

Developing a Credit Spread Trading Strategy using Tree-Based Machine Learning Methods

Jason Huang, Katie Li , Jeff Yang , Xuwei Yu

Abstract

We propose dollar-neutral trading strategies on Investment Grade Corporate Bond ETFs based on machine learning signals of future credit spread movements. Using various tree-based methods ranging from a single classification tree to a whole ensemble of black-box models, we aim to fit models that generalize well to out-of-sample data. We tune our models using a forward-chaining cross-validation scheme that is time-series appropriate, ensuring that we never use future information to predict the past. Our results indicate modest accuracy improvements relative to a baseline Logistic LASSO model. When backtesting the trading strategy on various Bond ETFs, we show it is possible to achieve very appealing returns for the risk taken over a three year time period in out-of-sample testing.

Contents

1	Introduction	1
2	Data and Methodology	2
2.1	Variables	2
2.2	Generated Features	2
2.3	Response Variable	3
2.4	Validation Strategy	3
3	Models	3
3.1	Classification Tree	3
3.2	Ensemble Methods	4
	Random Forest • Gradient Tree Boosting • XGBoost and Light-GBM	
4	Model Results	5
5	Trading Strategy	6
5.1	Methodology	6
5.2	Trading Strategy Results	6
5.3	Conclusion	7

1. Introduction

Credit spreads for corporate bonds have long been viewed as a fundamental dynamic, where investors are compensated for taking additional credit risk relative to default-free government securities. But, Longstaff (2005) [one] has argued against attributing default risk to existing spreads. Premiums demanded for high grade bonds cannot be explained solely by differences in default risk – many corporations hold capital reserves which more than cover for unlikely credit market events. We branch out from this fundamental view and explore

credit spreads as a picture of two moving components, focusing on factors affecting treasury notes and corporate bonds on individual bases. Our approach integrates both global and endogenous factors that encompass equity market risk and signal future corporate performance.

There has been growing popularity in utilizing machine learning methodologies for financial time series prediction. Various leading investment management firms have begun investing resources either in new hires or on data competitions. Two Sigma has recently finished hosting a stock prediction competition on Kaggle using news and market data. Citadel routinely hosts datathons across major US universities scouring for talent in statistics and machine learning. AQR has also started an in-house machine learning group in 2018.

The ability to accurately predict credit spreads using machine learning is extremely valuable for multiple reasons. Credit spreads tend to signal the health of an economy. When spreads are tight, it signals economic growth and prosperity. The opposite is true if wide credit spreads are observed. Having accurate predictions of credit spreads corresponds to possible trading strategies that can be quite profitable. Given predictions on whether credit spreads will widen or tighten, one can formulate a simple trading strategy as follows:

- For widening credit spreads, long the top end (BAML AAA corporate bond index) and short the bottom end (10 year treasuries).
- For tightening credit spreads, take opposite positions to the above

Our paper is structured as follows: in section 2, we introduce a wide range of features that have been collected and generated for this project. In sections 3 and 4, we describe the

machine learning models considered and report the results when applied on our data. In section 5, we expand on how we use the predictions to create a profitable trading strategy. Lastly, in section 6, we summarize our main findings.

2. Data and Methodology

2.1 Variables

To represent the credit spread, we used 10 year constant-maturity US treasuries and the ICE BAML US Corporate Effective Yield. The BAML index tracks the performance of US investment-grade rated corporate debt issued in the US market.

Because the credit spread consists of two moving parts, it is important to consider what factors cause each yield to move separately and what can cause the spread to narrow or widen. Therefore, we incorporated features that can roughly be characterized into four groups:

Macroeconomic The risk free rate can be viewed as a function of crucial economic variables, as lower maturity treasury notes and bills are disproportionately influenced by movements in the federal funds rate (FFR). The FFR is also determined by macroeconomic indicators paramount to FOMC decisions. The macroeconomic variables (retrieved from FRED database [13]) we consider are:

- *NAIRU* (non-accelerating inflation rate of unemployment) - a key metric that illustrates levels of economic activity and consumer spending
- *Unemployment rate* (including discouraged workers) - provides information about the state of the economy and ease of finding employment
- *Balance of Payments* - a measure of US export strength and weakness
- *Manufacturing* - captures the state of the US economy
- *Gold Prices* - often viewed as a safe alternative to treasuries for investors to flee towards in times of crisis

Interest Rate Because we are predicting corporate yield spreads, it is important to capture information from the term structure of interest rates. For example, an inverted yield curve can often signal impending market downturns. Thus, the term structure of US treasuries is incorporated using various “DGS” (constant maturity treasury yield) variables from the FRED database. A time series for each of 1, 2, 5, 7, 10, 20 year maturities are included to capture as much of the trends as possible. The 10 year treasury yield is inherently accounted for in our spread variables, but we also generate various lag features on these variables.

