

Transaction Cost Analytics for Corporate Bonds

Xin Guo *

Charles-Albert Lehalle †

Renyuan Xu ‡

December 10, 2021

Abstract

Electronic platform has been increasingly popular for executing large corporate bond orders by asset managers, who in turn have to assess the quality of their executions via Transaction Cost Analysis (TCA). One of the challenges in TCA is to build a realistic benchmark for the expected transaction cost and to characterize the price impact of each individual trade with given bond characteristics and market conditions.

Taking the viewpoint of retail investors, this paper presents an analytical methodology for TCA of corporate bond trading. Our analysis is based on the TRACE Enhanced dataset; and starts with estimating the initiator of a bond transaction, followed by estimating the bid-ask spread and the mid-price dynamics. With these estimations, the first part of our study is to identify key features for corporate bonds and to compute the expected average trading cost. This part is on the time scale of weekly transactions, and is by applying and comparing several regularized regression models. The second part of our study is using the estimated mid-price dynamics to investigate the amplitude of its price impact and the decay pattern of individual bond transaction. This part is on the time scale of each transaction of liquid corporate bonds, and is by applying a transient impact model to estimate the price impact kernel using a non-parametric method.

Our benchmark model allows for identifying abnormal transactions and for enhancing counterparty selections. A key discovery of our study is the price impact asymmetry between customer-buy orders and consumer-sell orders.

Keywords: Bond liquidity, transaction costs analysis, price impact, Enhanced TRACE, regression analysis, regularization method, data-driven decision making

*Department of Industrial Engineering and Operations Research, University of California, Berkeley, USA. Email: xinguo@berkeley.edu.

†Abu Dhabi Investment Authority (ADIA) and Imperial College London. Email: c.lehalle@imperial.ac.uk.

‡Epstein Department of Industrial and Systems Engineering, University of Southern California, US & Mathematical Institute, University of Oxford, UK. Email: renyuanx@usc.edu.

1 Introduction

Corporate bonds are critical to firm finance and play an important part in asset management [Nagel, 2016]. Different from equity shares which are mostly traded via order books available in multilateral trading facilities, corporate bonds are mainly traded via bilateral mechanisms [Fermanian et al., 2015] due to limited available electronic platforms [Linciano et al., 2014]. Even with the introduction of the TRACE reporting system in US in June 2002 and the establishment of MiFID 2 for electronic bond tradings in Europe in January 2018, bond trading remains far less transparent than equity trading [Bessembinder et al., 2008].

After the 2008 financial crisis, the macroprudential regulation requires more transparency of corporate bond trading to reduce information asymmetry [Hendershott and Madhavan, 2015] between intermediaries and their clients, leading to an increase in capital requirement and in turn preventing banks from taking large inventories as before [Wilson et al., 2014]. This lower inventories, combined with the requirement of more transparency, pushes banks and dealers towards flow driven business via electronification [Harris, 2015].

In this new trading environment, asset managers, lacking pricing tools and private databases enjoyed by maker-dealers, have to assess the quality of corporate bonds execution via Transaction Cost Analysis (TCA). Details of TCA are then shared with the portfolio managers of investment firm to review market liquidity and for allocation and hedging purposes [Albanese and Tompaidis, 2008].

TCA is difficult as there is a dire lack of benchmark [Collins and Fabozzi, 1991] for bond trading, unlike equity trading where the bid-ask spread is an obvious and easy choice for the benchmark. Instead, TCA needs to break down costs of a particular bond trading according to brokers from all possible execution venues in fragmented markets, including order books, requests-for-quotes, voice, dark pools, and block discovery mechanisms.

Our work. The goal of this paper is to establish a TCA benchmark in bond trading for retail investors. That is, we take the standpoint of an individual investor to evaluate the execution performance of each transaction.

Our TCA is based on the Enhanced TRACE dataset from 2015 to 2016. We assume that TCA consists of the bid-ask spread measuring the cost of illiquidity and the mid-price move measuring the impact of an individual trade.

Our analysis starts with a preliminary step of estimating the initiator of a bond transaction (Section 3.1). Initiator, currently missing from the TRACE database, indicates whether a given transaction is buyer-initiated or seller-initiated. This estimated initiator of each trade enables us to estimate the bid-ask spread and the mid-price dynamics (Section 3.1).

With this preliminary step, the first part of our study is to identify the most important features for corporate bonds and to compute the expected average trading cost (Section 4). This study is on the time scale of weekly transactions, and is carried out via comparing several regularized regression models including the two-step Lasso, the Ridge regression, and the two-step Elastic Net regression. The response variable in the regression analysis is the estimated bid-ask spread from the preliminary analysis.

Our regression approach manages to select features for corporate bonds that are consistent with existing works, including the volatility of the bond price, the number of years from the issue date, and the *activities* of the bond characterized by the number of trades and the traded amount (in dollars) per week. In addition, the number of trades and the traded amounts are found to play two opposite roles: the larger the amount traded in dollars, the smaller the bid-ask spread; the more trades (for the same amount in dollars), the larger the bid-ask spread. It is worth mentioning that the R^2 value obtained from our regression analysis ranges from 0.50 to 0.60, whereas the R^2 in existing works via regressions varies from 0.05 to 0.20 in [Hendershott and Madhavan, 2015], 0.30 to 0.50 in [Edwards et al., 2007], and 0.50 to 0.80 in [Dick-Nielsen et al., 2012].

The second part of our TCA is using the estimated mid-price dynamics to investigate the amplitude of its price impact and the price decay pattern of individual bond transaction. This study (detailed in Section 5) is on the time scale of each transaction of liquid corporate bonds. It is done by applying a transient impact model (TIM) to estimate the price impact kernel via a non-parametric method. The transient impact functions estimated in our study is found to share several important characteristics with those in the equity market:

- a price jump when the trade occurs,
- a price decay after the initial jump,

- and the stabilization at a “*permanent level*” higher than the initial price: this permanent impact can be interpreted as the informational content of the trade.

In addition, we discover an asymmetry in the amplitude of the initial price jump: *buy-initiated transactions with more instantaneous impact than sell-initiated transactions* on corporate bonds. Note that such an asymmetry, not present in the equity market, has also been reported in [Hendershott and Madhavan, 2015] and [Ruzza, 2016] for corporate bonds.

Existing works on TCA of corporate bonds. Empirical studies on transaction costs of corporate bonds are mostly from post-TRACE, as it was difficult for retail investors to obtain data in the pre-trace era. The post-TRACE trade reporting obligation started in US in July 2002. Because of the exogenous shock from the entry of TRACE, a number of earlier works [Goldstein et al., 2007, Ruzza, 2016, Bessembinder et al., 2008] focused on the early years of its introduction in order to identify the cost effect of this transparency. Another family of post-TRACE studied the influence of adopting electronic and multilateral trading [Hendershott and Madhavan, 2015] and decreasing borrowing costs from 2004 to 2007 [Asquith et al., 2013].

All these studies reached similar conclusions that the trading costs of corporate bonds decreased on average over the last twenty years. The main proxy for transaction costs adopted in these works was the (expected) bid-ask spread [Glosten and Milgrom, 1985], [Edwards et al., 2007]. Their main statistical approach was ordinary least square (OLS) regression to account for bond-specific or context-driven variations. The explanatory variables in these studies [Dick-Nielsen et al., 2012, Edwards et al., 2007, Goldstein et al., 2007, Eom et al., 2004] were the coupon, the maturity date, the number of years to maturity, the volatility, the risk-free rate, the expected recovery rate of the company, and the probability of default.¹ See Table 12 in Appendix A for a summary of the dataset and the years of bonds studied in these empirical analysis.

Existing works on asymmetric price impact. The price impact asymmetry between customer-initiated buy orders and customer-initiated sell orders in corporate bond market has been documented in the literature. For example, this asymmetry is reflected in the regression of Table IV of [Hendershott and Madhavan, 2015] since the coefficients of the buy and sell orders are not of the same amplitude for over-the-counter (OTC) trades (but not for electronic trades). Such an asymmetry was also reported in Figure 15 of [Mizrach, 2015] and Table 1 of [Ruzza, 2016]. The former plotted the yearly average price change after five trades from 2003 to 2015, with the impact of buys around 25% more than the impact of sells. The latter suggested that the average difference between the price of a transaction and the average price of the day is 56bp to 33bp for institutional buyers and -25bp to -21bp for institutional sellers on TRACE data from 2004 to 2012.

Our work is different from existing studies of average transaction costs with OLS, as we apply regularized regression models to select features within a broader class of candidate features. Unlike previous studies on the price impact asymmetry from a static point of view, we investigate the price impact curve of individual trade and analyze its asymmetry in a dynamic setting, characterizing both the amplitude of the price impact and the decay via TIM models. Moreover, the analysis and methodology presented in this paper are general and can be applied to conduct TCA on other datasets including Standard TRACE.

2 Data Processing and Bond Selection

Enhanced TRACE. TRACE, an acronym for *the Trade Reporting and Compliance Engine*, is the FINRA-developed mechanism that facilitates the mandatory reporting of over-the-counter secondary market transactions in eligible fixed income securities. TRACE database contains some useful (though limited) information and has been used for empirical studies by [Dick-Nielsen, 2014] and [Harris, 2015].

The main difficulty working with TRACE is the lack of information on the liquidity offer. For example, there are neither quotes, nor bid prices, nor ask prices. Instead, only final transactions are recorded, together with the type of the transaction: dealer-to-dealer, dealer-to-customer, or customer-to-customer. Therefore, besides

¹Note that [Biais and Green, 2019] did not perform any linear regression, but relied on descriptive statistics, probably due to the lack of explanatory variables available during this period.

TRACE, we also rely on Thomson Reuters to retrieve information on the bonds traded, such as the amount issued, the coupon rate, the sector information, rating information, and Libor and Overnight Indexed Swap rate. In addition, we obtain the outstanding amount through the Mergent Fixed Income Securities Database (FISD).

There are two types of TRACE datasets, Standard TRACE dataset and Enhanced TRACE dataset. Both TRACE datasets contain corporate bond transactions. The difference is that transactions are available on Standard TRACE with a delay of two weeks, with the volume of the transaction capped at 1MM for high yield bonds and 5MM for investment grade bonds. The Enhanced TRACE dataset has uncapped volumes, with transactions available with a delay of six months. There is a separate dataset provided by FINRA for monthly price, return, coupon, and yield information for all corporate bonds traded since July 2002.

Our study is primarily based on the Enhanced TRACE dataset. We use the nontruncated transaction volumes on Enhanced TRACE along with other information from FISD and Thomson Reuters to construct the estimation of the bid-ask spread (Section 3). It is worth noting that the Enhanced TRACE dataset and the Standard TRACE dataset yield insignificant differences in terms of the estimation of the expected bid-ask spread, as shown in Section C.4.

Data processing. The data used in our study is from January 1, 2015 to December 31, 2016, obtained from Wharton WRDS. During this period, there are 34,809,405 original trade reports, 390,193 reports of trade cancellations (approximately 1.1 percent of all original trade reports), 497,249 corrected trade reports (about 1.4 percent), and 28,005 reports of trade reversals. Trade reversals are transactions that have been changed after more than 20 days since they were initially recorded. Occasionally there are multiple correction records for the same original trade and cancel records that cancel previously corrected trades. There are 54,885 CUSIP²-days spread over 656 calendar days, many of which are weekends and holidays. The CUSIP-days are computed by counting all the trade days over all the CUSIP bonds.

In particular, for each transaction of a bond, one can recover from Enhanced TRACE the following information:

- t_k^b : the timestamp for the k th transaction of bond b ;
- P_k^b : the price of the k th transaction of bond b ;
- V_k^b : the volume of the k th transaction of bond b ;
- the side of the dealer-to-customer transaction: customer buy order or customer sell order.

The data cleaning procedure combines the approaches of [Dick-Nielsen, 2014] and [Harris, 2015], as detailed in Appendix B.2. In total, about 17.50% reports are filtered out from the original Enhanced TRACE dataset. Among all the remaining 28,719,813 records, 14,071,375 (49%) are dealer-to-customer trades and the remaining 14,648,438 (51%) are trades between dealers. These statistics are summarized in Table 15.

After the data cleaning, appropriate bond selection is necessary to facilitate the analysis of transaction costs, for both the regression and price impact analysis.

Bond selection for regression analysis There are two types of bonds for regression analysis, investment grade bonds and high yield bonds, which are picked from the standard universe of U.S. corporate bonds. The investment grade bonds are selected from iShares iBoxx Investment Grade Corporate Bond ETF, and the high yield bonds from the components of iShares iBoxx High Yield Corporate Bond ETF. There are 1,033 current holdings of the former, among which 538 bonds have more than one transaction recorded in Enhanced TRACE during the time period of Jan 1, 2015 to Dec 31, 2016. There are 1,575 current holdings of the latter, 1485 of which have transaction records during the same period. Moreover, there are 30 bonds that belong to both iShares iBoxx High Yield Corporate Bond ETF and iShares iBoxx Investment Grade Corporate Bond ETF. The rating levels of all these 30 bonds have been adjusted since issuance. Hence, there are a total of 1,993 bonds for the regression analysis. These selected bonds consist of 31.05 % of the total 14,071,375 customer-to-dealer

²CUSIP stands for *Committee on Uniform Securities Identification Procedures*.

reports from all bonds. Table 15 reports this selection as “Selection LR” and Table 1 reports the statistics of these selected bonds.

	Total	Investment Grade	High Yield	Bonds with rating changes
Number of trades	4,371,363	3,102,791	1,109,177	159,395
Number of customer buy	2,549,932	1,834,873	623,839	91,220
Number of customer sell	1,821,371	1,267,888	485,308	68,175
Total trading volume (billion)	3401.76	2438.26	850.46	113.04
Avg trading volume	778,202.926	785,834.74	766,775.08	709,163.43
Avg price	102.05	104.26	97.81	88.58
Prop volume of customer buy	55.3 %	56.8 %	51.7 %	58.2 %
Prop volume of customer sell	44.7 %	43.2 %	48.3 %	41.8 %

Table 1: Description of selected 1,993 bonds for regression (dealer-customer trades).

Note that our analysis throughout the paper focuses on trades between customers and dealers, which are statistically different from trades between dealers. The study for the latter requires the initiator analysis for dealer-to-dealer trades, which is infeasible to estimate given the current information from the database.

Bond selection for price impact analysis. Given that the calculation of price impact curves requires a higher trading frequency, out of all the 1,993 bonds for regression analysis, the top-200 traded bonds (in terms of number of transactions) are selected to compute price impact curves. See the “Selection for PI” step in Table 15 and the statistical summary of these 200 bonds in Table 2. Tables 1 and 2 show that these top-200 bonds account for 32% of the total number of transactions among the 1,993 bonds. Note that among the 30 bonds with a rating level adjustment, 13 belong to the top-200 traded bonds.

	Total	Investment Grade	High Yield	Bonds with rating changes
Number of trades	1,404,507	980,005	309,777	116,561
Number of customer buy	836,825	588,128	180,378	68,319
Number of customer sell	567,682	390,537	128,903	48,242
Total trading volume (billion)	605.5	430.83	107.70	66.98
Avg trading volume	431,119.48	440,217.83	348,241.14	574,636.26
Avg price	100.08	103.32	94.58	87.55
Prop volume of customer buy	52.2 %	53.12%	50.1 %	49.6 %
Prop volume of customer sell	47.8 %	46.88%	49.9 %	50.4 %

Table 2: Description of the selected 200 bonds for price impact analysis (dealer-customer trades).

3 Preliminary Analysis

Two key components for TCA are the bid-ask spread and the mid-price dynamics, for which it is necessary to identify the riskless-principle-trades (RPTs) and the initiator of a transaction.

3.1 RPT and Initiator

In the TRACE database, the information of the initiator, i.e., whether a given trade is a buyer-initiated or a seller-initiated, is missing. In fact there is a substantial fraction of transactions between dealers and customers [Harris, 2015] where the dealer has found two clients and put herself in between the transactions. These transactions are called *riskless principal trades* (RPT) since the dealer does not take any inventory risk by matching two clients. Consequently, it is not possible to recover the initiator of the RPT because there is no information on which of the two clients has initiated the trades.

