

A Time Series Approach to Explainability for Neural Nets with Applications to Risk-Management and Fraud Detection

Marc Wildi³ and Branka Hadji Misheva⁴

March 23, 2022

³ZHAW Zurich University of Applied Sciences, IDP; e-mail: wlmr@zhaw.ch

⁴ZHAW Zurich University of Applied Sciences, IDP; e-mail: hadj@zhaw.ch

Table of contents

- 1 Introduction
- 2 Problem 1: Random-Nets
- 3 Problem 2: Black-Box
 - Complexity vs. Interpretability
 - Problem 2: Deploying Explainability
 - Problem 2: Utility of XAI Methods for Finance
- 4 Solution Problem 2: XAI, Time Series and Neural Nets
- 5 Application of LPD to Crypto (BTC)
- 6 Application: XAI and Risk Management

Neural Nets: Forecasting

- ① Forecast performances: review of **international forecast competitions**
 - M1 (1982), M2 (1993), M3 (2000), NN3 (2007) and NN5 (2009) competitions: **classic linear** approaches **outperform** computationally intensive approaches (including neural nets)
 - M4 (2020) and M5 (2021): **hybrid approaches** based on a mix of ARIMA and neural nets **outperform** (hybrid approach won M4 competition)
- ② Accruing interest in **neural nets for forecasting** (in particular economic time series)

Neural Nets: Main Problems/Issues

- ❶ Problem 1: **random nets** (numerical computations, learning/estimation)
 - Numerical optimization algorithms do **not reach global** optimum (complexity)
 - Parameters are **not properly identified**: different nets lead to similar performances
 - Computationally intensive
- ❷ Problem 2: **Black-box** or how to relate 'output' to 'input'?
- ❸ Problem 3: **Overfitting** of richly parameterized nets (will be addressed in follow-up paper)

Problem 1: Random-Nets

- 1 Illustration of problem based on application of simple neural net to simple data
- 2 Simple data

$$y_t = x_t + \epsilon_{1t}$$

$$x_t = \epsilon_{2t}$$

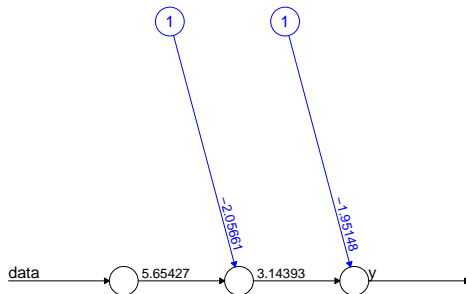
where ϵ_{it} are two iid Gaussian noise processes

- 3 Simple net: see next slide

Illustration Problem 1: Simple Net

- Non-linear function (data transformed to unit-interval)

$$o_t = \sigma \left(-1.951 + 3.14 \sigma(-2.057 + 5.65 x_t) \right)$$



Error: 0.899312 Steps: 126

Illustration Problem 1: Random Estimates

① Estimates

	neuralnet seed(1)	neuralnet seed(6)	own algo seed(1)
b1	-2.057	0.703	2.122
w1	5.654	-4.156	1.715
b2	-1.951	1.255	-40.445
w2	3.144	-5.146	42.696
mse	0.899	0.894	0.890

Table: Estimated parameters and MSE-performances for different random-seeds and optimization algorithms

Problem 1: Random Nets

- Optimization algorithm does not converge (in one of the simplest cases)
 - **Estimates inherit randomness** from random initialization of parameters (and data-generating process)
 - Problems: **inconsistency** (more data does not solve problem) and **un-identification** (different nets lead to similar performances)
- Summarize above findings by appellation **random nets**
- Increasing net-complexity (number of layers/neurons) magnifies the above problems
- **XAI for NN cannot rely on net parameters** (weights/biases)

Problem 2: Complexity vs Interpretability

- Tradeoff: more complex models might perform better in terms of predictions** but as a result are **less explainable** whereas simpler models rank high on interpretability but eventually lower on predictive accuracy (though remember competition outcomes)

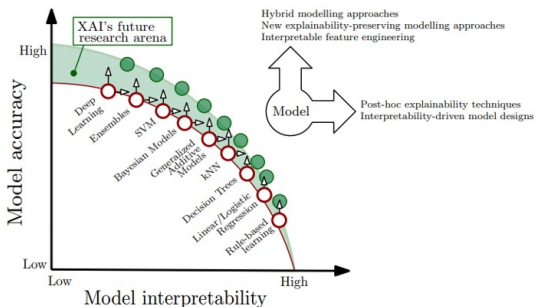


Figure: The trade-off. Image source: Arrieta et al. 2019

Problem 2: Why Do We Need Explainability?

