



Machine Learning on the OMXC25 Stocks

Bastian Bakkensen
Carlos Fernando Duarte Faurby
Daniel Hans Munk
Rasmus Nielsen

Date: 15-06-2022

UNIVERSITY OF COPENHAGEN



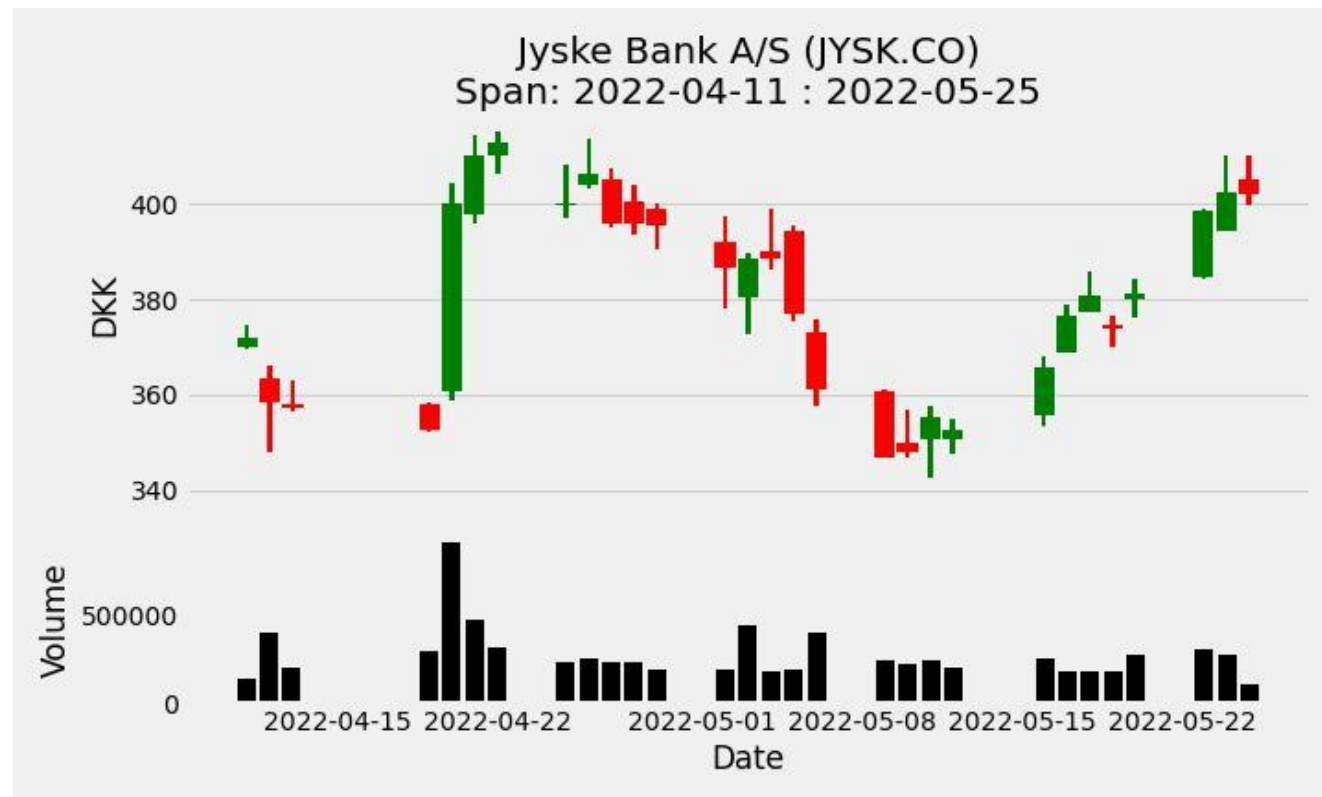
Objective:

With ~20 years of stock data from the OMX Copenhagen 25 index, what is the predictive ability of modern Machine Learning?

Data - the 'ohlcv' format

Ørsted
 DSV
 Carlsberg
 Novo Nordisk
 Chr. Hansen Holding
 Tryg
 Novozymes
 FLSmidth & Co.
 Genmab
 Rockwool
 Royal Unibrew
 Coloplast
 Vestas Wind Systems

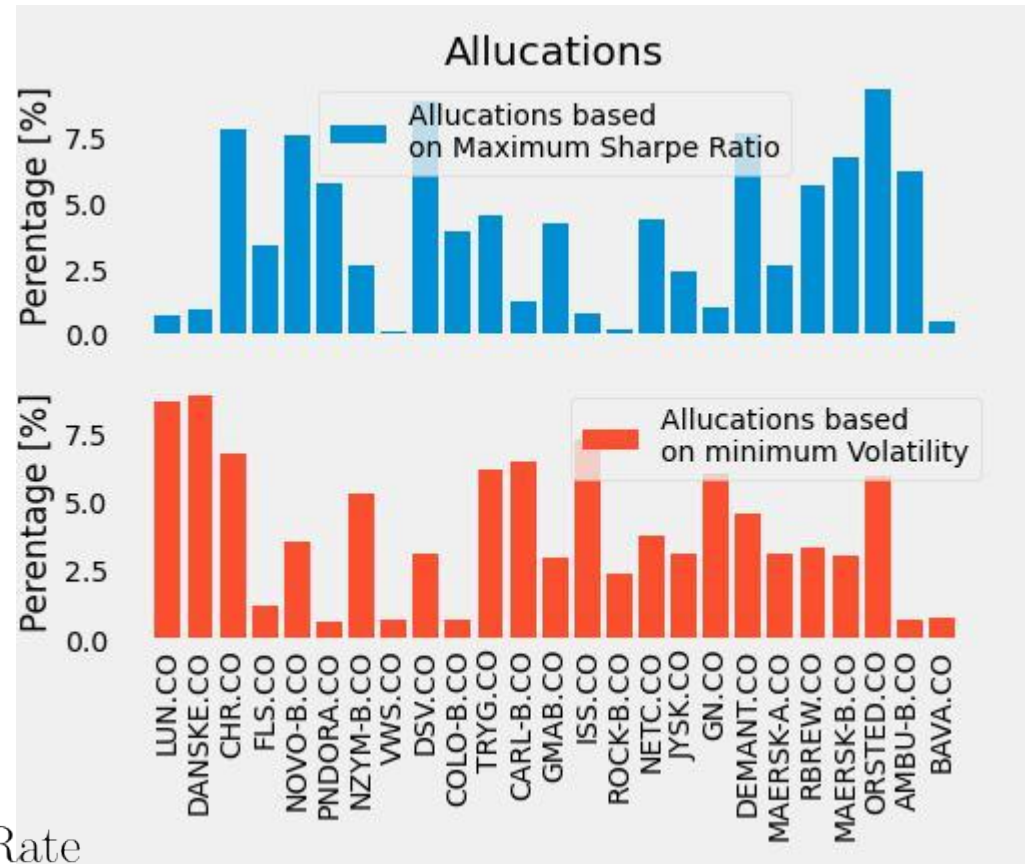
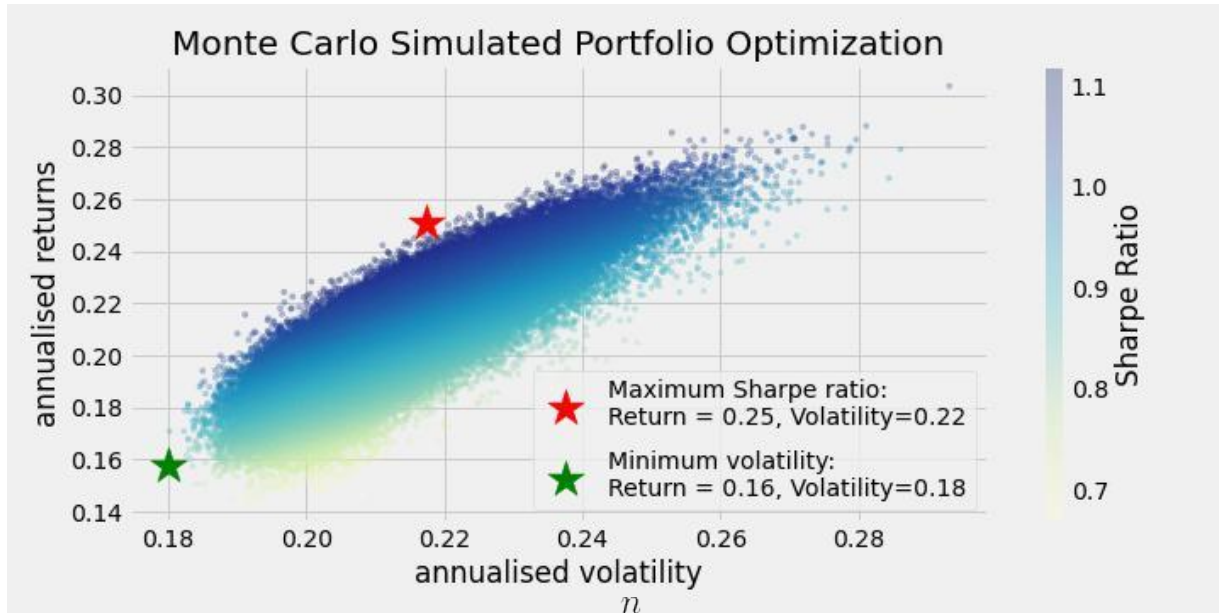
Pandora
 Ambu
 A.P. Møller-Mærsk (A)
 A.P. Møller-Mærsk (B)
 Demant
 ISS
 Danske Bank
 GN Store Nordic
 Netcompany Group
 Bavarian Nordic
 Jyske Bank
 Lundbeck



In hindsight it is easy!



Modern Portfolio Theory



$$Mean(\text{Portfolio Returns}) = \sum_{i=1}^n w_i r_i$$

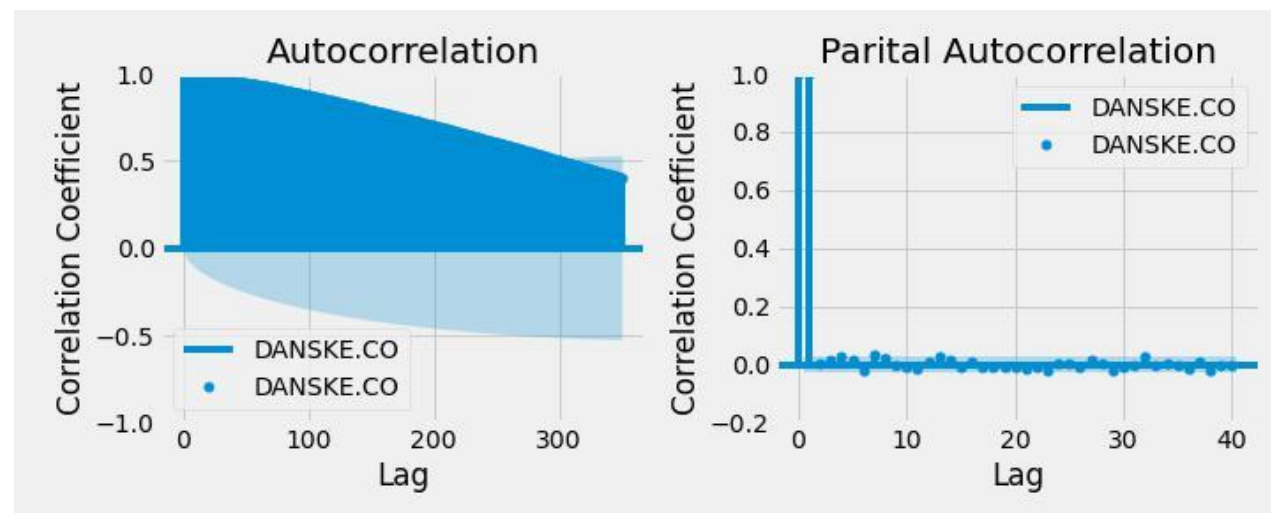
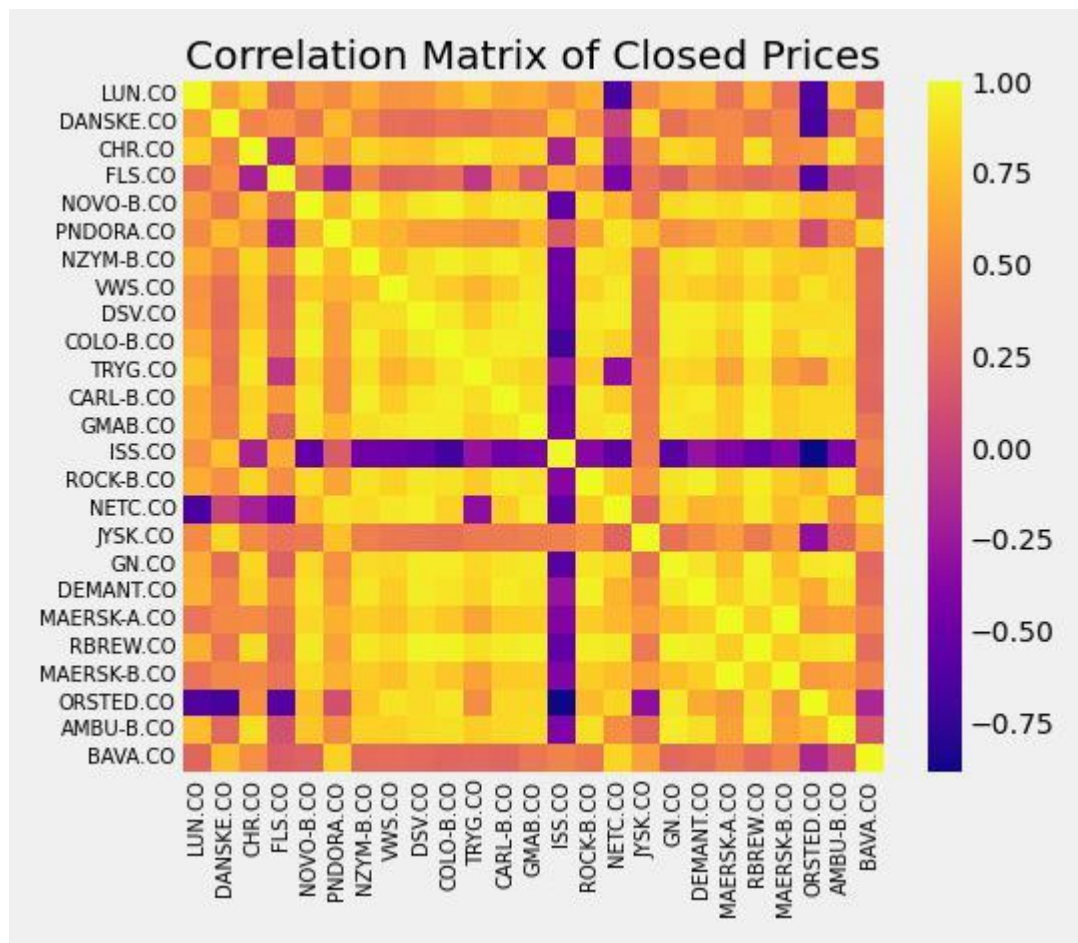
$$Var(\text{Portfolio Returns}) = \mathbf{w}^T \cdot (\mathbf{cov} \cdot \mathbf{w})$$

$$\text{Sharpe Ratio} = \frac{\text{Annual } Mean(\text{Portfolio Returns}) - \text{Risk Free Rate}}{\text{Annual } STD(\text{Portfolio Returns})}$$

$$\text{Minimum volatility} = \text{argmin}\{\text{Annual } STD(\text{Portfolio Returns})\}$$

$$\text{Maximum Sharpe Ratio} = \text{argmax}\{\text{Sharpe Ratio}\}$$

Correlations and Auto-Correlations



Autocorrelation:

No predictive ability after day 300

Technical Indicators

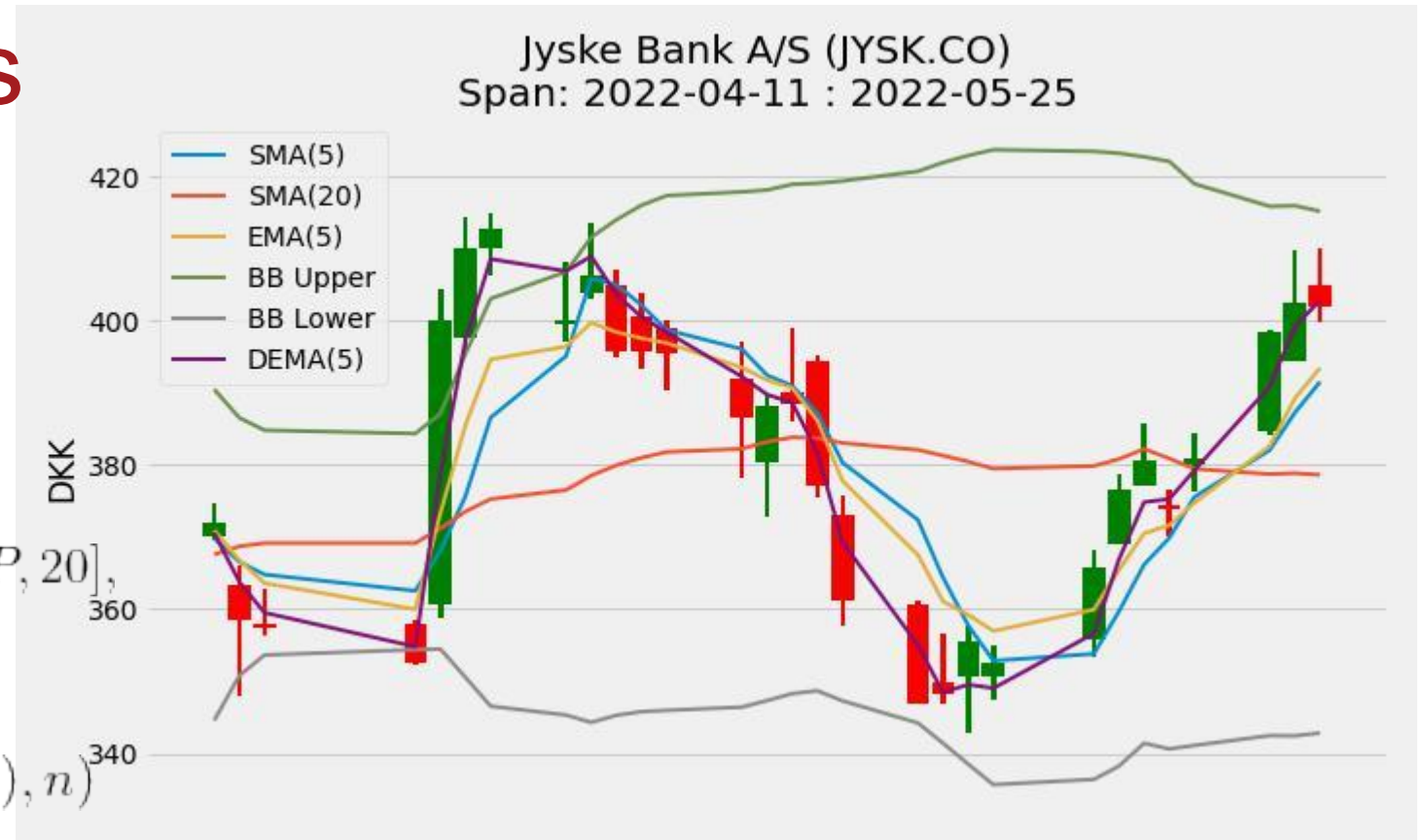
$$SMA(n) = \frac{A_1 + A_2 + \dots + A_n}{n}$$

$$\text{Bollinger Bands} = SMA(TP, 20) \pm 2\sigma[TP, 20],$$

$$TP = \frac{High + Low + Close}{3}$$

$$DEMA(n) = 2EMA(n) - EMA(EMA(n), n)$$

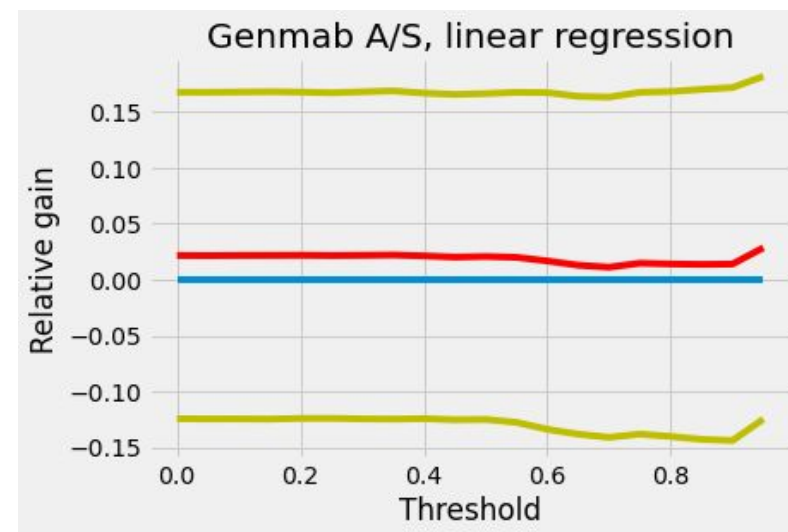
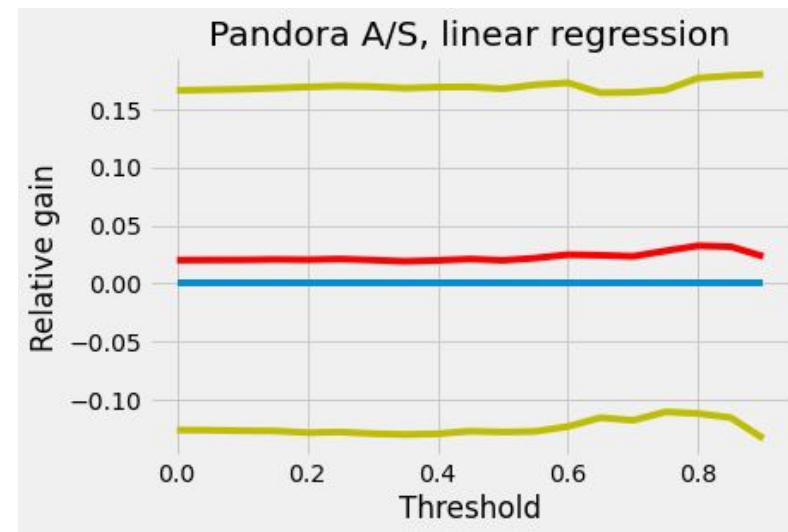
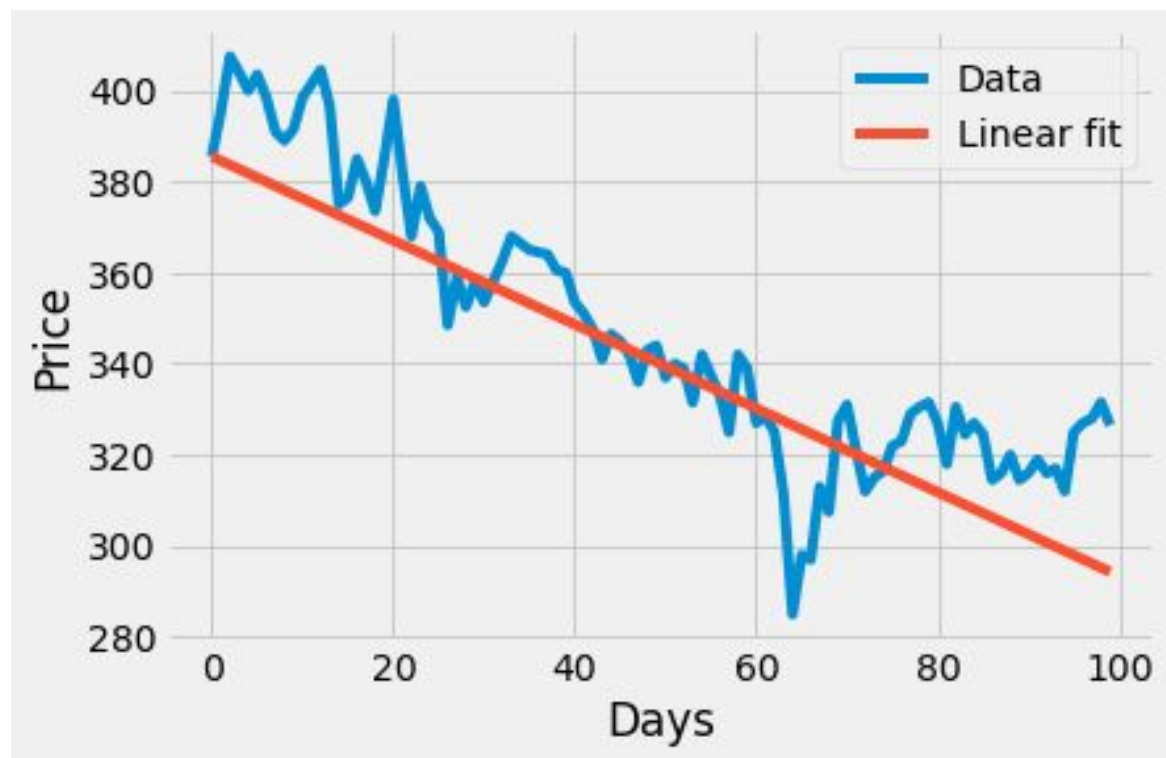
$$EMA(n) : EMA_{today} = Close_{today} \times \left(\frac{2}{1+n} \right) + EMA_{yesterday} \times \left(1 - \frac{2}{1+n} \right)$$