Our first step is therefore to identify these RPTs. See Table 13a and Table 13b for the statistics of RPTs and non-RPTs. Table 14a and Table 14b report the potential RPTs and non-RPT dealer-customer trades of the top-200 traded bonds. See also Appendix B.1 for a detailed literature review on RPT.³

After identifying and removing all the potential RPTs, we consider the transaction initiated by the client. We define the *sign of the transaction* ϵ_k^b as +1 (i.e., “buy”) if a client buys from a dealer and ϵ_k^b as -1 (i.e., “sell”) if a client sells to a dealer. When it is not possible to determine the sign of a trade as in the above RPT case, we assign ϵ_k^b to be zero.⁴

Figure 1 reports the auto-correlation of the order signs with lag 20. It suggests that with high confidence level, there is a persistent positive auto-correlation among order signs which delays very slowly. In comparison, [Bouchaud et al., 2009] showed that the sign of market orders on equity market is strongly correlated in time.

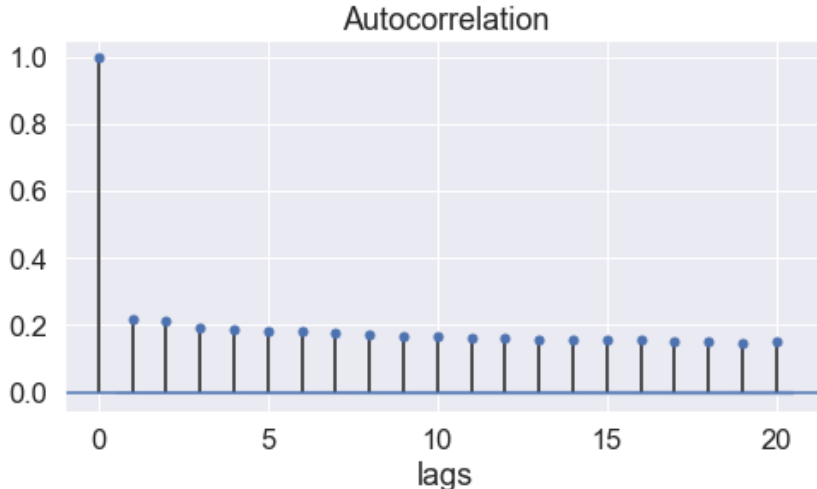


Figure 1: Auto-correlation of the signs.

3.2 Bid-ask spread and mid-price estimation

After identifying the initiator of each trade, we now analyze the two essential building blocks for TCA: the bid-ask spread and the mid-price dynamics.

To start, let us find two consecutive trades that have opposite signs $\epsilon_{k+1}^b = -\epsilon_k^b$ with $\epsilon_k^b \neq 0$ and are sufficiently close in time (i.e., $|t_{k+1}^b - t_k^b| < \Delta t$). Let us then define the estimate of bid-ask spread in absolute value as:

$$\psi_{k+1}^b := (P_{k+1}^b - P_k^b) \cdot \epsilon_{k+1}^b, \tag{1}$$

³Our percentage of RPTs is lower than that reported in [Harris, 2015], partly because of different datasets with different time periods. [Harris, 2015] used the Standard TRACE dataset from April 1, 2014 to March 31, 2015, where the markers (“1MM+” and “5MM+”) for larger trades assign the same value to many large trades. Finally, we only count the RPTs for a subset of bonds whereas [Harris, 2015] estimated the PRTs for a larger set of bonds.

⁴Note that there are other methods to estimate the sign of a trade when quote price is not available. See [Holthausen et al., 1987] for the tick test and [Lee and Ready, 1991] for the inverse tick test.

with P_k^b the price of the k -th transaction of bond b . We next estimate the mid-price at t_k as:

$$M_k^b := P_k^b - \epsilon_{k+1}^b \frac{\psi_{k+1}^b}{2}, \quad (2)$$

and define the bid-ask spread in basis point (relative value) by

$$s_{k+1}^b := \frac{\psi_{k+1}^b}{M_k^b} \times 10000. \quad (3)$$

Note that the choice of Δt is 5-minute, which is largely due to the low trading frequency of the corporate bond market.⁵ Consequently, only 15.6% of the transactions are used to calculate the bid-ask spread among bonds that are selected from Section 2.

We next check the reliability and stationarity of the estimated bid-ask spread.

Reliability of the estimated spread. We compare the estimated bid-ask spread with the one computed using bid and ask quotes provided by Composite Bloomberg Bond Trader (CBBT) for those bonds that are available in both the CBBT and Enhanced TRACE data sets. CBBT is a composite price based on the most relevant executable quotations on FIT, Bloomberg’s Fixed Income Trading platform. The CBBT pricing source provides average bid-ask prices based on executable quotes listed on Bloomberg’s trading platform. [Fermanian et al., 2016] used the CBBT data as a measure of bond liquidity. We only have access to quote price data from Bloomberg CBBT from June 1, 2015 to May 31, 2016 (12 months) for 2,361 investment grade bonds that belong to the iboxxIG universe, among which we have identified 1,401 bonds with records in both the Bloomberg CBBT database and the Enhanced TRACE subset.

Figure 2 below shows the plot for the empirical distribution of the spread from CBBT and the estimated spread from Enhanced TRACE for two arbitrarily chosen bonds, whose statistics are reported in Table 3. It is noticeable (and expected) that the CBBT spreads are larger than those estimated from real trades available in Enhanced TRACE. As [Fermanian et al., 2016] pointed out, CBBT bid-ask spread estimates are based on quotes, and not on real transactions. As a consequence they include quotes that are not attractive enough (i.e., not small enough) to trigger a transaction. Since the bid-ask spread is the first component of implicit transaction costs, trades occur when they are smaller than the average bid-ask spread.

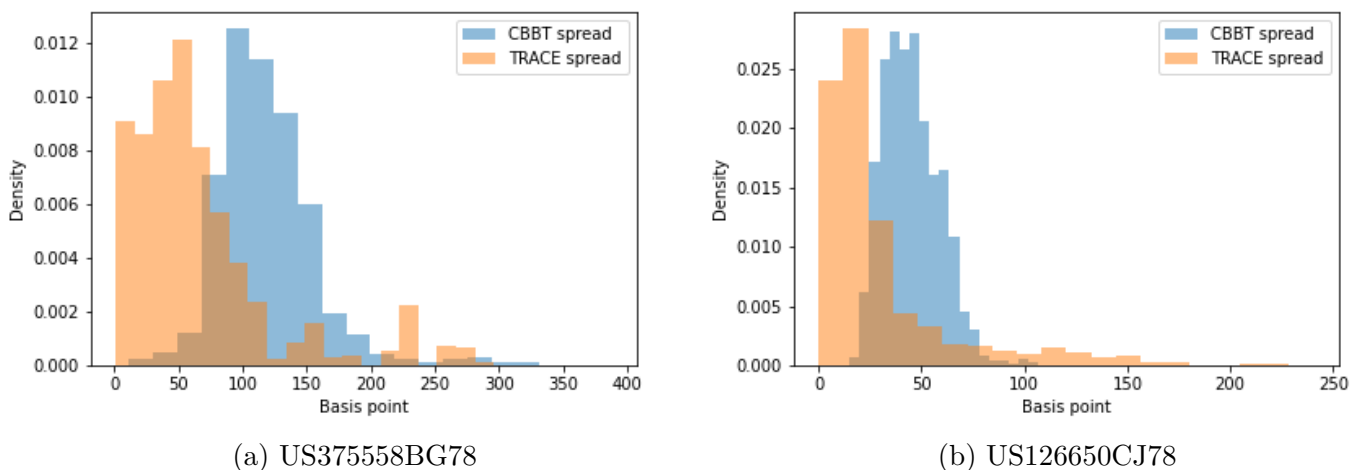


Figure 2: Empirical distributions of the spread (Eqn (3)).

⁵ $\Delta t < 1$ -minute is not realistic for the (illiquid) corporate bond market; meanwhile the choice of $\Delta t \geq 10$ -minutes is infeasible due to lack of data. Comparing results with $\Delta t = 4$ -minutes, $\Delta t = 6$ -minutes, and $\Delta t = 5$ -minutes leads to the choice of $\Delta t = 5$ -minute.

	Coupon	Amount outstanding (USD)	Average spread (bp)		Daily average number of updates	
			CBBT	TRACE	CBBT	TRACE
US375558BG78	4.6%	1,000,000,000	118.78	66.24	2336.49	1.63
US126650CJ78	2.8%	2,750,000,000	45.12	29.88	3406.20	11.50

Table 3: Spread comparison.

Stationarity of the bid-ask spread. We next check the consistency of the two approaches via a stationarity test on the ratio of the two estimates: the CBBT bid-ask spread and our trades-based estimates. Denote $s_{b,w}^{\text{CBBT}}$ as the average spread for bond b over the period w taken from Bloomberg CBBT and $s_{b,w}^{\text{TRACE}}$ as the average of the estimated bid-ask spread for bond b in period w from Enhanced TRACE (see Eqn (3)). Similarly, define $R_{b,w} = s_{b,w}^{\text{CBBT}}/s_{b,w}^{\text{TRACE}}$ the ratio between these two spreads for bond b over the period w .

First of all, note that the empirical estimate of the bid-ask spread using Enhanced TRACE transactions are smaller than the CBBT ones: the average ratio is between 0.9 and 1 and its median is between 0.7 and 1, as summarized in Table 4.

Month	1	2	3	4	5	6
Number of observations	1027	823	820	847	890	892
25%	0.7165	0.655	0.647	0.666	0.655	0.643
Median	0.971	0.931	0.899	0.900	0.852	0.873
75%	1.238	1.213	1.224	1.182	1.144	1.19
Mean	0.970	0.956	0.951	0.941	0.920	0.931
Month	6	7	8	9	10	12
Number of observations	939	914	1069	1062	1040	1047
25%	0.705	0.717	0.670	0.643	0.598	0.618
Median	0.963	1.000	0.893	0.830	0.779	0.834
75%	1.289	1.309	1.182	1.090	1.056	1.130
Mean	0.994	1.019	0.941	0.889	0.854	0.893

Table 4: Statistics of the ratios.

We also split the year 2016 into eleven groups of two consecutive months and check if this ratio is stationary from one period of the two months to the other. We use two tests for the stationarity of $R_{b,w}$. The first is the one-way ANOVA test and the second is the Kruskal-Wallis H-test. The former tests the stationarity of the mean and the latter tests the stationarity of the median. The mathematical formulations and definitions of the ANOVA test and the Kruskal-Wallis H-test are provided in Appendix C.1.

Month	1 and 2	2 and 3	3 and 4	4 and 5	5 and 6	6 and 7
ANOVA	4.59	0.10	0.39	1.71	0.52	15.32
(P-value)	0.032	0.751	0.533	0.191	0.472	0.000
H-test	3.564	0.163	0.154	1.890	0.296	15.16
(P-value)	0.038	0.686	0.6947	0.169	0.586	0.000
Month	7 and 8	8 and 9	9 and 10	10 and 11	11 and 12	
ANOVA	2.21	25.54	14.15	6.72	7.81	
(P-value)	0.136	0.000	0.0001	0.01	0.005	
H-test	2.30	23.78	12.03	8.54	5.82	
(P-value)	0.129	0.000	0.0005	0.003	0.015	

Table 5: Results of ANOVA and Kruskal-Wallis H-tests.

Table 5 summarizes the results of both the ANOVA test and the Kruskal-Wallis H-test. With a 99% confidence level, we accept (cannot reject) the null hypothesis (i.e., the ratios are stationary over time) in both the ANOVA test and the Kruskal-Wallis test for 7 of the total 11 comparisons. We will thus use this estimated bid-ask spread in all subsequent analyses because it can be operated over years of data using Enhanced TRACE, where CBBT is costly to obtain and linked to a private procedure owned by Bloomberg. Nevertheless, these stationarity tests imply that a large investor using CBBT estimates could rely on the methodology presented thereafter and apply a ratio to interpret our results in terms of “units” of CBBT.

We have found similar results in terms of reliability and stationarity for the estimated mid-price dynamics, with details skipped here to avoid repetition.

4 Regularized Regression Analysis for Bid-Ask Spread

In this section, we use regularized regressions to identify the key features that drive the bid-ask spread, which provides the estimated cost for investors who needs to move from one side (e.g., the buy side) to another side (e.g., the sell side).

We will exploit several regularized regression models including OLS, two-step Lasso, Ridge, and two-step Elastic Net regressions, along with a K -fold cross-validation method, to identify the most significant features and associated parameters for these models.

As illustrated in Section 2, there are total of 1,993 bonds selected for this regression analysis, along with a total 152,408 (weekly) samples processed from the Enhanced TRACE dataset during January 2015 and December 2016. Our regression analysis is performed on a weekly basis, with the weekly average bid-ask spread computed according to Eqn. (3) serving as the response variable.

4.1 Review of Methodologies

We start by reviewing the necessary notations and steps for the regression analysis that will be used throughout the paper.

OLS. OLS assumes that the regression function is in linear form. That is, given $\mathbf{Y} := (y_1, y_2, \dots, y_n) \in \mathbb{R}^n$ the vector of n observations of independent variables, and $\mathbf{X} := (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_{w-1})$ with covariates $\mathbf{1} \in \mathbb{R}^n$ and $\mathbf{x}_i \in \mathbb{R}^n$ ($i = 1, 2, \dots, w - 1$), OLS is to find:

$$\hat{\boldsymbol{\theta}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^w} \left\{ \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \right\}. \quad (4)$$

In an OLS, R^2 is used to measure the goodness of fit for the model. Meanwhile, an associated p -value indicates the significance level of the feature.

Two-step Lasso. The first step is to use Lasso regression to select the covariates by solving the following optimization problem:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^w} \left\{ \frac{1}{N} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \sum_{j=1}^{w-1} |\theta_j| \right\}. \quad (5)$$

Here a fixed constant λ , called the *tunable hyperparameter*, controls both the size and the number of coefficients: a higher value of λ leads to a smaller number of covariates in the linear model. In the second step, an OLS with only the selected covariates is applied [Belloni and Chernozhukov, 2013]. That is, given the Lasso estimator $\hat{\boldsymbol{\theta}}_l^\lambda$ in (5), the subsequent OLS refitting is to find $\bar{\boldsymbol{\theta}}_l^\lambda$ such that:

$$\bar{\boldsymbol{\theta}}_l^\lambda \in \arg \min_{\text{supp}[\boldsymbol{\theta}] = \text{supp}[\hat{\boldsymbol{\theta}}_l^\lambda]} \left\{ \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \right\}. \quad (6)$$

$\bar{\boldsymbol{\theta}}_l^\lambda$ is thus called the estimator for the LSLasso (least-squares Lasso), also known as post-Lasso. This two-step Lasso estimation procedure has been shown to produce a smaller bias than Lasso for a range of models [Belloni and Chernozhukov, 2013], [Lederer, 2013], and [Chételat et al., 2017].

Ridge regression. The penalty term in the Ridge regression is of the L_2 norm. That is, for a fixed hyperparameter λ , Ridge regression is to solve for:

$$\hat{\boldsymbol{\theta}}_r^\lambda \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^w} \left\{ \frac{1}{N} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \lambda \sum_{j=1}^{w-1} \theta_j^2 \right\}. \quad (7)$$

Two-step Elastic Net regression. Elastic Net (EN) regression, introduced in [Zou and Hastie, 2005], is a hybrid of Lasso and Ridge. That is, for a fixed hyperparameter (λ, α) with $\alpha \in [0, 1]$, EN is to solve for:

$$\hat{\boldsymbol{\theta}}_e^\lambda \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^w} \left\{ \frac{1}{N} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + \alpha \lambda \sum_{j=1}^{w-1} |\theta_j| + (1 - \alpha) \lambda \sum_{j=1}^{w-1} \theta_j^2 \right\}. \quad (8)$$

Note that the Lasso regression is recovered from Eqn. (8) by taking $\alpha = 1$ and the Ridge regression is recovered by taking $\alpha = 0$. Similar to the two-step Lasso, the second step of the Two-step Elastic Net is to fit an OLS with only the selected covariates. That is, given the EN estimator $\hat{\boldsymbol{\theta}}_e^\lambda$ in (8), the subsequent OLS refitting is to find $\bar{\boldsymbol{\theta}}_e^\lambda$ such that:

$$\bar{\boldsymbol{\theta}}_e^\lambda \in \arg \min_{\text{supp}[\boldsymbol{\theta}] = \text{supp}[\hat{\boldsymbol{\theta}}_e^\lambda]} \left\{ \|\mathbf{Y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \right\}. \quad (9)$$

Cross-validation. In all three regularized regression models, the selection of hyperparameters is by the standard K -fold cross-validation approach to improve the predictive power of the model. That is, the dataset is randomly divided into K subsets. Each time, one of the K subsets is used as the validation set and the remaining $K - 1$ subsets form a training set. In this approach, every data point is in a validation set exactly once and in a training set $K - 1$ times. The variance of the resulting estimate is reduced as K increases.