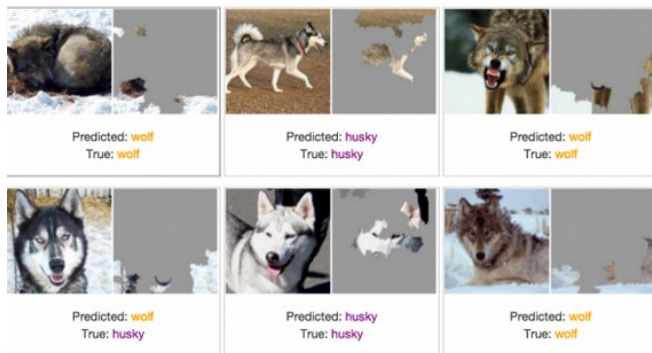
Let's consider the 'Husky vs Wolf' experiment results.



- The classifier makes one mistake!

Problem 2: Why Do We Need Explainability?

Next, we investigate which features drive the classification.



- The decision is based on whether there are white patches on the image!

Problem 2: Why Do We Need Explainability?

- In practice, it is essential to verify that for a given task, the high measured **accuracy** is actually the result of **proper problem representation** and that the model is capturing **relevant dependencies** between the observed features.
- This ensures **trust** in the system.
- Moreover, **regulation** demands it.

No black box excuses – explainability/traceability of models is necessary and can improve the analysis process | It is the responsibility of supervised firms to ensure that BDAI-based decisions can be explained and are understood by third-party experts. Supervisory authorities take a critical view of models that are categorised purely as black boxes. New approaches allow firms using such models to at least gain some insight into how these models work and identify the reasons behind decisions. In addition, a better understanding of models provides an opportunity to improve the analysis process – allowing, for instance, the responsible units in the supervised firm to identify statistical problems.

Figure: Extract: Bafin AI and Big Data Report 2020

Deploying Explainability: Classic XAI Methods

- To ensure trust in modelling, in most cases **post-hoc explainability** is required.
- Post-hoc explainability techniques aim to provide understandable information about **how output** (predictions) is generated **from input**.
- In general, methods are considered in view of two main criteria (Linardatos et al. 2021):
 - the **type of algorithm** on which they can be applied (model-specific vs. model-agnostic)
 - the **unit being explained** (if the method provides an explanation which is instance-specific then this is a **local explainability** technique and if the method attempt to explain the behavior of the entire model, then this is a **global explainability** technique).



Figure: Machine Learning Interpretability Techniques

Problem 2: Deploying Explainability: LIME

- LIME, short for locally interpretable model agnostic explanations, is an explanation technique which aims to identify an interpretable model over the representation of the data that is locally faithful to the classifier
- More formally, the explanations provided by LIME is obtained by following :

$$\xi(x) = \underset{g \in G}{\operatorname{argmin}} L(f, g, \pi_x) + \Omega(g) \quad (1)$$

g : An explanation considered as a model

$g \in G$: class of potentially interpretable models such as linear models and decision trees

$f : R^d \rightarrow R$: The main classifier being explained

$\pi_x(z)$: Proximity measure of an instance z from x

$\Omega(g)$: A measure of complexity of the explanation $g \in G$

- The goal is to minimize the locality aware loss L without making any assumptions about f .