Full List of Technical Indicators

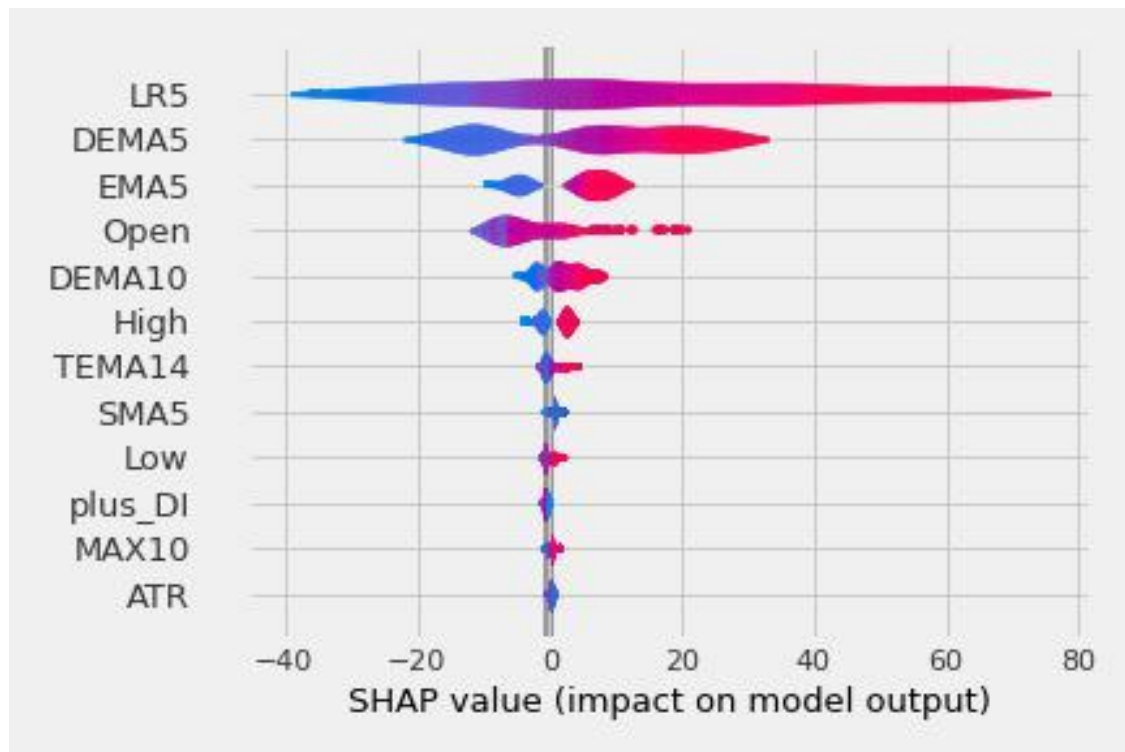
Acceleration Bands (ABANDS)	Linear Regression Intercept (LRI)	Parabolic Sar (SAR)
Accumulation/Distribution (AD)	Linear Regression Slope (LRM)	Simple Moving Average (SMA)
Average Directional Movement (ADX)	Moving Average Convergence	Standard Deviation (STDDEV)
Adaptive Moving Average (AMA)	Divergence (MACD)	Stochastic (STOCH)
Absolute Price Oscillator (APO)	Max (MAX)	Stochastic Fast (StochF)
Aroon (AR) Aroon Oscillator (ARO)	Money Flow Index (MFI)	T3 (T3)
Average True Range (ATR)	Midpoint (MIDPNT)	Triple Exponential Moving Average
Volume on the Ask (AVOL)	Midprice (MIDPRI)	(TEMA)
Volume on the Bid and Ask (BAVOL)	Min (MIN)	Triangular Moving Average (TRIMA)
Bollinger Band (BBANDS)	MinMax (MINMAX)	Triple Exponential Moving Average
Band Width (BW)	Momentum (MOM)	Oscillator (TRIX)
Commodity Channel Index (CCI)	Normalized Average True Range (NATR)	Time Series Forecast (TSF)
Chande Momentum Oscillator (CMO)	On Balance Volume (OBV)	TT Cumulative Vol Delta (TT CVD)
Double Exponential Moving Average	Price Channel (PC)	Ultimate Oscillator (ULTOSC)
(DEMA)	Percent Price Oscillator (PPO)	Volume At Price (VAP)
Directional Movement Indicators (DMI)	Price Volume Trend (PVT)	Volume (VOLUME)
Exponential Moving Average (EMA)	Rate of Change (ROC)	Volume Delta (Vol Δ)
Fill Indicator (FILL)	Rate of Change (ROC100)	Volume Weighted Average Price
Ichimoku (ICH)	Rate of Change (ROCP)	(VWAP) Williams % R (WillR)
Keltner Channel (KC)	Rate of Change (ROCR)	Weighted Moving Average
Linear Regression (LR)	Relative Strength Indicator (RSI)	(WMA)
Linear Regression Angle (LRA)	Session Volume (S_VOL)	Welles Wilder's Smoothing Average
		(WWS)

Linear regression

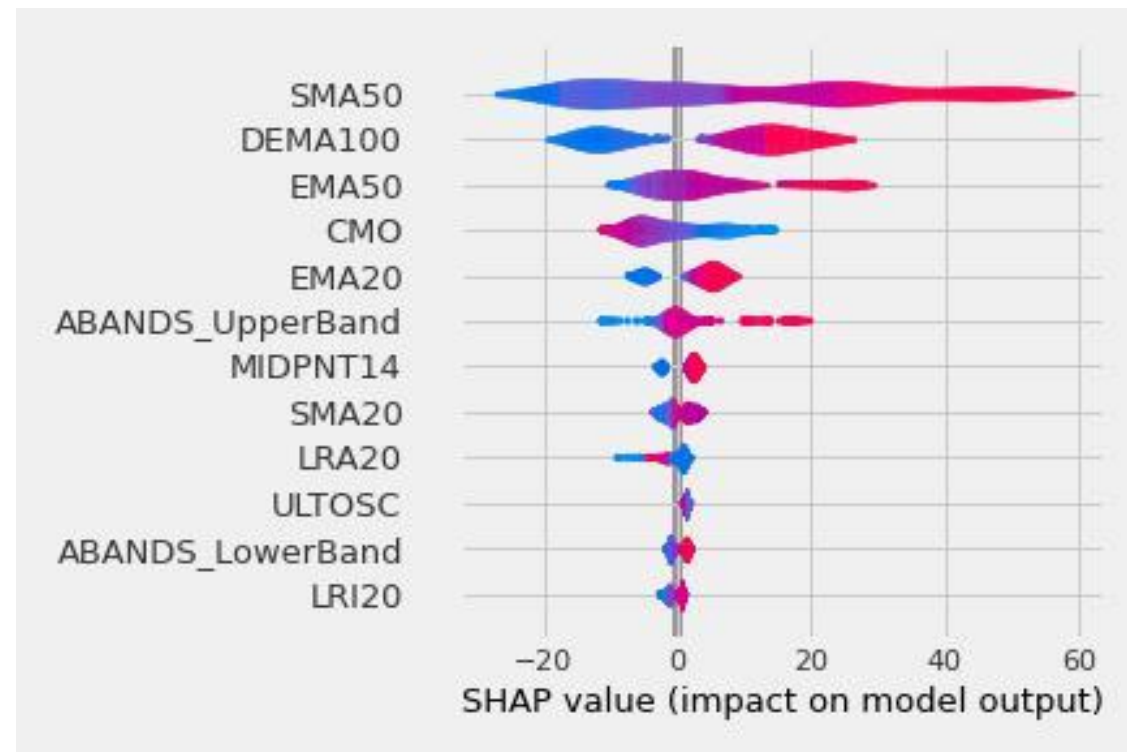


SHAP with LightGBM (Technical Indicators)

Predicting 1 day ahead (Jysk)

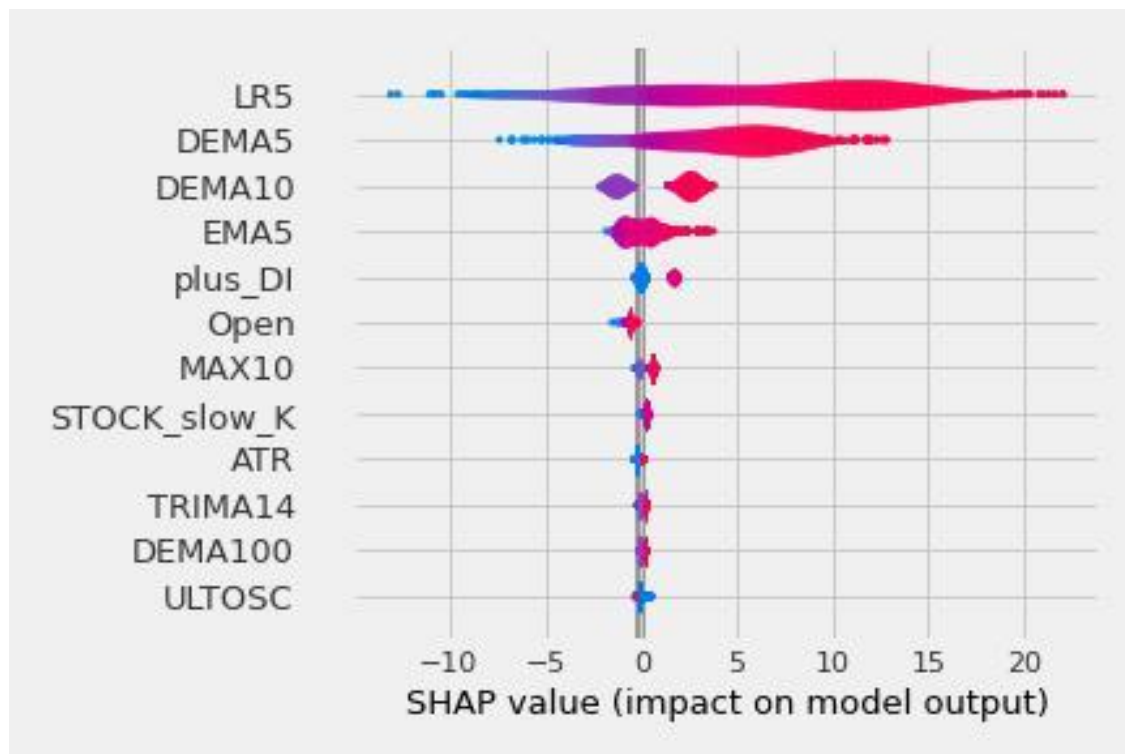


Predicting 20 day ahead (Jysk)

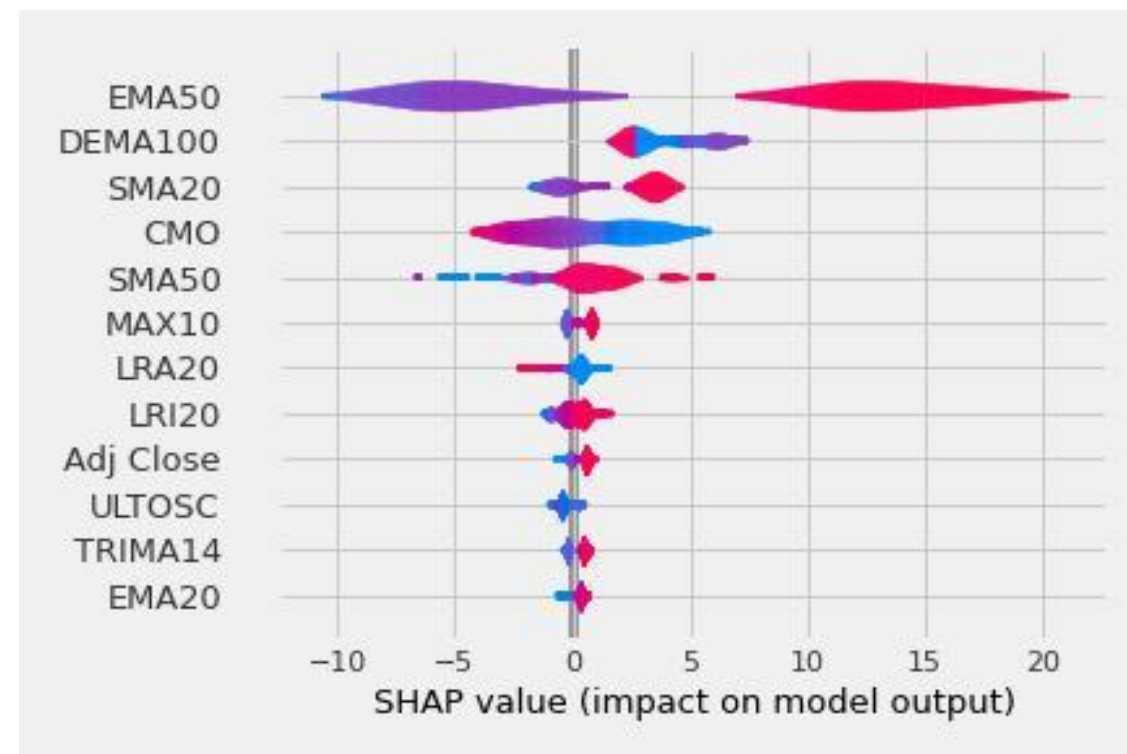


SHAP, Danske Bank

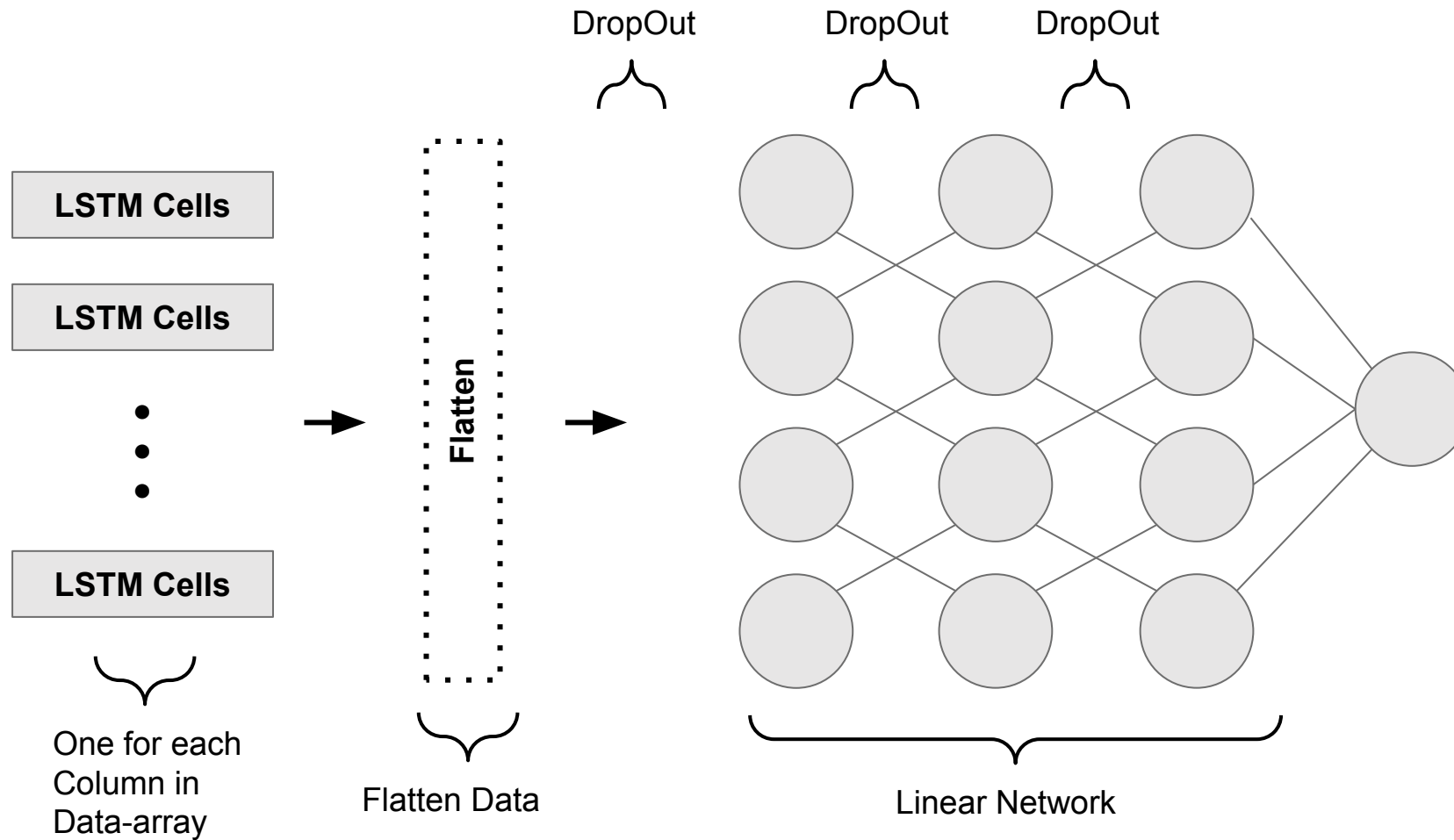
Predicting 1 day ahead (Danske Bank)



Predicting 20 day ahead (Danske Bank)



PyTorch LSTM Network



Params:

- LogLoss
- Cross-validation
- Shuffle
- Batches

Output:
Probability of price
increasing in m days

Metrics

Testing Metrics

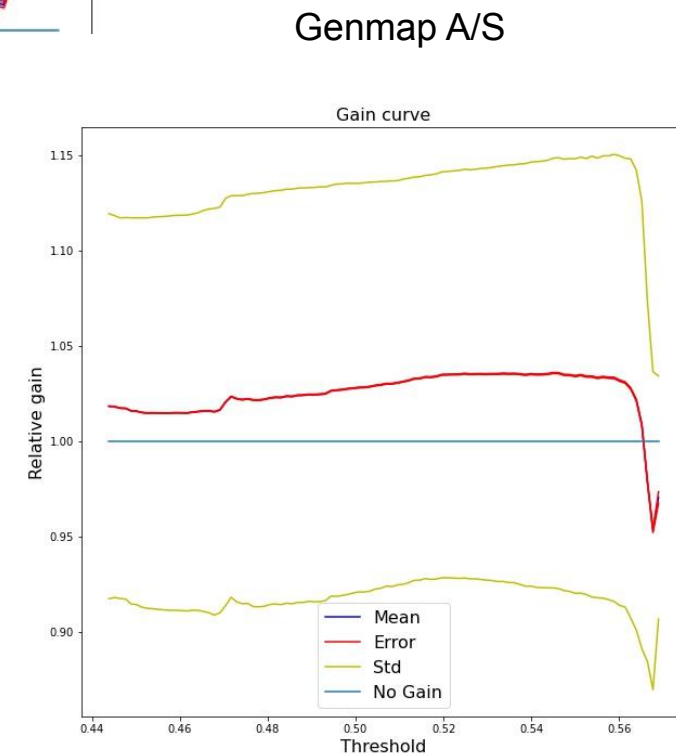
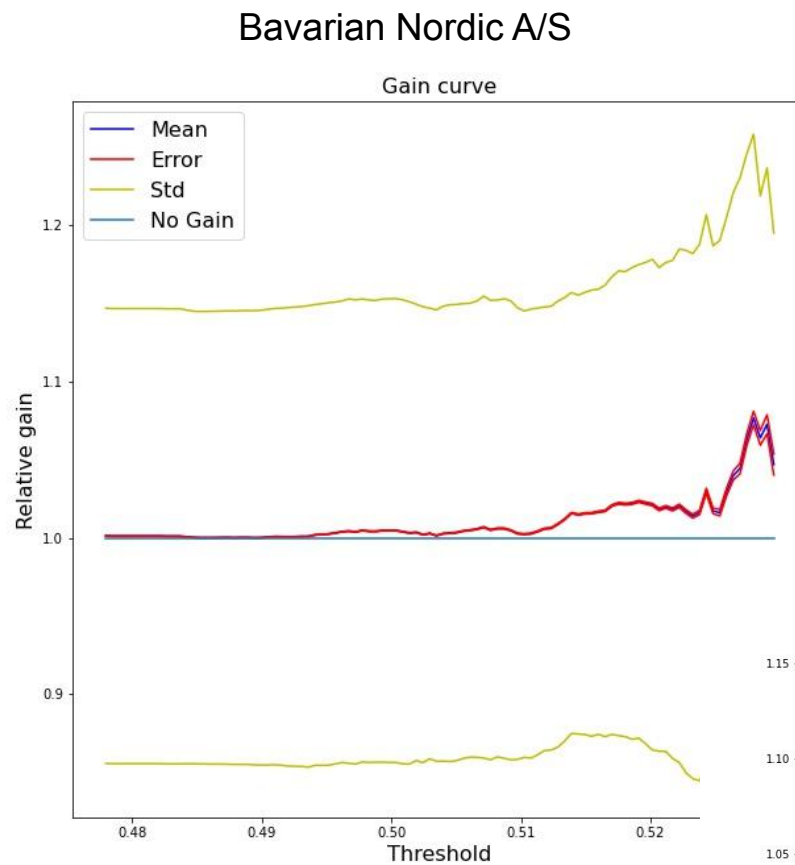
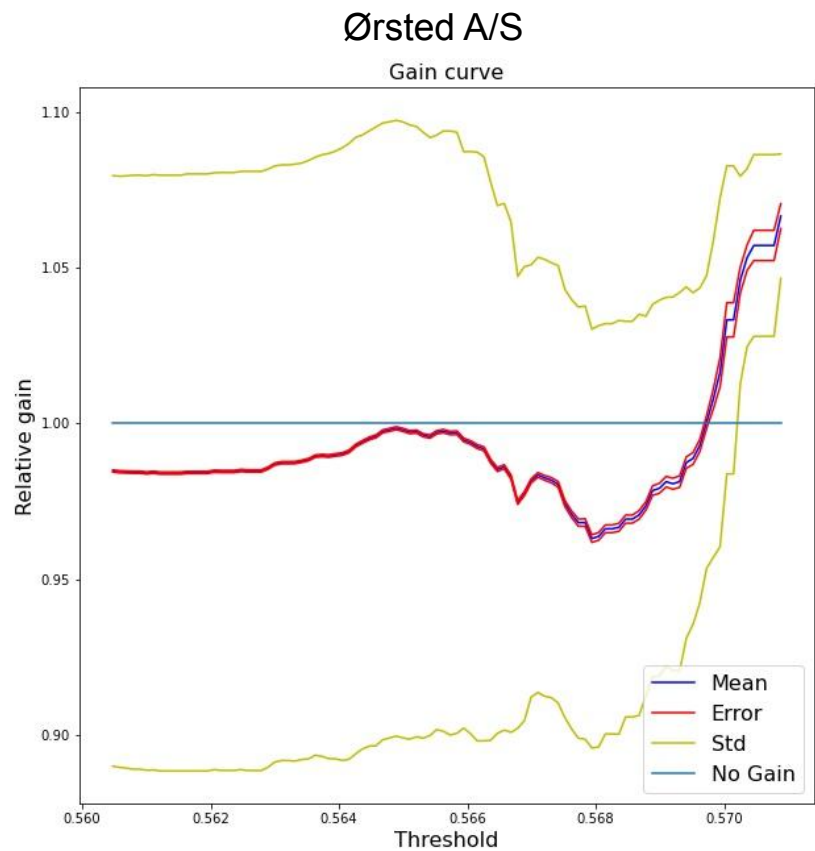
- LogLoss
- AUC
- Accuracy
- MaxGain
- MaxGain Probability

