

# AI による数法則発見の時系列データへの 拡張と金融データへの応用

Modeling and Visualization of Social Reality  
Using Latent Profile Analysis and Number Law Discovery Methods  
for Evidence-Based Policy Making

蒲田 涼馬 (Ryoma Gamada)  
u455007@st.pu-toyama.ac.jp

富山県立大学大学院 工学研究科 電子・情報工学専攻  
情報基盤工学講座

N212, 09:30-10:00 Tuesday, February 13, 2024.

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

情報技術の発達により、社会における様々なデータを観測・収集することが可能に

→ 経済分析においても将来予測などの研究が急速に発展.  
しかし要因分析に関する研究はそれほど進んでいるとは言えない.

経済に影響を与える要因を分析する研究

因果探索による要因分析.  
シンボリック回帰を用いた要因分析.

本研究

シンボリック回帰を用いて分析を行う.

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

### なぜシンボリック回帰？

経済分野では，原因と結果の間に成り立つ関係性が重要

→ 複数の要因が複雑に影響しあうため，因果探索では具体的にどのような絡み合って影響を与えるかを詳細に分析できない．

### アプローチ

公開されている金融データ，経済データ，市場間データを用いて分析を行い，データ間の関係性を数理モデルによって表す．

数理モデルの例

$$(\text{データ A}) = 2.0 \cdot (\text{データ B}) + 1.0 \cdot (\text{データ C}) - 1.0 \cdot (\text{データ D})$$

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## 公開されているデータ

XMMT5 や Investing.com, 日本銀行時系列サイトで様々なデータが公開されている.

Table 1: 公開されている様々な経済データ

データ項目	
為替レート	金利
コモディティ価格	エネルギー価格
マネーストック	ボラティリティ指数
出来高	スプレッド
株価指数	ニュース

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## 目的

時系列経済データを用いて、時系列を考慮した数法則の発見を行いデータ間の関係性をモデル化する手法を提案する。

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## 使用する手法の概要

機械学習を用いたシンボリック回帰手法である「End to end symbolic regression with Transformers」を拡張させ、時系列を考慮した分析を行う。可読性と得られる情報量を重視し、人間が式を見ることである程度その式が何を表しているのかをわかるレベルのものを生成させる。

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## End to end Symbolic Regression with Transformers の概要

### データから数式を自動発見する深層学習アプローチ

従来のシンボリック回帰の課題: 計算コストが非常に高い.

Transformer アプローチの着想:

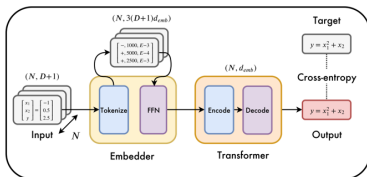
数式は、演算子、定数、変数といった要素が並んだシーケンスとして表現する.

例)  $y = x + 2 \cdot \sin(z) \rightarrow + \times * 2 \sin z$  Transformer はシーケンスデータの複雑なパターン学習と高速なシーケンス生成能力を持つため、数式発見に応用できる.

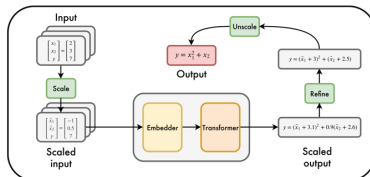
# Transformer によるシンボリック回帰 (前準備)