Problem 2: Deploying Explainability: SHAPLEY

- The Shapley value is the average marginal contribution of a feature value across all possible coalitions.
- Given a model:

$$f(x_1, x_2, x_3 \dots x_n) \quad (2)$$

- The marginal contribution $\Delta_v(i, S)$ a feature i :

$$\Delta_v(i, S) = v(S \cup i) - v(S) \quad (3)$$

- Let Π be the set of permutations of the integers up to N and given $\pi \in \Pi$ let $S_{i,\pi} = j : \pi(j) < \pi(i)$ are the players preceding player i in π , then:

$$\phi_v(i) = \frac{1}{N!} \sum_{\pi \in \Pi} \Delta_v(i, S_{i,\pi}) \quad (4)$$

Deploying Explainability: Utility of Classical XAI Methods for Finance

- Classical approaches and their current implementation are **not tailored for time series** and in particular for **financial data** (subject to trends, vola-clusters,...).
- **Key limitation** of many classical methods is the fact that they **ignore feature dependence** which is a defining property of financial data.
- Specifically, perturbation-based methods are fully dependent on **the ability to perturb samples in a meaningful way**. In the context of financial data:
 - if **features are correlated**, the artificial coalitions created will lie **outside of the multivariate joint distribution** of the data,
 - if the data are independent, coalitions can still be meaningless;
 - generating artificial data points through random replacement disregards the time sequence hence producing **unrealistic values** for the feature of interest.

Problem 2: Example

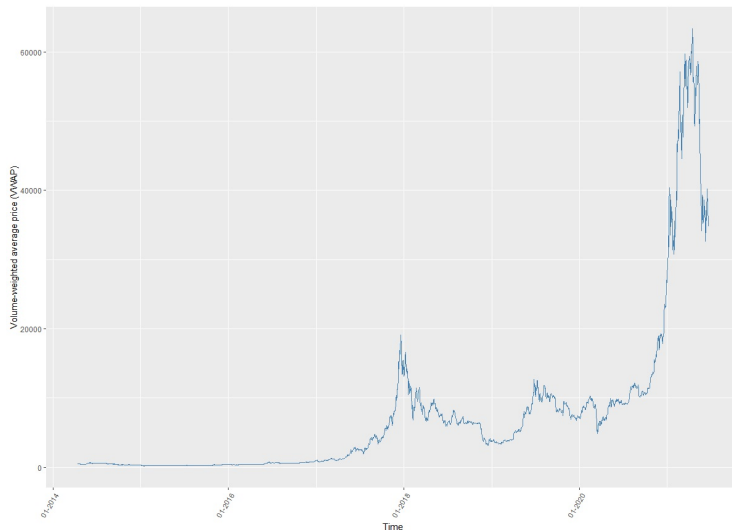


Figure: BTC Prices (correct ordering)

Problem 2: Example

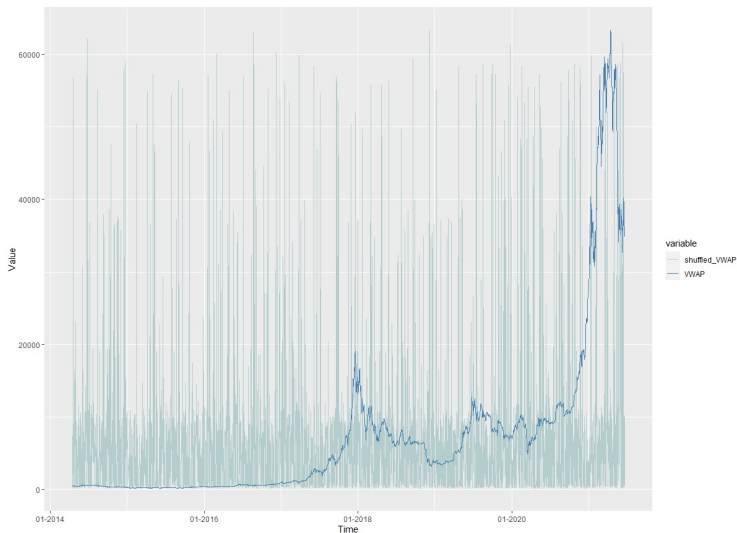


Figure: BTC Prices (correct ordering vs shuffled values)