Out-of-sample test. With the hyperparameters selected from the cross-validation step, the coefficients for a regression model are estimated using the training and validation datasets. The performance of the regression model is then evaluated with the test dataset. For financial applications, the time period of the test dataset needs to be after those of the training and validation datasets to ensure the information adaptiveness.

Given a test dataset $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$ with size m and $\tilde{\mathbf{Y}} = (\tilde{y}_1, \dots, \tilde{y}_m)$, the following relative error function is used as the criterion to measure the performance:

$$\text{Relative err} = \frac{1}{m} \sum_{j=1}^m \frac{|\tilde{y}_j - \hat{y}_j|}{|\tilde{y}_j|}, \quad (10)$$

where \tilde{y}_j is the true label and \hat{y}_j is the label predicted from the regression model for test sample j ($j = 1, 2, \dots, m$).

4.2 Features for Regression Analysis

The features in the regression consist of two categories. One category concerns bond information, including time to maturity date, time since issued date, coupon rate, amount outstanding, and duration. The other category focuses on trade information including average transaction price, volatility, proportion of customer-buys (sells), LIBOR-OIS rate, and the 5-year treasury rate during the given week. More specifically, we consider:

- **Volatility:** calculated from the trade price. For bond b , assume there are n trades in week w . Recall P_j^b as the trade price of the j^{th} transaction ($j = 0, 1, 2, \dots, n$) of bond b . Denote the log return $r_i^b = \log(\frac{P_i^b}{P_{i-1}^b})$ ($i = 1, 2, \dots, n$) and the average return $\bar{r}^b = \sum_{i=1}^n r_i^b / n$. Then the volatility in week w is:

$$\sigma_b = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (r_i^b - \bar{r}^b)^2} \cdot 100.$$

Notice that n may vary from bond to bond and from week to week.

- **Number of trading days:** the number of days that bond b is traded during the week.
- **Log(zero trade days):** the log of the number of days that bond b is not traded during the week.
- **Proportion of buy/sell number:** estimated by counting the number of customer-buy orders and the number of customer-sell orders in week w and calculating the proportion of buys and sells for each bond b .
- **Proportion of buy/sell volume:** estimated by taking the total volume (in dollars) for customer-buy orders and customer-sell orders in week w and calculating the proportion of buys and sells for each bond b .
- **Trading activity:** the log of the *number of trades* in the week.
- **Total volume:** the weekly total trading volume in dollars of both customer-dealer trades and dealer-dealer trades.
- **Average price:** the weekly average trade price in dollars.
- **Coupon:** annual coupon payments paid by the issuer relative to the bond's face or par value. The coupon rate is the yield the bond paid on its issue date. This yield changes as the value of the bond changes, thus giving the bond's yield to maturity.
- **Duration:** an approximation of a bond's price sensitivity to changes in interest rates which is defined as:

$$D^b = \sum_t \frac{PV(C_t^b)}{\sum_t PV(C_t^b)} \times t$$

for bond b , where C_t^b is the cash flow on date t , $PV(C_t^b)$ is its present value (evaluated at the bond's yield), and $\sum_t PV(C_t^b)$ is the total present value of the cash flow, which is equal to the bond's current price.

- **Years to maturity:** the time to maturity date calculated in years.
- **Years since issuance:** the time since issued date counted in years.

- **Amount outstanding:** the principal amount outstanding of a bond; sometimes referred to as the notional amount.
- **Turnover:** the volume of bonds traded relative to the total volume of outstanding bonds. The inverse of the turnover can be interpreted as the average holding time of the bond. For instance, a turnover of one implies an average holding time of about two weeks.
- **LIBOR-OIS rate:** London inter-bank offer rate (LIBOR) is the rate at which banks indicate they are willing to lend to other banks for a specified term of the loan: Overnight indexed swap (OIS) rate is the rate on a derivative contract on the overnight rate. The term LIBOR-OIS spread is assumed to be a measure of the health of banks because it reflects the default risk associated with lending to other banks. In this analysis, the 1-month LIBOR-OIS rate is used to indicate the bank health condition over time.
- **Indicator of high yield (HY) or investment grade (IG) bond:** indicator of the bond .
- **Indicator of different sectors:** including nine different sectors such as basic materials sector (S1), communications sector (S2), consumer & cyclical sector (S3), consumer & non-cyclical sector (S4), energy sector (S5), financial sector (S6), industrial sector (S7), technology sector (S8), and utilities sector (S9).

Table 6 provides descriptive statistics of these response variables and features.

	Mean	std.	q-25%	Median	q-75%
Bid-ask spread (bp)	68.79	69.1	23.14	45.82	90.91
Volatility	0.785	0.582	0.388	0.646	1.039
Trading activity	1.28	0/35	1.04	1.26	1.51
Log(\$ traded volume)	6.90	0.63	6.50	6.94	7.33
Nbe trading days	4.41	0.80	4.00	5.00	5.00
Log(zero trade days)	0.12	0.18	0.00	0.00	0.30
Prop. nbe buy	0.45	0.18	0.33	0.45	0.57
Prop. nbe sell	0.55	0.18	0.43	0.55	0.67
Prop. \$ buy	0.49	0.25	0.32	0.50	0.66
Prop. \$ sell	0.51	0.25	0.34	0.51	0.69
Avg. price	103.69	10.93	100.06	102.87	107.93
Coupon	4.69	1.78	3.38	4.65	5.93
Duration	5.72	4.21	2.75	6.0	8.0
Years to maturity	8.23	7.91	3.0	6.0	8.0
Years since issuance	3.99	2.93	2.0	3.0	5.0
Turnover ($\times 10^{-2}$)	1.80	6.10	0.30	0.80	1.90
LIBOR-OIS	0.21	0.09	0.14	0.20	0.25

	Mean	std.		Mean	std.
High yield	0.27	0.44	Invest. grade	0.73	0.44
S: Basic Material	0.03	0.18	S: Communications	0.15	0.38
S: Consumer, Cyclical	0.11	0.31	S: Consumer, Non-cyclical	0.14	0.35
S: Energy	0.16	0.37	S: Financial	0.27	0.44
S: Industrial	0.07	0.25	S: Technology	0.07	0.25
S: Utilities	0.01	0.11			

Table 6: Statistics of the response variable and the features.

Hyperparameter selection. Specific to the regression models aforementioned, denote μ as the parameter for one of the regression models (for example, $\mu = (\lambda_e, \alpha)$ for EN). Partition in log-scale is used for m different

hyperparameter values for $\boldsymbol{\mu}$ and the training dataset is divided into K folds for cross-validation. For each leave-out fold i , $R_i^2(\boldsymbol{\mu})$ is computed with regression coefficients calculated using the other $K - 1$ folds. Hence for each λ , there is an empirical distribution of $\tilde{R}^2(\boldsymbol{\mu}) = \{R_i^2(\boldsymbol{\mu}), i = 1, 2, \dots, K\}$. Denote $\widehat{R}^2(\boldsymbol{\mu})$ and $\sigma_{R^2}(\boldsymbol{\mu})$ as the mean and standard deviation of the empirical distribution with parameter $\boldsymbol{\mu}$, and define the confidence interval by:

$$\mathfrak{J}_1(\boldsymbol{\mu}) = \left[\widehat{R}^2(\boldsymbol{\mu}) - \frac{\sigma_{R^2}(\boldsymbol{\mu})}{\sqrt{K}}, \widehat{R}^2(\boldsymbol{\mu}) + \frac{\sigma_{R^2}(\boldsymbol{\mu})}{\sqrt{K}} \right]. \quad (11)$$

Then $\boldsymbol{\mu}$ is picked such that the number of $\tilde{R}^2(\boldsymbol{\mu})$ in $\mathfrak{J}_1(\boldsymbol{\mu})$ is maximized. Moreover, define:

$$\mathfrak{J}_2(\boldsymbol{\mu}) = \left[\widehat{R}^2(\boldsymbol{\mu}) - \sigma_{R^2}(\boldsymbol{\mu}), \widehat{R}^2(\boldsymbol{\mu}) + \sigma_{R^2}(\boldsymbol{\mu}) \right]. \quad (12)$$

Note that $\mathfrak{J}_2(\boldsymbol{\mu})$ in (12) is a relaxation of $\mathfrak{J}_1(\boldsymbol{\mu})$ in (11). When the number of $\{R_i^2(\boldsymbol{\mu})\}$ is not sensitive to $\boldsymbol{\mu}$ in $\mathfrak{J}_1(\boldsymbol{\mu})$, one can compare $\mathfrak{J}_2(\boldsymbol{\mu})$ instead.

The training dataset consists of data from January 1, 2015 to December 31, 2016, and the test dataset (for out-of-sample performance) consists of data from January 1, 2017 to March 31, 2017.

4.3 Features Identified by Regression Analysis

In this section, we present the results from the two-step Lasso, the Ridge, and the EN regressions, including the most significant features and associated parameters identified by these models.

4.3.1 Benchmark model: OLS

We first summarize the result from the benchmark method OLS. As seen in Table 8, all but two of the estimated coefficients are statistically significant at any reasonable level of significance. The two exceptions are *year to maturity* and *turnover*. Moreover,

- The coefficients of *Prop number of buys* and *Prop number of sells* have the same sign but different values. The coefficient of *Prop number of buys* is roughly one third of the coefficient of *Prop number of sells*. Similarly, both of the coefficients of *Prop buy volume* and *Prop sell volume* are positive. The coefficient of *Prop buy volume* is roughly half of the coefficient of *Prop sell volume*. This shows the asymmetry between customer buy orders and customer sell orders. This is consistent with numerous studies, i.g., [Fermanian et al., 2016], suggesting that dealers offer tighter quotes for larger trades than for smaller ones.
- *Avg price* has a small effect on the bid-ask spread.
- The indicators of different sectors have different coefficients, but the overall values are small.
- The *Log(Total volume)* coefficient is negative as expected. With value -21.4028 , the estimated coefficient implies that a increase of 10,000 in trade size would make a retail-size trade into a large institutional-size trade and would reduce the bid-ask spread by 100 basis points.
- The *Indicator of investment grade bonds* coefficient is negative and the *Indicator of high yield bonds* coefficient is positive. This is consistent with the well-documented empirical findings: larger spreads for high yield bonds and smaller spreads for investment grade bonds.

4.3.2 Features Identified from Two-Step Lasso

We then present the results from a class of two-step Lasso parameterized by different $\boldsymbol{\mu} = \lambda_l$ and discuss how to select the best λ_l with cross-validation.

In this analysis, 20 different values of $\boldsymbol{\mu} = \lambda_l$ are picked with a partition in the range of $[10^{-1}, 10^3]$ in the log scale. Note that the ranges of hyperparameters are different for two-step Lasso, Ridge, and two-step EN,

shown in the figures of cross-validation scores (Figures 8, 9, and 10). The range is selected according to the sensitivity of the model with a larger prior partition grid.

Figure 8 shows the 25%, 50%, and 75% percentiles of out-of-sample \tilde{R}^2 with different λ_l values. One can see that all three 25%, 50% and 75% curves decrease fast before $\lambda_l^* = 2.98$ and tend to be flat after λ_l^* . Also, both $\mathfrak{J}_1(\lambda_l)$ and $\mathfrak{J}_2(\lambda_l)$ are large when $\lambda_l = \lambda_l^*$. Hence, λ_l^* is a good choice of the regularization level. Table 16 shows λ_l 's along with $\mathfrak{J}_1(\lambda_l)$ and $\mathfrak{J}_2(\lambda_l)$, in which the number of \tilde{R}^2 are the largest, respectively. Table 9 shows the features selected from the first step of the two-step Lasso, with corresponding parameters $\lambda_l^* = 1.13, 2.98, \text{ and } 7.85$, respectively. It also shows the models from the OLS regression in the second step of the two-step Lasso. For instance, in Model L2 of Table 9 with $\lambda_l^* = 2.98$, the model is of the form with four features such that:

$$\begin{aligned} \text{Bid-ask spread} = & \quad 83.48 \times \text{Volatility} + 39.07 \times \text{Trading activity} \\ & - 19.45 \times \text{Log(Total volume)} + 0.23 \times \text{Issued years} + 86. \end{aligned} \quad (13)$$

In addition, as seen from Table 9:

- The coefficient of *Volatility* is positive with value 83.48. This is consistent with existing theoretical and empirical studies in market microstructure in that a higher return volatility is predicted to lead to decreased liquidity (e.g., [Stoll, 1978]).
- The coefficient of *Issued years* is positive with value 0.23, which means a newly issued bond will have a small bid-ask spread. This is consistent with the work of [Konstantinovskiy et al., 2016], which argued that recent and large issues are cheaper to trade than seasoned and small ones.
- The number of trades per day N and the trade volume V (in dollars) suggest a joint impact of order $\log(N/\sqrt{V})$ on the bid-ask spread. Section 4.4 provides a detailed analysis of this relationship.
- Finally, $\lambda_l = 7.85$ leads to the features of *Volatility* and *Issued year* in model L3. Compared to model L2 with four features and $R^2 = 51.5\%$, R^2 in EN3 drops to 43.22%. In our view, the model with four features, significantly reduced from the original 26 features, is preferable to EN3.

It is worth noting that our findings are supported by previous studies. For instance, [Chacko et al., 2005] found that credit quality, the age of a bond, the size of a bond issue, the original maturity value of a bond at issuance date, and provisions such as a call, put, or convertible options all have strong impacts on liquidity. [Choi and Huh, 2019] showed that trade size and maturity date are important features to understand in the bid-ask spread.

4.3.3 Features Identified from Ridge

We next discuss the results from a group of Ridge regression methods parameterized by different λ_r values and analyze how to select the best λ_r with cross-validation.

In this analysis, 20 different values of λ_r are chosen in the range of $[10^2, 10^8]$ with uniform partition in the log scale. Figure 9 shows the 25%, 50% and 75% percentiles of out-of-sample \tilde{R}^2 with different λ_r values. One can see that all three 25%, 50%, and 75% curves start to decrease at $\lambda_r^* = 1.27 \cdot 10^6$. Hence $1.27 \cdot 10^6$ is a good choice for the regularization level.

Table 17 shows $CI(\lambda_r)$ and $CI_2(\lambda_r)$ for different values of λ_r , in which the number of \tilde{R}^2 are the largest, respectively. Table 10 shows the results of Ridge regressions with parameters $\lambda_r^* = 1.62 \cdot 10^4, 6.95 \cdot 10^4, \text{ and } 1.27 \cdot 10^6$.

The analysis by the Ridge regression is consistent with the findings from the two-step Lasso. In particular,

- When λ_r goes up, the coefficients of the following features go to 0 very fast: *Indicator functions of different sectors*, *Proportion of buy (or sell) volumes (or numbers)*, *Turnover*, and *Number of trading days*. Note that from Table 10 these features are also excluded from Model L3 and L4 of Table 9, which means that results from these two approaches are consistent.

- When λ_r takes a large value $6.95 \cdot 10^4$, the *Volatility*, *Issued years*, *Trading activity* and *Log(Total volume)* are still significant. This is also consistent with the findings from Lasso in Table 9.
- Both two-step Lasso and Ridge regressions point to the significance of the time value and the special structure of bonds. The variable *years since issuance* is significant in two of the two-step Lasso models, L1 and L3, and all three Ridge regression models.