Additionally, we look at the TED spread, which is the difference between the three month treasury and three month LIBOR rates. This captures additional macroeconomic and

interest rate information from international fixed income markets. The LIBOR and US treasuries are often used for pricing of more complex financial derivative products, and thus the difference of the two is an important predictive variable.

Foreign Exchange Prevailing economic theory such as Uncovered Interest Rate Parity suggests there should be an empirical relationship between exchange rates and interest rates. Given exchange rate fluctuations, foreign investors will see different levels of attractiveness in investing in US corporate bonds. We select exchange rates based on countries that the US trades heavily with.

Stock Market Kwan (1996) [2] and Gebhardt et al. (2005) [8] suggest a strong relationship between past stock market returns and future bond returns for a given company. Thus, we include various features generated from the S&P index [14] to proxy for this relationship in the aggregate market. We have also incorporated the daily volume of the S&P as a measure of aggregate liquidity in equity markets, and this may potentially be correlated with aggregate corporate bond market liquidity. Finally, we add the VIX, a proxy for aggregate market volatility.

2.2 Generated Features

Our lagged features are constructed using three rolling windows - 5, 20, and 60 trading days, corresponding to approximately 1 week, 1 month, and 3 months. For each rolling window, we compute simple moving averages, exponential moving averages, and standard deviations for the chosen features. We also process sharpe and tail ratio features for the spread variable using a 60 day rolling window. The calculations are as follows:

Let X_t denote a variable at time t . For $w \in \{5, 20, 60\}$

Simple Moving Average

$$MA = \frac{\sum_{i=1}^w X_{t-i}}{w}$$

Exponential Moving Average

$$EMA = \sum_{i=1}^w \left(1 - \frac{2}{w+1}\right)^{i-1} X_{t-i}$$

For y_t denote the spread at time t .

Sharpe Ratio of Spread

$$Sharpe = \frac{\frac{1}{60} \sum_{i=1}^{60} y_{t-i}}{\sqrt{\frac{1}{60-1} \sum_{i=1}^{60} (y_{t-i} - \frac{1}{60} \sum_{i=1}^{60} y_{t-i})^2}}$$

Tail Ratio of Spread (60 day window)

$$Tail = \frac{95\text{th percentile value}}{|5\text{th percentile value}|}$$

$$\begin{aligned}
 & -\rho(X_t^{\text{Open}}, X_{t-7}^{\text{Close}}, 10) \\
 \text{(a) Alpha 6} & \\
 & \begin{cases} -\phi(X_t^{\text{Close}} - X_{t-7}^{\text{Close}}, 60) * \text{sgn}(X_t^{\text{Close}} - X_{t-7}^{\text{Close}}), & \mu(X_t^{\text{Volume}}, 20) < X_t^{\text{Volume}} \\ -1 & \text{otherwise} \end{cases} \\
 \text{(b) Alpha 7} & \\
 & \begin{cases} X_t^{\text{Close}} - X_{t-1}^{\text{Close}}, & \min(X_t^{\text{Close}} - X_{t-1}^{\text{Close}}, 5) > 0 \\ X_t^{\text{Close}} - X_{t-1}^{\text{Close}}, & \max(X_t^{\text{Close}} - X_{t-1}^{\text{Close}}, 5) < 0 \\ -(X_t^{\text{Close}} - X_{t-1}^{\text{Close}}) & \text{otherwise} \end{cases} \\
 \text{(c) Alpha 9} & \\
 & -(X_t^{\text{Close}} - X_{t-1}^{\text{Close}}) * \text{sgn}(X_t^{\text{Volume}} - X_{t-1}^{\text{Volume}}) \\
 \text{(d) Alpha 12} &
 \end{aligned}$$

Figure 1. Alpha Signals Used

Investors tend to view bonds as substitute investments when the stock market is performing poorly. By generating additional features from the S&P variables, we extract market trends/signals and incorporate this information into our models. In particular, we implement real quant trading alpha signals from Kakushadze (2016) [3] on the S&P variables along with the same lag features previously described.

