not always the best choice, in fact if we assume the probability distribution of a word sequence as:

$$P(w_{1:T} | W_0) = \prod_{t=1}^{T} P(w_t | w_{1:t-1}, W_0)$$

 W_0 being the initial context word sequence

Greedy Search





We obtain:

P("The","nice","woman") = 0.2P("The","dog","has") = 0.36

Beam search is one of the techniques to mitigate this problem, reducing the risk of missing hidden high probability word sequences by keeping the most likely n_{beams} of hypotheses at each time step.

Temporal Fusion Transformer





Real world time series datasets consist of many type of input features:

- Past observed inputs
- Apriory known future inputs
- Static features

TFT (Temporal fusion transformer) [3] is a novel interpetrable forecasting deep neural network model from Google AI.

Architecture





TFT is based on sequential processing using both:

- Recurrent layers, to capture local dynamics
- Attention layers, to capture long-term context

GRN layer







GRN layer





Gated Residual Network (**GRN**) is a building block of TFT. The GRN takes in a primary input a and an optional context vector c

$$\begin{aligned} \mathsf{GRN}_{\omega}(\boldsymbol{a},\boldsymbol{c}) &= \mathsf{LayerNorm} \ \left(\boldsymbol{a} + \mathrm{GLU}_{\omega}\left(\eta_{1}\right)\right) \\ \eta_{1} &= \boldsymbol{W}_{1,\omega}\eta_{2} + \boldsymbol{b}_{1,\omega} \\ \eta_{2} &= \mathrm{ELU}\left(\boldsymbol{W}_{2,\omega}\boldsymbol{a} + \boldsymbol{W}_{3,\omega}\boldsymbol{c} + \boldsymbol{b}_{2,\omega}\right) \end{aligned}$$

where:

$$\mathsf{GLU}_{\omega}(\gamma) = \sigma \left(\mathbf{W}_{4,\omega} \gamma + \mathbf{b}_{4,\omega} \right) \odot \left(\mathbf{W}_{5,\omega} \gamma + \mathbf{b}_{5,\omega} \right)$$

ELU is the Exponential Linear Unit activation function. $\sigma(.)$ is the Sigmoid activation function.

Variable Selection Network







Variable Selection Network





Variable selection also allows TFT to remove any unnecessary noisy inputs which could negatively impact performance.

Being $\xi_t^{(j)}$ the embedding of the *j*-th variable at time *t* and $\Xi_t = \left[\xi_t^{(1)^T}, \ldots, \xi_t^{(m_\chi)^T}\right]^T$ a vector of past inputs at time *t*:

$$oldsymbol{v}_{\chi_t} = \mathsf{Softmax}\left(\mathsf{GRN}_{v_\chi}\left(\Xi_t, oldsymbol{c}_s
ight)
ight)$$

 c_s is a context vector and v_{χ_t} is a vector of variable selection weights.

Variable Selection Network





At each time step, an additional layer of non-linear processing is employed by feeding each $\xi_t^{(j)}$ through its own GRN:

$$\tilde{\xi}_{t}^{(j)} = \mathsf{GRN}_{\tilde{\xi}(j)}\left(\xi_{t}^{(j)}\right)$$

Where each GRN share weights across all time steps t. Finally input variables are weighted and combined:

$$\widetilde{\xi}_t = \sum_{j=1}^{m_\chi} \mathsf{v}_{\chi_t}^{(j)} \widetilde{\xi}_t^{(j)}$$

Sequence to sequence layer





The encoding and decoding part is done using LSTM layers, feeding $\xi_{t-k:t}$ into the encoder and $\xi_{t+1:t+\tau_{max}}$ into the decoder.

At inference time a vector of zeros is the input of the LSTM decoder.

Interpretable multi-head attention





This layer is the same multi-head attention proposed in Transformer architecture.