The difference between Lasso and Ridge regression: *Avg price* is not significant in all three two-step Lasso models, whereas it is significant in all three Ridge models. This inconsistency is expected because of the collinearity among features. When features are correlated, Lasso tends to select one feature from a group of correlated features while Ridge tends to penalize the group of correlated features towards the same coefficients [Zou and Hastie, 2005]. Indeed, as shown in Models R1, R2 and R3, the coefficients of *Prop number of buys* and *Prop number of sells* have the same value but different signs; the coefficients of *Prop buy volume* and *Prop sell volume* also have the same value but different signs. Additionally, the reappearance of *Avg price* in Model EN3 is due to this group effect as well.

4.3.4 Features Identified from Two-Step EN

Finally, we discuss the results from the two-step EN models parameterized by different λ_e values and analyze how to select the best λ_e with cross-validation.

Figure 10 shows the 25%, 50%, and 75% percentiles of out-of-sample \tilde{R}^2 with different λ_e values given different $\alpha = 0.2, 0.5$ and 0.8 . When $(\alpha, \lambda_e) = (0.5, 10^3)$ and $(\alpha, \lambda_e) = (0.8, 10^3)$, more than 170 empirical \tilde{R}^2 falls into \mathfrak{J}_2 . This is because the hyperparameter over-penalizes the model such that all the coefficients and the empirical \tilde{R}^2 are all nearly zero. Therefore, these sets of hyperparameters should be excluded.

Instead, for the analysis the following parameters are selected $(\alpha, \lambda_e) = (0.5, 0.774)$, $(\alpha, \lambda_e) = (0.8, 2.15)$, and $(\alpha, \lambda_e) = (0.5, 129)$. Parameter $(\alpha, \lambda_e) = (0.8, 2.15)$ leads to the following set of features: *Volatility*, *Number of trades*, *Log(Total volume)* and *Issued year*. This is consistent with the feature selection in the two-step Lasso model L2. $(\alpha, \lambda_e) = (0.5, 129)$ leads to the features: *Volatility* and *Average price* in model EN3. Compared with model EN2 with four features and $R^2 = 51.5\%$, the R^2 in EN3 drops to 42.62%. Similar to the argument for L3, model EN2 with four features is superior to EN3, in our view.

Through all three different regression models, the consensus is that *Volatility* and *Issued years* are important features.

4.3.5 Out-of-sample Performance

It is well recognized out-of-sample forecast performance is generally considered more trustworthy than in-sample performance [Alpaydin, 2020]. The latter can be more sensitive to outliers and data mining, while the former tends to better reflect the information available to the forecaster in “real time”.

In this subsection, we test the out-of-sample performance of all regression models on an unseen dataset during the period of January-March 2017. The distributions of the relative errors (defined in (10)) are provided in Figure 3 and the mean of the relative errors are recorded in Table 7. For the OLS model, errors smaller than 0.3 account for more than 70% of the testing dataset. Similar results hold for the two-step LASSO and two-step EN for which errors smaller than 0.3 account for more than 60% of the testing dataset. The errors are bigger for Ridge, as expected, since the coefficients are biased.

The value of R^2 in Table 8 indicates that the OLS model can explain around 60% of the variance hence its out-of-sample performance is acceptable. Similar results are obtained from the two-step LASSO model EN1 and the two-step LASSO model L2. Ridge regression model R3 has a larger mean relative error since the coefficients are biased with the L2 penalty term, which is also expected.

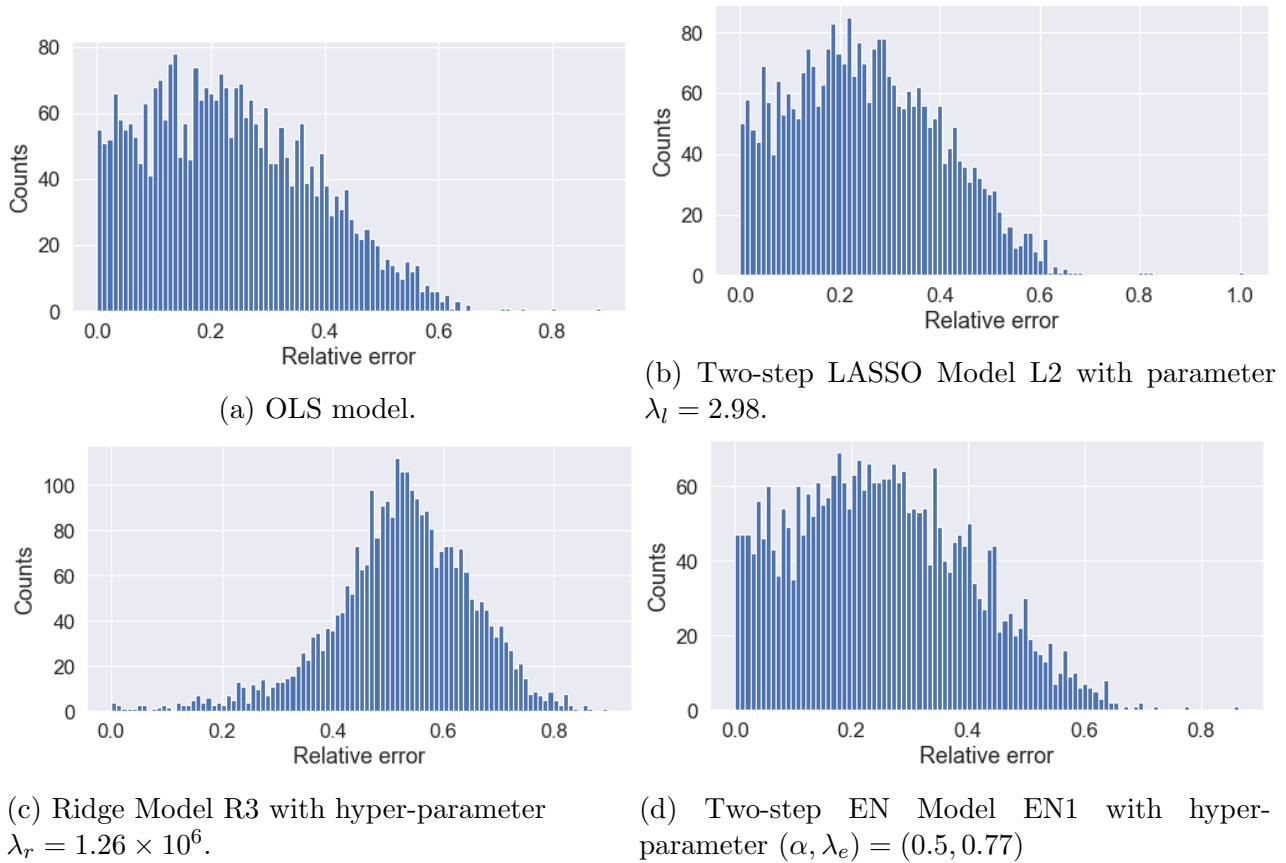


Figure 3: Out-of-sample test on data during January-March 2017.

Model	OLS	L2	R3	EN1
Mean relative error	0.222	0.249	0.522	0.252

Table 7: Mean relative error of out-of-sample test on data during January-March 2017.

4.4 Summary: Bid-Ask Spread of Corporate Bonds

Different linear regressions are performed for feature selection and for estimating the bid-ask spread using two kinds of variables: one describing the bond and the other characterizing the market. Tables 8-11 summarize the results, with comparisons to the benchmark OLS approach.

These regressions allow for computing an “expected bid-ask spread” for a given week, which can be used as a benchmark cost for TCA. In particular,

- *Volatility* is an important feature, as expected by both empirical observations and theory: the larger the volatility, the larger the bid-ask spread. It has been observed in practice that an increase of 5% in volatility, that is 1/2 of its standard deviation in our dataset, corresponds to an increase of the bid-ask spread by 25 basis points, which is around one third of its standard deviation.
- *Number of trades per day* N and *Traded volume* V (in dollars) are both important variables (in log units), with coefficients suggesting $\log(N/\sqrt{V})$ being the feature impacting the bid-ask spread in the basis points,⁶ implying that:

⁶The ridge regression suggests the feature being in the form of $\frac{N}{\sqrt{V}}$, with an addition term of the average price. Possible explanation: penalization in the Ridge regression tends to avoid large coefficients in the regression.

- for a given *trading activity* N , the larger the traded volume, the smaller the bid-ask spread (in basis points);
- for a given traded volume in dollars, the lower the average trade size (i.e., the more trades), the larger the bid-ask spread.

Our result is compatible with the documented stylized fact that for corporate bonds: small trade size obtain a worse bid-ask spread than large trades [Fermanian et al., 2016].

- The value of the coupon and the duration of the corporate bond play a small role in the formation of the bid-ask spread (both with a positive coefficient).
- Last but not least, the *Number of years to maturity* and the *Years since issuance* are selected by our robust regressions. Keep in mind these two variables are linked, via the maturity of the bond, thanks to the relation: $Year\ to\ maturity = Maturity - Years\ since\ issuance$. Hence naturally, the coefficient of *year to maturity* is negative while the one of *years since issuance* is positive: the further away from the maturity, the smaller the bid-ask spread (in basis points).

This could support the market folklore that there is only a short time period after the issuance when corporate bond trading is not too expensive on secondary markets.

Other variables appearing in the OLS are not robust enough to be selected by penalized regressions. Removing these 17 variables from the regression only reduces the R^2 from around 0.55 to around 0.50, a minor price to pay for the increased robustness. It is worth noticing that the R^2 of all these OLS and regularized regressions are around 50%, which is in line with the best results obtained in the literature: [Dick-Nielsen et al., 2012] obtained R^2 between 0.50 and 0.80, while the R^2 of other studies were far below 0.50 (see Section 1). The out-of-sample performance in Section 4.3.5 further indicates the promise of the regression models.

In addition, the regression results (Table 20 for OLS and Table 19 for two-step LASSO) between the Enhanced TRACE and Standard TRACE datasets confirm similar performance in terms of R^2 for these two datasets. Therefore, one could choose either dataset for the bid-ask spread estimation, depending on the time-lag and the accuracy of the trade volume.

5 Price Impact Analysis

After analyzing the average TCA on the weekly basis (Section 4), we now move on to the second part of TCA analysis. We will focus on the individual trade, and study the amplitude of its price impact and the price impact decay after the transaction for liquid corporate bonds. The goal of the price impact analysis is to determine the necessary set of events to fit the mid-price dynamics and to understand how the impact of each type of event decays over time.

As a benchmark comparison, recall several stylized facts on equity market from [Bouchaud et al., 2009]:

- buy trades on average push the price up and sell trades on average drive the price down;
- the impact curve as a function of the volume of the trade is strongly concave. In other words, large volumes impact the price only marginally more than small volumes;
- the sign of market orders is strongly autocorrelated in time.

To see whether these facts hold for corporate bonds, we apply a transient price impact model to estimate the price impact amplitude and decay pattern. Since the trading frequency is much lower on corporate bonds than on equities, it is more appropriate to use the “event time” in our transient impact model instead of the chronological time used in other non-parametric price impact models for equity markets [Biais et al., 2016, Van Kervel and Menkveld, 2019].

Our first attempt is to model naively the mid-price by a single-event TIM (TIM1) (Section 5.2). Using the signature plot as a metric for the goodness-of-fit shows that TIM1 is not sufficient to describe the mid-price movements for corporate bonds. Meanwhile, statistical evidence implies an asymmetry between the price

impacts from customer-buy orders and customer-sell orders, as detailed in Section 5.3.1. This statistical evidence motivates us to propose a TIM model with two types of events (TIM2) (Section 5.3.2): customer-buy orders and customer-sell orders. The signature plot indicates good performance of this improved TIM framework.

5.1 Review: Transient Price Impact Models

Let us first review the classic transient impact model (TIM) following [Bouchaud et al., 2009] with the “event time”. Assume Π is a set of event-types considered on the market and the mid-price M_k^b of a corporate bond b follows [Eisler et al., 2012, Taranto et al., 2018, Lehalle and Laruelle, 2018]:

$$M_k^b = \sum_{k'=k}^{-\infty} \sum_{\pi \in \Pi} \left(G_\pi^b(k - k') \mathbf{1}(\pi_k^b = \pi) (V_{k'}^b)^\alpha \epsilon_{k'}^b + \eta_{k'}^b \right) + M_{-\infty}^b, \quad (14)$$

where $\epsilon_k^b \in \{-1, +1\}$ is the sign of the k -th trade, estimation of which has been detailed in Section 3.1, V_k^b is the volume of the k th trade, π_k^b is the type of the k -th trade, α is a power index, $G_\pi^b(\delta k)$ is a decaying kernel of type π event, η^b is a noise, and $M_{-\infty}^b$ is a initial value for the mid-price. η is a random change of the fair price independent of ϵ and it is assumed to be i.i.d.

Note that G_π^b is typically an exponential or a power law (i.e., $G_\pi^b(\delta t) \propto \exp(-\lambda \delta t)$ or $(1 + \delta t)^{-\gamma}$) [Eisler et al., 2012, Taranto et al., 2018, Lehalle and Laruelle, 2018]. $G_\pi^b(\delta k)$ can be interpreted as the response function *per bond* when $\alpha = 1$; $G_\pi^b(\delta k)$ can be understood as the response function *per order* when $\alpha = 0$ and the volume is ignored. In equity markets, there has been empirical evidence showing that $\alpha \approx 0.1$.

5.2 First Attempt: Single-Event Transient Impact Model

In this section, we will show that the naive TIM1 model (i.e., Eqn (14) with one type of events) does not fit the price impact curves for cooperate bonds.

To see this, note that the mid-price dynamics of bond b under TIM1 are:

$$M_k^b = \sum_{k'=k}^{-\infty} \left(G^b(k - k') (V_{k'}^b)^\alpha \epsilon_{k'}^b + \eta_{k'}^b \right) + M_{-\infty}^b, \quad (15)$$

where $\epsilon_k^b \in \{-1, +1\}$ is the sign of the k th trade as estimated in Section 3.1, V_k^b is the volume of the k th trade, α is a power index, $G^b(\delta k)$ is a decaying kernel of the bond b mid-price, η^b is noise and $M_{-\infty}^b$ is an initial value for the mid-price, and η^b is a random change of the fair price independent of ϵ^b and is assumed to be i.i.d.

Under (14), the change of the mid-price can be written as:

$$R_k^b(1) := M_{k+1}^b - M_k^b = G^b(0) (V_{k+1}^b)^\alpha \epsilon_{k+1}^b + \eta_{k+1}^b + \sum_{j=0}^{\infty} \underbrace{(G^b(j+1) - G^b(j))}_{\Delta_1 G^b(j)} \cdot (V_{k-j}^b)^\alpha \cdot \epsilon_{k-j}^b. \quad (16)$$

Consequently, we can check the values of $S^b(l) = \mathbb{E} [R_k^b(1) \epsilon_{k-l+1}^b]$ and $C^b(n) = \mathbb{E} [(V_{t+n}^b)^\alpha \epsilon_{t+n}^b \epsilon_t^b]$ and obtain:

$$S^b(l) = G^b(0) C^b(l) + \sum_{j=0}^{+\infty} \Delta_1 G^b(j) \cdot C^b(l - j - 1). \quad (17)$$

If we only focus on the first N transaction in the calculation of the response function, then (17) can be written in the following matrix format:

$$\underbrace{\begin{pmatrix} S^b(1) - G^b(0)C^b(1) \\ S^b(2) - G^b(0)C^b(2) \\ \vdots \\ S^b(L) - G^b(0)C^b(L) \end{pmatrix}}_{=: \bar{S}^b(L)} = \underbrace{\begin{bmatrix} C^b(0) & C^b(-1) & C^b(-2) & \cdots & C^b(-N+1) \\ C^b(1) & C^b(0) & C^b(-1) & \cdots & C^b(-N) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ C^b(L-1) & \cdots & \cdots & \cdots & C^b(L-N) \end{bmatrix}}_{=: \bar{C}^b(N,L)} \underbrace{\begin{pmatrix} \Delta_1 G^b(0) \\ \Delta_1 G^b(1) \\ \vdots \\ \Delta_1 G^b(N-1) \end{pmatrix}}_{=: \bar{G}^b(N)}. \quad (18)$$

Note that it suffices to estimate $S^b(l)$ and $C^b(n)$ for different values of l and n , and the initial value $G^b(0)$. Afterwards, an estimator for $\overline{G}^b(N)$ can be constructed using (18) such that it follows:

$$\widehat{\overline{G}^b(N)} = \widehat{\overline{C}^b(N, L)}^{-1} \cdot \widehat{\overline{S}^b(L)}, \quad (19)$$

with $\widehat{\overline{C}^b(N, L)}$ and $\widehat{\overline{S}^b(L)}$ the estimates of $\overline{C}^b(N, L)$ and $\overline{S}^b(L)$, respectively.