Solution: Explainability (X-) Function

- Net-parameters cannot be straightforwardly interpreted (**random-nets**)
- Emphasize **input-output** relations instead: **infinitesimal** changes
- **Partial derivative** of
 - Net-output with respect to explanatory (input) variables: $\frac{\partial o_t}{\partial x_{it}}$
 - X-function $xf(o_t)$ of net-Output with respect to input: $\frac{\partial xf(o_t)}{\partial x_{it}}$
- Interpretability/XAI: observe sensitivities (gradient/derivative of X-function) over time

Simple X-Function: Identity and LPD

- Represent net-output o_t as time-varying 'linear' regression

$$\begin{aligned} o_t &= b_t + w_{1t}x_{1t} + w_{2t}x_{2t} + \dots + w_{pt}x_{pt} \\ y_t &= o_t + \epsilon_t \end{aligned}$$

with target y_t and explanatory input-data $x_{1t}, x_{2t}, \dots, x_{pt}$

- $\text{LPD}_t := (b_t, w_{1t}, w_{2t}, \dots, w_{pt})$
- b_t intercept, $w_{1t}, w_{2t}, \dots, w_{pt}$ **time-dependent 'sensitivities'** of output with respect to input: **partial effects**
 - Linear regression: $b_t = b_0, w_{1t} = w_1, w_{2t} = w_2, \dots, w_{pt} = w_p$ fixed (not time dependent)
- **Non-stationarity of data-generating** process can be tracked by **non-linearity of net**, see crypto-example

Simple X-Function: LPD and QPD

- **X-function: LPD** i.e. time-dependent linear replication of net
 - Sensitivity of LPD: **QPD** (quadratic parameter data)
 - **Second order derivative** or Hessian of net-output for each time point $t = 1, \dots, T$
 - Dimension: $T * p * p$ -dimensional array
- **QPD** can track periods of **departure from linearity**
 - Large $|\text{QPD}_t|$ means high non-linearity at time point t

Derivation LPD: Forward Propagation

- First hidden layer

$$\mathbf{dA}_t^{f(1)} = \left(\mathbf{W}^{(1)'} \cdot \mathbf{A}_t^{(1)} \cdot (\mathbf{e} - \mathbf{A}_t^{(1)}) \right)' \quad (5)$$

- Following hidden/output layers

$$\mathbf{dA}_t^{f(k)} = \left(\left(\mathbf{dA}_t^{f(k-1)} \mathbf{W}^{(k)} \right)' \cdot \mathbf{A}_t^{(k)} \cdot (\mathbf{e} - \mathbf{A}_t^{(k)}) \right)' \quad (6)$$

Derivation LPD: Backward Propagation

- Output layer

$$d\mathbf{A}_t^{\mathbf{b}^{(p)}} = \mathbf{A}_t^{(p)} \cdot (\mathbf{e} - \mathbf{A}_t^{(p)}) \quad (7)$$

- Previous hidden layers

$$d\mathbf{A}_t^{\mathbf{b}^{(k)}} = \mathbf{W}^{(k+1)} d\mathbf{A}_t^{\mathbf{b}^{(k+1)}} \cdot \mathbf{A}_t^{(k)} \cdot (\mathbf{e} - \mathbf{A}_t^{(k)}) \quad (8)$$

- Input layer

$$d\mathbf{A}_t^{\mathbf{b}^{(0)}} = \mathbf{W}^{(1)} d\mathbf{A}_t^{\mathbf{b}^{(1)}} \quad (9)$$

Derivation QPD

- Output layer

$$\mathbf{dda}_t^{(p)} = \mathbf{A}_t^{(p)} \cdot (\mathbf{e} - \mathbf{A}_t^{(p)}) \cdot (\mathbf{e} - 2\mathbf{A}_t^{(p)}) \quad (10)$$

$$\mathbf{ddA}_t^{(p)} = \left(\mathbf{W}^{(p)'} \mathbf{dA}_t^{f^{p-1}'} \cdot \mathbf{dda}_t^{(p)} \right)' \quad (11)$$