The alpha signals use the following functions:

- $\mu(x, t)$ - the trailing mean function looking t days back for variable x .
- $\rho(x, y, t)$ - the rolling correlation looking t days back for variables x and y .
- $\phi(x, t)$ - the time series rank function, where the values of x over the t day period are ranked and the function returns the rank of the current observation
- $\min(x, t)$ - the historical min in the last t days for variable x
- $\max(x, t)$ - the historical max in the last t days for variable x

Given Open, Close and Volume for some underlying X_t at time t , we compute the following alpha signals:

Alpha 6 - Figure 1 (a)

Alpha 7 - Figure 1 (b)

Alpha 9 - Figure 1 (c)

Alpha 12 - Figure 1 (d)

A summary of all of the variables can be found in Table 1.

2.3 Response Variable

Instead of predicting the exact spread on the next trading day, we decide to predict whether the spread will increase or decrease over the next day. Thus, we treat the problem as binary classification.

2.4 Validation Strategy

Our data ranges from 1998 to 2019. We hold out the last three years of data as a test set. With the rest of the data, we use cross-validation to evaluate the predictive power of various models and fine-tune each model's hyper-parameters. Traditional cross-validation techniques such as K-fold do not typically work well on time-series, as the data is often serially correlated. Additionally, data leakage must be prevented.

We use a rolling cross-validation technique [4] that ensures each training set consists only of observations occurring prior to any observation in the test set. No future information is used to train the models. Figure 2 shows how we partition the data. The cross validation accuracy is computed by averaging across all folds.

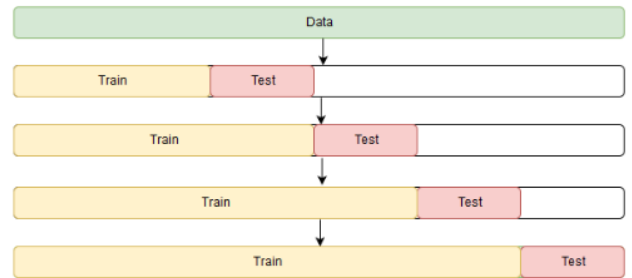


Figure 2. Cross-Validation Scheme for Time Series

3. Models

The goal of a tree-based model is to predict the value of a target variable by learning simple decision rules inferred from the data features. Tree-based methods are flexible, relatively interpretable, and have great predictive power. Because of these aspects, we decided to compare a series of tree-based models ranging in their black-box nature.

3.1 Classification Tree

Classification trees take a top-down, greedy approach to searching for the optimal split at each step of the tree. Using the

Table 1. Variables
 (*) indicates that MA, EWA, STD were generated on the variable

Type of Variable					
Spread	Macroeconomic	Interest Rate	FX	Stock Market	Generated
Risk-Free Rate	NAIRU	DSG (*)	USD/CAD (*)	S&P Open (*)	Alpha 6
BAML index	Balance of Payments	(1/2/5/7/10/20 yr constant maturity)	USD/EURO (*)	S&P Close (*)	Alpha 7
Spread Lag 1	Manufacturing	TED Rate (*)	USD/JPY (*)	S&P Adj Close	Alpha 9
Spread Lag 2	Unemployment Rate		USD/GBP (*)	S&P Volume (*)	Alpha 12
Spread Lag 3	Gold (*)			S&P High (*)	Spread Sharpe Ratio(*)
				S&P Low (*)	Spread Tail Ratio (*)
				VIX (*)	

CART algorithm, we select the feature and threshold that yield the largest information gain at each node.

Let data at node m be represented by Q . For each candidate split $\theta = (j, t_m)$, where j is the feature and t_m is the threshold, partition the data in to $Q_{left}(\theta)$ and $Q_{right}(\theta)$ based on the threshold.