To evaluate the model and quantify the price diffusion for different lags, define the *signature plot* [Bouchaud et al., 2009] as below:

$$D^b(l) = \frac{1}{l} \mathbb{E} \left[(M_{t+l}^b - M_t^b)^2 \right]. \quad (20)$$

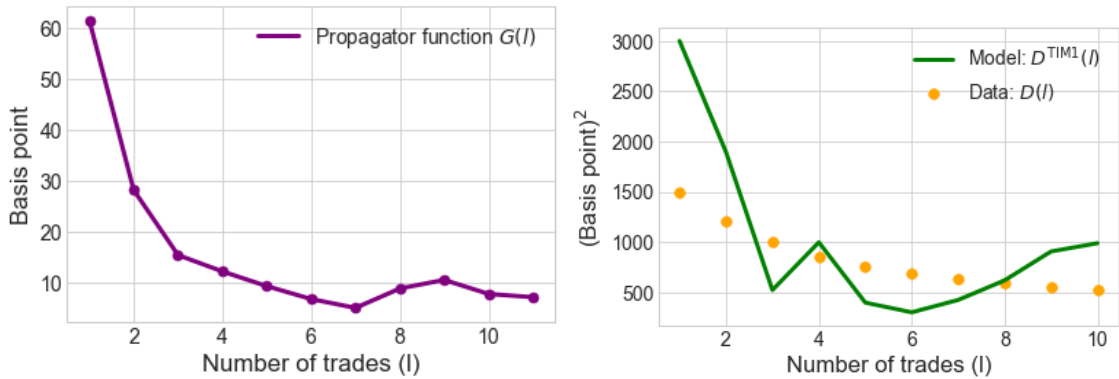
For the TIM1 model, the approximated signature plot follows:

$$D_{\text{TIM1}}^b(l) = \frac{1}{l} \sum_{0 \leq n < l} \left(G^b(l-n) \right)^2 + \frac{1}{l} \sum_{n > 0} \left(G^b(l+n) - G^b(n) \right)^2 + 2\Phi^b(l) + D_{\text{const}}^b, \quad (21)$$

where D_{const}^b is some constant and $\Phi^b(l)$ is the correlation-induced contribution to the price diffusion:

$$\begin{aligned} l \Phi^b(l) &= \sum_{0 \leq n < m < l} G^b(l-n) G^b(l-m) C^b(m-n) \\ &+ \sum_{0 \leq n < m} [G^b(l+n) - G^b(n)] [G^b(l+m) - G^b(m)] C^b(m-n) \\ &+ \sum_{0 \leq n < l} \sum_{m > 0} G^b(l-n) [G^b(l+m) - G^b(m)] C^b(m+n). \end{aligned}$$

Experiment set-up. We fit the TIM1 model (18)-(19) with $L = N = 10$ for the top-200 bonds (described in Section 2). In order to facilitate the comparison of different bonds, we calculate the propagator functions for relative price changes with $\alpha = 0.0$.⁷



(a) Propagator function G for TIM1 model. (b) Signature plots: data vs. TIM1 model

Figure 4: Fitted TIM1 model and the goodness-of-fit (aggregation over 200 bonds and $\alpha = 0$).

From Figure 4, observe that:

- unlike equity markets where the decay of the propagator function is *slow* for both large tick stocks and small tick stocks [Taranto et al., 2018], the decay of the propagator functions is *fast* in cooperate bond markets and $G(l) \approx 10(bp)$ when $l \geq 10$ (Figure 4a).
- since the signature plot serves as a metric to evaluate the fitted models, we observe from Figure 4b that $D_{\text{TIM1}}(l)$, the signature calculated from the TIM1 Model, does not fit well with the signature plot $D(l)$ calculated from the data. It appears that $D_{\text{TIM1}}(l)$ overestimates the signatures for small l .

⁷The results with small α are qualitatively the same.

5.3 Modified TIM Model and Asymmetric Price Impact

5.3.1 Statistical Evidence of Asymmetric Price Impact

We next present some statistical evidence on the asymmetry of price impacts. This study then motivates us to consider a two-type event model treating customer-buy and customer-sell orders separately.

To start, we adopt one-sided spread to test if there is any difference between the buy-side liquidity and the sell-side liquidity [Choi and Huh, 2019]:

$$\begin{aligned} \text{spread}_B &= \frac{\text{traded price} - \text{reference price}}{\text{reference price}} \mathbf{1}(\text{buy order}), \\ \text{spread}_S &= \frac{\text{reference price} - \text{traded price}}{\text{reference price}} \mathbf{1}(\text{sell order}). \end{aligned} \quad (22)$$

For each customer trade, we calculate its reference price as the volume-weighted average price of inter-dealer trades larger than \$100,000 in the same bond-day, excluding inter-dealer trades executed within 15 minutes. spread_B and spread_S are calculated at the bond-day level by taking the volume-weighted average of trade-level spreads. The average buy spread is 44.52 (bp) and the average sell spread is 38.74 (bp) across all 1,993 liquid bonds. Afterwards, we perform a t -test with the null-hypothesis that the buy spread and the sell spread have the same sample mean. The null-hypothesis is rejected with a p-value smaller than 1%, indicating that the buy spread is different from the sell spread. Meanwhile, we also perform t -tests for individual bonds – 1,215 out of 1,993 bonds have p-values smaller than 5% indicating different buy spread and sell spread distributions.

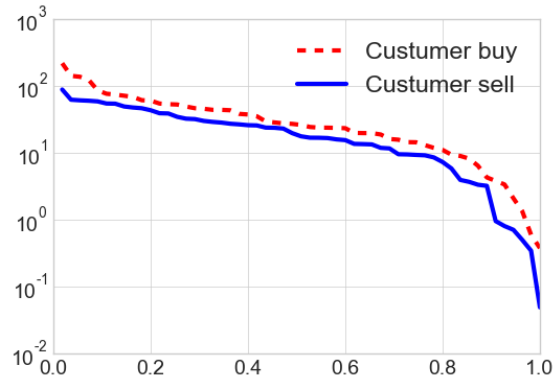


Figure 5: Rank frequency plot of buy-spread (spread_B) and sell-spread (spread_S) of all 1,993 bonds.

The rank frequency plot of buy-spread (spread_B) and sell-spread (spread_S) is visualized in Figure 5. Note that a rank-frequency distribution is a discrete form of a quantile function (inverse cumulative distribution) in reverse order, giving the size of the element at a given rank. From Figure 5, we observe that the distributions of the buy-spread and sell-spread have different tail behaviors.

5.3.2 Modified TIM Models and Estimation of Asymmetric Price Impact

The preliminary analysis of price impact asymmetry from 5.3.1 motivates us to propose a modified TIM model in which customer-buy orders and customer-sell orders are treated as different events in the calculation of propagator functions [Bouchaud et al., 2009, Eisler et al., 2012, Eisler and Bouchaud, 2016, Taranto et al., 2018, Schneider and Lillo, 2019]. See also [Jurksas et al., 2021] on liquidity spill-overs in sovereign bond market which estimate the price impact curves for buy and sell orders separately.

This model is inspired by [Taranto et al., 2018] where events with small trades and large trades are treated differently. Here we assume that there are two types of events $\Pi := \{+1, -1\}$ with +1 denoting the customer-buy orders and -1 denoting the customer-sell orders. The calculation of the propagator function is similar to (19) and as detailed in Appendix D.1. See also [Lehalle and Laruelle, 2018, Appendix A.12] for a more detailed discussion.

Experiment set-up. We fit the TIM2 model (18)-(19) with $L = N = 10$ for the top-200 bonds (described in Section 2). Similar as before, we calculate the propagator functions for relative price changes in order to make different bonds comparable. In the estimation, we take $\alpha = 0.0$ and the qualitative results for small α are similar.

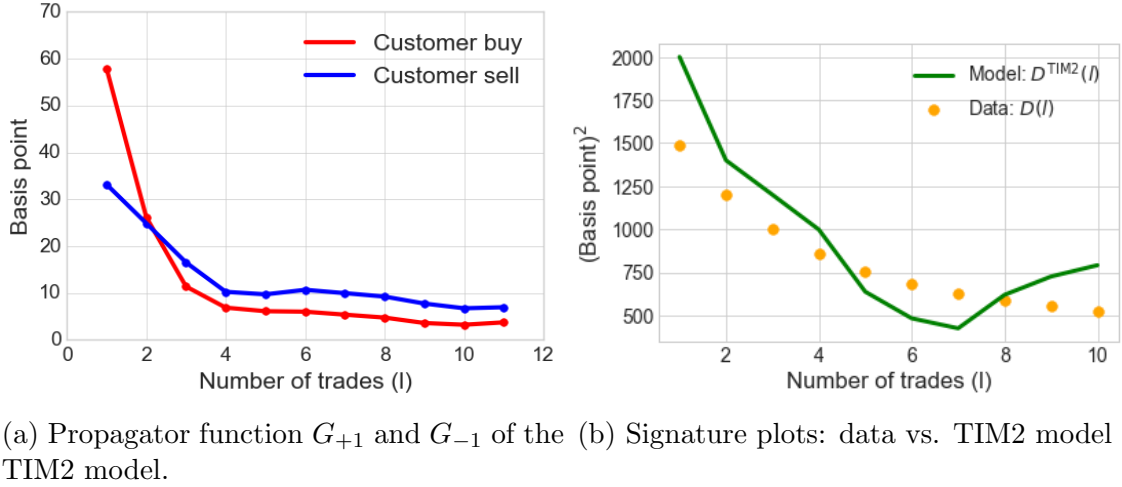


Figure 6: Fitted TIM2 model and the goodness-of-fit ($\alpha = 0$ and aggregation over 200 selected bonds).

One can observe the following from Figure 6: First, customer-buy orders have *larger* price impacts than customer-sell orders for the first few trade-times $l = 1, 2$ (Figure 6a). Second, the decay of the propagator function for customer-buy orders are slightly *faster* than customer-sell orders (Figure 6a). Moreover, comparing Figures 4b and 6b, we see the TIM2 model fits better with the signature plot calculated from the data. This implies that the TIM2 model with customer-buy orders and customer-sell orders being treated differently is better than the single-event model TIM1.

Figure 7 suggests heterogeneity among bonds in terms of the size of the market impact, the different impacts between buy-market orders and sell-market orders, and the shape of the decay. In addition, for newly issued bonds (i.e., Wells Fargo 94974BFY1 4.1%) or bonds that are close to maturity (i.e., Transocean 893830AS8 6.0%), the difference between the customer-buy propagator function and the customer-sell propagator function is *larger* than for the bonds that are in the middle of their life-time (i.e., Goldman Sachs 38141GGQ1 5.25%).

It is worth pointing out that price impact models based on no-arbitrage considerations in equity markets require the price impact to be symmetric [Huberman and Stanzl, 2004] and [Gatheral, 2010]. This no-arbitrage condition does not hold in bond markets due to the less liquidity and more fragmentation in bond market. Our discovery of the asymmetric price impacts indicates possible arbitrage opportunities in the secondary OTC market for corporate bonds.

Summary. We propose to use propagator functions to measure the *price impact of each single trade* for corporate bonds that are liquid enough. Our analysis finds two characteristics of the price impact of corporate bonds:

- The *asymmetry between buying and selling trades*. The mid price moves triggered by a trade on a corporate bond are larger for buying transactions than those for selling ones. In terms of TCA, it means that the asset manager has to respect such an asymmetry and take it into consideration during the evaluation of the counterparty dealers.
- Decay in price impact curves, similar to the one identified in equity markets [Eisler et al., 2012, Taranto et al., 2018]. The price impact curve consists of a jump corresponding to the adverse selection suffered by the dealer, followed by a decay stabilizing the price at the level of the permanent market impact.

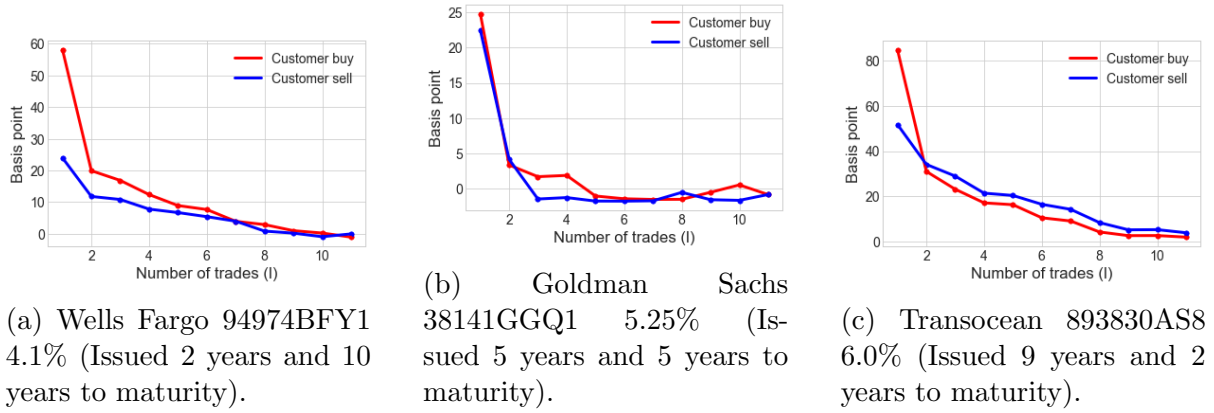


Figure 7: Heterogeneity among bonds ($\alpha = 0$).

6 Conclusion

This paper established a TCA benchmark in bond trading for retail investors and asset managers. It consists of (a) estimating the expected average cost on a weekly basis via regularized regression analysis and (b) investigating the amplitude of price impact and the price impact decay for each trade of liquid corporate bonds via TIM model.

The most important features identified in the regression analysis are volatility, trading activity, $\log(\text{total volume})$, and issued years. Meanwhile, asymmetry is discovered between buying and selling trades: mid-price moves triggered by a trade on corporate bonds are larger for buying than those for selling.

Our study suggests the following approach for TCA in practice:

1. For all corporate bonds of interest, asset managers first compute an expected bid-ask spread given the characteristics of the bond and market conditions using one of the regression approaches proposed in Section 4, and using either the Standard TRACE or the Enhanced TRACE datasets for bid-ask spread approximation.
2. This reference bid-ask spread can be used to benchmark the bid-ask spread obtained while requesting for quotes from counterparties. It can also be used to *score* all the obtained trades during the week.
3. Worst trades can be qualitatively evaluated using the average price impact curves obtained in Section 5.3. More specifically, if a trade has price impact larger than the curve showed in Figure 7, then it can be identified as a “worst trade” and the asset manager can conduct further analysis on the counterparty.

	Estimate	Standard error	t-value
Volatility	77.7272***	0.316	151.687
Number of trade days	-3.6648***	0.203	-18.014
Prop number of buys	10.3189***	0.689	14.971
Prop number of sells	32.4532***	0.700	46.333
Trading activity	46.3169***	0.531	87.162
Prop volume sell \$	16.3523***	0.608	26.893
Prop volume buy \$	26.4198***	0.642	41.155
Log(total volume)	-21.4028***	0.272	-78.757
Avg price	-0.1175***	0.017	-6.723
Coupon	-0.4707***	0.120	-3.914
Duration	1.4168***	0.148	9.596
Years to maturity	-0.0656	0.076	-0.858
Years since issuance	1.2552***	0.065	19.359
Turnover	-1.7717	2.260	-0.784
LIBOR-OIS	34.0088***	1.432	23.754
Indicator of high yield bonds	26.7859***	0.595	45.039
Indicator of investment grade bonds	15.9861***	0.604	26.461
Indicator of basic materials sector	8.9009***	0.661	13.462
Indicator of communications sector	5.1285***	0.374	13.695
Indicator of consumer, cyclical sector	4.0261***	0.421	9.570
Indicator of consumer, non-cyclical sector	4.9727***	0.378	13.140
Indicator of energy sector	3.8824***	0.362	10.731
Indicator of financial sector	4.0544***	0.325	12.460
Indicator of industrial sector	2.8149***	0.485	5.806
Indicator of technology sector	4.4085***	0.489	9.017
Indicator of utilities sector	4.5836***	1.116	4.107
Constant	42.7720***	1.112	38.477
N	152,408		
R ²	55.4	%	

Standard errors in parenthesis. Significance levels: * p<0.1, ** p<0.05, *** p<0.01. Two-tailed test. Source: TRACE Enhanced (2015-2016).