8/21

## Transformer によるシンボリック回帰



Training



Inference

## Embedder

入力データをトークン化する。

埋め込みルックアップテーブルを使ってトークンをベクトルに変換する (512 次元)。

ベクトルを FFN に入力し、短いベクトルに圧縮する。

この処理をすべてのデータ点に対して行う。

これを Transformer 本体に渡す。



## Transformer のメカニズム

Transformer はエンコーダーとデコーダーから成り、Attention メカニズムが核を担っている。

### ■ Attention メカニズム

シーケンス内の各要素が、ほかのどの要素に注意を向けるべきかを動的に判断し、その重要度について重みづけを行う。

$$Scores = QK^T \quad (1)$$

$$ScaledScore = \frac{QK^T}{\sqrt{d_k}} \quad (2)$$

$$AttentionWeight = softmax(\frac{QK^T}{\sqrt{d_k}}) \quad (3)$$

$$Attention = softmax(\frac{QK^T}{\sqrt{d_k}})V \quad (4)$$

ここで  $Q$  はクエリ行列,  $K$  はキー行列,  $V$  は値,  $d_k$  はキーの次元を意味する。

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## ■ Feed-Forward Network

ベクトルを圧縮する際やベクトルを拡張する際に使われる。  
Transformer モデルでは両方の使われ方がされている。

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (5)$$

ここで  $x$  は前のステップからの入力,  $W_1, W_2, b_1, b_2$  はモデルが学習する重みとバイアスを意味する。

### エンコーダー

入力されたデータ点を読み込み, それらがどんな関数の特徴を持っているかを分析する。

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## エンコーダー

Embedder によってベクトル化された  $N$  個のデータ点を受け取る。  
関連性の計算 (self-attention) を行い、各データ点がほかのすべてのデータ点とどの程度関連しているかを計算する。  
FFN によって関連性を加味した情報をさらに深く処理する。  
関連性の計算と FNN による処理を 4 層繰り返す。  
入力データ全体の特徴を要約した情報をデコーダーに渡す。

## エンコーダーの self-attention

入力された全データ点の、相互の関連性の強さを計算する。  
自分自身を含むすべての入力データ点を参照する。  
これによって各データ点のベクトルにデータセット全体における他の点との関係性を埋め込む。

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## デコーダー

エンコーダーからの情報と数式の一部を受け取る (初回は開始トークンのみを受け取る).

自己参照 (初回は開始トークンのみ, 2 回目以降はトークン列を受け取る)

エンコーダーからの要約情報に注目し, 作る関数を理解する.

上記 2 角情報をもとに, 数式の次に来るべきトークンを予測する.

自己参照から予測までの流れを終了トークンが出力されるまで 1 トークンずつ繰り返す.

完成した数式を出力する.

## デコーダーの自己参照

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}} + M)V \quad (6)$$

注目しているトークンとほかの全トークンとの関連度を計算する.

マスクを利用して未来のトークンに関する部分を  $-\infty$ , それ以外の部分を 0 にする.

これに softmax 関数を適用することでモデルが過去と現在のトークンのみを参照して次のトークンを予測するようにする.

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## クロスエントロピーによる数式の評価

$$H(p, q) = - \sum_i p(x_i) \log q(x_i) \quad (7)$$

ここで  $p(x_i)$  は正解の確率分布,  $q(x_i)$  はモデルが予測した確率分布を意味する.

これを計算することで生成した数式の評価を行い, これを高めていくように学習していったモデルを作る.

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## 今回やったこと

説明変数を正規化し数値実験を行った.

説明変数を絞りこむことによってどのような結果 (同定式と精度) が得られるのかを検証した.

別の手法の検討と E2ESR についての勉強を行った.  
新規性を考えた.

## 使用したデータ

項目：10 種類の経済データ

対象年：2015 年 1 月 1 日から 2024 年 12 月 31 日までの 10 年間

データ数：土日祝日を除いた日足で 2355 のデータ (前半 8 割で学習, 2 割でテスト)

目的変数：USDJPY の  $t$  のときの値

説明変数：Table2 の  $t-1$  のときの値

前は正規化をせずに実験を行ったが、今回は正規化を行った結果を求めた。

Table 2: 数値実験に用いたデータ

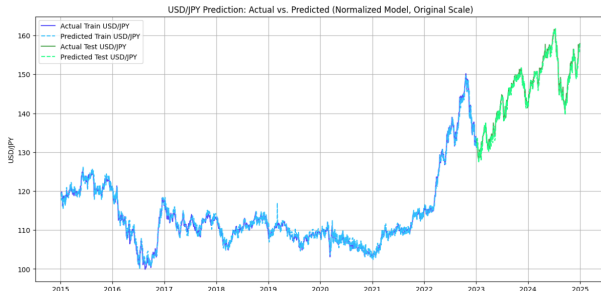
データ項目	
SP500 価格	日経平均株価
日米の金利差 10 年	日本 10 年国債
米国 10 年国債	日米の金利差 2 年
オイル価格	金価格
USDJPY	VIX 指数

## 実験結果

USDJPY を求めた結果を以下に示す.