InterpretableMultiHead $(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \tilde{\boldsymbol{H}} \boldsymbol{W}_{H}$

But each head shares the same value weights and the heads are aggregated with the average instead of concatenation

$$ilde{H}=1/H\sum_{h=1}^{m_{H}}$$
 Attention $\left(oldsymbol{Q}oldsymbol{W}_{Q}^{(h)},oldsymbol{K}oldsymbol{W}_{K}^{(h)},oldsymbol{V}oldsymbol{W}_{V}
ight)$

 W_V are value weights shared across all heads.

Quantile Regression





TFT generates prediction intervals on top of point forecasts. This is achieved by the simultaneous prediction of various percentiles (e.g. 10th, 50th and 90th) at each time step.

The prediction intervals are learned by jointly minimizing the quantile loss, summed across all quantile outputs:

$$QL(oldsymbol{y}, \hat{oldsymbol{y}}, oldsymbol{q}) = oldsymbol{q}(oldsymbol{y} - \hat{y})_+ + (1 - oldsymbol{q})(\hat{oldsymbol{y}} - oldsymbol{y})_+$$

Where $(.)_{+} = max(0,.)$

Quantile Regression





$$\begin{split} \mathcal{L}(\Omega, \boldsymbol{W}) &= \\ &= \sum_{y_t \in \Omega} \sum_{q \in \mathcal{Q}} \sum_{\tau=1}^{\tau_{\max}} \frac{QL(y_t, \hat{y}(q, t-\tau, \tau), q)}{M \tau_{\max}} \end{split}$$

Where Ω is the domain of training data containing *M* samples, *W* represents the weights of TFT and $Q = \{0.1, 0.5, 0.9\}$ is the set of output quantiles.

N-Beats





N-Beats [4] is a block-based deep neural network model for univariate time series forecasting.

The main key principles of this architecture are:

- Being simple and generic
- Not rely on time-series-specific feature engineering or input scaling

Input



The model takes a time serie $(x_0, x_1, ..., x_t)$ as input, called **Lookback Period** and predict $(x_{t+1}, ..., x_{t+h})$ called **Forecast Period**.

h is the forecast window length.



Basic block





The basic building block consists of a 4-layers FC (fully connected) network with RELU nonlinearities.

It is then divided into two parts:

- 2-layers FC to predict the **lookback period**, so basically it learns to reconstruct the input.
- 2-layers FC to predict the forecast period.

Stack basic blocks





The basic blocks are organized into stacks using doubly residual stacking principle. The **lookback period output** is used to compute the residual and becomes the input of the next basic block.

 $\mathsf{x}_\ell = \mathsf{x}_{\ell-1} - \widehat{\mathsf{x}}_{\ell-1}$

Where ℓ indicate the $\ell\text{-th}$ block

The **forecast period outputs** from each block are simply aggregated:

$$\widehat{y} = \sum_{\ell} \widehat{y}_{\ell}$$

Combined stacks





Stacks are combined in a higher-level structure. Each stack ouputs a **residual vector** (this vector represents learnings not learnt by the stack) and a **forecast vector**.

The residual vector of each stack becomes the input of the next one and the forecast vectors are aggregated in the same way seen before.



The paper also propose a different configuration of the architecture, augmented with certain inductive biases, to be interpretable.

• **Trend model** A typical characteristic of trend is that most of the time it is a monotonic function, or at least a slowly varying function. In order to mimic this behaviour we constrain the last layer of the basic block:

$$\widehat{\mathbf{y}}_{\boldsymbol{s},\ell}^{tr} = \boldsymbol{T} \theta_{\boldsymbol{s},\ell}^{f}$$

where $\theta_{s,\ell}^f$ are polynomial coefficients predicted by a FC network of layer ℓ of stack s and $T = [1, t, ..., t^p]$ is the matrix of powers of t.



The paper also propose a different configuration of the architecture, augmented with certain inductive biases, to be interpretable.