Table 8: OLS regression: the impact on bid-ask spread

	Model OLS All features (Benchmark)	Model L1 $\lambda_l = 1.13$	Model L2 $\lambda_l = 2.98$	Model L3 $\lambda_l = 7.85$
Volatility	77.7272*** (0.316)	80.7218*** (0.253)	83.5410*** (0.221)	86.5025*** (0.237)
Number of trade days	-3.6648*** (0.203)			
Prop number of buys	10.3189*** (0.689)			
Prop number of sells	32.4532*** (0.700)			
Trading activity	46.3169*** (0.531)	42.1450*** (0.451)	39.0703*** (0.433)	
Prop volume sell \$	16.3523*** (0.608)			
Prop Volume buy \$	26.4198*** (0.642)			
Log(total volume)	-21.4028*** (0.272)	-21.1246*** (0.252)	-19.4449*** (0.243)	
Avg price	-0.1175*** (0.017)			
Coupon	-0.4707*** (0.120)	1.2071*** (0.089)		
Duration	1.4168*** (0.148)	0.2453*** (0.035)		
Years to maturity	-0.0643 (0.076)	-0.1078 (0.075)		
Years since issuance	1.2576*** (0.065)	-0.0216 (0.052)	0.2336*** (0.046)	1.0396*** (0.011)
Turnover	-1.7717 (2.260)			
LIBOR-OIS	34.0088*** (1.432)			
Indicator of high yield bonds	26.7859*** (0.595)			
Indicator of investment grade bonds	15.9861*** (0.604)			
Indicator of basic materials sector	8.9009*** (0.661)			
Indicator of communications sector	5.1285*** (0.374)			
Indicator of consumer, cyclical sector	4.0261*** (0.421)			
Indicator of consumer, non-cyclical sector	4.9727*** (0.378)			
Indicator of energy sector	3.8824*** (0.362)			
Indicator of financial sector	4.0544*** (0.325)			
Indicator of industrial sector	2.8149*** (0.485)			
Indicator of technology sector	4.4085*** (0.489)			
Indicator of utilities sector	4.5836*** (1.116)			
Constant	42.7720	86.1737	85.6982	-2.9270
N	152,408	152,408	152,408	152,408
R ²	55.4 %	52.8 %	51.5 %	43.22 %

Standard errors in parenthesis. Significance levels: * p<0.1, ** p<0.05, *** p<0.01. Two-tailed test. Source: TRACE Enhanced (2015-2016).

Table 9: Two-step Lasso regression table: the impact on bid-ask spread (in bp)

	Model OLS All features (Benchmark)	Model R1 $\lambda_r = 1.62 \times 10^4$	Model R2 $\lambda_r = 6.95 \times 10^4$	Model R3 $\lambda_r = 1.27 \times 10^6$
Volatility	77.7272*** (0.316)	82.0769*** (0.522)	84.0327*** (0.364)	75.1199*** (0.380)
Number of trade days	-3.6648*** (0.202)	1.1389*** (0.339)	1.2833*** (0.241)	0.1705 (0.249)
Prop number of buys	10.3189*** (0.689)	-3.5702*** (1.148)	-1.1382 (0.817)	-0.0734 (0.844)
Prop number of sells	32.4532*** (0.700)	3.5702*** (1.167)	1.1382 (0.830)	0.0734 (0.858)
Trading activity	46.3169*** (0.531)	14.2859*** (0.885)	4.3609*** (0.630)	0.2794 (0.651)
Prop volume sell \$	16.3523*** (0.608)	-0.7126 (1.013)	-0.0237 (0.720)	0.0060 (0.745)
Prop volume buy \$	26.4198*** (0.642)	0.7126 (1.069)	0.0237 (0.761)	-0.0060 (0.786)
Log(total volume)	-21.4028*** (0.272)	-10.5886*** (0.453)	-4.5781*** (0.322)	-0.4096 (0.333)
Avg price	-0.1175*** (0.017)	-0.1313*** (0.029)	-0.0896*** (0.021)	-0.2174*** (0.021)
Coupon	-0.4707*** (0.120)	0.0114 (0.200)	-0.0264 (0.142)	0.2436 (0.147)
Duration	1.4168*** (0.1480)	1.7553*** (0.246)	0.0044*** (0.003)	0.2660 (0.181)
Years to maturity	-0.0643 (0.076)	-0.6328*** (0.127)	-0.5433*** (0.091)	0.0510 (0.094)
Years since issuance	1.2576*** (0.065)	1.0926*** (0.108)	1.0705*** (0.077)	0.6555*** (0.079)
Turnover	-1.7717 (2.260)	-0.0963 (3.764)	-0.0893 (2.677)	-0.0120 (2.768)
LIBOR-OIS	34.008 *** (1.432)	2.5574 (2.384)	0.6844 (1.696)	0.0386 (1.753)
Indicator of high yield bonds	26.7859*** (0.595)	1.8302* (0.990)	0.5512 (0.716)	0.0478 (0.728)
Indicator of investment grade bonds	15.9861*** (0.604)	-1.8302* (1.006)	-0.5512 (0.012)	-0.0478 (0.740)
Indicator of basic materials sector	8.9009*** (0.661)	0.6401 (1.101)	0.1724 (0.783)	0.0125 (0.810)
Indicator of communications sector	5.1285*** (0.374)	0.7563 (0.624)	0.3777 (0.444)	0.0373 (0.459)
Indicator of consumer, cyclical sector	4.0261*** (0.421)	-0.5567 (0.701)	-0.2340 (0.498)	-0.0256 (0.515)
Indicator of consumer, non-cyclical sector	4.9727*** (0.378)	-0.4812 (0.630)	-0.2448 (0.448)	-0.0371 (0.463)
Indicator of energy sector	3.882 *** (0.362)	-1.3953** (0.603)	-0.7044* (0.429)	-0.0304 (0.433)
Indicator of financial sector	4.0544*** (0.325)	0.1134 (0.542)	0.0911 (0.385)	-0.0054 (0.398)
Indicator of industrial sector	2.8149*** (0.485)	-0.1199 (0.807)	0.0657 (0.574)	0.0130 (0.594)
Indicator of technology sector	4.4085*** (0.489)	1.0098 (0.814)	0.4649 (0.579)	0.0346 (0.599)
Indicator of utilities sector	4.5836*** (1.116)	0.0335 (1.859)	0.0114 (1.322)	0.001 (1.367)
Constant	42.7720	0.000	0.000	0.000
N	152,408	152,408	152,408	152,408
R ²	55.4 %	53.5 %	52.0 %	50.0 %

Standard errors in parenthesis. Significance levels: * p<0.1, ** p<0.05, *** p<0.01. Two-tailed test. Source: TRACE Enhanced (2015-2016).

Table 10: Ridge regression table: the impact on bid-ask spread (in bp)

	Model OLS All features (Benchmark)	Model EN1 $(\alpha, \lambda_e) = (0.5, 0.774)$	Model EN2 $(\alpha, \lambda_e) = (0.8, 2.15)$	Model EN3 $(\alpha, \lambda_e) = (0.5, 129)$
Volatility	77.7272*** (0.316)	77.8015*** (0.300)	83.5410*** (0.221)	87.5398*** (0.253)
Number of trade days	-3.6648*** (0.203)	-3.3792*** (0.200)		
Prop number of buys	10.3189*** (0.689)	20.4300*** (0.725)		
Prop number of sells	32.4532*** (0.700)	36.5251 (0.732)		
Trading activity	46.3169*** (0.531)	45.9632*** (0.519)	39.0703*** (0.433)	
Prop volume sell	16.3523*** (0.608)			
Prop Volume buy	26.4198*** (0.642)			
Log(total volume)	-21.4028*** (0.272)	-21.6955*** (0.252)	-19.4449*** (0.243)	
Avg price	-0.1175*** (0.017)	-0.1083*** (0.016)		0.1264*** (0.014)
Coupon	-0.4707*** (0.120)			
Duration	1.4168*** (0.148)	1.5875*** (0.146)		
Years to maturity	-0.0643 (0.076)	-0.1680 (0.074)		
Years since issuance	1.2576*** (0.065)	0.9346*** (0.056)	0.2336*** (0.046)	
Turnover	-1.7717 (2.260)			
LIBOR-OIS	34.0088*** (1.432)			
Indicator of high yield bonds	26.7859*** (0.595)	33.1587*** (0.650)		
Indicator of investment grade bonds	15.9861*** (0.604)	23.7964*** (0.625)		
Indicator of basic materials sector	8.9009*** (0.661)			
Indicator of communications sector	5.1285*** (0.374)			
Indicator of consumer, cyclical sector	4.0261*** (0.421)			
Indicator of consumer, non-cyclical sector	4.9727*** (0.378)			
Indicator of energy sector	3.8824*** (0.362)			
Indicator of financial sector	4.0544*** (0.325)			
Indicator of industrial sector	2.8149*** (0.485)			
Indicator of technology sector	4.4085*** (0.489)			
Indicator of utilities sector	4.5836*** (1.116)			
Constant	42.7720	56.9551	85.6982	-12.7213
N	152,408	152,408	152,408	152,408
R ²	55.4 %	53.8 %	51.5 %	42.62 %

Standard errors in parenthesis. Significance levels: * p<0.1, ** p<0.05, *** p<0.01. Two-tailed test. Source: TRACE Enhanced (2015-2016).

Table 11: Two-step EN regression table: the impact on bid-ask spread (in bp)

References

- [Albanese and Tompaidis, 2008] Albanese, C. and Tompaidis, S. (2008). Small transaction cost asymptotics and dynamic hedging. *European Journal of Operational Research*, 185(3):1404–1414.
- [Alpaydin, 2020] Alpaydin, E. (2020). *Introduction to machine learning*. MIT press.
- [Asquith et al., 2013] Asquith, P., Au, A. S., Covert, T., and Pathak, P. A. (2013). The market for borrowing corporate bonds. *Journal of Financial Economics*, 107(1):155–182.
- [Belloni and Chernozhukov, 2013] Belloni, A. and Chernozhukov, V. (2013). Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547.
- [Bessembinder et al., 2008] Bessembinder, Hendrik, Maxwell, and William (2008). Markets transparency and the corporate bond market. *The Journal of Economic Perspectives*, 22(2):217–234.
- [Bessembinder et al., 2006] Bessembinder, H., Maxwell, W., and Venkataraman, K. (2006). Market transparency, liquidity externalities, and institutional trading costs in corporate bonds. *Journal of Financial Economics*, 82(2):251–288.
- [Biais et al., 2016] Biais, B., Declerck, F., and Moinas, S. (2016). Who supplies liquidity, how and when? Technical report, BIS Working Paper.
- [Biais and Green, 2019] Biais, B. and Green, R. (2019). The microstructure of the bond market in the 20th century. *Review of Economic Dynamics*, 33:250–271.
- [Bouchaud et al., 2009] Bouchaud, J.-P., Farmer, J. D., and Lillo, F. (2009). How markets slowly digest changes in supply and demand. In *Handbook of financial markets: dynamics and evolution*, pages 57–160. Elsevier.
- [Chacko et al., 2005] Chacko, G., Mahanti, S., Mallik, G., and Subrahmanyam, M. G. (2005). The determinants of liquidity in the corporate bond markets: An application of latent liquidity.
- [Chakravarty and Sarkar, 2003] Chakravarty, S. and Sarkar, A. (2003). Trading costs in three US bond markets. *The Journal of Fixed Income*, 13(1):39–48.
- [Chételat et al., 2017] Chételat, D., Lederer, J., Salmon, J., et al. (2017). Optimal two-step prediction in regression. *Electronic Journal of Statistics*, 11(1):2519–2546.
- [Choi and Huh, 2019] Choi, J. and Huh, Y. (2019). Customer liquidity provision: Implications for corporate bond transaction costs. *Available at SSRN 2848344*.
- [Chordia et al., 2005] Chordia, T., Sarkar, A., and Subrahmanyam, A. (2005). An empirical analysis of stock and bond market liquidity. *Review of Financial Studies*, 18(1):85–129.
- [Collins and Fabozzi, 1991] Collins, B. M. and Fabozzi, F. J. (1991). A methodology for measuring transaction costs. *Financial Analysts Journal*, 47(2):27–36.
- [Dick-Nielsen, 2014] Dick-Nielsen, J. (2014). How to clean Enhanced TRACE data. *Available at SSRN 2337908*.
- [Dick-Nielsen et al., 2012] Dick-Nielsen, J., Feldhütter, P., and Lando, D. (2012). Corporate bond liquidity before and after the onset of the subprime crisis. *Journal of Financial Economics*, 103(3):471–492.
- [Edwards et al., 2007] Edwards, A. K., Harris, L. E., and Piwowar, M. S. (2007). Corporate bond market transaction costs and transparency. *The Journal of Finance*, 62(3):1421–1451.
- [Eisler and Bouchaud, 2016] Eisler, Z. and Bouchaud, J.-P. (2016). Price impact without order book: A study of the otc credit index market. *Available at SSRN 2840166*.

- [Eisler et al., 2012] Eisler, Z., Bouchaud, J.-P., and Kockelkoren, J. (2012). The price impact of order book events: market orders, limit orders and cancellations. *Quantitative Finance*, 12(9):1395–1419.
- [Eom et al., 2004] Eom, Y. H., Helwege, J., and Huang, J.-Z. (2004). Structural models of corporate bond pricing: An empirical analysis. *Review of Financial Studies*, 17(2):499–544.
- [Fermanian et al., 2016] Fermanian, J.-D., Guéant, O., and Pu, J. (2016). The behavior of dealers and clients on the european corporate bond market: the case of multi-dealer-to-client platforms. *Market Microstructure and Liquidity*, page 1750004.
- [Fermanian et al., 2015] Fermanian, J.-D., Guéant, O., and Rachez, A. (2015). *Agents’ Behavior on Multi-Dealer-to-Client Bond Trading Platforms*.
- [Friewald and Nagler, 2014] Friewald, N. and Nagler, F. (2014). Dealer inventory and the cross-section of corporate bond returns. *Social Science Research Network Working Paper Series*.
- [Gatheral, 2010] Gatheral, J. (2010). No-dynamic-arbitrage and market impact. *Quantitative finance*, 10(7):749–759.
- [Glosten and Milgrom, 1985] Glosten, L. R. and Milgrom, P. R. (1985). Bid, ask and transaction prices in a specialist market with heterogeneously informed traders. *Journal of Financial Economics*, 14(1):71–100.
- [Goldstein et al., 2007] Goldstein, M. A., Hotchkiss, E. S., and Sirri, E. R. (2007). Transparency and liquidity: A controlled experiment on corporate bonds. *Review of Financial Studies*, 20(2):235–273.
- [Harris, 2015] Harris, L. (2015). Transaction costs, trade throughs, and riskless principal trading in corporate bond markets. *Social Science Research Network Working Paper Series*.
- [Hendershott and Madhavan, 2015] Hendershott, T. and Madhavan, A. (2015). Click or call? auction versus search in the over-the-counter market. *The Journal of Finance*, 70(1):419–447.
- [Holthausen et al., 1987] Holthausen, R. W., Leftwich, R. W., and Mayers, D. (1987). The effect of large block transactions on security prices: A cross-sectional analysis. *Journal of Financial Economics*, 19(2):237–267.
- [Huberman and Stanzl, 2004] Huberman, G. and Stanzl, W. (2004). Price manipulation and quasi-arbitrage. *Econometrica*, 72(4):1247–1275.
- [Jurksas et al., 2021] Jurksas, L., Teresiene, D., and Kanapickiene, R. (2021). Liquidity spill-overs in sovereign bond market: An intra-day study of trade shocks in calm and stressful market conditions. *Economies*, 9(1):35.
- [Konstantinovskiy et al., 2016] Konstantinovskiy, V., Ng, K. Y., and Phelps, B. D. (2016). Measuring bond-level liquidity. *Journal of Portfolio Management*, 42(4):116.
- [Lederer, 2013] Lederer, J. (2013). Trust, but verify: benefits and pitfalls of least-squares refitting in high dimensions. *arXiv preprint arXiv:1306.0113*.
- [Lee and Ready, 1991] Lee, C. M. and Ready, M. J. (1991). Inferring trade direction from intraday data. *The Journal of Finance*, 46(2):733–746.
- [Lehalle and Laruelle, 2018] Lehalle, C.-A. and Laruelle, S. (2018). *Market microstructure in practice*. World Scientific.
- [Linciano et al., 2014] Linciano, N., Fancello, F., Gentile, M., and Modena, M. (2014). The liquidity of dual-listed corporate bonds. empirical evidence from italian markets. Technical report, CONSOB. italy14bonds.
- [Mizrach, 2015] Mizrach, B. (2015). Analysis of corporate bond liquidity. Technical report, FINRA.
- [Nagel, 2016] Nagel, J. (2016). Electronic trading in fixed income markets. Technical report, NBIS.