Impurity measures how well the two response classes are separated. Ideally, we want nodes to perfectly separate the 0 and 1 class; in this case, the impurity would be 0. For a given candidate split, the impurity of the node is:

$$G(Q, \theta) = \frac{n_{left}}{N_m} H(Q_{left}(\theta)) + \frac{n_{right}}{N_m} H(Q_{right}(\theta))$$

Then, parameters are selected to minimize impurity.

$$\theta^* = \operatorname{argmin}_{\theta} G(Q, \theta)$$

In this fashion, we continue splitting the tree until the maximum depth is reached. To prevent overfitting, we fine-tune the *maximum depth* and *minimum sample size* of a split and leaf to encourage shorter trees. Single decision trees can be highly unstable so we incorporate a *minimum impurity decrease* parameter, only allowing a node to be split if the reduction in impurity is at least as large as the parameter.

3.2 Ensemble Methods

A single classification tree is prone to overfitting and bias. To improve on the classification tree, we apply a series of ensemble methods that combine individual trees and incorporate features such as bagging and boosting.

Bagging Bagging methods seek to build several independent estimators and then average their predictions. This reduces the variance of the combined estimator.

3.2.1 Random Forest

A random forest has several key differences from a decision tree classifier. Firstly, when splitting a node, we no longer select from the best possible split among all features. Instead, we choose a random subset of features and pick the best split from that smaller group. Moreover, each tree is built from a bootstrap sample of the training set, with data points drawn with replacement. Because of the randomness introduced, individual trees are less correlated with each other.

The individual classifiers are combined by averaging their probabilistic predictions.

We tune the same parameters as the classification tree, but add an additional parameter, *number of estimators*, that controls the number of trees in the forest. Increasing the number of trees helps train the classifier, but the generalization ability decreases past a certain point.

Boosting Boosting methods try to reduce the bias of the combined estimator by combining several weak models.

3.2.2 Gradient Tree Boosting

The gradient boosting classifier is an additive model using classification trees of a fixed size as weak learners $h_m(x)$. Thus we can denote our model $F(x)$ as a weighted sum of M weak learners:

$$F_M(x) = \sum_{m=1}^M \gamma_w h_m(X)$$

When the $m - th$ learner is added to the model, it becomes

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

The new learner tries to minimize the loss L , given the previous ensemble of trees $F_{m-1}(x)$:

$$h_m = \operatorname{argmin}_h \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + h(x_i))$$

This optimization problem can be solved using gradient descent, and works for arbitrary differentiable loss functions. For our classification problem, we use the negative binomial log-likelihood loss function. The main parameters we tune are the *number of learners* as well as the *learning rate*. The learning rate is a regularization parameter and scales back the step length in gradient descent. There is a trade-off between number of learners and the learning rate, as a smaller learning rate typically requires more weak learners before being able to predict well. Overall, this class of models has good predictive power and is robust to outliers in the output space.

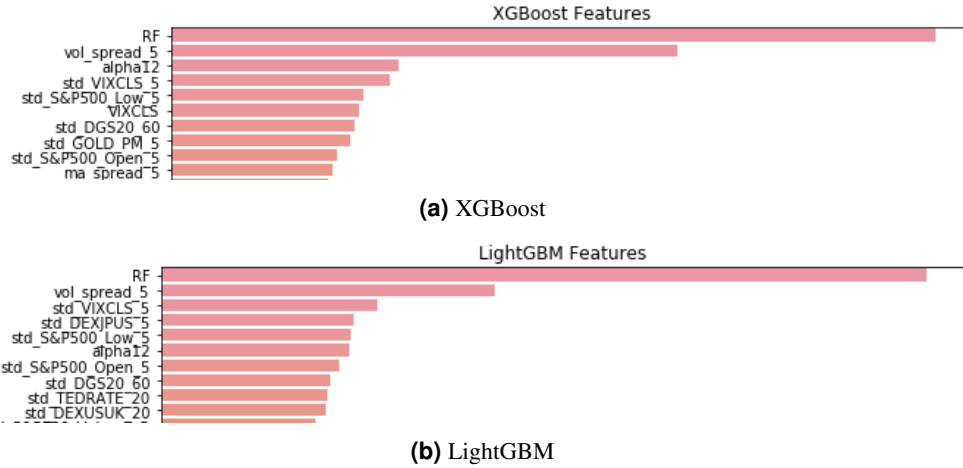


Figure 3. Top 10 features for XGBoost and LightGBM Models

3.2.3 XGBoost and LightGBM

XGBoost and LightGBM are the two most black-box models we consider. These algorithms have dominated Kaggle competitions over the past few years because of their incredible predictive power.