- Previous hidden layers

$$\begin{aligned} \mathbf{dda}_t^{(k)} &= \mathbf{A}_t^{(k)} \cdot (\mathbf{e} - \mathbf{A}_t^{(k)}) \cdot (1 - 2\mathbf{A}_t^{(k)}) \\ \mathbf{ddA}_t^{(k)} &= \left(\mathbf{W}^{(k+1)} \mathbf{ddA}_t^{(k+1)'} \cdot \mathbf{A}_t^{(k)} \cdot (\mathbf{e} - \mathbf{A}_t^{(k)}) \right)' \quad (12) \\ &\quad + \left(\mathbf{W}^{(k+1)} \mathbf{dA}_t^{b^{(k+1)}} \cdot \mathbf{dda}_t^{(k)} \cdot \mathbf{W}^{(k)'} \mathbf{dA}_t^{f^{k-1}'} \right)' \quad (13) \end{aligned}$$

Derivation QPD

- First hidden layer

$$\begin{aligned} \mathbf{dda}_t^{(1)} &= \mathbf{A}_t^{(1)} \cdot (\mathbf{e} - \mathbf{A}_t^{(1)}) \cdot (1 - 2\mathbf{A}_t^{(1)}) \\ \mathbf{ddA}_t^{(1)} &= \left(\mathbf{W}^{(2)} \mathbf{ddA}_t^{(2)'} \cdot \mathbf{A}_t^{(1)} \cdot (\mathbf{e} - \mathbf{A}_t^{(1)}) \right)' \quad (14) \end{aligned}$$

$$+ \left(\mathbf{W}^{(2)} \mathbf{dA}_t^{(2)} \cdot \mathbf{dda}_t^{(1)} \cdot \mathbf{W}^{(1)'} \right)' \quad (15)$$

- Input layer

$$\mathbf{ddA}_t^{(0)} = \left(\mathbf{W}^{(1)} \mathbf{ddA}_t^{(1)'} \right)' = \mathbf{W}^{(1)} \mathbf{ddA}_t^{(1)'} \quad (16)$$

Derivation of Generic Differentiable X-Functions

- **Generic X-function:** straightforwardly derivation from **LPD**
(modify the above expression at **output layer** only)

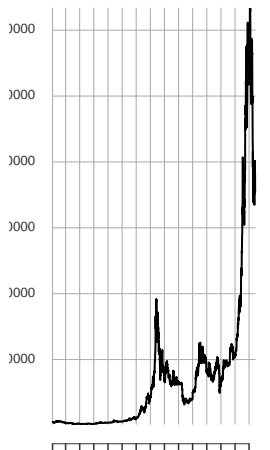
$$\mathbf{d}(\mathbf{x}\mathbf{f} \circ \mathbf{A}^{\mathbf{f}})_t^{(p)} := \left(\nabla \mathbf{x}\mathbf{f}(\mathbf{A}_t^{(p)}) \cdot \mathbf{d}\mathbf{A}_t^{\mathbf{f}(p)} \right)'$$

Alternative X-Functions

- X-function is **MSE-criterion**: track **importance** of input data over time
- X-function is **trading performance** (Sharpe ratio, max-drawdown,...): track effect of input data on trading performances over time
- ...