## 決定係数

- トレーニングデータについての決定係数 : 0.9890
- トレーニングデータについての RMSE : 0.957
- test データについての決定係数 : 0.9796
- test データについての RMSE : 1.1425



はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに



はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

$$\begin{aligned}
 & -0.10040714384915 \cdot x_0 - 0.000198358776965863 \cdot x_7 + 0.951613884352736 \cdot x_8 + 5.3976 \\
 & 0.011125679 + \frac{-2.21694913680499e-5 \cdot x_0 - 0.000115127670935178 \cdot x_1 + 0.043894457}{2668756 \cdot x_2 + 333.516649543171 \cdot x_3 - 130.630846400977 \cdot x_4 + (0.301880750900265} \\
 & 0.677 \\
 & - 0.00456788346207082 \cdot x_4) \cdot (1.85830014281174e-5 \cdot x_1 + 11.5184087304255 \cdot x_2 - (4 \\
 & .35176808096492 \cdot x_2 + 6.64848083125045) \cdot (0.00108863330413424 \cdot x_0 - 0.054506250 \\
 & 8133566) - 199.764377322288) + 249.687084647262
 \end{aligned}$$

簡略後の式

別添 pdf

## 実験 2 の結果と考察

18/21

### 実験結果

USDJPY を求めた結果を以下に示す。  
重回帰の重みが 0.01 以下になったものを除外して実験した。

### 決定係数

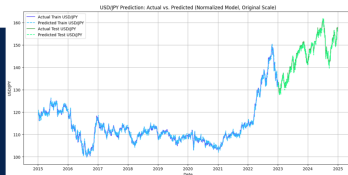
- トレーニングデータについての決定係数 : 0.9890
- トレーニングデータについての RMSE: 0.957
- test データについての決定係数 : 0.9796
- test データについての RMSE: 1.1425

--- 特徴量の重み分析 (閾値: 0.01) ---

特徴量名	重み (絶対値)	除外/保持
SP500	0.000056	除外
Nikkei225	0.000014	除外
金利差_3P-US_2Y	0.019127	保持
3P_10Y_Yield	0.223148	保持
US_10Y_Yield	0.153652	保持
金利差_3P-US_10Y	0.010400	保持
WTI_Oil	0.003235	除外
Gold	0.000164	除外
VIX	0.000013	除外
USD/JPY_lag1	0.985967	保持

元の特徴量数: 10  
除外後の特徴量数: 5  
除外された特徴量数: 5  
保持された特徴量: ["金利差\_3P-US\_2Y", "3P\_10Y\_Yield", "US\_10Y\_Yield", "金利差\_3P-US\_10Y", "USD/JPY\_lag1"]

重回帰の結果



波形

$$\frac{1}{100} \left( \frac{USD/JPY_{lag1}}{100} + 0.019127 \left( \frac{3P\_US\_2Y}{100} + 0.223148 \left( \frac{3P\_10Y\_Yield}{100} + 0.153652 \left( \frac{US\_10Y\_Yield}{100} + 0.010400 \left( \frac{3P\_US\_10Y}{100} + 0.985967 \left( \frac{USD/JPY_{lag1}}{100} \right) \right) \right) \right) \right) \right)$$

得られた式

## 前回聞かれたこと

これは人が作ったものをオンラインで使っているだけなのではないか  
→オンラインで使っているわけではないが、事前学習済みのモデルをダウンロードして使っているだけだった。  
ただ 1 から学習させてモデルを作成するシステムもリポジトリにはあったため学習させることも可能  
ただ計算時間はかなりかかりそう。

## 新規性をどうするか

- ダウンロードした事前学習モデルでファインチューニングを行うことで時系列に対応、可読性を向上させる手法を考える。
- 一から学習させるためにプログラムを組みなおし時系列を考慮、可読性も向上させる。
- 非定常性を考慮した分析を行うことで金融時系列により強いモデルを構築する (上記両方の場合)

はじめに

統計データの特徴  
と研究の概要

End to end  
symbolic  
regression with  
transformers

今回やったこと

数値実験並びに  
考察

おわりに

## 可読性の向上

可読性の高さを決める要因として以下が考えられている。

- ノード数, 深さ, 演算子の数と種類, 項の数, コード長

## 可読性向上についての研究

- 損失関数に複雑度を組み込んで複雑になりすぎず, 精度もある程度維持させる.
- 短い数式や簡潔な式に高いスコアをつけ, 学習させる.

## まとめ

手法の学習アルゴリズムについての勉強をし、また前回の実験を正規化した条件で行った。

変数の絞り込みによって精度と可読性がどのように変わってくるのか、簡単な実験を行って検証してみた。

プログラムを書き換え、一から学習を行えるように試みた (途中)

- 正規化を行うことで以前は軽い重みだった変数の重みが重くなり、式がより複雑になった。
- 変数を絞り込むことで式は簡略化され、見やすくなった。
- プログラムの書き換えの方はまだ終わっていない。

## 今後の展望

- まずは手法について完璧に理解する。
- 一から学習させて実行してみる。(従来研究よりは少ないデータ量)
- 今回は  $t-1$  からの影響を考えたが、もっと