- [Ruzza, 2016] Ruzza, A. (2016). Agency issues in corporate bond trading. Technical report, SSRN.
- [Schneider and Lillo, 2019] Schneider, M. and Lillo, F. (2019). Cross-impact and no-dynamic-arbitrage. *Quantitative Finance*, 19(1):137–154.
- [Schultz, 2001] Schultz, P. (2001). Corporate bond trading costs: A peek behind the curtain. *The Journal of Finance*, 56(2):677–698.
- [Stoll, 1978] Stoll, H. R. (1978). The supply of dealer services in securities markets. *The Journal of Finance*, 33(4):1133–1151.
- [Taranto et al., 2018] Taranto, D. E., Bormetti, G., Bouchaud, J.-P., Lillo, F., and Tóth, B. (2018). Linear models for the impact of order flow on prices. i. history dependent impact models. *Quantitative Finance*, 18(6):903–915.
- [Van Kervel and Menkveld, 2019] Van Kervel, V. and Menkveld, A. (2019). High-frequency trading around large institutional orders. *Journal of Finance*.
- [Wilson et al., 2014] Wilson, D., Trivedi, K., Weisberger, N., Karoui, L., Timcenko, A., Ursua, J., Cole, G., and Yin, S. (2014). The state of play in the leveraged finance market: Ok for now. Technical Report 33, Global Economics Weekly.
- [Zou and Hastie, 2005] Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. In *Journal of the Royal Statistical Society, Series B*, pages 301–320.

A Literature review

Reference	Dataset(s) Name(s)	Period covered
[Schultz, 2001]	CAI	1995-1997
[Chakravarty and Sarkar, 2003]	NAIC	1995-1997
[Chordia et al., 2005]	TAQ, ISSM (Nyse), GovPX	1991-1998
[Bessembinder et al., 2006]	NAIC + TRACE	2001 ; 2002
[Goldstein et al., 2007]	TRACE	2002-2004
[Biais and Green, 2019]	Nyse Archives	1926-1930; 1943-1948
[Dick-Nielsen et al., 2012]	TRACE	2003-2005
[Asquith et al., 2013]	TRACE	2004-2007
[Friewald and Nagler, 2014]	Labelled TRACE	2003-2013
[Mizrach, 2015]	TRACE	2002-2007
[Hendershott and Madhavan, 2015]	TRACE + MarketAxes	2011-2011
[Ruzza, 2016]	TRACE	2004-2012

Table 12: List of empirical papers on transaction costs of corporate bonds

B Data Processing

B.1 Assigning a sign to a trade and identifying RPT

To estimate the sign of transactions, we will first reproduce the essentials of preprocessing to identify such RPTs in [Harris, 2015]. We identify potential RPTs as pairs of sequentially adjacent trades of the same size for which one trade is a customer trade. To find these trades in the Enhanced TRACE data, we first identify all size runs (sequences) of two or more trades of equal size. Next, for each size run, we consider which trades, if any, consist of a pair of trades in a potential RPT. We identify potential RPTs if one trade of two adjacent trades within a size run is a dealer trade with a customer, or if both trades in an adjacent pair are customer trades *and* the dealer both buys and sells. We identify the first such pair as a potential RPT, and then continue searching the size run for any additional pairs that do not involve trades already identified as being part of a potential RPT. [Harris, 2015] found that the RPT rate is above 42% and 41% of customer trades appear to be RPTs. The RPT rate for our entire Enhanced TRACE data set is 23.9%. Moreover, Table 13a shows we found 21.8% RPTs.

	Total	Dealer-customer	Dealer-dealer
Total number	9,413,109	4,523,268	4,889,841
Number of RPT	2,052,644	1,145,127	907,517
Percentage of RPT	21.8%	25.3%	18.5%

(a) Statistics of potential RPTs for selected 1,993 bonds.

	Total	Customer-buy	Customer-sell
Total number	3,378,141	1,921,608	1,456,533
Number percentage	100%	57%	43%
Total volume	3.10×10^{12}	1.73×10^{12}	1.37×10^{12}
Volume percentage	100%	55.8%	44.2%

(b) Statistics of non-RPT dealer-customer trades.

Table 13: Statistics of selected 1,993 bonds for the BA-spread regression.

	Total	Dealer-customer	Dealer-dealer
Total number	3,251,042	1,387,290	1,819,731
Number of RPT	783,022	414,503	365,459
Percentage of RPT	24.1%	29.9%	20.1%

(a) Distribution of potential RPTs.

	Total	Customer-buy	Customer-sell
Total number	972,787	566,232	406,555
Number percentage	100%	58%	42%
Total volume	5.12×10^{11}	2.68×10^{11}	2.45×10^{11}
Volume percentage	100%	52.3%	47.7%

(b) Distribution of non-RPT Dealer-customer trades.

Table 14: Statistics of selected 200 bonds for the price impact analysis.

B.2 Data Filtering

The data cleaning procedure combines the approaches in [Dick-Nielsen, 2014] and [Harris, 2015], with the following steps:

1. Remove canceled trades and apply corrections to ensure that only trades that are actually settled are accounted for. After the removal of canceled trades and canceled corrections records, there are 32,931,539 trades.
2. Remove the transactions reported by agents as both principal and agent in the dealer-to-dealer transactions report to FINRA (see [Dick-Nielsen, 2014]). As a result, 2,095,934 (6.36%) of the reports are removed, with 30,835,605 reports remaining after this step.
3. Remove the transactions on unusual trading days such as weekends and holidays. Thus 5,753 (0.02%) records are removed, with 30,829,852 reports left after this step.

4. Exclude all trade reports with an execution time outside of the normal 8:00AM to 5:15PM ET trading hours. Therefore 745,619 (2.4%) are removed, with 30,084,233 reports remaining after this step.
5. Remove all irregular trades with sales condition codes that indicate late reports, late reports after market hours, weighted average price trades, or trades with special price flags. As a result, 583,157 (1.9%) reports are removed, with 29,501,076 reports left after this step.
6. Remove trade reports with a price below 10. This price filter step excludes 217,321 (0.74%) of the remaining trades, with 29,283,755 reports left after this step.
7. Select reports classified as *corporate bonds* in the dataset. Remove those reports with sub-product indicators such as Mortgage Backed Securities Transactions. Consequently 563,942 (1.94%) of the remaining reports are filtered out, with 28,719,813 reports left.

	Removal (nbe)	Removal (pct)	Number left
Step 1			
Keep settled trades	1,877,866	5.4%	32,931,539
Step 2			
Keep trades reported by dealers	2,095,934	6.36%	30,835,605
Step 3			
Keep business days	5,735	0.02%	30,829,853
Step 4			
Keep opened hours	745,619	2.4%	30,084,233
Step 5			
Keep regular trades	583,157	1.9%	29,501,076
Step 6			
Keep compatible prices	217,321	0.074%	29,182,755
Step 7			
Keep bonds only	563,942	1.94%	28,719,813
Selection for LR			
For bid-ask spread regression	–	–	4,371,363
Selection for PI			
For price impact curves	–	–	1,404,507

Table 15: Data filtering procedure.

C Statistical Test and Regression Analysis

C.1 ANOVA test and Kruskal-Wallis H-test

Suppose there are W groups of observations. (In our example, $W = 6$.) There are n_w observations in group w and the total number among all groups is n . Within each group, $w = 1, 2, \dots, W$ and the observations are denoted as $y_{w,1}, \dots, y_{w,n_w}$ with sample size n_w . Denote $\bar{y}_w = \frac{\sum_{i=1}^{n_w} y_{w,i}}{n_w}$ as the sample mean in group w and $\bar{y} = \frac{\sum_{w=1}^W \sum_{i=1}^{n_w} y_{w,i}}{n}$ as the sample mean of all observations.

One-way ANOVA test. A one-way ANOVA test is applied to samples from two or more groups, possibly with differing sizes. In a one-way ANOVA test, the formula for the F-ratio is $F = \frac{MS_B}{MS_W}$, where $MS_B = \frac{\sum_{j=1}^J n_j (\bar{y}_j - \bar{y})^2}{n-1}$ is the between-group mean square value and $MS_W = \frac{\sum_{j=1}^J \sum_{i=1}^{n_j} (y_{j,i} - \bar{y}_j)^2}{n(n-1)}$ is the within-group mean square value.

C.2 KS test

Denote by $F(x) = \mathbb{P}(X_1 \leq x)$ a cumulative density function of a true underlying distribution of the data and define an *empirical* cumulative density function by $F_n(x) = \mathbb{P}_n(X \leq x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$. Then $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0$.

Suppose that the first sample X_1, X_2, \dots, X_m of size m has a cumulative distribution function (CDF) $F(x)$ and the second sample Y_1, Y_2, \dots, Y_n of size n has a CDF $G(x)$. Suppose that one wants to test

$$H_0 : F = G \quad \text{vs.} \quad H_1 : F \neq G.$$

Let $F_m(x)$ and $G_n(x)$ be their respective empirical CDFs, then $D_{mn} = \left(\frac{mn}{(m+n)} \right)^{\frac{1}{n}} \sup_x |F_m(x) - G_n(x)|$ satisfies the following property of convergence in the distribution:

$$\mathbb{P}(|D_{mn}| < t) \rightarrow H(t) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2i^2 t},$$

where $H(t)$ is the CDF of the KS distribution.

Kruskal-Wallis H-test. The Kruskal-Wallis H-test is a non-parametric version of ANOVA. The test works on two or more independent samples, which may have different sizes. The mathematical formula for H-statistic is

$$H = \frac{12}{n(n+1)} \sum_{j=1}^W \frac{T_j^2}{n_j} - 3(n+1),$$

where T_j is the sum of ranks in the j^{th} group.

C.3 Cross-validation results

Lasso.

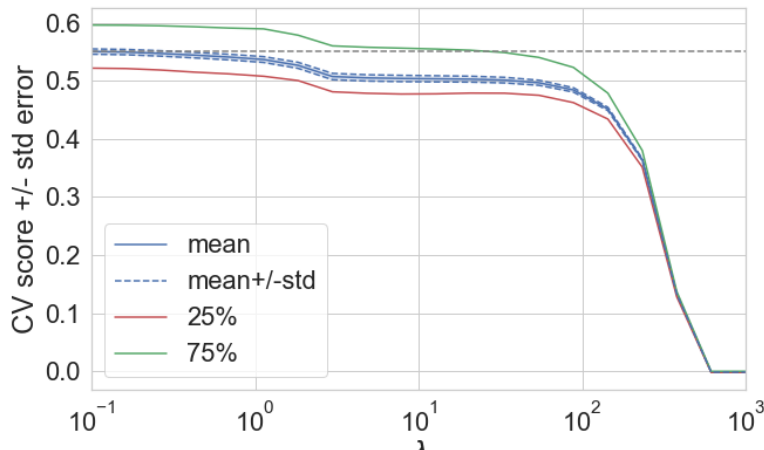


Figure 8: Cross-validation score for Lasso.

λ_l	1.0×10^{-1}	1.62×10^{-1}	2.64×10^{-1}	4.28×10^{-1}	6.95×10^{-1}
\widehat{R}^2	0.550	0.549	0.547	0.544	0.541
Number in \mathfrak{J}_1	15	12	12	14	15
Number in \mathfrak{J}_2	152	153	155	156	157
λ_l	1.13	1.83	2.98	4.83	7.85
\widehat{R}^2	0.527	0.517	0.507	0.505	0.484
Number in \mathfrak{J}_1	18	16	16	16	14
Number in \mathfrak{J}_2	158	160	161	161	162
λ_l	1.27×10	2.07×10	3.36×10	5.46×10	8.86×10
\widehat{R}^2	0.484	0.483	0.481	0.477	0.464
Number in \mathfrak{J}_1	14	16	13	16	16
Number in \mathfrak{J}_2	160	161	161	160	158
λ_l	1.44×10^2	2.34×10^2	3.79×10^2	6.16×10^2	1.00×10^3
\widehat{R}^2	0.442	0.364	0.133	-0.001	-0.001
Number of \tilde{R}^2 in $CI(\lambda_l)$	17	13	16	8	8
Number of \tilde{R}^2 in $CI_2(\lambda_l)$	155	149	149	175	175

Table 16: Number of \tilde{R}^2 in the confidence interval for Lasso.

Ridge.

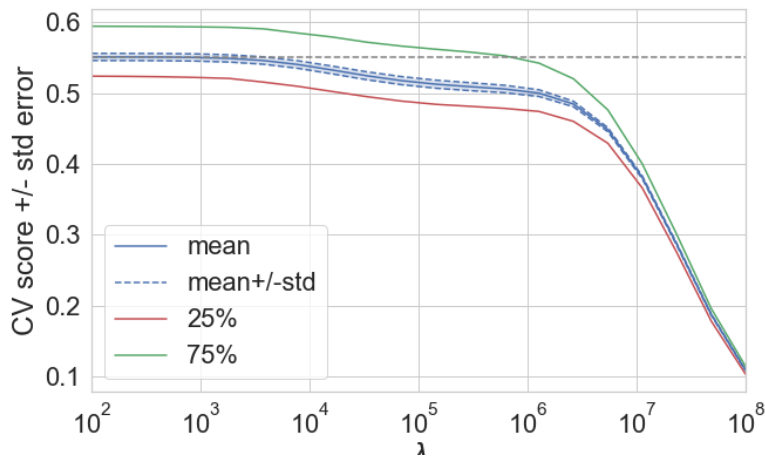


Figure 9: Cross-validation score for Ridge.

λ_r	1.00×10^2	2.07×10^2	4.28×10^2	8.85×10^2	1.83×10^3
\widehat{R}^2	0.551	0.551	0.551	0.550	0.549
Number in \mathfrak{I}_1	9	9	11	9	11
Number in \mathfrak{I}_2	168	167	166	166	166
λ_r	3.79×10^3	7.84×10^3	1.62×10^4	3.36×10^4	6.95×10^4
\widehat{R}^2	0.546	0.540	0.533	0.524	0.518
Number in \mathfrak{I}_1	11	14	15	13	11
Number in \mathfrak{I}_2	168	169	170	172	173
λ_r	1.44×10^5	2.98×10^5	6.16×10^5	1.27×10^6	2.64×10^6
\widehat{R}^2	0.513	0.509	0.506	0.500	0.485
Number in \mathfrak{I}_1	11	11	13	17	11
Number in \mathfrak{I}_2	172	171	171	167	164
λ_r	5.46×10^6	1.23×10^7	2.34×10^7	4.83×10^7	1.00×10^8
\widehat{R}^2	0.448	0.381	0.287	0.188	0.109
Number in \mathfrak{I}_1	8	8	6	7	7
Number in \mathfrak{I}_2	159	149	143	136	136

Table 17: Number of \widehat{R}^2 in the confidence interval for Ridge regression.

EN.