XGBoost stands for Extreme Gradient Boosting. XGBoost implements gradient boosted decision trees, but with deep considerations in terms of systems optimization in order to provide scalable, portable, and accurate predictions. Besides traditional gradient boosting, the algorithm includes stochastic gradient boosting with sub-sampling at the row, column, and column per split levels, and regularized gradient boosting with both L1 and L2 regularization.

The training objective is to minimize the following objective function:

$$Obj = L + \Omega$$

L is the loss function, which in this case is binary classification log loss:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Ω is the regularization term, which controls the complexity of the model.

$$\Omega = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

where T is the number of leaves and w_j^2 is the score on the j^{th} leaf.

LightGBM is another black-box gradient boosting method. It boasts a reduced cost of calculating many of the tree computations and overall memory usage. While traditional tree-based methods are top-down (growing by level), LightGBM grows trees leaf-wise, choosing the leaf with the maximum reduction in loss.

Because neither model appears to have an advantage over the other, we decided to combine LightGBM and XGBoost into one classifier. We hope that by averaging the predictions for these two models, we can get a lower-variance, more stable estimator.

4. Model Results

Cross validation and out-of-sample accuracy for all of the models can be found in table 2.

We use a regularized logistic regression model as a baseline. Because of the large number of features we have included in the model, using a LASSO penalty helps with feature selection and with preventing overfitting. This results in a simple but effective baseline model, that yielded a 59.6% test set accuracy and test set log-loss of 0.696.

The decision tree model performed surprisingly well for a relatively simple model, though it did not beat the logistic LASSO in either CV or test accuracy. Cross validation scores for Random Forest improved on the baseline; however, it did not perform well on the test set which could be a sign of overfitting.

We created a "ensemble model" that averages the predictions for decision tree, gradient boosting, and XGBoost & LightGBM. The various gradient boosting models and the ensemble model performed the best, with very good accuracy and much lower test set log-loss.

As we may want to know which of our covariates contribute the most to model performance, we have shown the top 10 features for the LightGBM and XGBoost models in Figure 3. In both models, the Risk Free Rate, volatility of the spread over 5 days, alpha12, standard deviation of the VIX, standard deviation of the S&P500, and standard deviation of DGS 20 over 60 days were among the top 10 features. In XGBoost, the standard deviation of gold over 5 days was in the top 10, while for LightGBM the standard deviation of the TED rate over 20 days was in the top 10. This shows that changes in the spread are extremely affected by volatility features.

Table 2. Model Results

Result	Logistic LASSO (Base Line)	Decision Tree (1)	Random Forest (2)	Gradient Boosting (3)	XGBoost & LightGBM (4)	Ensemble (1),(3) and (4)
CV Accuracy	Fold 1: 0.533 Fold 2: 0.600 Fold 3: 0.556 Fold 4: 0.605 Fold 5: 0.585	Fold 1: 0.488 Fold 2: 0.593 Fold 3: 0.537 Fold 4: 0.606 Fold 5: 0.570	Fold 1: 0.558 Fold 2: 0.672 Fold 3: 0.597 Fold 4: 0.674 Fold 5: 0.701	Fold 1: 0.587 Fold 2: 0.625 Fold 3: 0.552 Fold 4: 0.652 Fold 5: 0.657	Fold 1: 0.580 Fold 2: 0.620 Fold 3: 0.547 Fold 4: 0.617 Fold 5: 0.622	Fold 1: 0.513 Fold 2: 0.614 Fold 3: 0.541 Fold 4: 0.615 Fold 5: 0.597
CV Logloss	Fold 1: 0.840 Fold 2: 0.666 Fold 3: 1.015 Fold 4: 0.664 Fold 5: 0.667	Fold 1: 2.644 Fold 2: 0.674 Fold 3: 0.800 Fold 4: 0.675 Fold 5: 0.680	Fold 1: 0.665 Fold 2: 0.620 Fold 3: 0.656 Fold 4: 0.621 Fold 5: 0.590	Fold 1: 0.667 Fold 2: 0.636 Fold 3: 0.671 Fold 4: 0.635 Fold 5: 0.624	Fold 1: 0.671 Fold 2: 0.651 Fold 3: 0.682 Fold 4: 0.651 Fold 5: 0.648	Fold 1: 0.687 Fold 2: 0.651 Fold 3: 0.687 Fold 4: 0.648 Fold 5: 0.648
Test Accuracy	0.596	0.589	0.536	0.595	0.607	0.605
Test Logloss	0.696	0.682	0.688	0.676	0.673	0.673