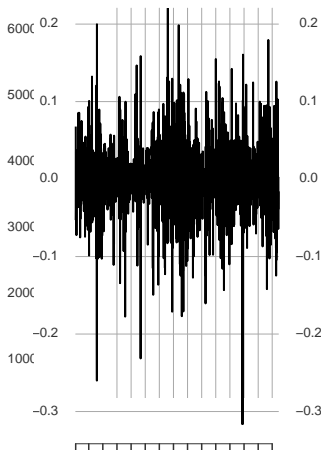
Application to BTC Crypto-Currency

BTC 2014-04-15 / 2021-06-20



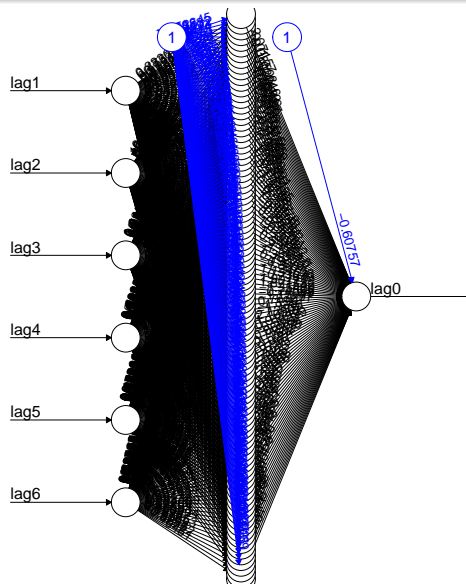
Apr 15 2014 Oct 01 2017 Oct 01 2020

Log-returns 2014-04-15 / 2021-06-20



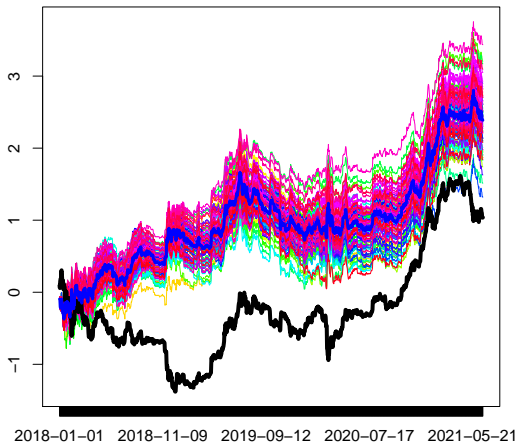
Apr 15 2014 Oct 01 2017 Oct 01 2020

Net-Architecture



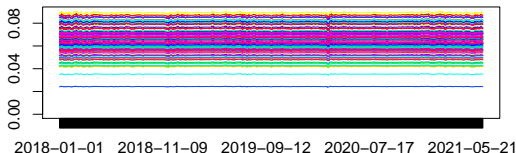
Random Nets and Random Performances (trading based on sign-rule)

g-perf (sign-rule): 100 random nets (colored) vs. buy-and-hold

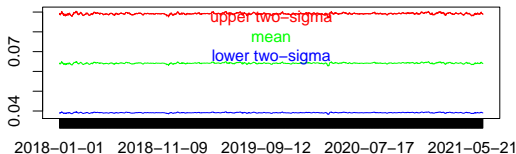


Mean and 2-Sigma Bands of Lag-One BTC (of LPDs of Random Nets)

LPDs of 100 random nets for explanatory 6



Mean and two-sigma band of above LPD realizations



LPD: Mean, Standard Deviation and Synthetic t-Test

- Mean and standard deviations of LPDs over 100 random-nets:

$$\overline{\text{LPD}}_t := \frac{1}{100} \sum_{i=1}^{100} \text{LPD}_{ti}$$

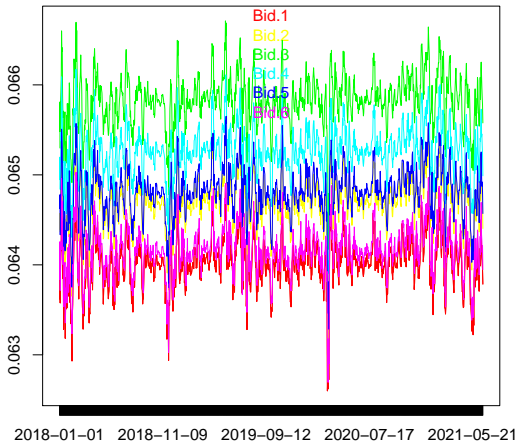
$$\sigma_t := \frac{1}{100} \sum_{i=1}^{100} (\text{LPD}_{ti} - \overline{\text{LPD}}_t)^2$$