$$\text{Gain}(t) = \text{mean}(\text{Return}\{P(\text{UpTick}) > t\})$$

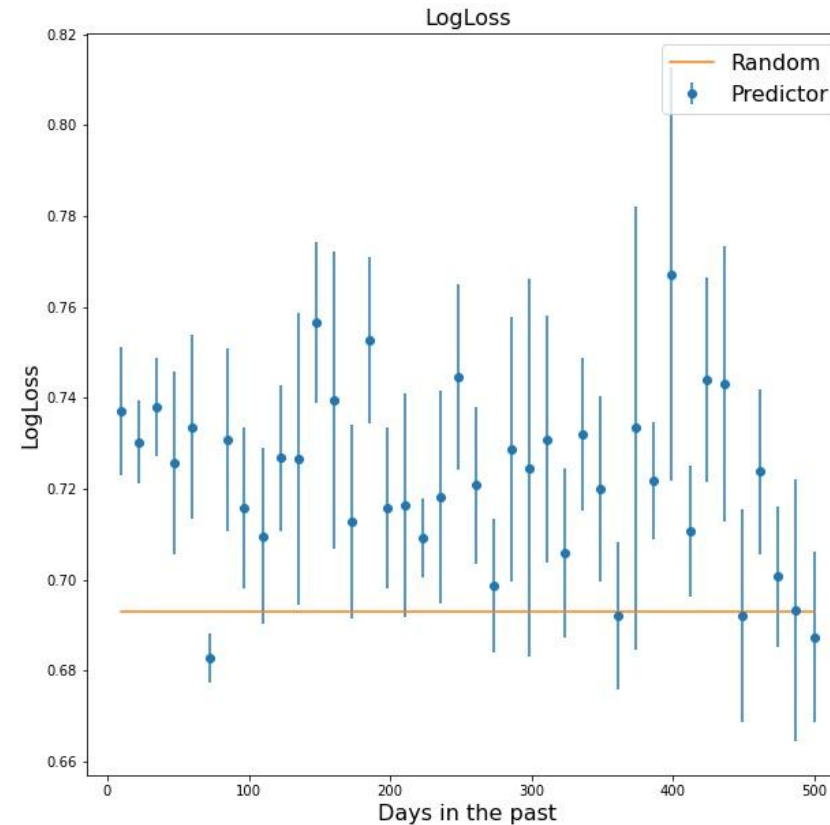
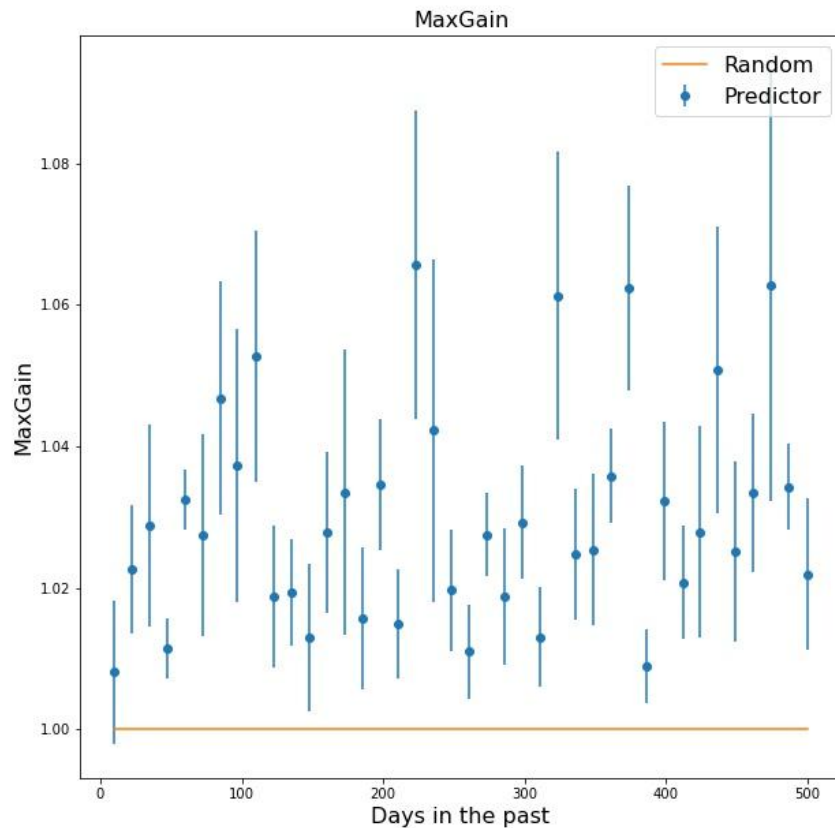
$$\text{MaxGain} = \text{Max}\{\text{Gain}(t), t\}$$

$$\text{MaxGainProbability} = \text{Max}\left\{\frac{\text{Gain}(t) - 1}{\text{STD}(\text{Gain}(t))}, t\right\}$$

Gain curves



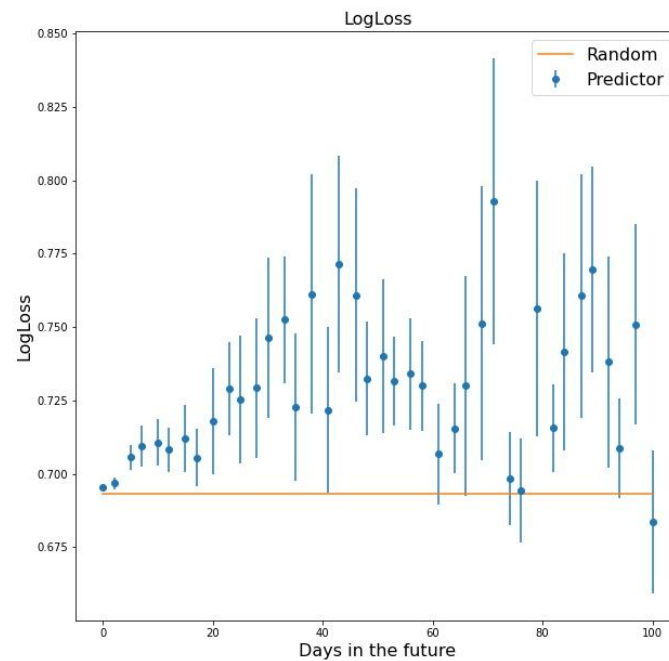
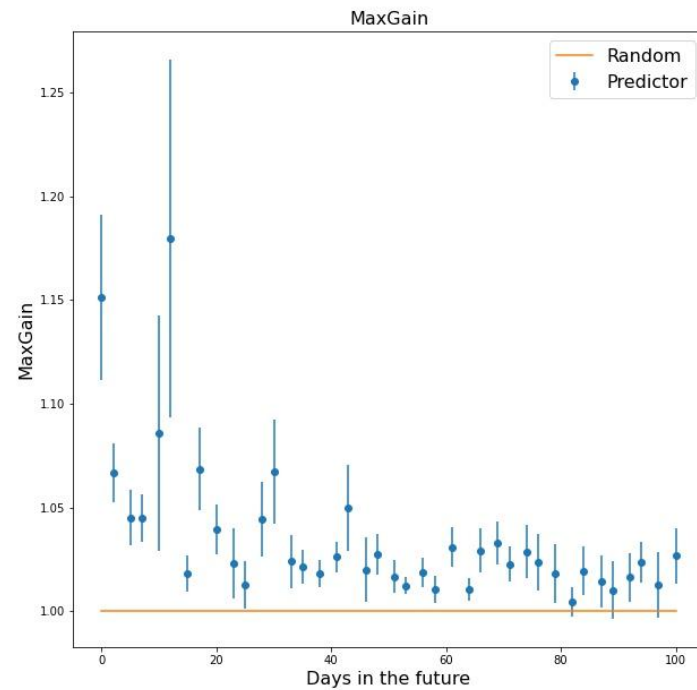
Parameter optimization: Days in the past



Nr. of days in the past to base prediction on

Used 5-fold cross validation to get errors

Parameter optimization: Days in the future

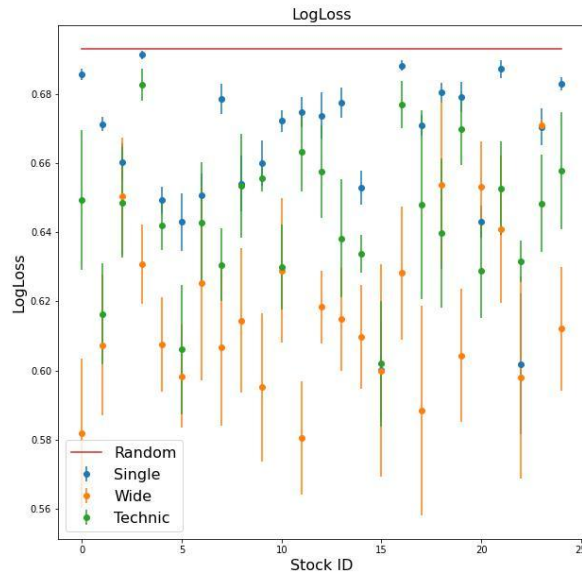


Nr. of days into the future to predict

Used 5-fold cross validation to get errors

Data types

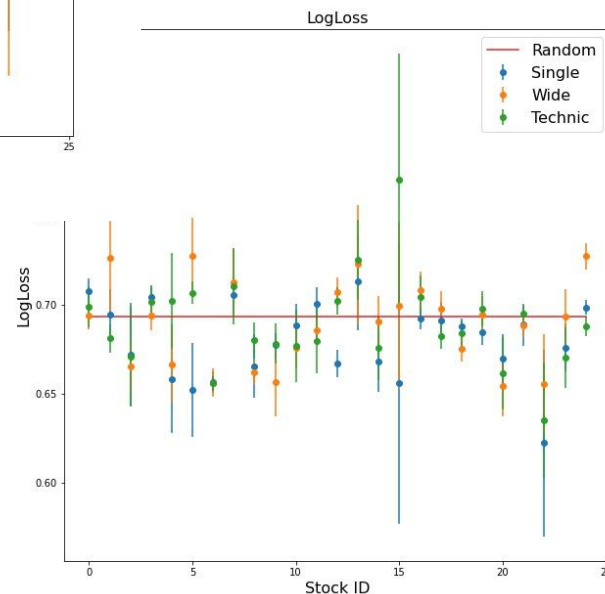
Shuffled data



Wide is best for shuffled data

Essentially the same for ordered data, technical indicators are a bit better

Ordered data



Different data for prediction:

- Single: Stock data for the stock to predict
- Wide: Stock data for all stocks to predict one of them
- Technic: Stock data and technical indicators for the stock to predict

Failed attempts

Training on all stocks

idea:

- More data to learn trends from
- Better chance of recognizing critical events

Reality:

- Stocks are too different

Using technical indicators for all stocks to predict one stock value

Idea:

- Using all stocks worked well
- Technical indicators are useful

Reality:

- Too much data ($68 \times 25 \times 20 = 34000$ inputs to linear network): Unable to train

Evaluation

Model	Performance	Training
LSTM network	Good performance	Slow and sometimes unreliable
Linear regression	Slightly better than random	Super fast
Support vector machine	Random	Fast
EchoState network	Bad (did not work)	Fast

More things to be done..

- SHAP values on PyTorch
- Get errors on model via bootstrapping
- Get optimal threshold for all stocks
- Optimize Hyperparameters
- Incorporate Stock-Portfolio Allocations on the go
- Shadow-Trading with optimal model

Appendix

Details and workflow of the projects

The following were done to obtain the objective: Make an ML algorithm to get as much return on the OMXC25 Stocks as possible:

We tried the following:

Making an LSTM-NN which could predict if future stocks would increase in value. Here we first found the number of days into the future that we wanted to predict as well as how many days in the past should be used for each prediction. Then we used the standard OHLCV data for each stock, compared that to using all data (WIDE) from all of the OMXC25 index. We then calculated +60 Technical Indicators which we also used in addition.

The idea was that if the model is sure (above some threshold) that a stock will go up, then buy and sell later. This buy-and-hold strategy was our way of checking how much “money” we have earned.

There are loads of studies on stock market behavior, and using Modern Portfolio Theory (MPT) we cross-referenced what we found with the ML to see if our machine learning algorithm actually made sense.

This was the case since MPT told us, that all stocks in the OMXC25 index were correlated (except ISS) and that made sense when the WIDE dataset were giving better result than just the pure OHLCV data. In addition the MPT told us that Ørsted is a good choice for good returns, both with the SHARPE-ratio strategy and with the minimum volatility strategy. This was true when looking at the result for MaxGain graph for Ørsted which was really good.

We also tried to look at a “simple” linear fit-model where we used a linear fit to predict future values.

The Dataset were really large when adding technical indicators and it was hard to tell which of them were the best. There we used a lightGBM model to see the best performing technical indicators. However we wanted to do it on PyTorch, but this wouldn't work.

Another thing we have worked on is to incorporate Echo-States into the algorithm, however here we have a lot of problems and it was hard to get any good predictions.

Roc- and gain curves

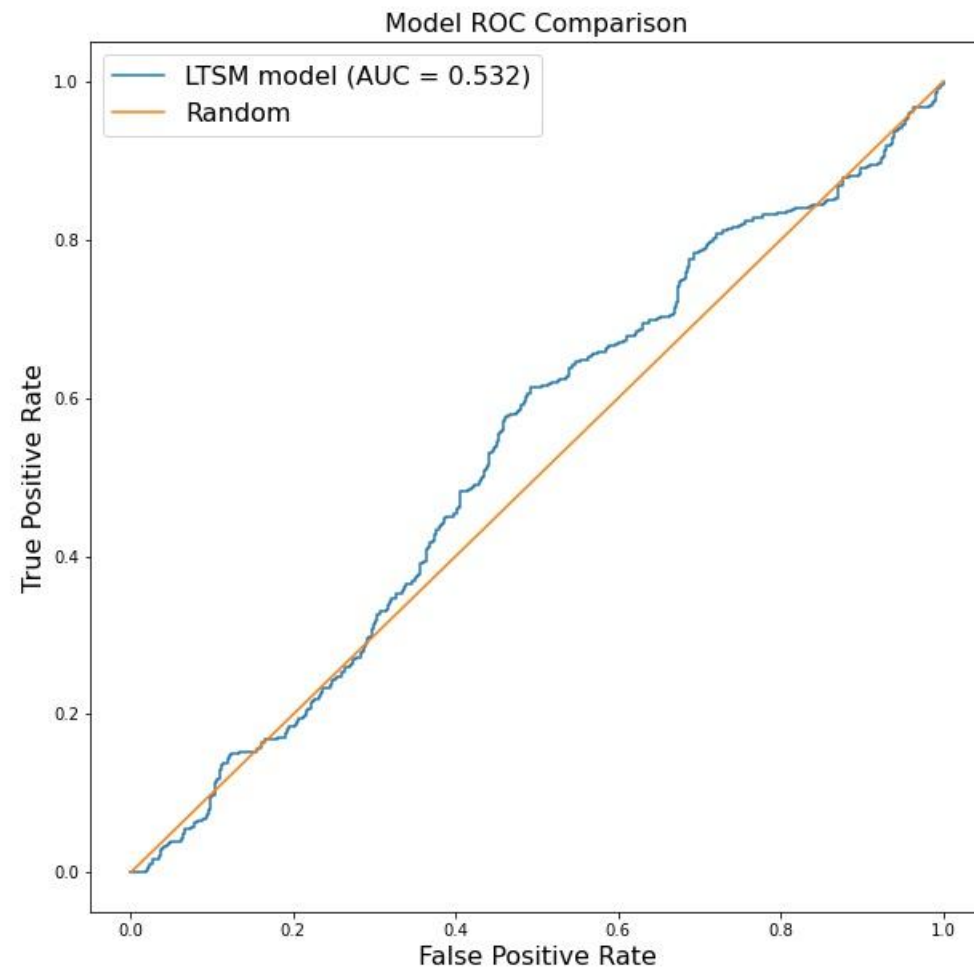
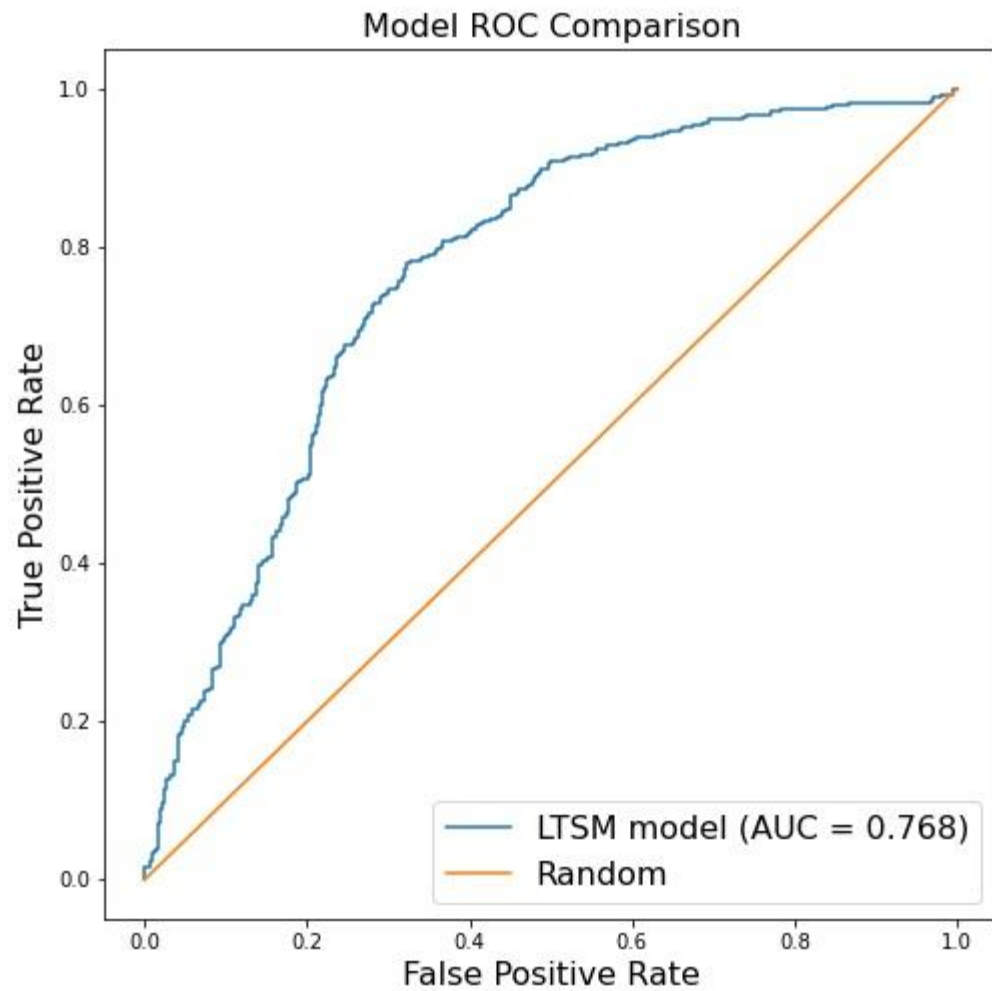
Left: Shuffled data