Accuracy is not a great measurement of predictive performance in a finance setting, because we are more interested in making good bets sized by the level of expected risk. Thus, looking at the log-loss on the test set may be a better indicator of model performance. As we will soon expand upon, we are very interested in trading performance which additionally depends upon the uncertainty of predictions. The log-loss measure better captures this dynamic by considering the deviation of predicted probabilities from the true label. Selecting models based on this metric likely leads to better risk-adjusted performance (returns relative to volatility).

5. Trading Strategy

5.1 Methodology

We employ ETFs that track various bond indices. For the top leg of the spread, we use QLTA (For predictions from the AAA, AA and A indices) and PBBBX (For predictions from the BBB index). For the bottom leg, we use IEF. Predictions for yield directions are inversely related to bond prices and correspondingly also inversely related to our ETF prices. If we predict spreads will widen, we short the top leg ETF and long the bottom leg ETF. If we predict spreads will tighten, we long the top leg ETF and short the bottom leg ETF. The nominal amounts are scaled such that the strategy will be self financing on any given day and we use notional amounts \$1 on both the long and short side.

Our predictions come from probability estimates which allows us to refine the strategy further. Instead of predicting a binary class, our models will return a probability score p , which we will use to calculate a confidence measure $2 * p - 1$. If this value has absolute value near 1, we are very confident in our predictions of increase or decrease. However, if the value is near 0, we are not very confident in our predictions.

We multiply this confidence measure to the \$1 amounts on the long and short trades to further scale our trading strategy. The profit of the trading strategy is computed in the following way:

Let T_t be the price of an ETF for the top leg and B_t be the price of an ETF for the bottom leg. Let X_t denote the value of the trading strategy at time t . Consider just two time periods, where we enter a position at time t and exit at time $t + 1$ for demonstration purposes:

$$C_t = 2 * p_t - 1$$

$$\Delta_{T_t} = \frac{1}{T_t}$$

$$\Delta_{B_t} = \frac{1}{B_t}$$

$$X_t = \begin{cases} (T_t * \Delta_{T_t} - B_t * \Delta_{B_t}) * C_t & p_t > 0.5 \\ (B_t * \Delta_{B_t} - T_t * \Delta_{T_t}) * C_t & p_t \leq 0.5 \end{cases}$$

$$\Rightarrow X_{t+1} = \begin{cases} (T_{t+1} * \Delta_{T_t} - B_{t+1} * \Delta_{B_t}) * C_t & p_t > 0.5 \\ (B_{t+1} * \Delta_{B_t} - T_{t+1} * \Delta_{T_t}) * C_t & p_t \leq 0.5 \end{cases}$$

Then, profit $\pi_{t+1,t}$ between trading days can be easily computed on a daily level as $X_{t+1} - X_t$, for all days t we in the test set. Since the strategy is self financing, and at most can be long and short \$1, the daily profits in fact are simply the daily returns of the strategy. We can thus get annualized estimates of average return and volatility as follows (for T being a set of times that we consider in the test set):

$$\mu_{\text{annualized}} = 252 * \frac{1}{|T|} \sum_{t \in T} \pi_{t+1,t}$$

$$\sigma_{\text{annualized}} = \sqrt{252 * \frac{1}{|T| - 1} \left(\pi_{t+1,t} - \frac{1}{|T|} \sum_{t \in T} \pi_{t+1,t} \right)^2}$$

5.2 Trading Strategy Results

Table 3 details the trading strategy results. From our modes, there is a clear consensus that predictions from A rated bonds do not perform well with the trading strategy that uses the