- **Importance of input variables:** significance based on synthetic t-Test

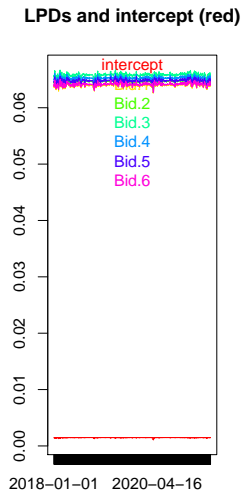
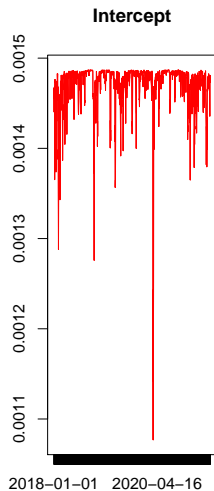
$$\mathbf{t}_t := \frac{\overline{\text{LPD}}_t}{\sigma_t}$$

Mean LPD of All Input Variables

100 sample LPD of 100 random nets: explanatory (colored) and i



Mean Intercept and Mean LPD of All Input Variables



Findings 1: First Order Linear Approximation

- Resolve **problem 2** (black-box): **LPD**
- Resolve **problem 1** (randomn-nets): **aggregate mean-LPD**
- Findings: **first order linear approximation** (explainability)
 - (Mean-/consensus) Neural net is close to an unassuming **equally-weighted MA(6)**-filter

$$\bar{o}_t = \frac{1}{100} \sum_{i=1}^{100} o_{t,i} = \overline{\text{LPD}}_t \left(\begin{array}{c} 1 \\ \mathbf{x}_t \end{array} \right) \approx 0.0015 + 0.065 \sum_{j=1}^6 x_{jt}$$

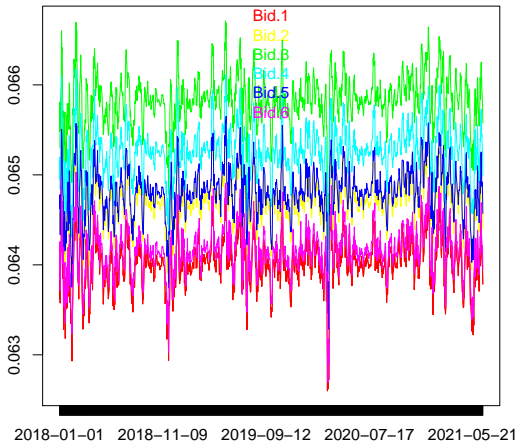
see also Bundi,Wildi (2019)

Findings 1: First Order Linear Approximation

- We can gain '**trust**' in the consensus (mean) net output by relating it to a **simple forecast heuristic**
- Unfortunately, there does not seem to be great benefits of using a complex net-architecture in this case (recall forecast competitions)
- But maybe there is more to extract from the LPD than a 'first order' approximation: let's look at the second order departures from linearity

Mean LPD of All Input Variables

100 sample LPD of 100 random nets: explanatory (colored) and i



Findings 2: Second-Order Non-Linearity

- **Second-order non-linearity:** departure from constant LPD-levels
- **Second-order non-linearity:** stationary-process **strongly correlated** across input variables and across random realizations of LPD

	Lag 1	Lag 2	Lag 3	Lag 4	Lag 5
Cor. across mean LPDs	1.000	0.980	0.980	0.980	0.980
Cor. of random LPDs	0.980	0.990	0.990	0.990	0.990

Table: Correlations

- **Common second-order non-linearity factor:** pervades all dimensions of LPD
- Conjecture: **second-order non-linearity factor** is driven by **common data (-generating process)**

Risk-Management: Exploit Second-Order Non-Linearity Factor

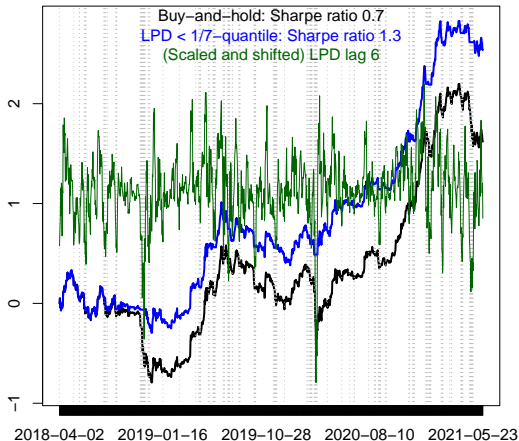
- Idea Risk-Management (and Fraud Detection with slight alterations):
 - **Identify 'uncertain times' or 'unusual states'** of the market in real-time
 - **Downsize** market exposure during **uncertain times/unusual states**
- Identification of **uncertain times** and **unusual states**
 - Unusually **weak LPD** (weak dependence structure)
 - Unusually **strong QPD** (strong non-linearity)
 - Uncertainty: **wide/expanded sigma-bands** (standard deviation) of random LPD (random nets are inconclusive about dependence structure)