- Training is done on the entire dataset for random points used for testing

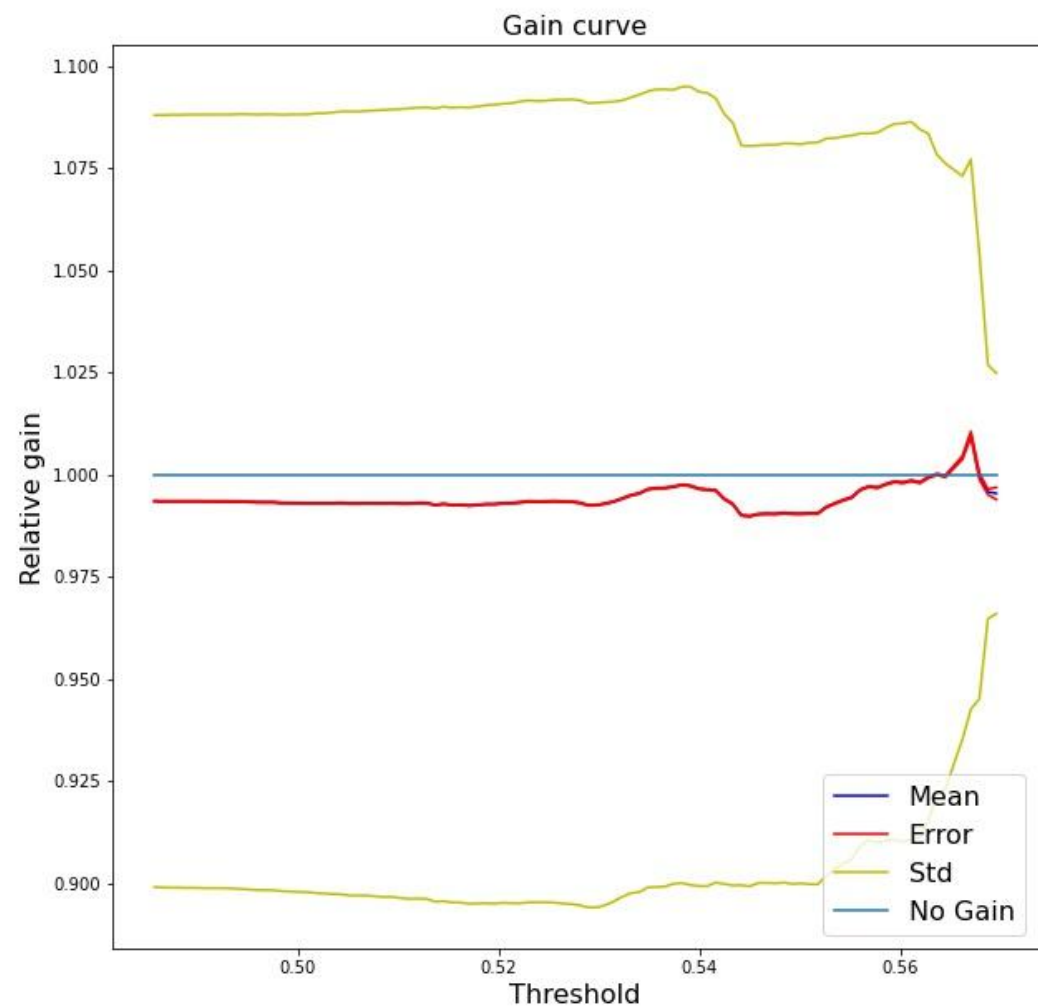
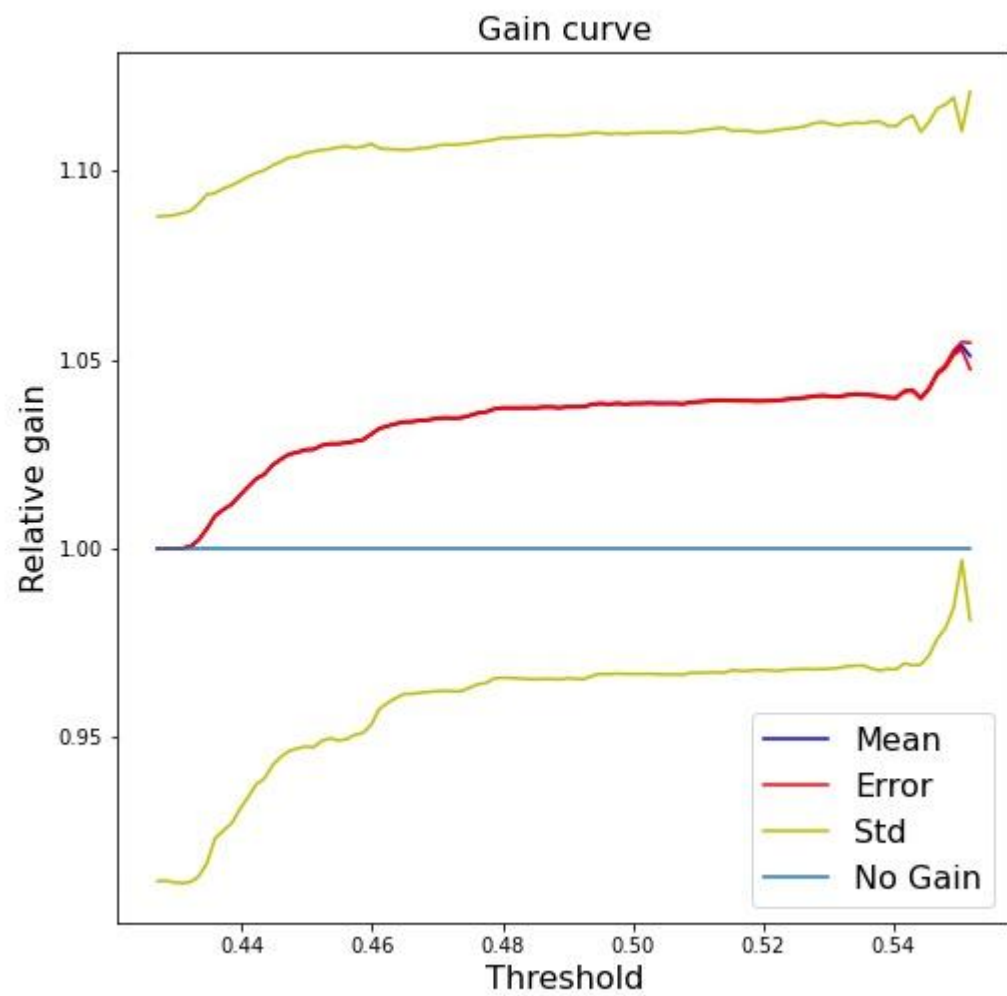
Right: Ordered data

- Training is done up until some time, testing is done on the rest of the dataset

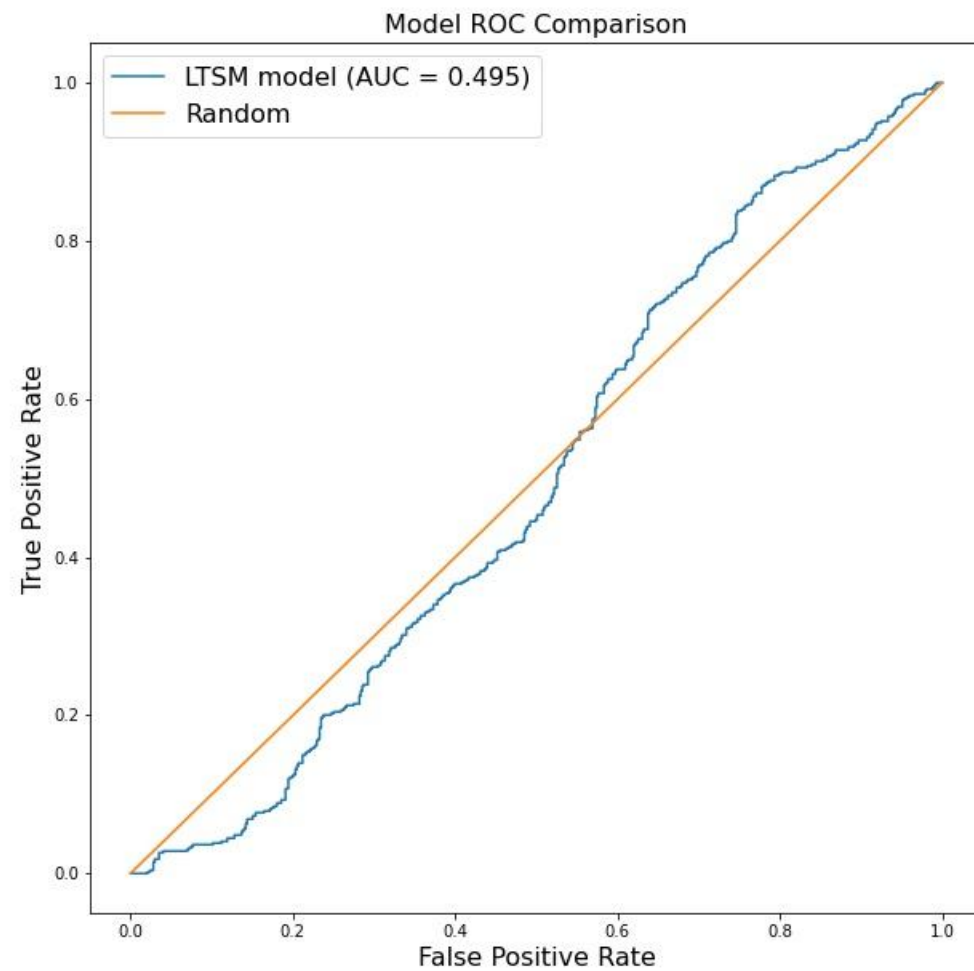
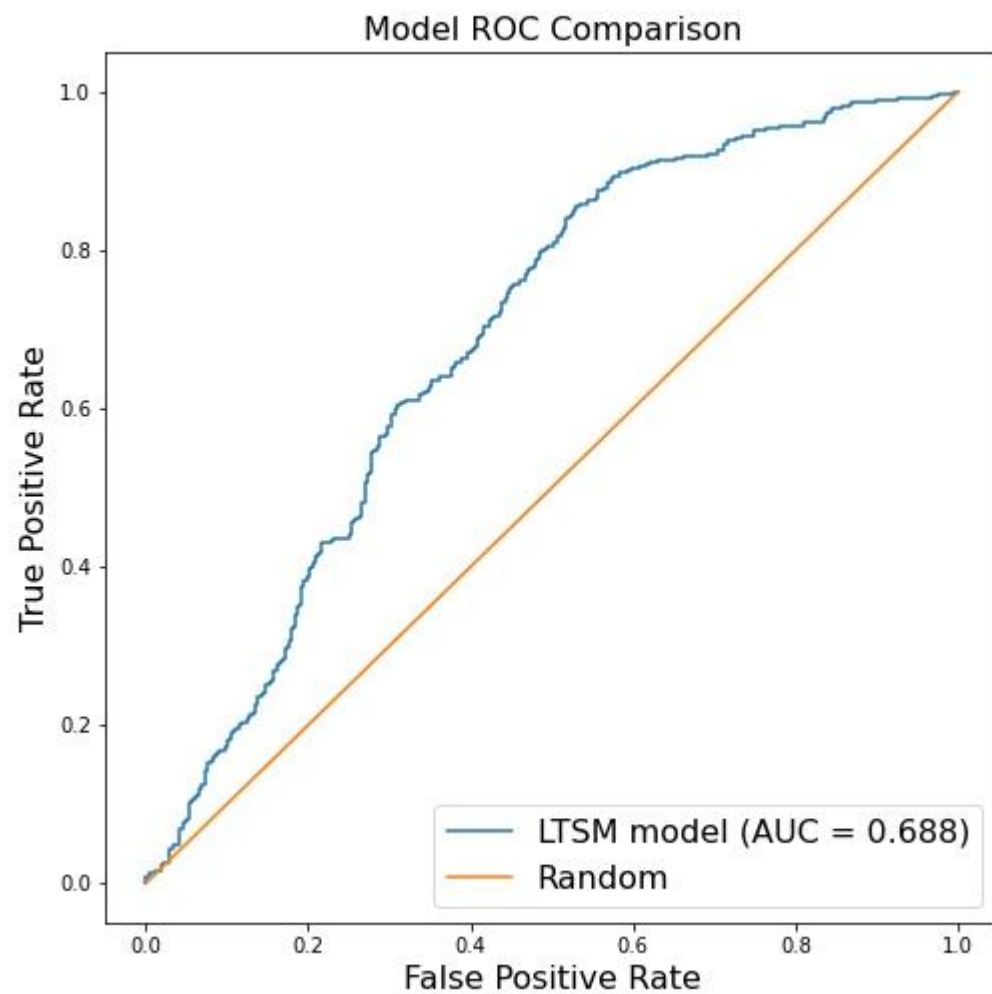
H. Lundbeck A/S



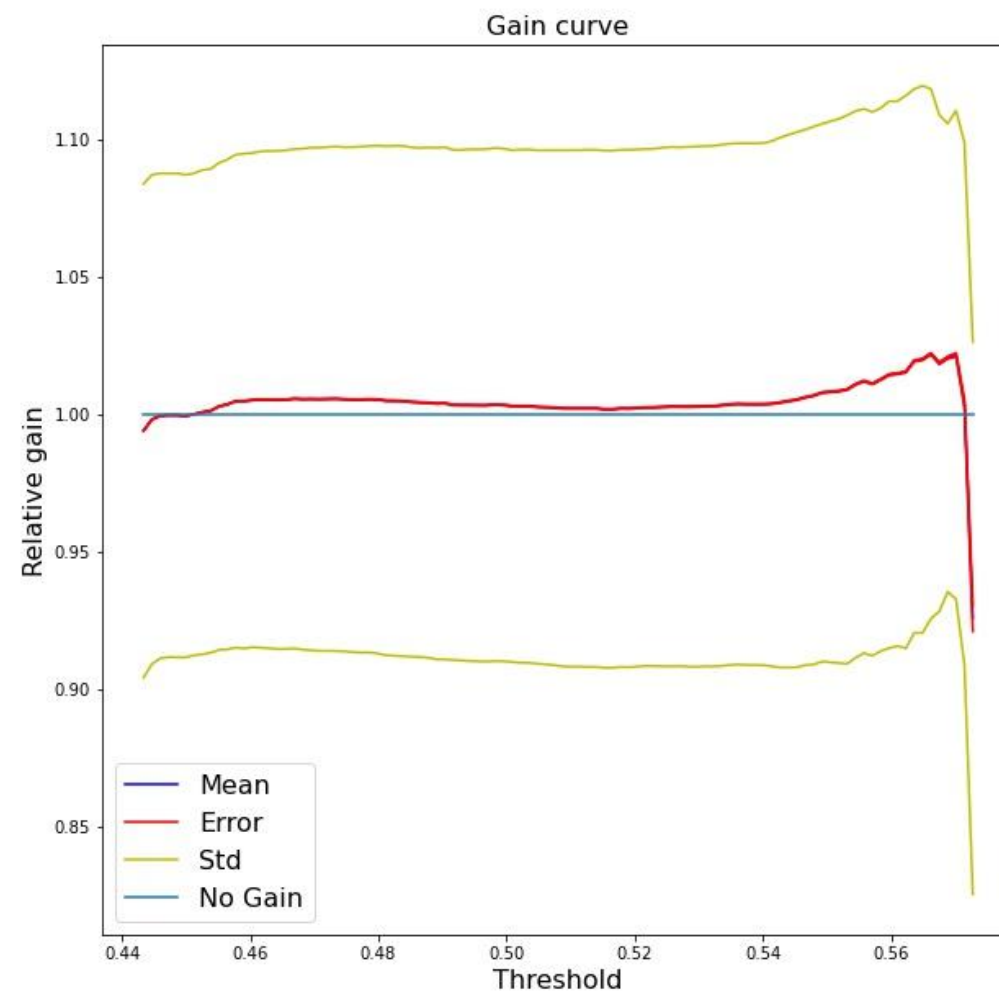
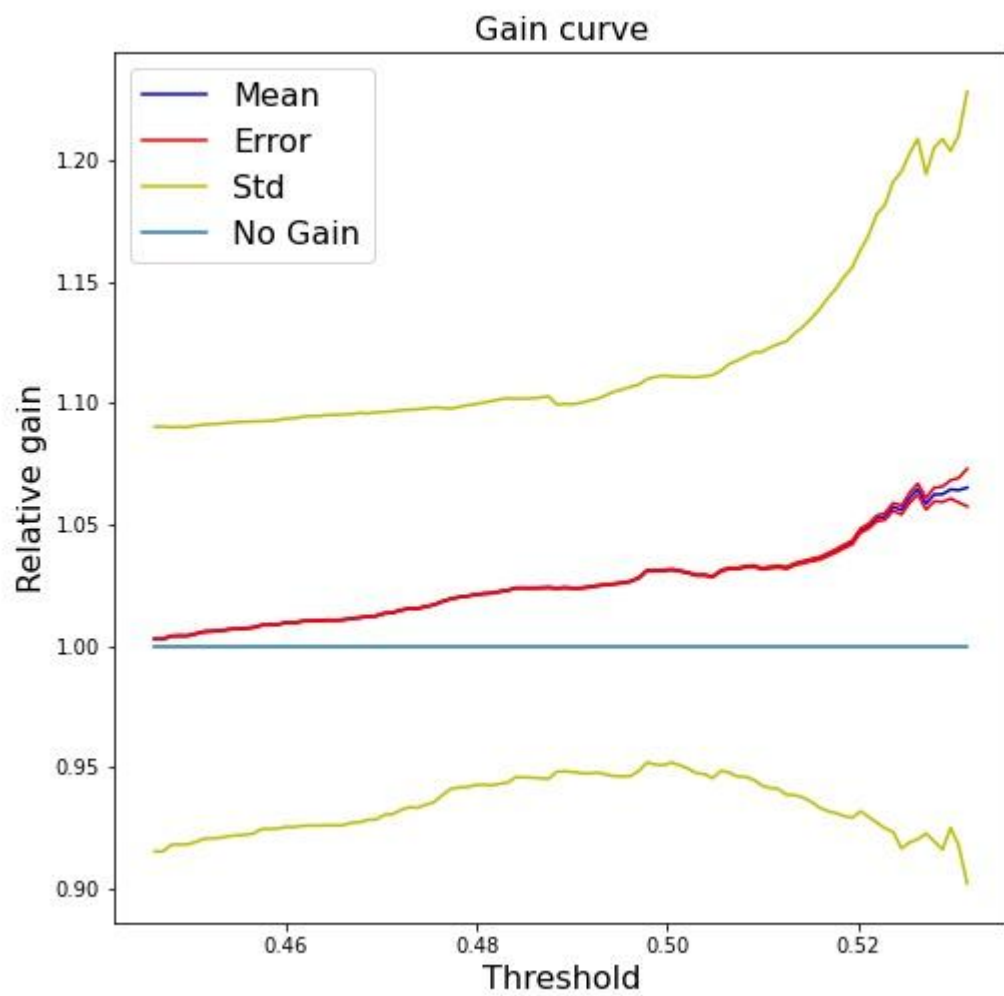
H. Lundbeck A/S



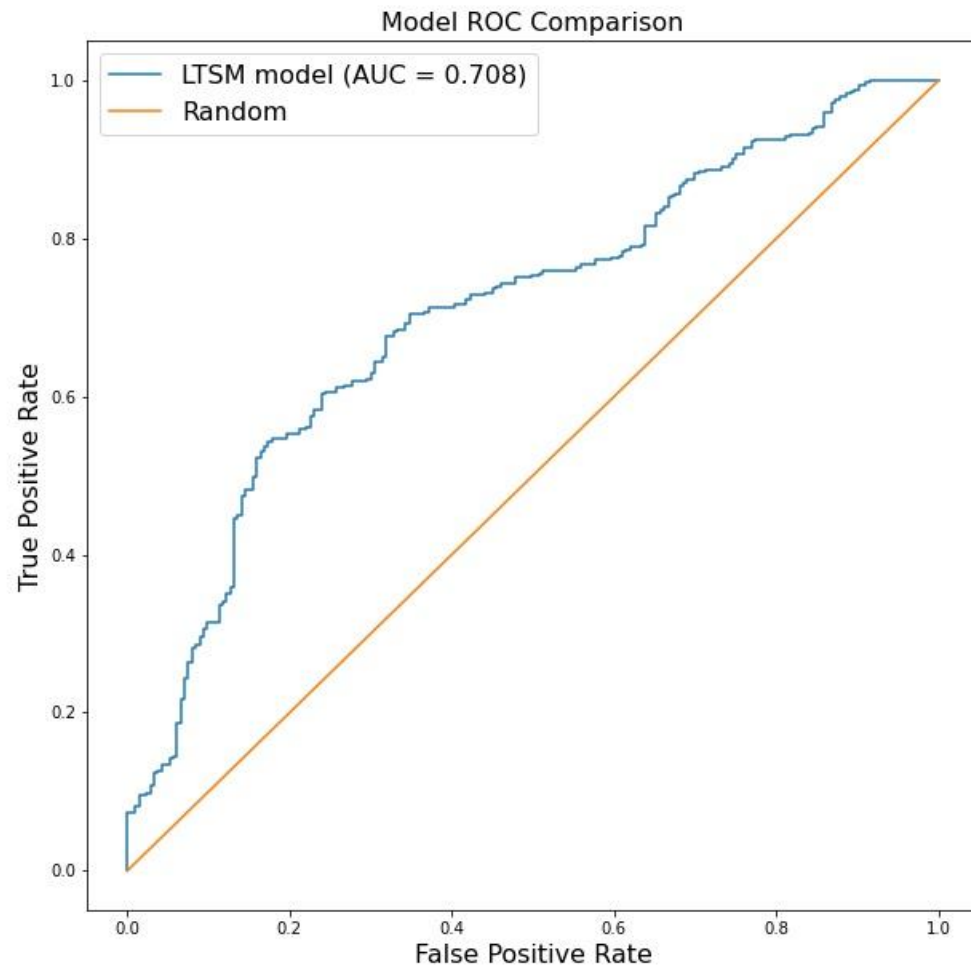
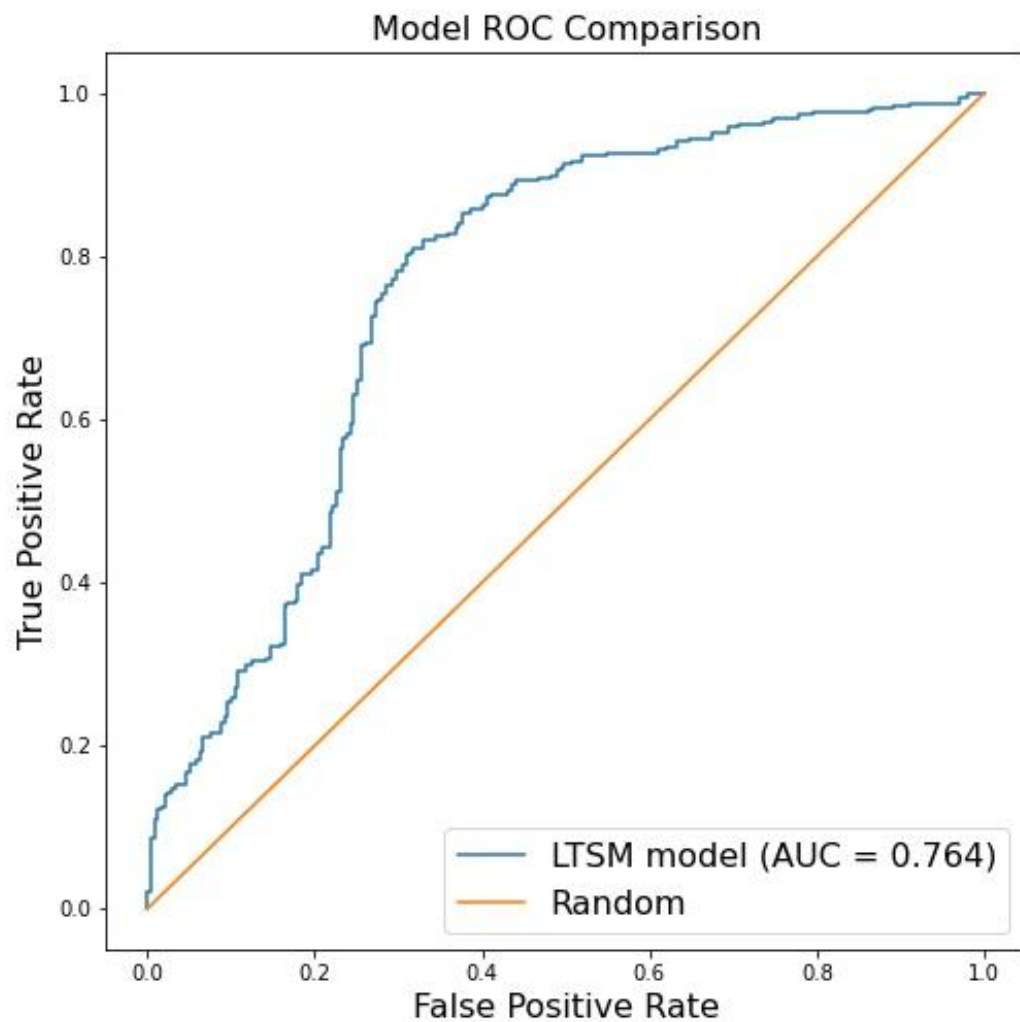
Danske Bank A/S