Table 3. Trading Strategy Results

Classifier	Spread Type	Annualized Return	Annualized Volatility
Gradient Boosting	AAA	17.14%	10.39%
Gradient Boosting	AA	19.35%	10.38%
Gradient Boosting	A	8.80%	10.48%
Gradient Boosting	BBB	25.77%	12.50%
Decision Tree	AAA	11.96%	10.42%
Decision Tree	AA	11.96%	10.42%
Decision Tree	A	4.04%	10.49%
Decision Tree	BBB	18.42%	12.55%
XGBoost&LightGBM	AAA	16.47%	10.40%
XGBoost&LightGBM	AA	16.43%	10.40%
XGBoost&LightGBM	A	2.58%	10.49%
XGBoost&LightGBM	BBB	24.68%	12.50%
Ensemble	AAA	18.65%	10.38%
Ensemble	AA	18.65%	10.38%
Ensemble	A	2.26%	10.50%
Ensemble	BBB	28.06%	12.48%

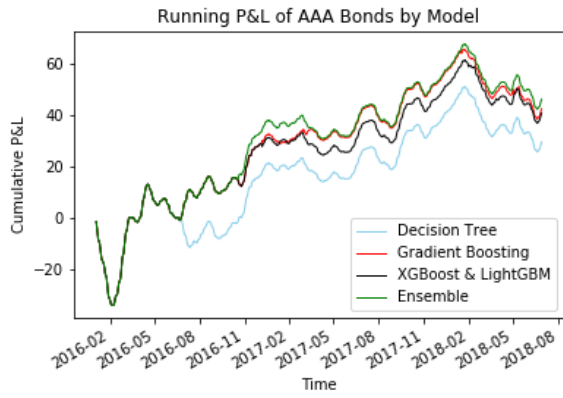
ETF QLTA (which is comprised of AAA to A rated bonds). However, the other spread types perform extremely well. The Decision Tree performed the worst, but this is to be expected as it is the simplest model we considered. Surprisingly, Gradient Boosting by itself performed exceptionally well, even with using predictions from the A rated spread. Additionally, a striking feature is that predictions on the A rated spreads all have volatility around 10.4 to 10.5 percent, and for the BBB rated spread a volatility of around 12.5 percent. The ensemble model does exceptionally well for the BBB spread predictions, and this strategy in particular produces a very appealing risk to reward profile.

In Figure 4 on the next page, we show a time series of the trading performance on the test set dates. The ensemble method is superior in the AAA, AA, and BBB bonds. However, Gradient Tree Boosting does the best for the A bonds and all of the models are much more volatile for these bonds.

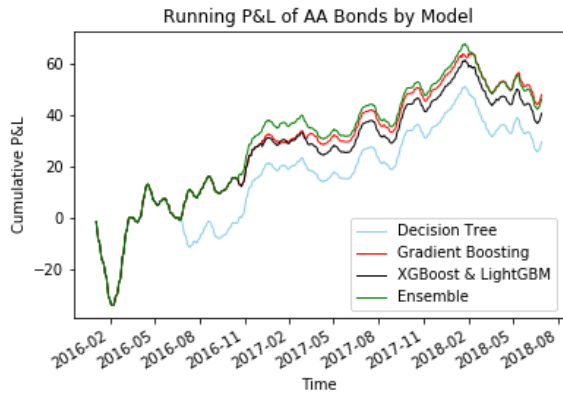
5.3 Conclusion

We investigated various tree-based machine learning methodologies to predict yield curve spread movements, and examined the profitability of trading strategies based on these predictions. We find that while it is extremely difficult to get a very high out of sample accuracy, it is possible to extract extremely lucrative and profitable trading strategies from modestly accurate predictions. By scaling predictions by the confidence of our predictions, we are able to achieve great risk to return profiles trading the ETFs: QLTA, PBBBX and IEF. By constructing long short portfolios that are self-financing, and scaling the amount by our predictive confidence, we can trade on a daily timescale by making bets on spreads widening or tightening. In the future, to get even better results, we believe researchers can investigate at the individual bond level and consider a set of liquid instruments. By building a model