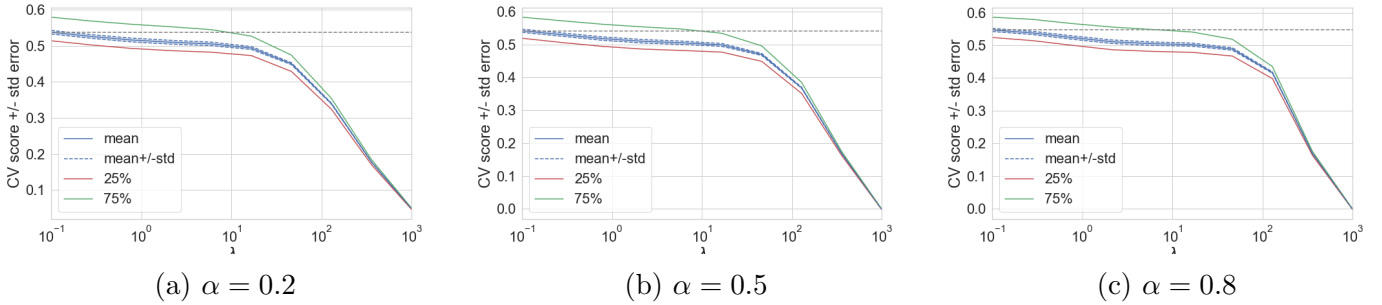


Figure 10: Cross-validation score for EN.

$\alpha = 0.2, \lambda_e =$	1.0×10^{-1}	2.78×10^{-1}	7.74×10^{-1}	2.15	5.99
\widehat{R}^2	0.537	0.526	0.516	0.510	0.505
\mathfrak{I}_1	15	16	15	13	12
\mathfrak{I}_2	161	160	160	160	159
$\alpha = 0.2, \lambda_e =$	1.68×10	4.64×10	1.29×10^2	3.59×10^2	1.00×10^3
\widehat{R}^2	0.494	0.450	0.340	0.178	0.048
\mathfrak{I}_1	9	2	16	11	13
\mathfrak{I}_2	152	150	133	137	142
$\alpha = 0.5, \lambda_e =$	1.0×10^{-1}	2.78×10^{-1}	7.74×10^{-1}	2.15	5.99
\widehat{R}^2	0.542	0.530	0.518	0.5010	0.505
\mathfrak{I}_1	16	16	17	11	12
\mathfrak{I}_2	160	162	163	161	160
$\alpha = 0.5, \lambda_e =$	1.68×10	4.64×10	1.29×10^2	3.59×10^2	1.00×10^3
\widehat{R}^2	0.498	0.470	0.368	0.169	-0.001
\mathfrak{I}_1	13	9	17	12	16
\mathfrak{I}_2	158	154	154	138	174
$\alpha = 0.8, \lambda_e =$	1.0×10^{-1}	2.78×10^{-1}	7.74×10^{-1}	2.15	5.99
\widehat{R}^2	0.547	0.538	0.524	0.508	0.504
\mathfrak{I}_1	16	16	15	17	12
\mathfrak{I}_2	159	160	160	162	159
$\alpha = 0.8, \lambda_e =$	1.68×10	4.64×10	1.29×10^2	3.59×10^2	1.00×10^3
\widehat{R}^2	0.502	0.487	0.417	0.170	-0.001
\mathfrak{I}_1	13	10	9	12	16
\mathfrak{I}_2	160	158	147	139	174

Table 18: Number of \widehat{R}^2 in confidence interval for EN.

C.4 Comparison between the Enhanced TRACE and Standard TRACE datasets

The comparison between the Enhanced TRACE and Standard TRACE datasets for OLS and two-step LASSO is summarized here.

Table 20 compares the OLS regression with data in Enhanced TRACE and Standard TRACE from the same period January 01, 2015-December 31, 2016. It shows similar R^2 values, with 55.4% for Enhanced TRACE and 54.4% for Standard TRACE. The relative difference between the regression coefficients is small except for the following features.

- *Prop Volume sell \$* and *Prop Volume buy \$*: For Enhanced TRACE, the coefficient of *Prop Volume buy \$* is 50% larger than that of *Prop Volume sell \$*. This relationship is reversed for Standard TRACE. Further studies of the distribution of capped transactions show that 65% of these transactions are customer sell orders. This contributes to the changes of these two features in the Standard TRACE as well as the difference in regression coefficients. See Figure 11.
- *Turnover* and *Years to maturity*: The coefficients of these two features are much bigger for Standard TRACE. This difference is tolerable as neither of two features is significant ($p \geq 0.1$).

Comparing the two-step LASSO model on these two datasets (see Table 19) confirms similar conclusions: the outstanding features are consistent and the regression coefficients are compatible.

	Enhanced TRACE	Standard TRACE
Volatility	80.7218*** (0.253)	72.96 *** (0.21)
Trading activity	42.1450*** (0.451)	45.46 *** (0.496)
Log(Total Volume)	-21.1246*** (0.252)	-19.39 *** (0.261)
Years since issuance	0.25 *** (0.052)	0.32 *** (0.041)
Constant	86.1737	63.25
N	152,408	130,716
R^2	52.8 %	52.5 %

Standard errors in parenthesis. Significance levels: * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Two-tailed test.

Source: Enhanced TRACE (2015-2016).

Table 19: Two-step Lasso regression table: the impact on bid-ask spread (in bp).

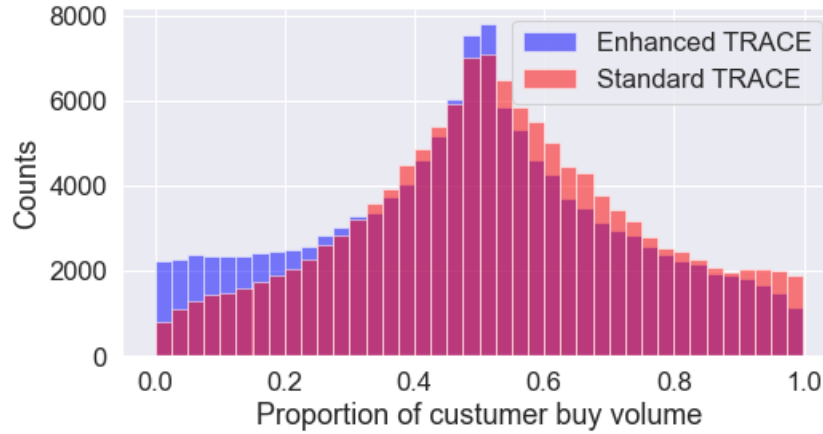


Figure 11: Proportion of customer buy orders in dollars (mean 0.48 for Enhanced TRACE and mean 0.52 ofr Standard TRACE).

	TRACE Enhanced		Standard TRACE	
Volatility	77.73	***	66.4	***
Number of trade days	-3.66	***	-2.85	***
Prop number of buys	10.32	***	9.93	***
Prop number of sells	32.45	***	28.68	***
Trading activity	46.32	***	51.27	***
Prop volume sell \$	16.35	***	23.18	***
Prop Volume buy \$	26.42	***	15.43	***
Log(total volume)	-21.40	***	-20.62	***
Avg price	-0.12	***	-0.14	***
Coupon	-0.47	***	-0.42	***
Duration	1.42	***	1.98	***
Years to maturity	-0.07		-0.37	
Years since issuance	1.26	***	1.33	***
Turnover	-1.77		-23.54	
LIBOR-OIS	34.01	***	30.62	***
Indicator of high yield bonds	26.79	***	23.49	***
Indicator of investment grade bonds	15.99	***	13.12	***
Indicator of basic materials sector	8.90	***	8.88	***
Indicator of communications sector	5.13	***	4.03	***
Indicator of consumer, cyclical sector	4.03	***	3.41	***
Indicator of consumer, non-cyclical sector	4.97	***	4.40	***
Indicator of energy sector	3.88	***	3.49	***
Indicator of financial sector	4.05	***	2.99	***
Indicator of industrial sector	2.815	***	3.16	***
Indicator of technology sector	4.41	***	4.06	***
Indicator of utilities sector	4.58	***	4.29	***
Constant	42.77	***	35.80	***
N	152,408		130,716	
R ²	55.4	%	54.4	%

Standard errors in parenthesis. Significance levels: * p<0.1, ** p<0.05, *** p<0.01. Two-tailed test. Source: TRACE Enhanced (2015-2016).

Table 20: OLS regression: Comparison between TRACE Enhanced and Standard TRACE

D Additional Details of the Price Impact Analysis

D.1 Estimation of TIM2 Model

For simplicity, we omit the subscript b (or bond b) in the derivation here. Assume that there are two types of events $\Pi := \{+1, -1\}$ with $+1$ denoting customer-buy orders and -1 denoting customer-sell orders. In this case, the mid-price dynamics (14) leads to the following expression for mid-price changes:

$$\begin{aligned}
R_k(1) := M_{k+1} - M_t &= \sum_{\pi \in \Pi} G_{\pi}(0) I(\pi_{k+1} = \pi) V_{k+1}^{\alpha} \epsilon_{k+1} + \eta_{k+1} \\
&+ \sum_{j=0}^{\infty} \sum_{\pi' \in \Pi} \underbrace{(G_{\pi'}(j+1) - G_{\pi'}(j))}_{\Delta_1 G_{\pi'}(j)} I(\pi_{k-j} = \pi') \epsilon_{k-j} V_{k-j}^{\alpha}. \tag{23}
\end{aligned}$$

As a consequence, we can write the conditional response functions and response correlation matrix as:

$$S_\pi(l) = \mathbb{E}[R_k \cdot \epsilon_{k-l+1} | \pi_{k-l+1} = \pi] = \frac{\mathbb{E}[R_k \cdot \epsilon_{k-l+1} I(\pi_{k-l+1} = \pi)]}{\mathbb{P}(\pi)}, \quad (24)$$

$$C_{\pi, \pi'}(n) = \mathbb{E}[\epsilon_t \epsilon_{t+n} V_{t+n}^\alpha | \pi_t = \pi, \pi_{t+n} = \pi'] = \frac{\mathbb{E}[\epsilon_t I(\pi_t = \pi) \cdot V_{t+n}^\alpha \epsilon_{t+n} I(\pi_{t+n} = \pi')]}{\mathbb{P}(\pi_t = \pi, \pi_{t+n} = \pi')}, \quad (25)$$

for $1 \leq l \leq L$ and $-N \leq n \leq L$. Then the response function can be written as, for $\pi = +1$ or -1 :

$$S_\pi(l) = \sum_{\pi' \in \Pi} \mathbb{P}(\pi' | \pi) G_{\pi'}(0) C_{\pi, \pi'}(l) + \sum_{j=0}^{+\infty} \sum_{\pi' \in \Pi} \mathbb{P}(\pi' | \pi) \Delta_1 G_{\pi'}(j) \cdot C_{\pi, \pi'}(l - j - 1). \quad (26)$$

Denote $\tilde{C}_{\pi, \pi'}(l) = \mathbb{P}(\pi_{t+l} = \pi' | \pi_t = \pi) C_{\pi, \pi'}(l)$ and $\tilde{S}_\pi(l) = S_\pi(l) - \sum_{\pi' \in \Pi} G_{\pi'}(0) \tilde{C}_{\pi, \pi'}(l)$, then:

$$\underbrace{\begin{pmatrix} \tilde{S}_{+1}(1) \\ \tilde{S}_{+1}(2) \\ \vdots \\ \tilde{S}_{+1}(L) \\ \tilde{S}_{-1}(1) \\ \tilde{S}_{-1}(2) \\ \vdots \\ \tilde{S}_{-1}(L) \end{pmatrix}}_{\tilde{\mathcal{S}}(L)} = \underbrace{\begin{bmatrix} \tilde{C}_{+1,+1}(0) & \tilde{C}_{+1,+1}(-1) & \cdots & \tilde{C}_{+1,+1}(-N+1) & \tilde{C}_{+1,-1}(0) & \tilde{C}_{+1,-1}(-1) & \cdots & \tilde{C}_{+1,-1}(-N+1) \\ \tilde{C}_{+1,+1}(1) & \tilde{C}_{+1,+1}(0) & \cdots & \tilde{C}_{+1,+1}(-N+2) & \tilde{C}_{+1,-1}(1) & \tilde{C}_{+1,-1}(0) & \cdots & \tilde{C}_{+1,-1}(-N+2) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \tilde{C}_{+1,+1}(L-1) & \cdots & \cdots & \tilde{C}_{+1,+1}(-N+L) & \tilde{C}_{+1,-1}(L-1) & \cdots & \cdots & \tilde{C}_{+1,-1}(-N+L) \\ \tilde{C}_{-1,+1}(0) & \tilde{C}_{-1,+1}(-1) & \cdots & \tilde{C}_{-1,+1}(-N+1) & \tilde{C}_{-1,-1}(0) & \tilde{C}_{-1,-1}(1) & \cdots & \tilde{C}_{-1,-1}(N-1) \\ \tilde{C}_{-1,+1}(1) & \tilde{C}_{-1,+1}(0) & \cdots & \tilde{C}_{-1,+1}(-N+2) & \tilde{C}_{-1,-1}(1) & \tilde{C}_{-1,-1}(0) & \cdots & \tilde{C}_{-1,-1}(-N+2) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \tilde{C}_{-1,+1}(L-1) & \cdots & \cdots & \tilde{C}_{-1,+1}(-N+L) & \tilde{C}_{-1,-1}(L-1) & \cdots & \cdots & \tilde{C}_{-1,-1}(-N+L) \end{bmatrix}}_{\tilde{\mathcal{C}}(N,L)} \underbrace{\begin{pmatrix} \Delta_1 G_{+1}(0) \\ \Delta_1 G_{+1}(1) \\ \vdots \\ \Delta_1 G_{+1}(N-1) \\ \Delta_1 G_{-1}(0) \\ \Delta_1 G_{-1}(1) \\ \vdots \\ \Delta_1 G_{-1}(N-1) \end{pmatrix}}_{\tilde{\mathcal{G}}(N)}, \quad (27)$$

where we have:

$$\begin{aligned} S_1(l) &= \mathbb{E}[R_k | \epsilon_{k-l+1} = 1], \quad S_{-1}(l) = -\mathbb{E}[R_k | \epsilon_{k-l+1} = -1], \\ \tilde{C}_{1,-1}(n) &= -\frac{\mathbb{P}(\pi_t = 1, \pi_{t+n} = -1) \mathbb{E}[V_{t+n}^\alpha | \pi_t = 1, \pi_{t+n} = -1]}{\mathbb{P}(\pi_t = 1)}, \\ \tilde{C}_{-1,1}(n) &= -\frac{\mathbb{P}(\pi_t = -1, \pi_{t+n} = 1) \mathbb{E}[V_{t+n}^\alpha | \pi_t = -1, \pi_{t+n} = 1]}{\mathbb{P}(\pi_t = -1)}, \\ \tilde{C}_{-1,-1}(n) &= \frac{\mathbb{P}(\pi_t = -1, \pi_{t+n} = -1) \mathbb{E}[V_{t+n}^\alpha | \pi_t = -1, \pi_{t+n} = -1]}{\mathbb{P}(\pi_t = -1)}, \\ \tilde{C}_{1,1}(n) &= \frac{\mathbb{P}(\pi_t = 1, \pi_{t+n} = 1) \mathbb{E}[V_{t+n}^\alpha | \pi_t = 1, \pi_{t+n} = 1]}{\mathbb{P}(\pi_t = 1)}. \end{aligned}$$

The signature plot for TIM2 model [Eisler et al., 2012] can be similarly defined as:

$$\begin{aligned} l D_{\text{TIM2}}(l) &= \sum_{0 \leq n < l} \sum_{+1} G_{+1}(l-n)^2 P(+1) + \sum_{n > 0} \sum_{+1} [G_{+1}(l+n) - G_{+1}(n)]^2 P(+1) \\ &+ 2 \sum_{0 \leq n < n' < l+1, -1} G_{+1}(l-n) G_{-1}(l-n') C_{+1,-1}(n'-n) \\ &+ 2 \sum_{0 < n < n' < l+1, -1} [G_{+1}(l+n) - G_{+1}(n)] [G_{-1}(l+n') - G_{-1}(n')] C_{+1,-1}(n-n') \\ &+ 2 \sum_{0 \leq n < l} \sum_{n' > 0} \sum_{+1, -1} G_{+1}(l-n) [G_{-1}(l+n') - G_{-1}(n')] C_{-1,+1}(n'+n). \end{aligned} \quad (28)$$