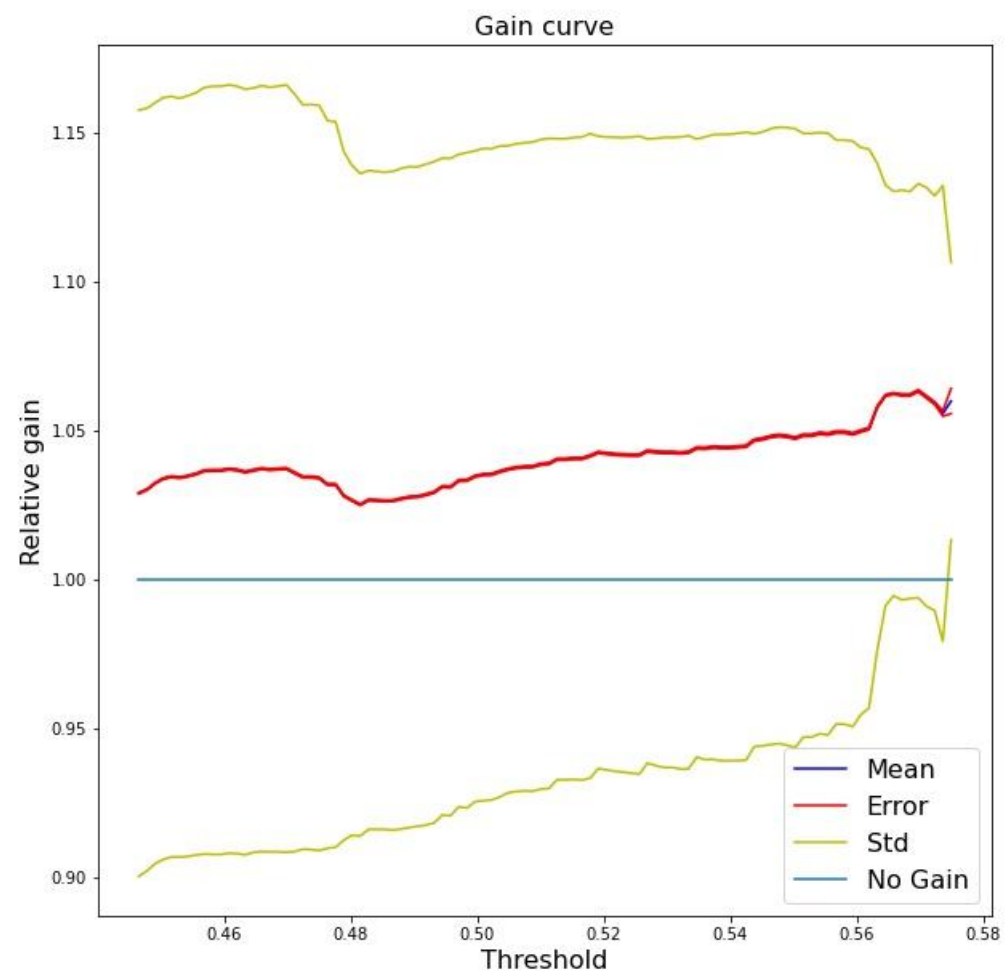
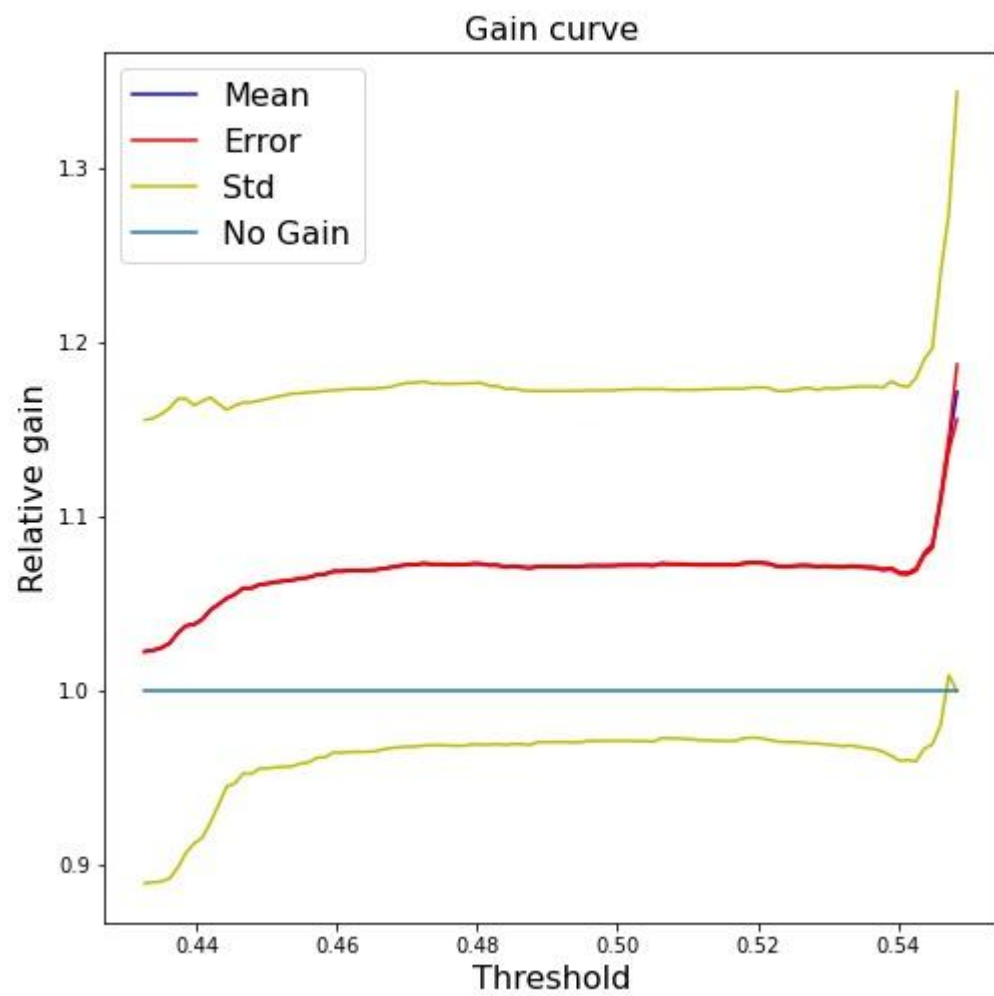
Danske Bank A/S



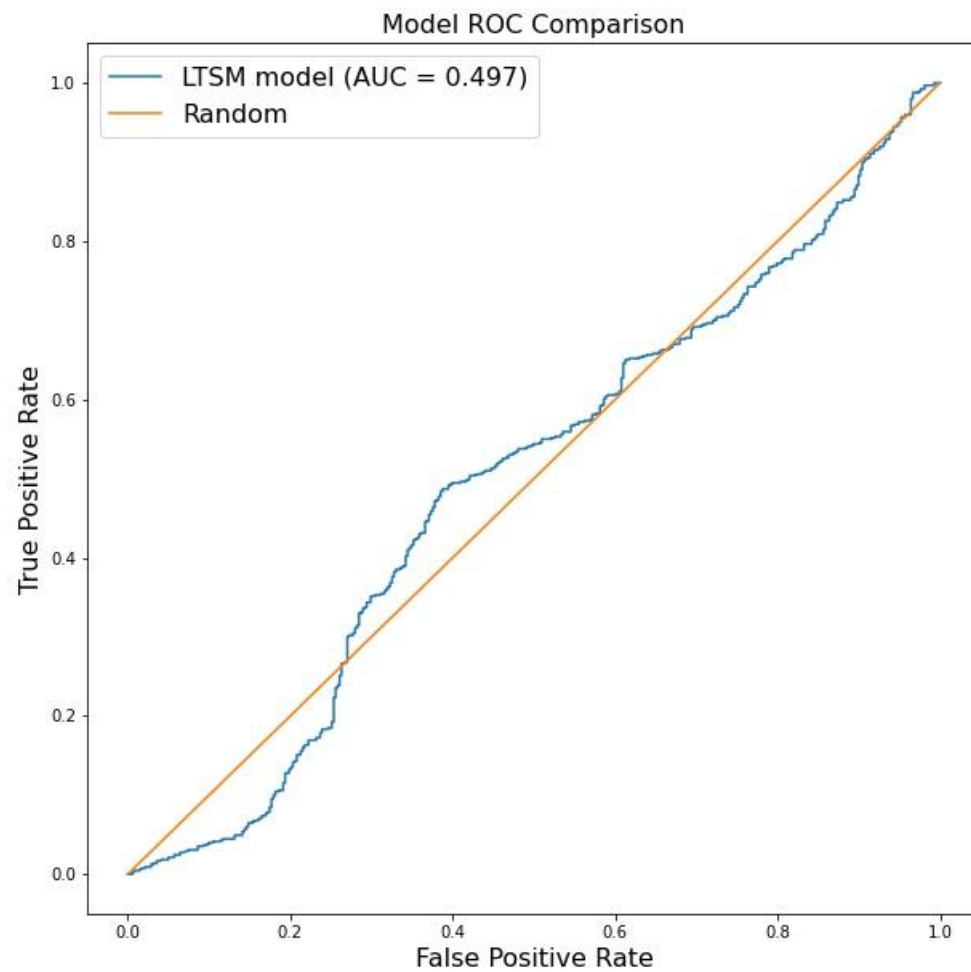
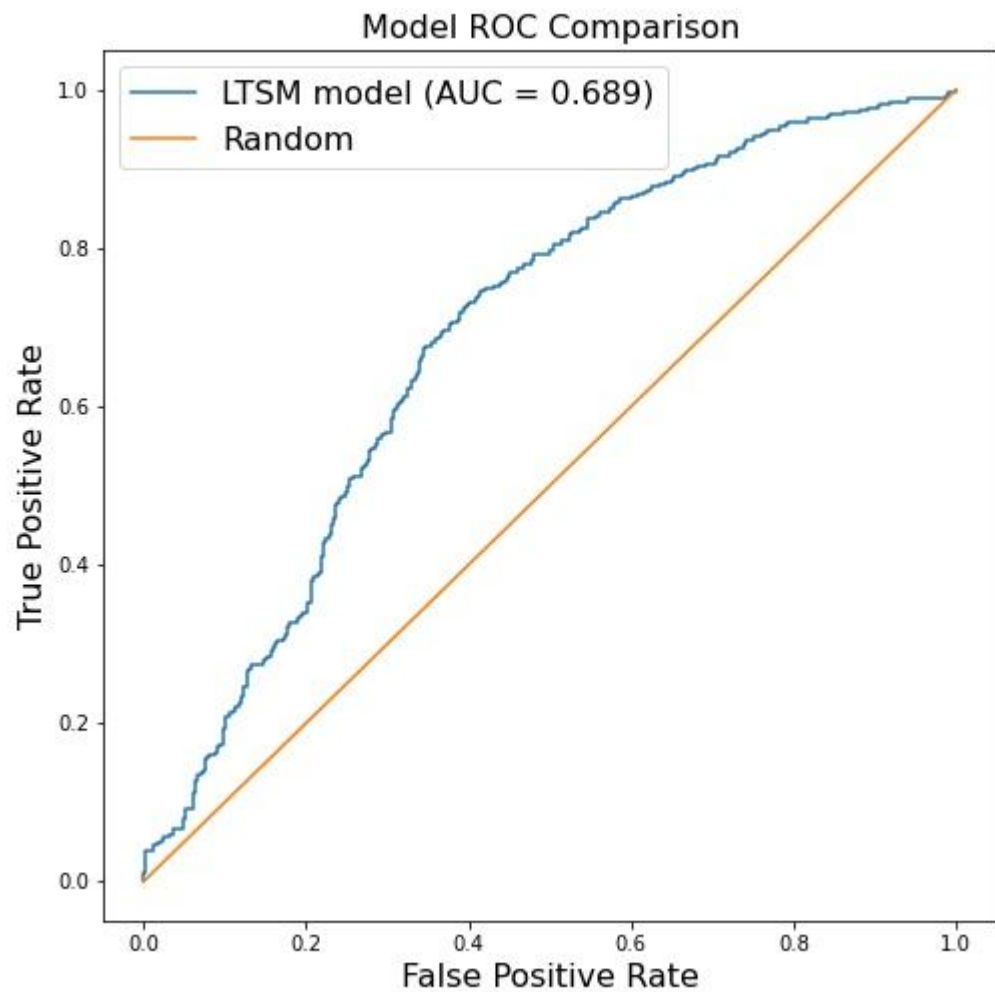
Pandora A/S



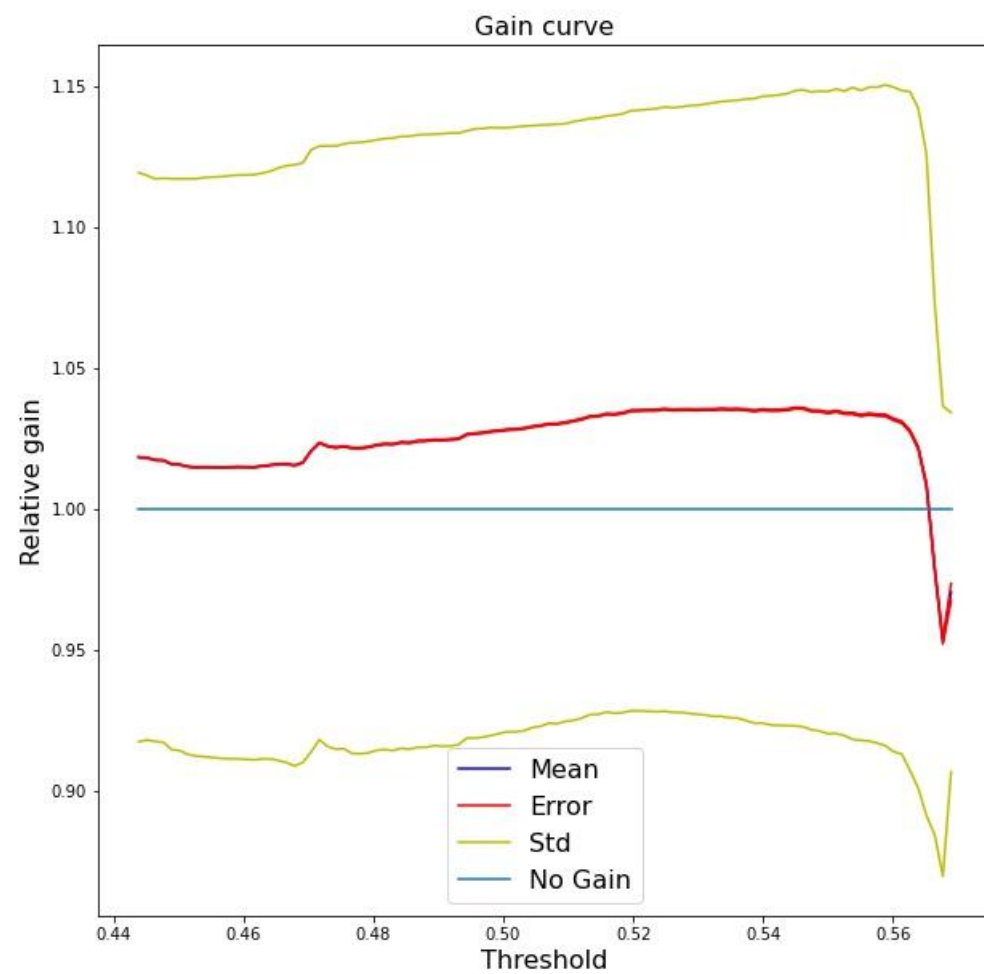
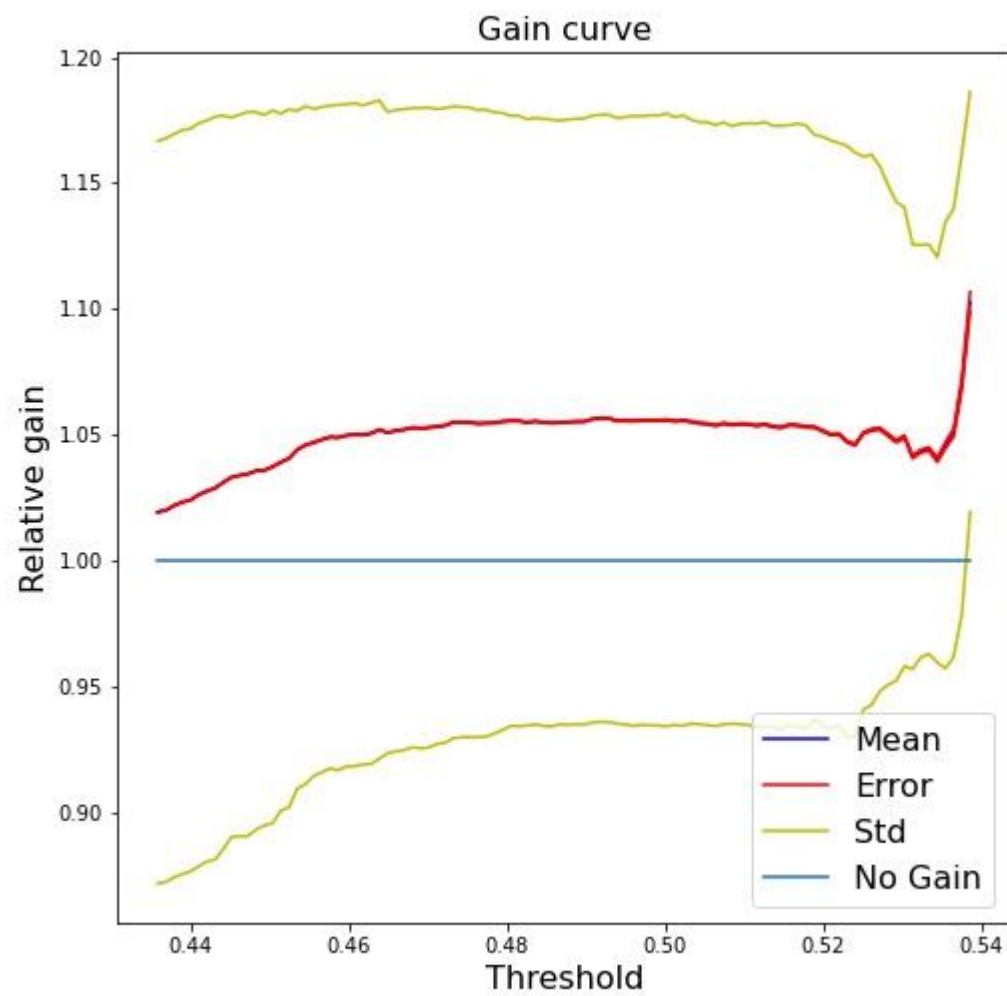
Pandora A/S