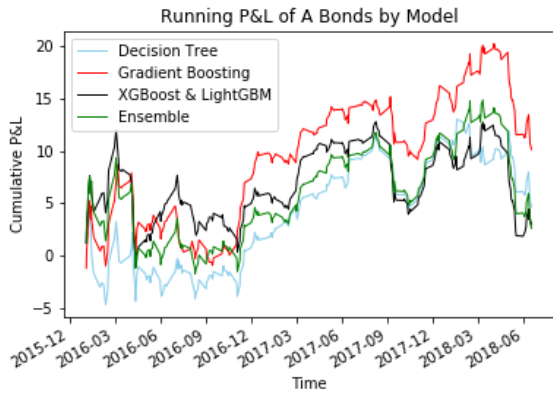
at the individual bond level, more data points will be available to train and test on, yielding even more robust results. With our conservative approach to cross validation, we still believe that our results generalize well out of sample, as shown by the performance of the trading strategies we tested.



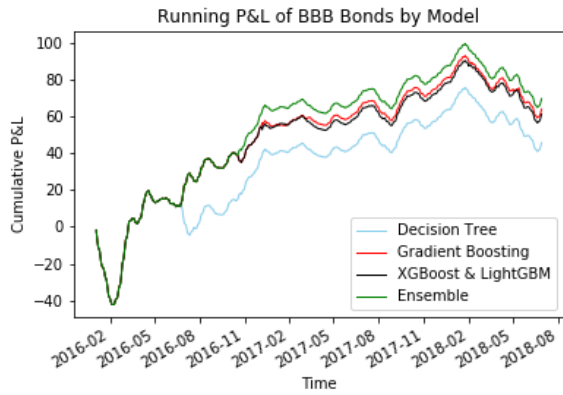
(a) AAA Bonds



(b) AA Bonds



(c) A Bonds



(d) BBB Bonds

Figure 4. Cumulative Trading P&L for Bond ETF Strategy

References

- [1] E. Neis F. A. Longstaff S. Mithal. *Corporate Yield Spreads: Default Risk or Liquidity? New Evidence from the Credit Default Swap Market*. 2005. URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.337.1116&rep=rep1&type=pdf>.
- [2] S. H. Kwan. *Firm-specific information and the correlation between individual stocks and bonds*. 1996. URL: https://www.sciencedirect.com/science/article/pii/S0304405X95008364?fbclid=IwAR1wdRP-nw6jyEBh-p050QIfBAR4qEq96okQrem37A%5C_HhjcdM-WGK6ul9hM.
- [3] Z. Kakushadze. *101 Formulaic Alphas*. 2015. URL: <https://arxiv.org/pdf/1601.00991.pdf>.
- [4] R. Hyndman. *Cross-validation for time series*. 2016. URL: <https://robjhyndman.com/hyndsight/tscv/>.
- [5] *Ensemble Methods*. URL: <https://scikit-learn.org/stable/modules/ensemble.html>.
- [6] M. Hauskrecht. *Decision Trees*. URL: <https://people.cs.pitt.edu/~milos/courses/cs2750-Spring03/lectures/class19.pdf>.
- [7] J. Dunnmon S. Ganguli. *Machine Learning Methods for Better Models for Predicting Bond Prices*. URL: <https://arxiv.org/pdf/1705.01142.pdf>.
- [8] B. Swaminathan W. Gebhardt S. Hvidkjaer. *The Cross-Section of Expected Corporate Bond Returns: Betas or Characteristics?* 2000. URL: http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.198.5763&rep=rep1&type=pdf&fbclid=IwAR1EC4--LI9ziz0EhJE36qT0m51Ux6751N35E2KeJ6rGWCvs%5C_JxC2SOHcZU.
- [9] *LightGBM*. URL: <https://lightgbm.readthedocs.io/en/latest/>.
- [10] J. Brownlee. *A Gentle Introduction to XGBoost for Applied Machine Learning*. 2016. URL: <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/>.
- [11] *Decision Trees*. URL: <https://scikit-learn.org/stable/modules/tree.html#tree>.
- [12] D. Xiu S. Gu B. Kelly. *Empirical Asset Pricing via Machine Learning*. 2018. URL: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3159577.
- [13] Federal Reserve Bank of St. Louis and US. data retrieved from FRED, <https://fred.stlouisfed.org/>. 2018.
- [14] Yahoo Finance. data retrieved from <https://finance.yahoo.com/quote/%5ESPX?p=%5E^SPX&.tsrc=fin-srch>. 2018.