• **Seasonality model** Typical characteristic of seasonality is that it is a regular, cyclical, recurring fluctuation. In order to mimic this behaviour we constrain the last layer of the basic block:

$$\widehat{\mathbf{y}}_{\boldsymbol{s},\ell}^{\mathsf{seas}} = \boldsymbol{S} \theta_{\boldsymbol{s},\ell}^{f}$$

where $\theta_{s,\ell}^f$ are polynomial coefficients predicted by a FC network of layer ℓ of stack s and $S = [1, \cos(2\pi t), \ldots, \cos(2\pi \lfloor H/2 - 1 \rfloor t)), \sin(2\pi t), \ldots, \sin(2\pi \lfloor H/2 - 1 \rfloor t))]$ is the matrix of sinusoidal waveforms.

Unsupervised Anomaly Detection on Time Series





An **anomaly**, or an **outlier**, is a data point which is significantly different from the rest of the data. [2]

Generally, the data in most applications is created by one or more generating processes that reflect the functionality of a system. When the underlying generating process behaves in an unusual way, it creates outliers.

Unsupervised Anomaly Detection on Time Series



• Classical approach:

Detecting an anomaly based on **PCA** has been widely used in practice because the only information needed is a dataset describing the normal process operation. The most common PCA-based monitoring statistic is **Hotelling's** T^2 :

$$\Gamma^2 = \sum_{i=1}^{n.comp} rac{t_i^2}{s_i^2}$$

Where s_i is the i-th standard deviation of the t_i score. For testing new data, when the value of **T** exceeds a threshold value then is a fault.



• Classical approach:

Q-statistic denotes the change of the events that are not explained by the model of principal components. It is a measure of the difference, or residual between a sample and its projection into the model.

$$Q_i = \tilde{\pmb{x}}_i \tilde{\pmb{x}}_i^{\mathrm{T}} = \pmb{x}_i \left(\mathsf{I} - \mathsf{P}\mathsf{P}^{\mathsf{T}} \right) x_i^{\mathrm{T}}$$

Where **P** is the matrix of principal components obtained using **PCA**. When a vector of new data is available, the Q-statistic is calculated and compared with a threshold value, if the confidence limit is violated then a fault is declared.

Unsupervised Anomaly Detection on Time Series





• Forecasting Based Approaches:

In this methodology, a prediction is performed with a forecasting model (such as a deep learning model) for the next time period.

When a new forecasted value is out of confidence interval, or exceeds a certain threshold, the sample is flagged as anomaly.



These slides heavily borrow from a number of awesome sources. I'm really grateful to all the people who take the time to share their knowledge on this subject with others.

In particular:

- Stanford CS224n: Natural Language Processing with Deep Learning http://web.stanford.edu/class/cs224n/
- Understanding 1D and 3D Convolution Neural Network https://towardsdatascience.com/ understanding-1d-and-3d-convolution-neural-network-keras-9d8f76e29610
- The illustrated Transformer http://jalammar.github.io/illustrated-transformer/

Credits II



• The Positional Encoding https:

//kazemnejad.com/blog/transformer_architecture_positional_encoding/

- Hugging Face https://huggingface.co/
- N-BEATS https://kshavg.medium.com/ n-beats-neural-basis-expansion-analysis-for-interpretable-time-series-for

References I



[1] S. Bai, J. Z. Kolter, and V. Koltun.

An empirical evaluation of generic convolutional and recurrent networks for sequence modeling.

arXiv:1803.01271, 2018.

- [2] A. Bl'azquez-Garc'ia, A. Conde, U. Mori, and J. Lozano. A review on outlier/anomaly detection in time series data. *ArXiv*, abs/2002.04236, 2020.
- [3] B. Lim, N. Loeff, S. Arik, and T. Pfister.

Temporal fusion transformers for interpretable multi-horizon time series forecasting. 2020.

References II



 B. N. Oreshkin, D. Carpov, N. Chapados, and Y. Bengio.
 N-beats: Neural basis expansion analysis for interpretable time series forecasting. In International Conference on Learning Representations, 2020.

[5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin.

Attention is all you need.

In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008. Curran Associates, Inc., 2017.