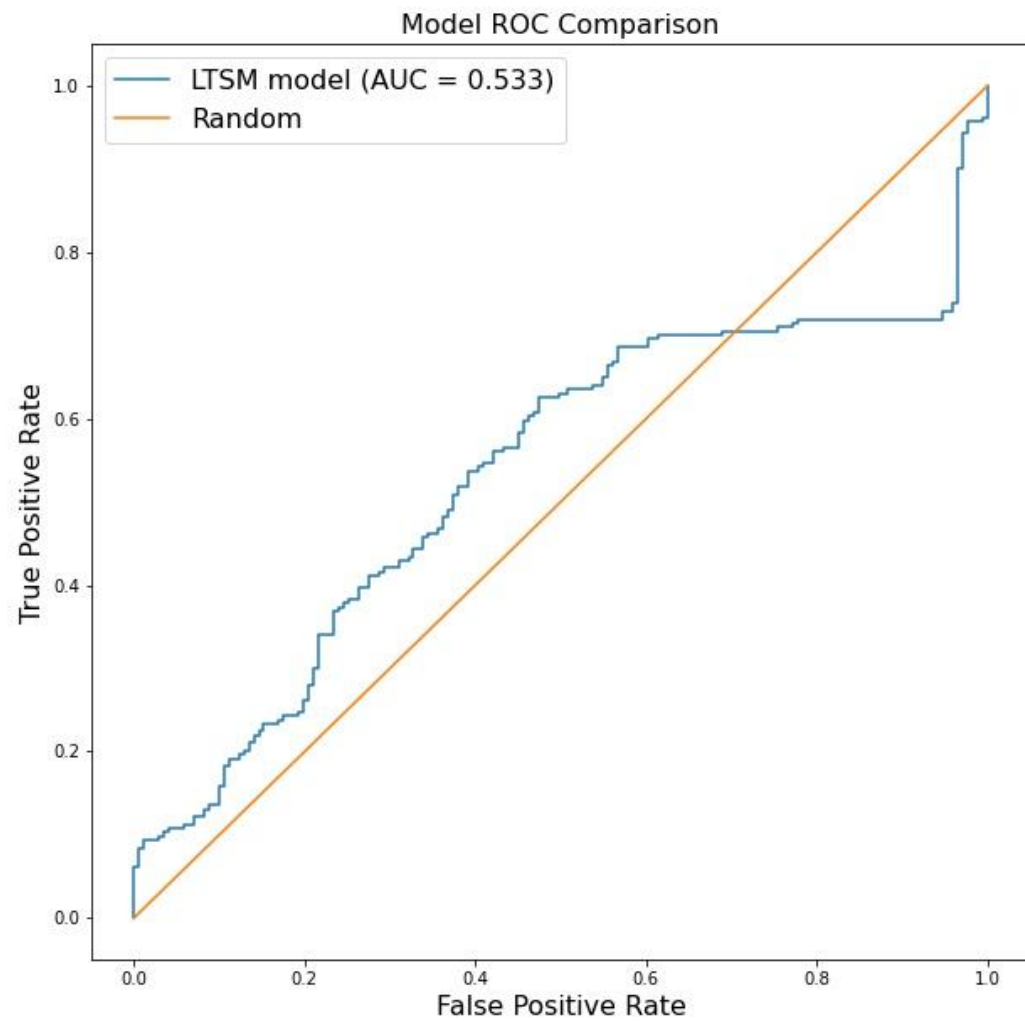
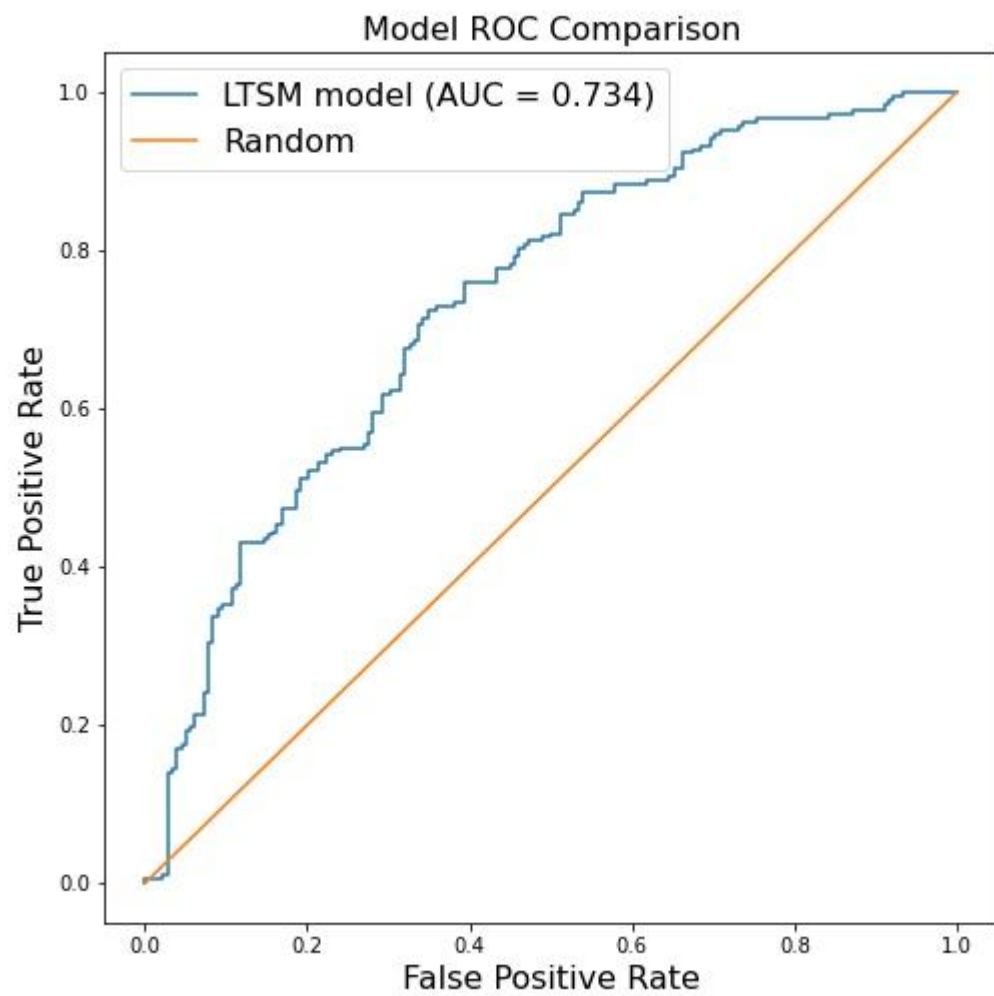
Genmab A/S



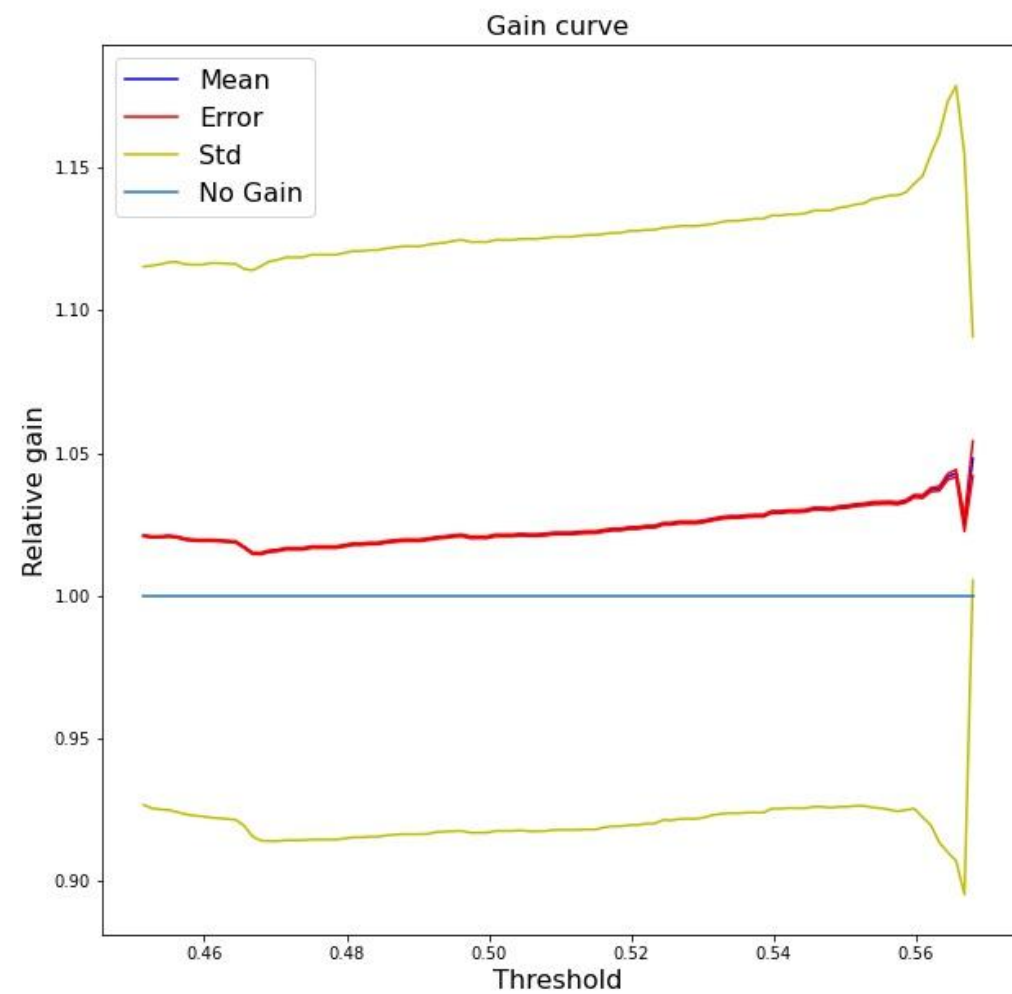
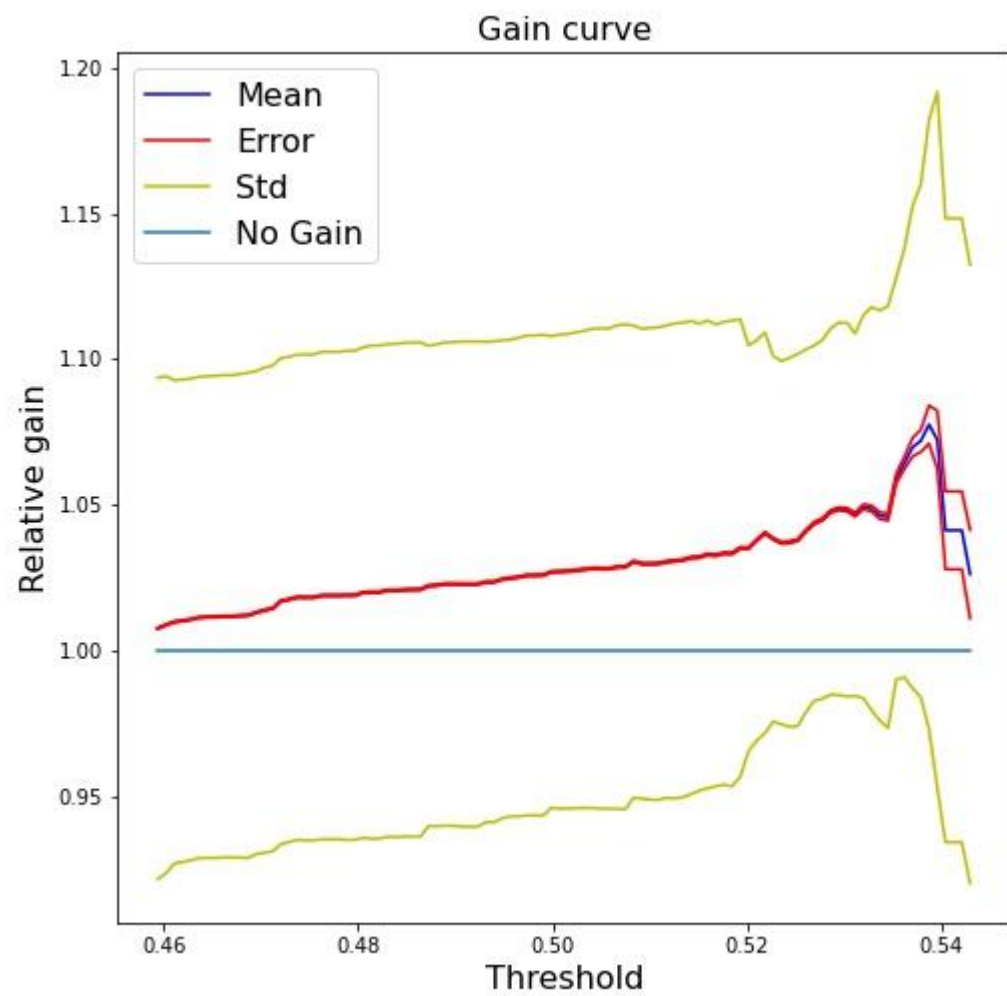
Genmab A/S



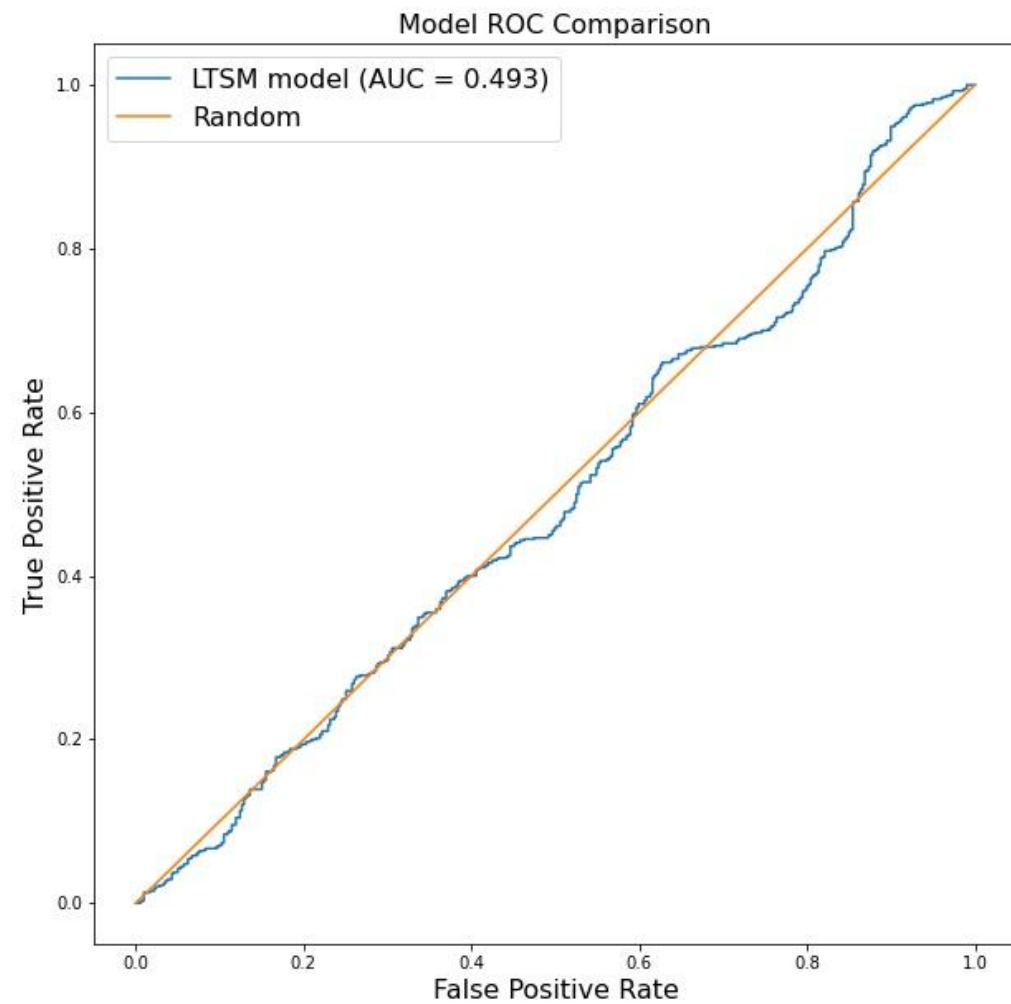
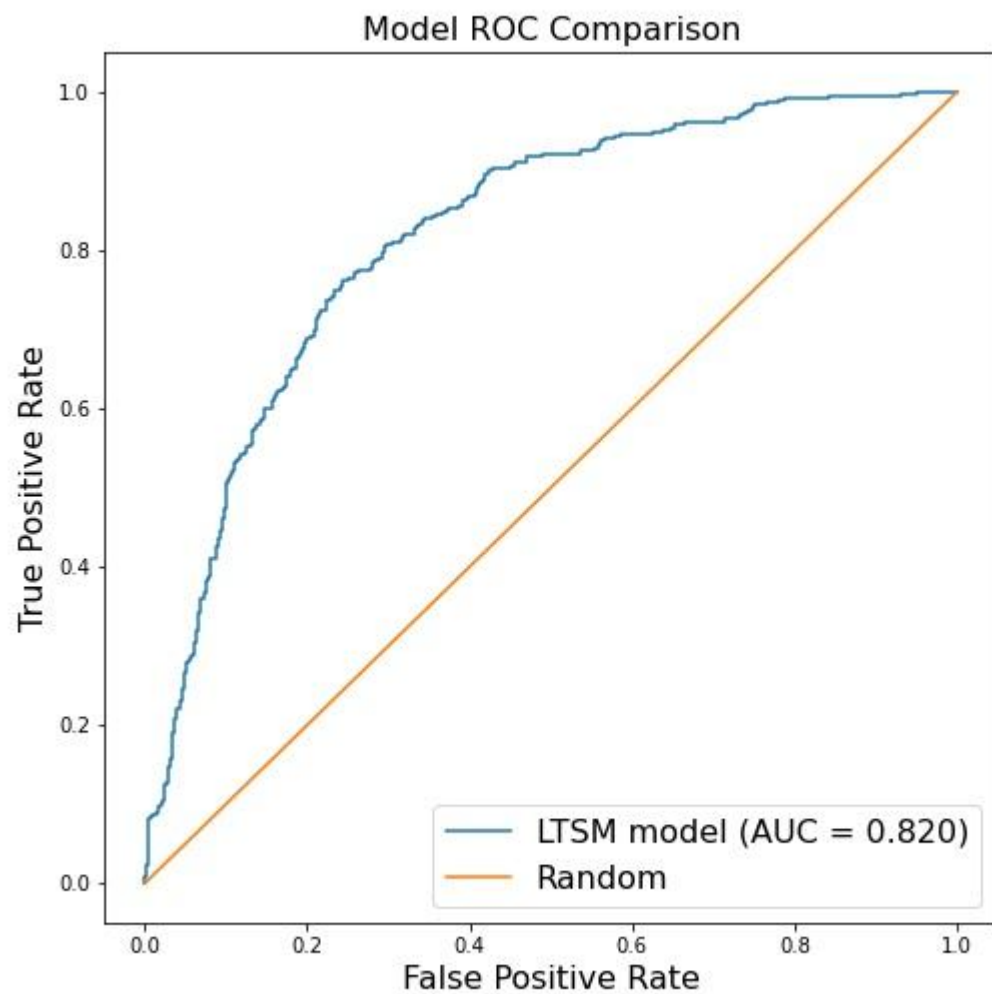
ISS A/S



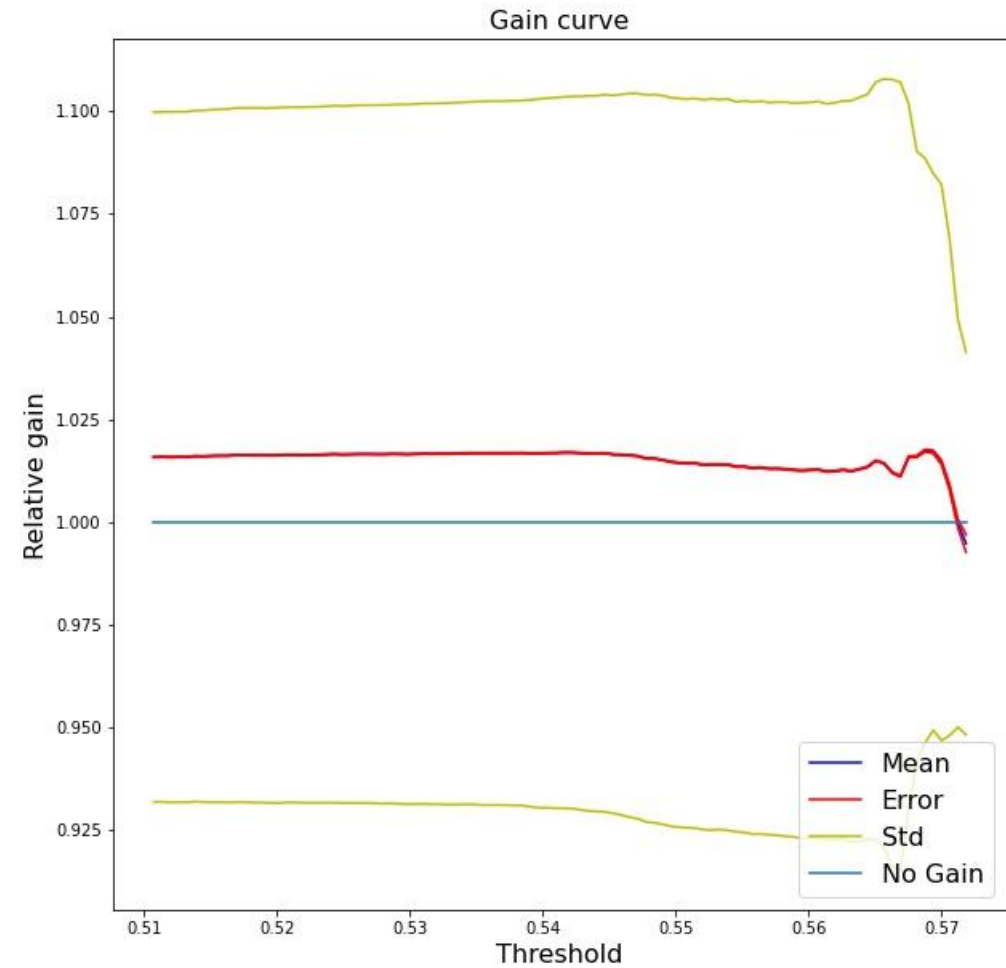
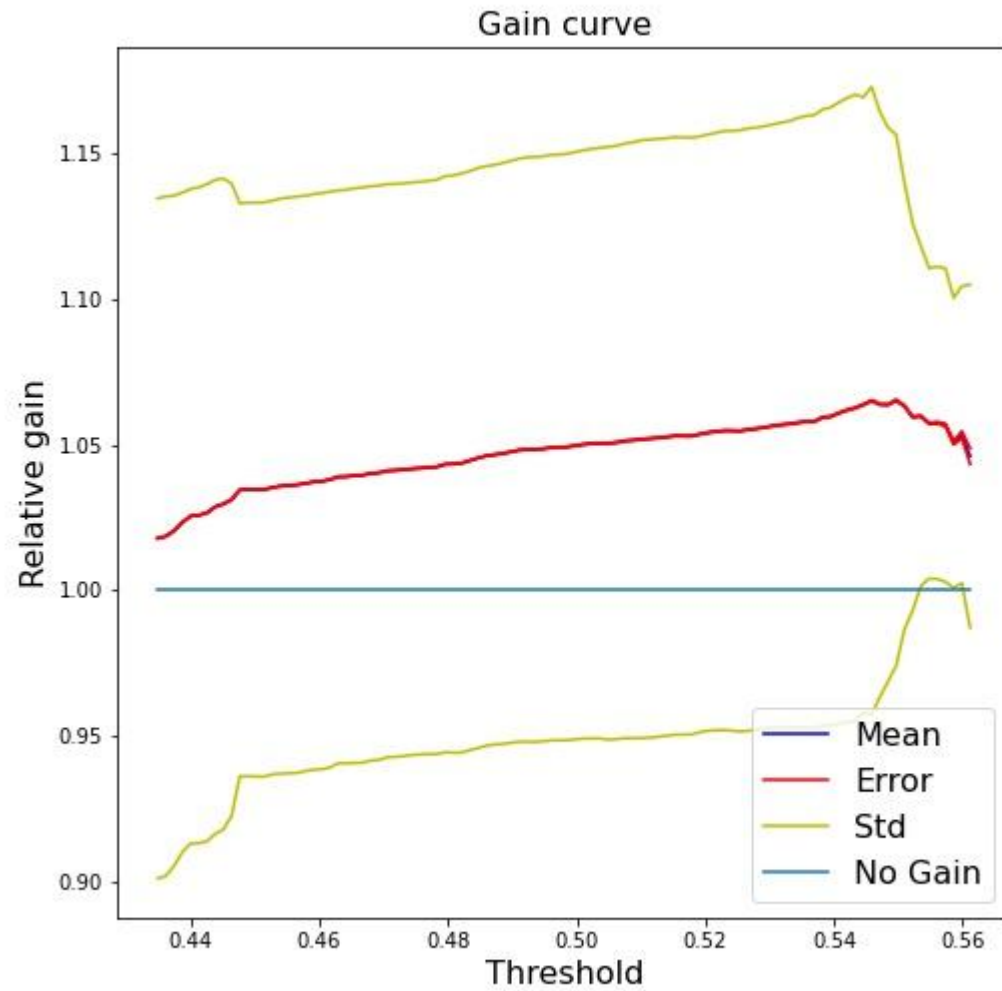
ISS A/S