Unusually Weak LPD: Operationalization

- Assumption: market-exit signals have a probability of $1/7$ (**one per week** in the mean)
- **Market exit** (cash-position) is triggered if **LPD drops below $1/7$ -quantile**.
- **Quantile** is based on a rolling-window of length **1 quarter** (100 days)

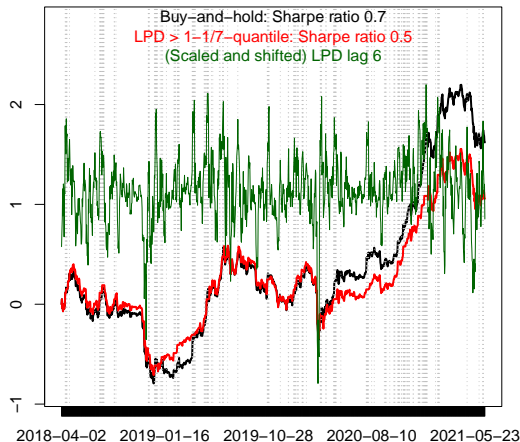
Market-Exit: Critical Time-Points (Unusually Weak LPD)

y-and-hold (black) vs. LPD-down (blue): $< 1/7$ -quantile, length 9

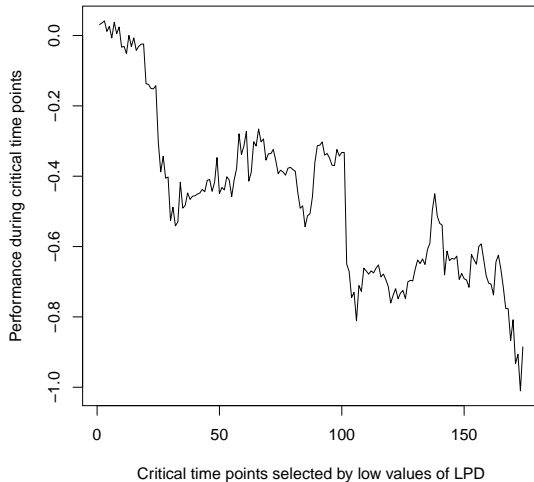


Market-Exit: Un-Critical Time-Points (Unusually Strong LPD)

Buy-and-hold (black) vs LPD-up (red): $> 1-1/7$ quantile, length 90,



Performance During Critical Time-Points



Analysis Critical Time Points

	Next day's mean return
Today's LPD: Weak dependence	-0.509%
Today's LPD: Normal dependence	0.244%
Today's LPD: Strong dependence	0.305%
All time points	0.142%

Findings

- **Concomitance of negative (next day's) BTC-growth and of (today's) weak data dependence (LPD)**
 - Panic or herding

Summary

- **Problem 1** (random nets): use aggregate consensus outputs/forecasts
- **Problem 2** (explainability)
 - Input/output relations: **time-dependent regression**
 - X-functions: **LPD** (regression), **QPD** (strength of non-linearity), **MSE** (importance of data), **customized** (for example trading performances)

Summary

- Application of LPD to BTC
 - Interpretability: **first-order linear approximation** (simple forecast heuristic)
 - Departure from linearity: **common second-order non-linearity factor**
- Addressing **Risk-management** and **fraud-detection** with XAI-tool
 - Exploit **second-order non-linearity factor**
 - Identify **uncertain times** and **unusual states**: small LPD or large QPD or wide σ -bands
 - Down-size market exposure (RM) or analyse for eventual fraud

Summary

- Extension of original XAI-prospect: inferring a **better understanding of the data** from a **better understanding of the model**