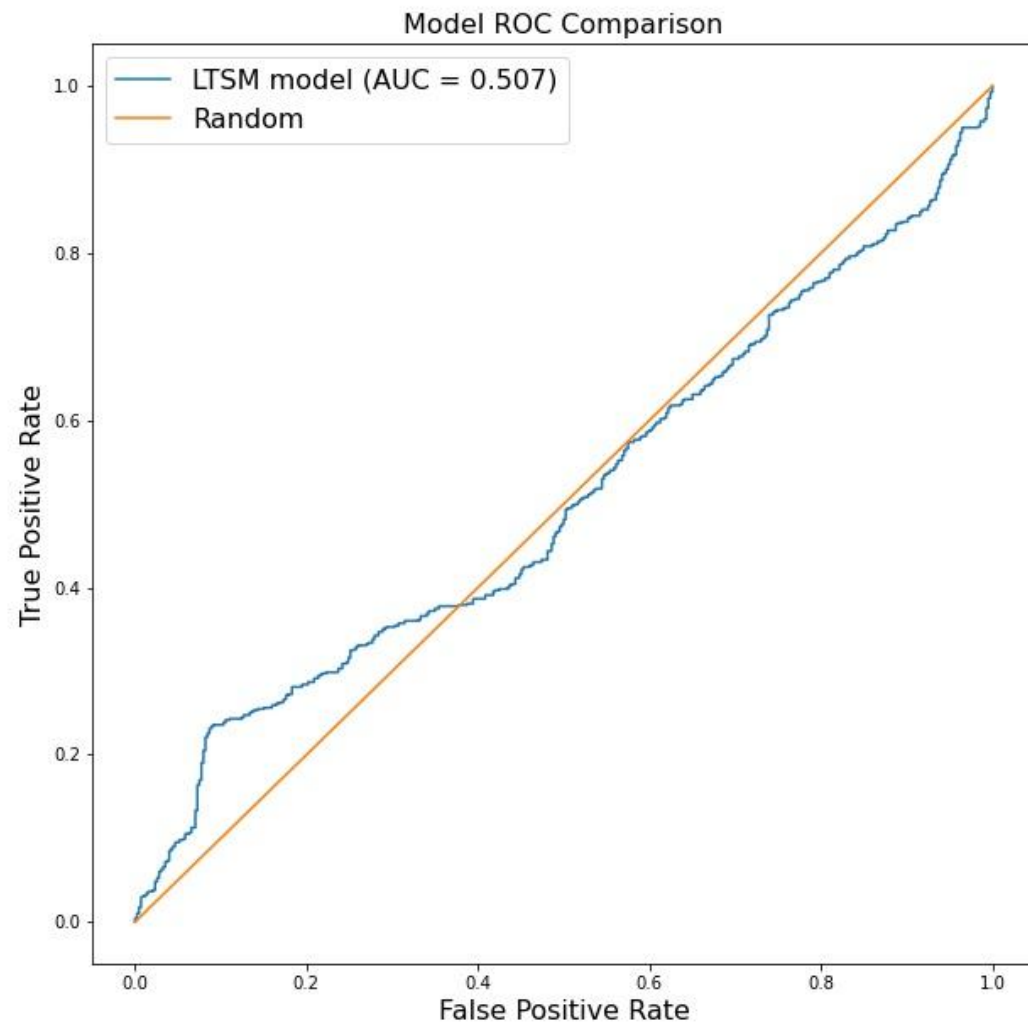
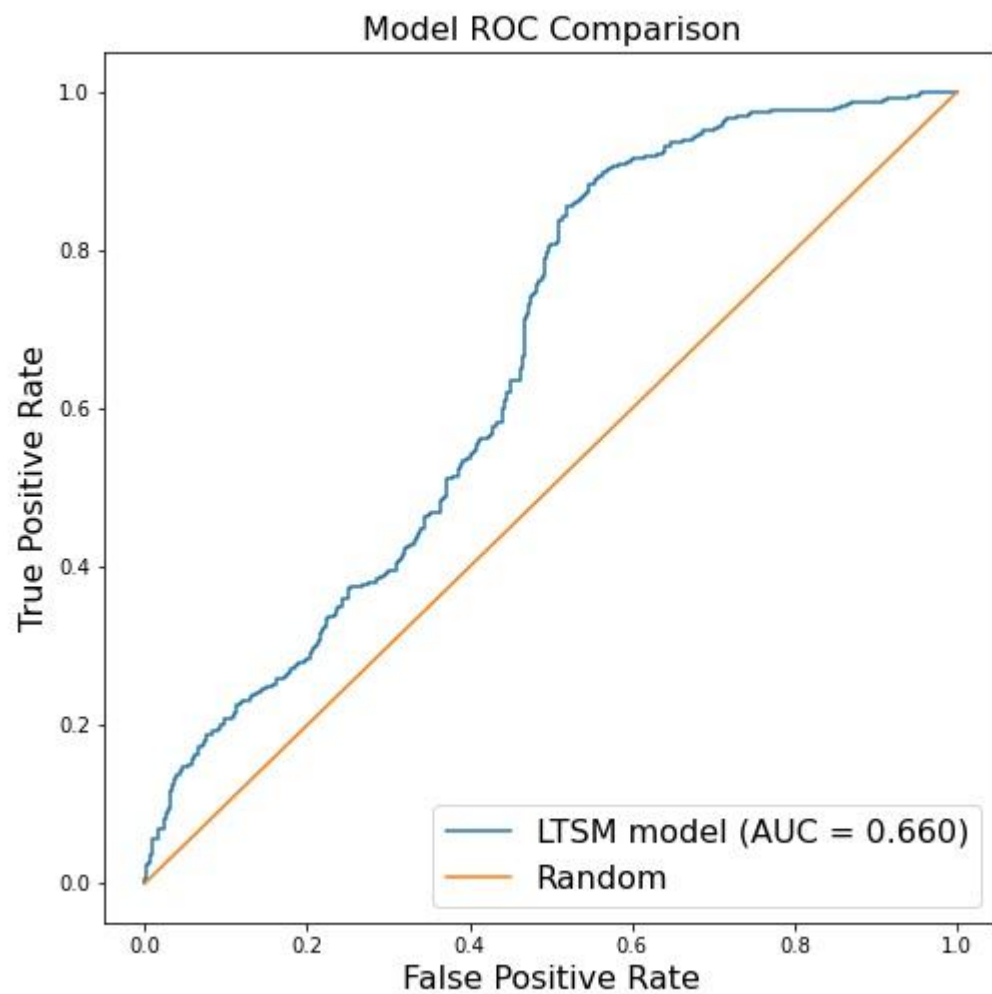
Royal Unibrew A/S



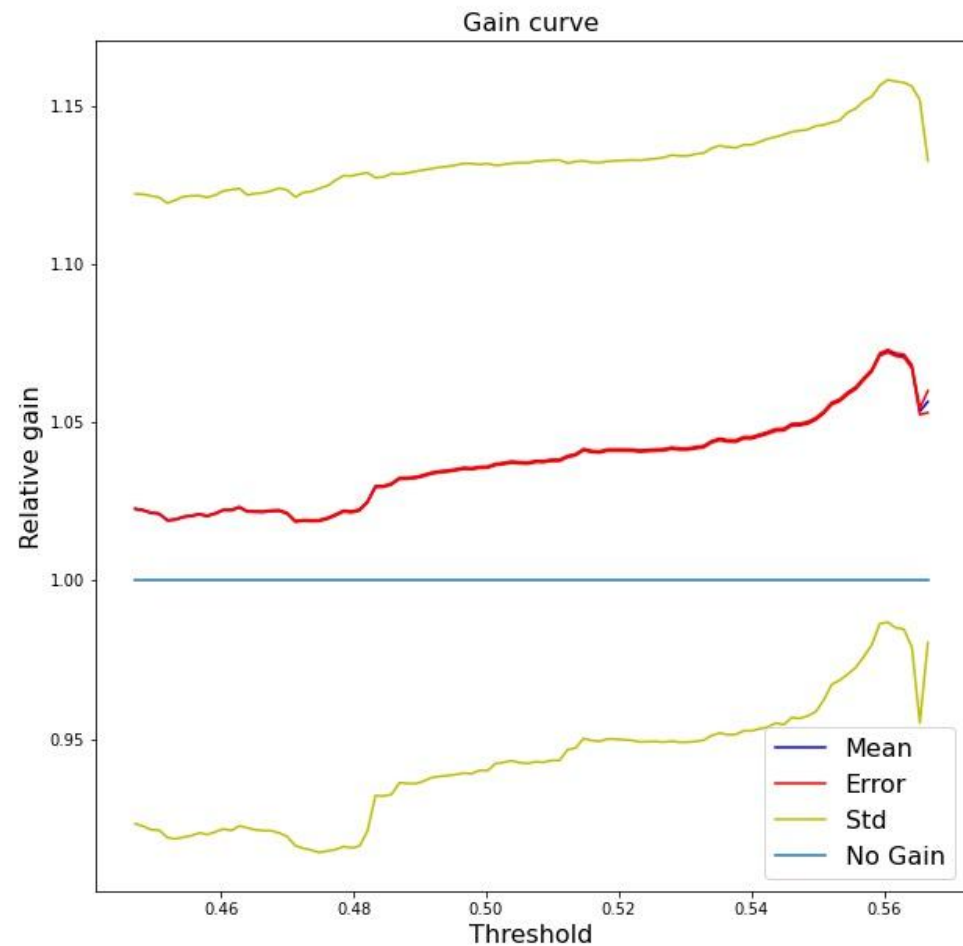
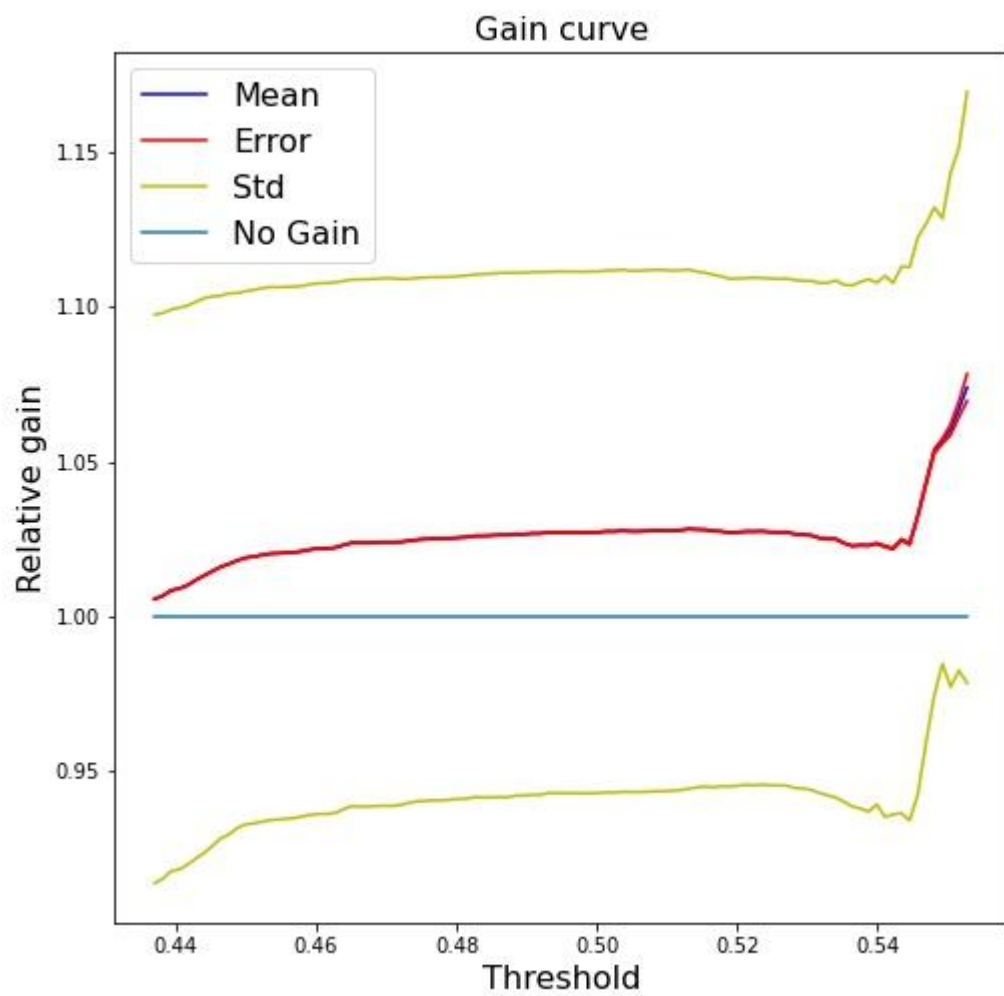
Royal Unibrew A/S



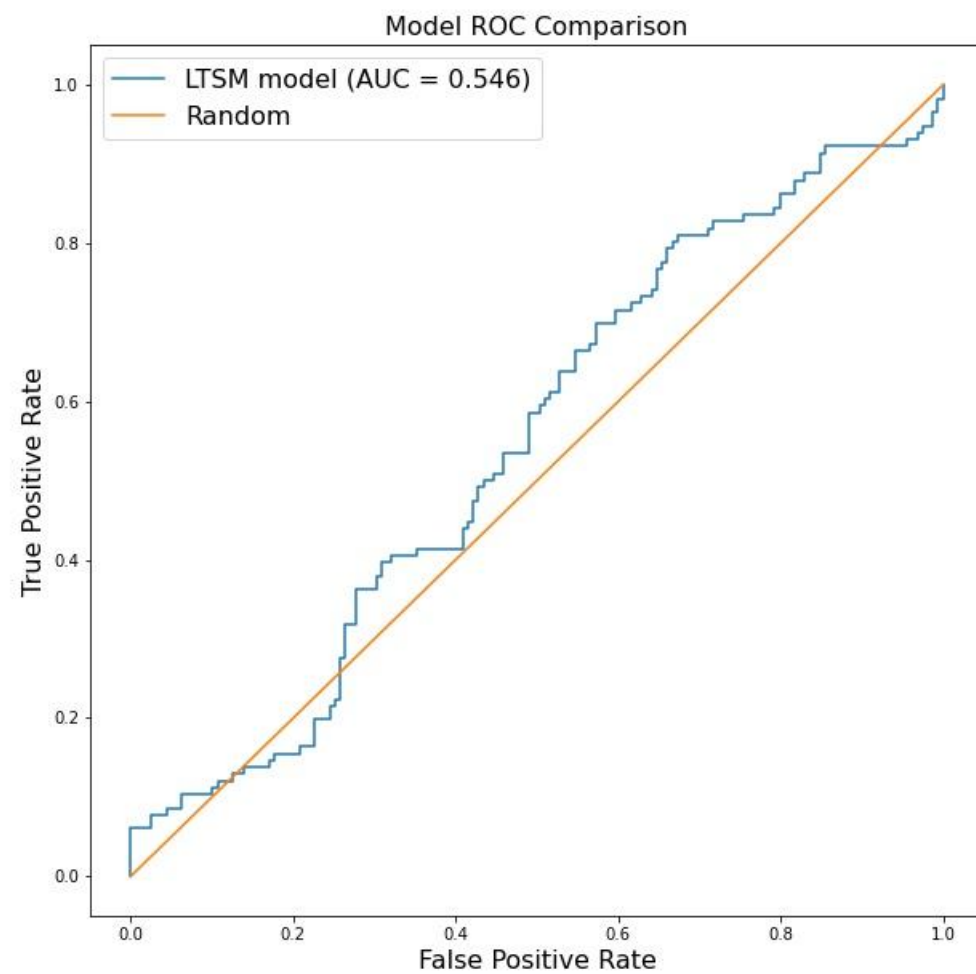
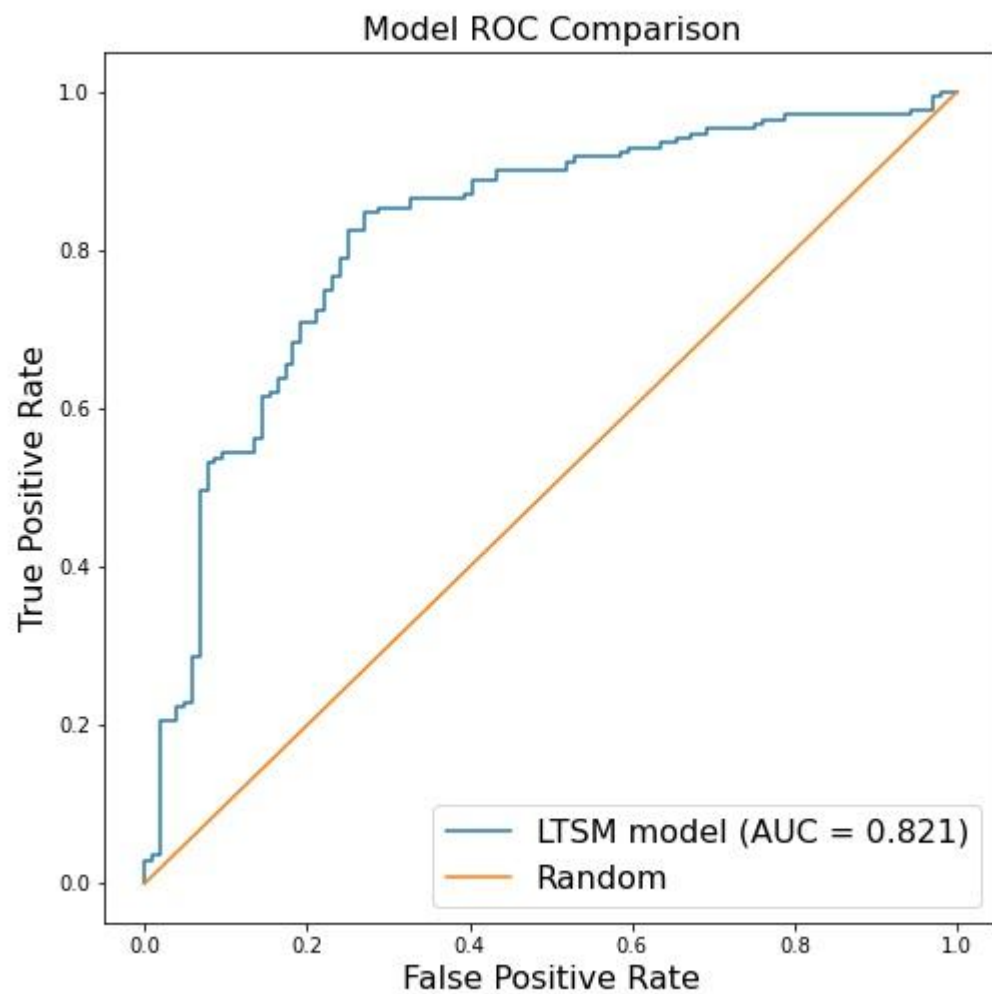
A.P. Møller - Mærsk A/S



A.P. Møller - Mærsk A/S

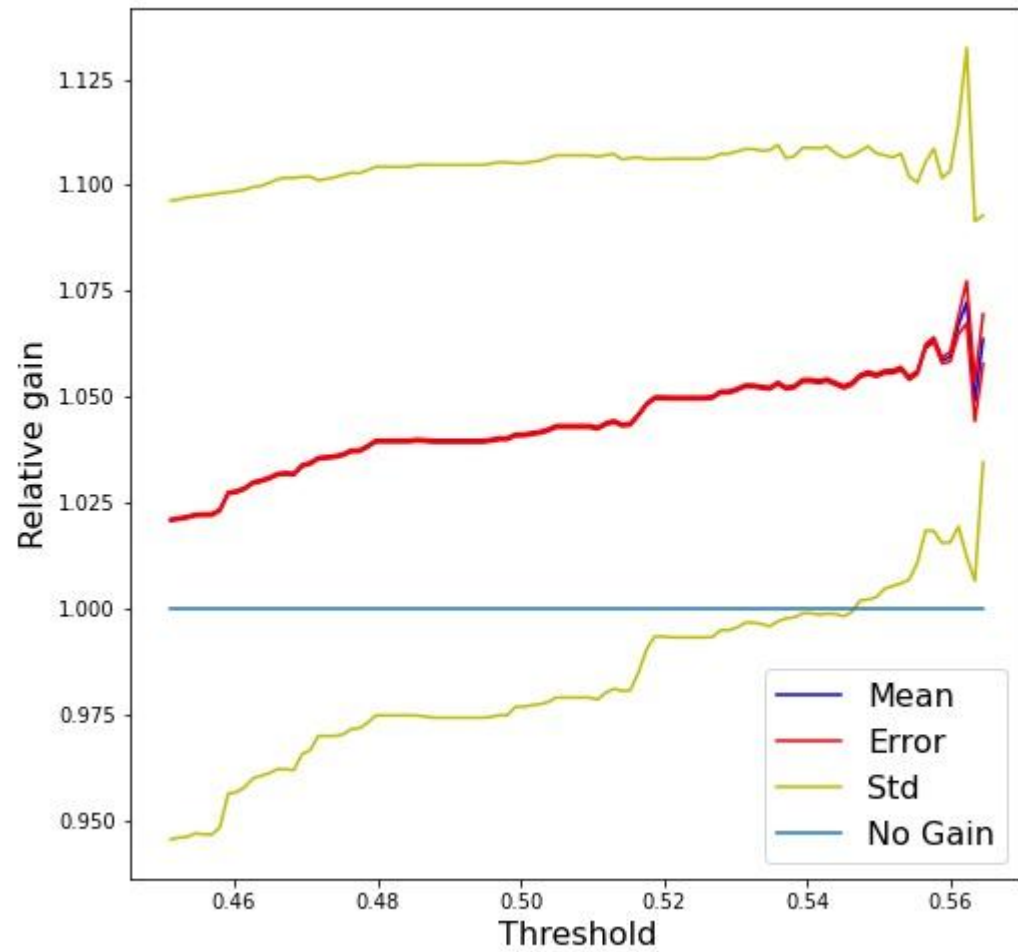


Ørsted A/S

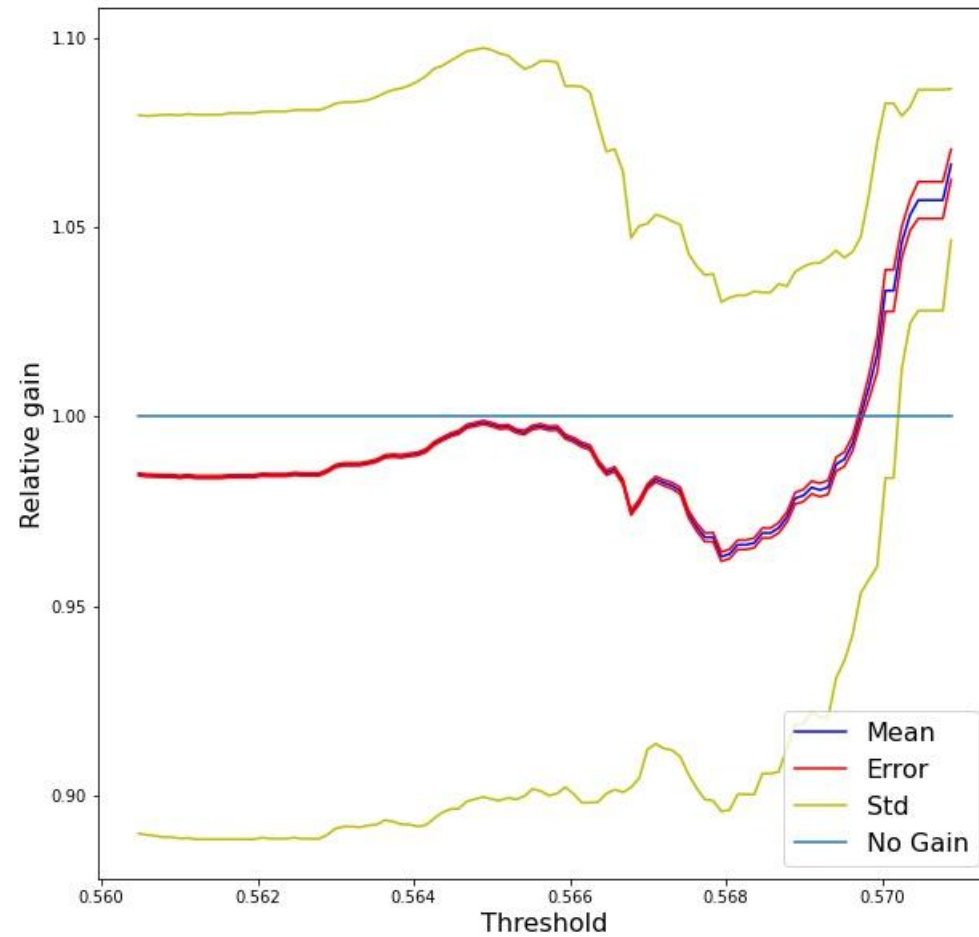


Ørsted A/S

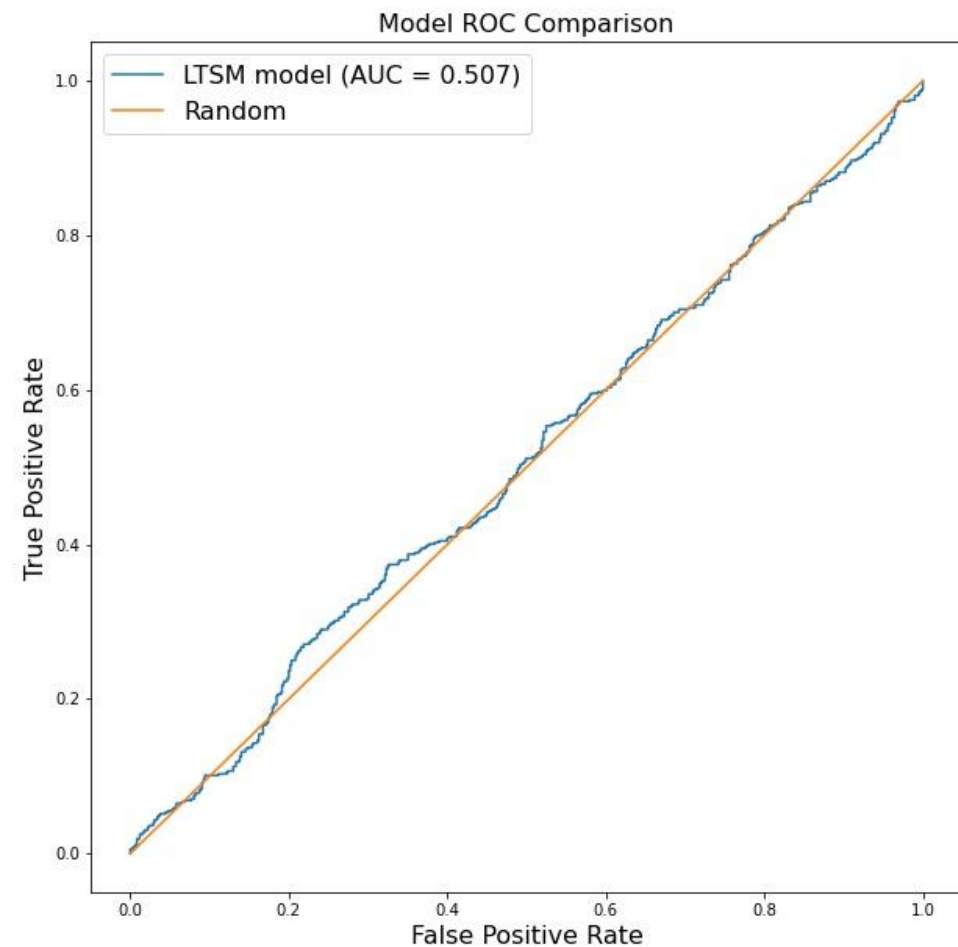
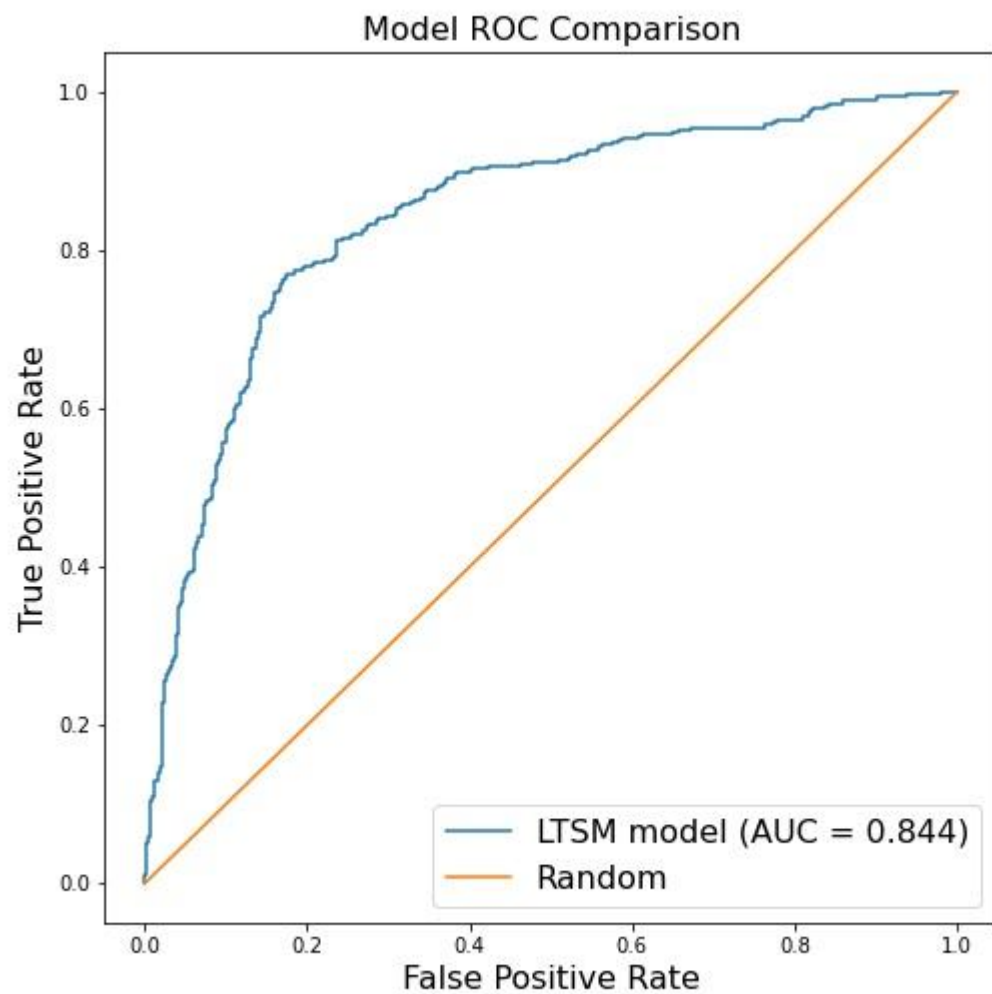
Gain curve



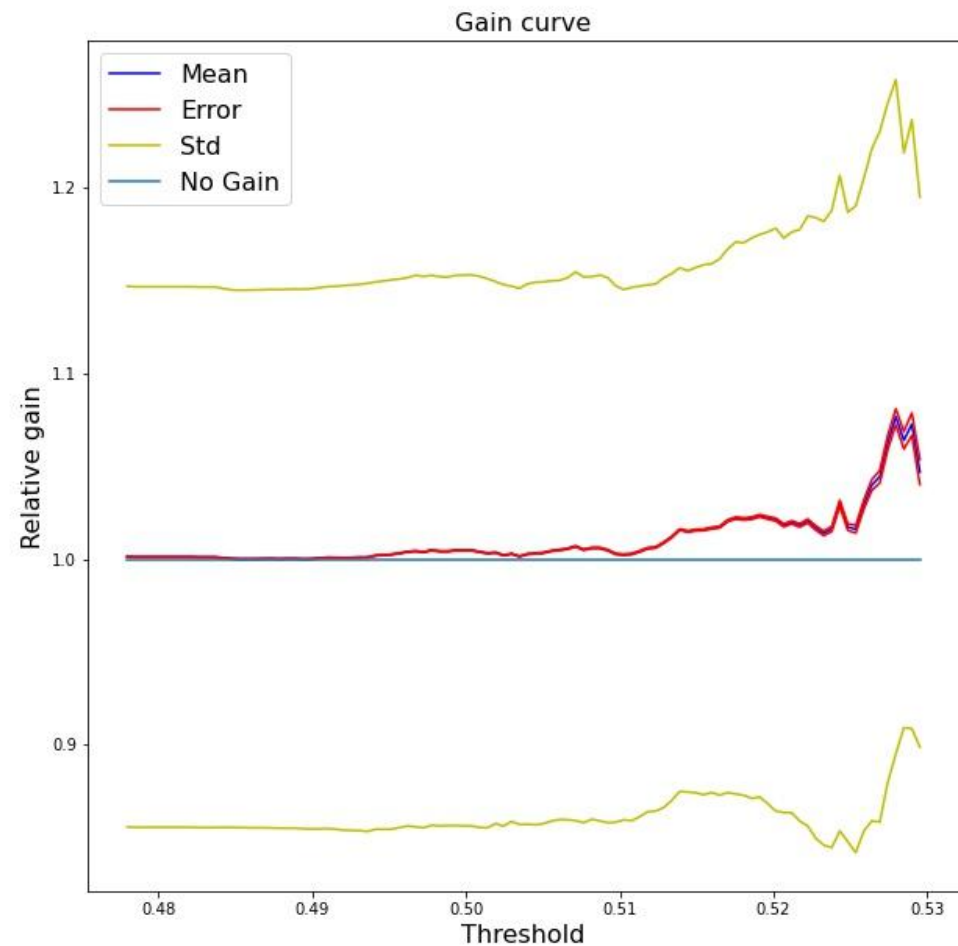
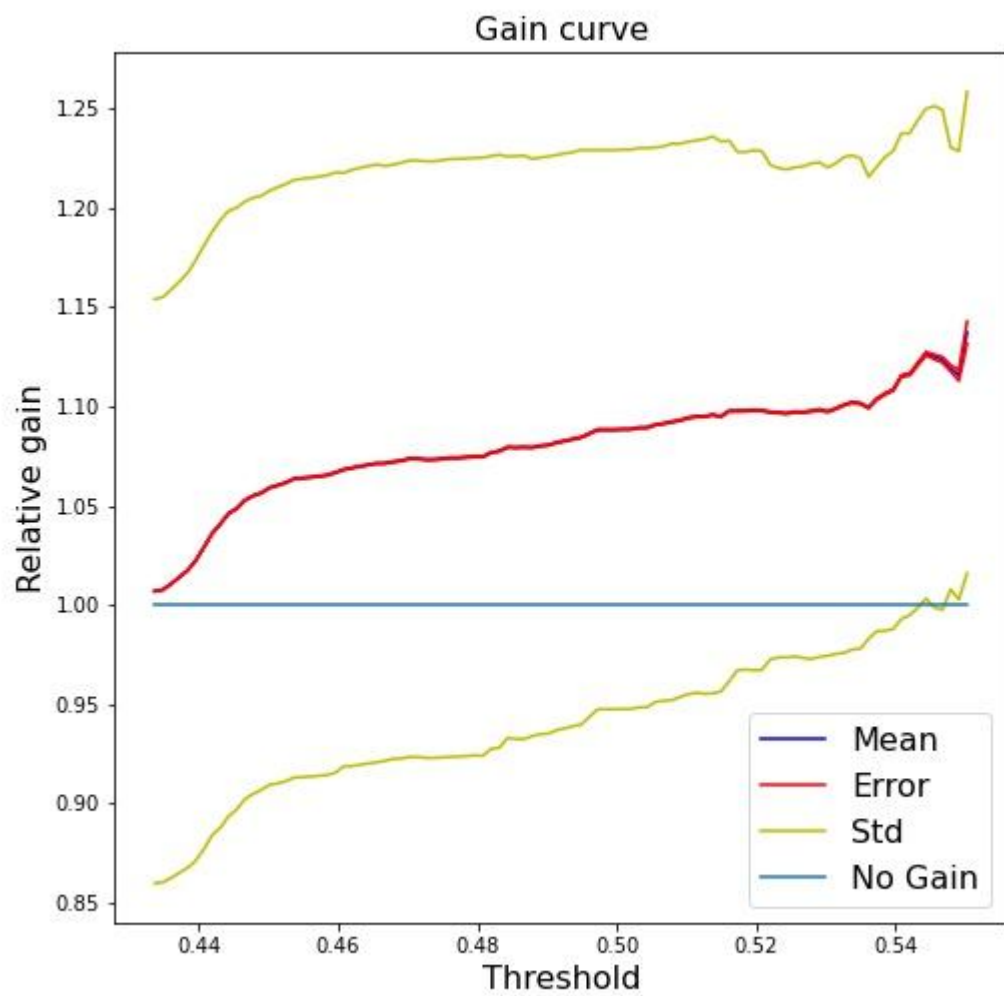
Gain curve



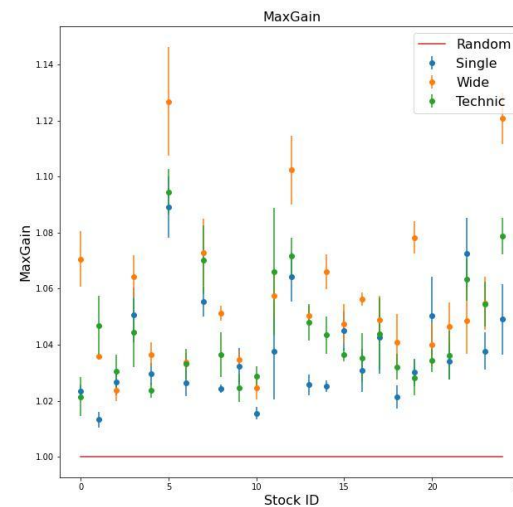
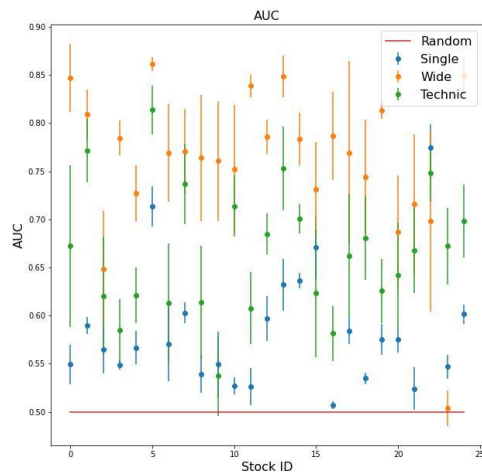
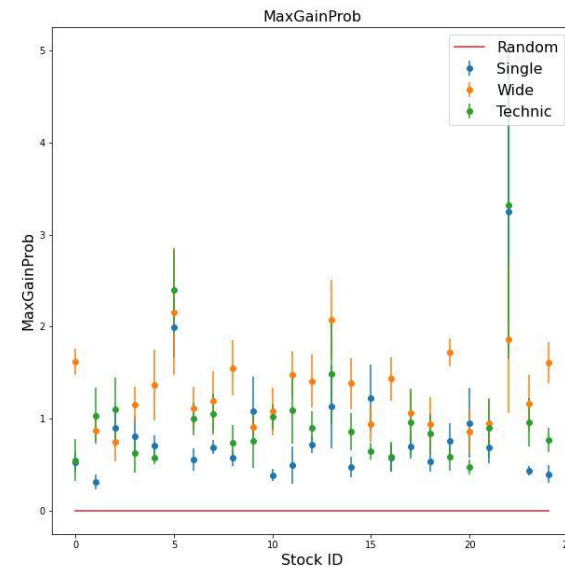
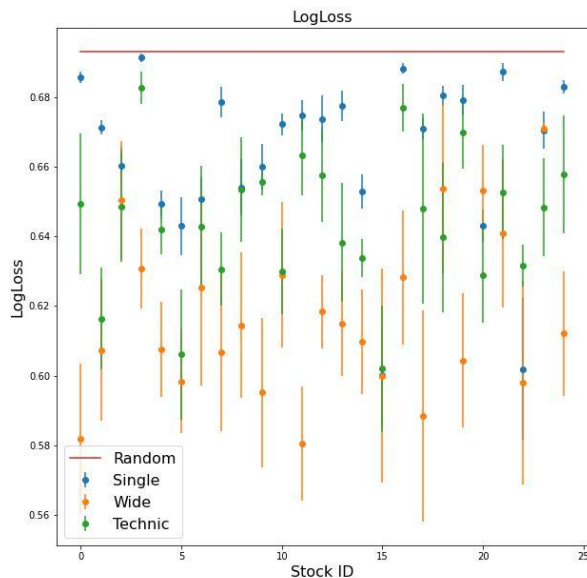
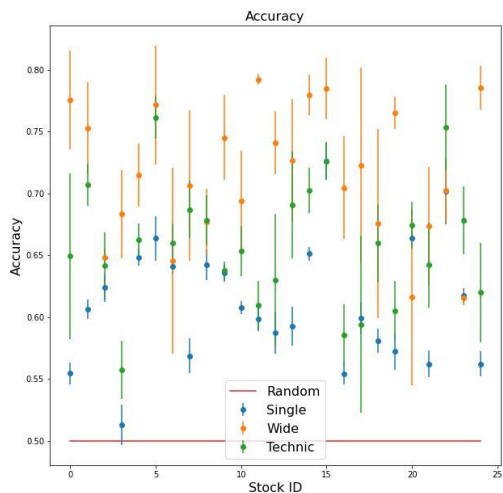
Bavarian Nordic A/S



Bavarian Nordic A/S



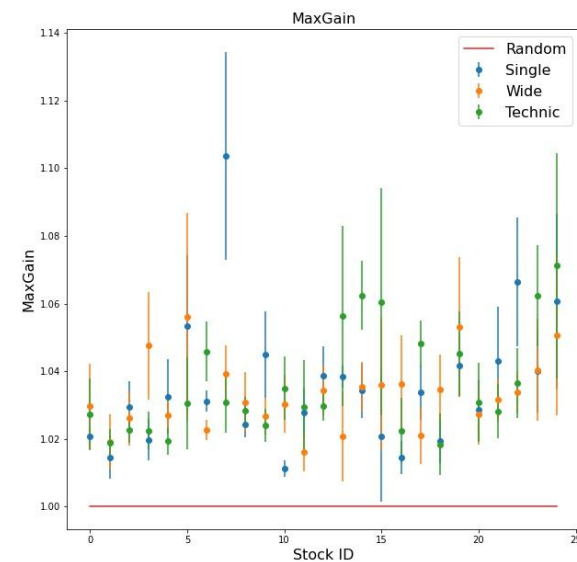
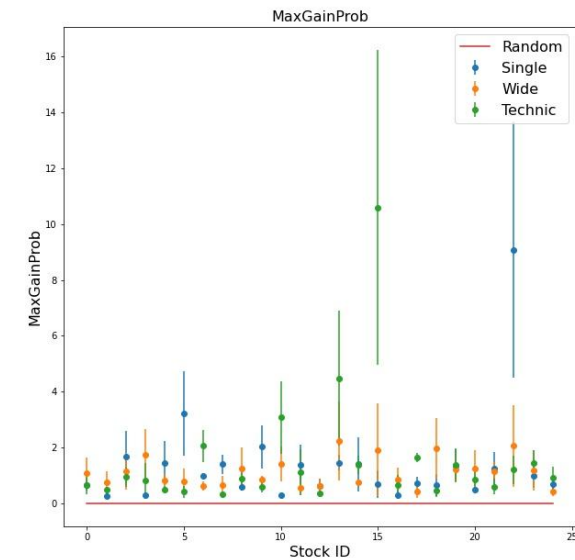
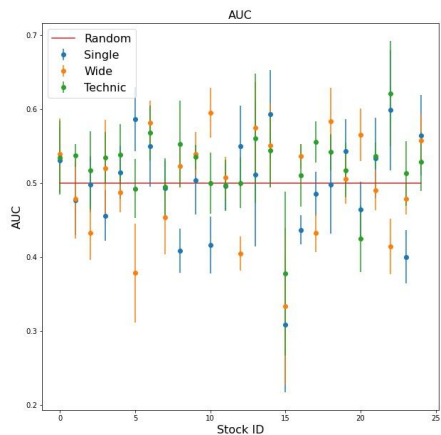
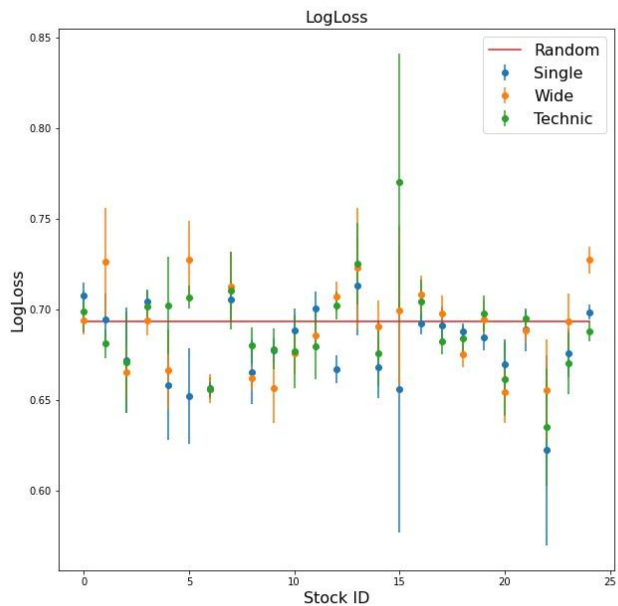
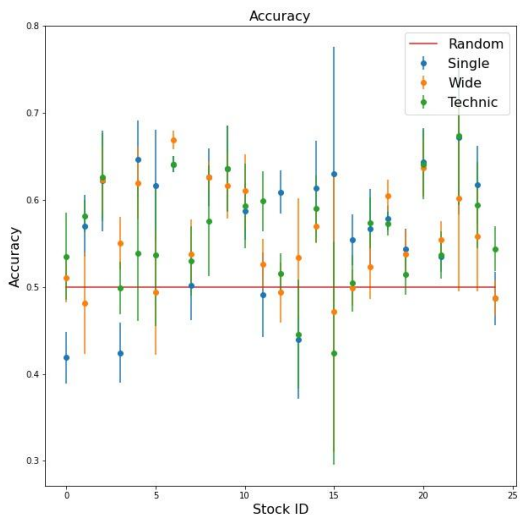
Shuffled LSTM evaluation



5-fold cross validation to get errors
 Lines indicate a completely random model

Wide most is best, followed by technical indicators and single mode is worst

Ordered LSTM evaluation



5-fold cross validation to get errors
 Lines indicate a completely random model

Wide most is best, followed by technical indicators and single mode is worst

Support vector machine

