

Supporting Information

How Complexity and Uncertainty Grew with Algorithmic Trading

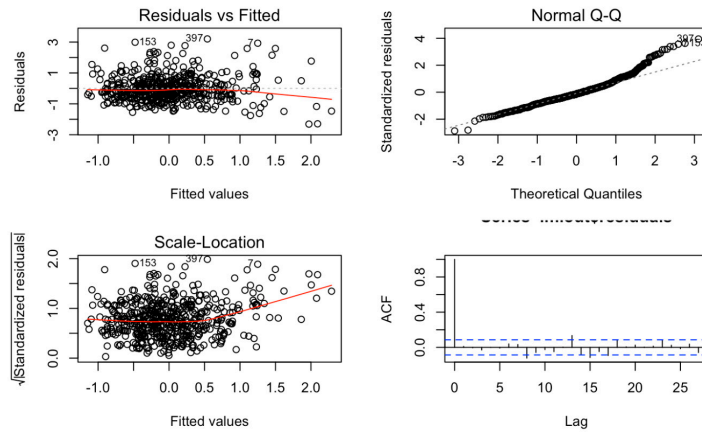
Martin Hilbert ^{1,*} and David Darmon ²

¹ Communication, Computational Social Science; University of California, Davis, CA 95616, USA; hilbert@ucdavis.edu

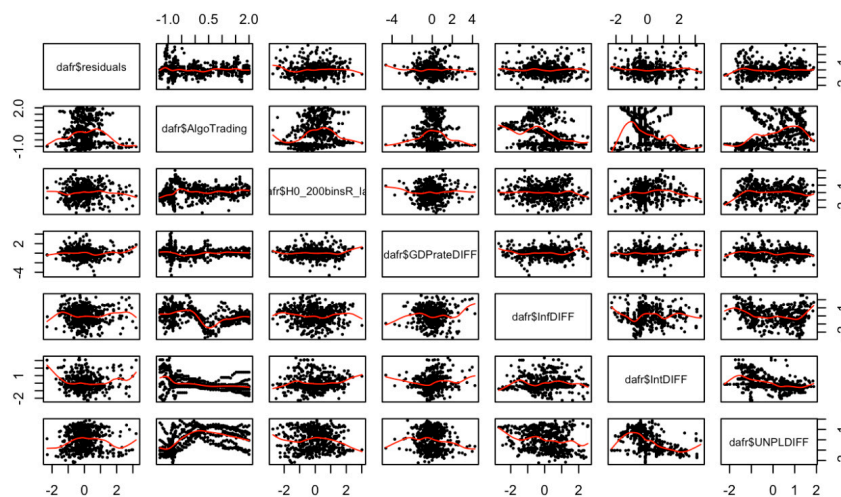
² Department of Mathematics; Monmouth University, West Long Branch, NJ 07764, USA; ddarmon@monmouth.edu

SI.1. Assumptions of Linear Regression

We ran a series diagnostic analysis on the residuals of our multiple linear regressions and feel comfortable with the multiple linear regression with normal noise model. We found that the most worrisome violation stems from the normality of the residuals.

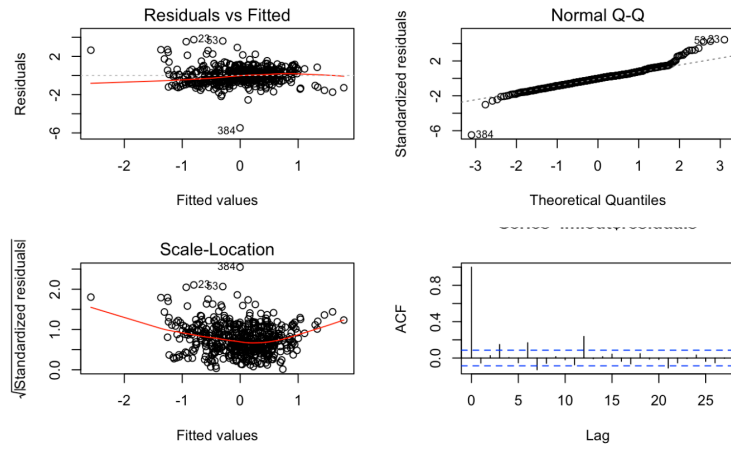


(a)

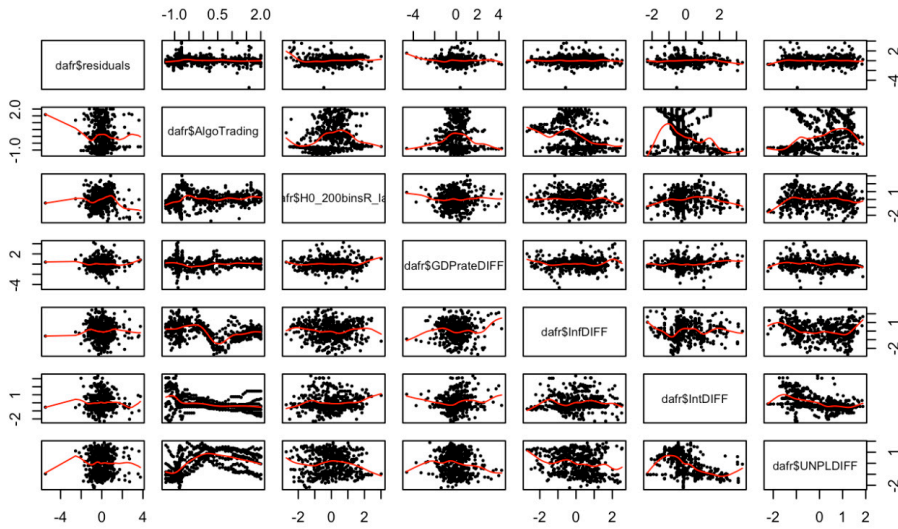


(b)

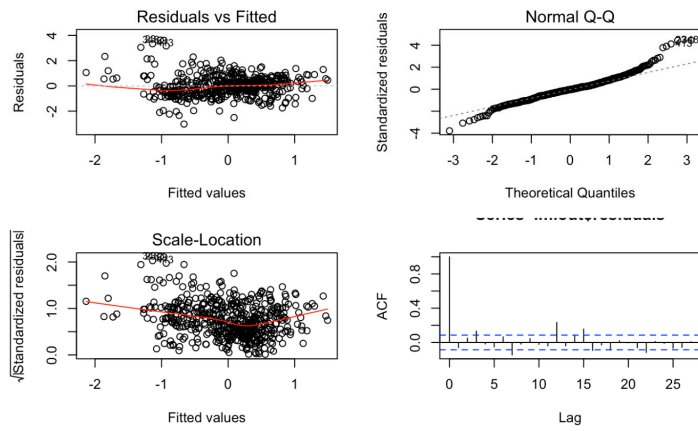
Figure S1. Cont.



(c)

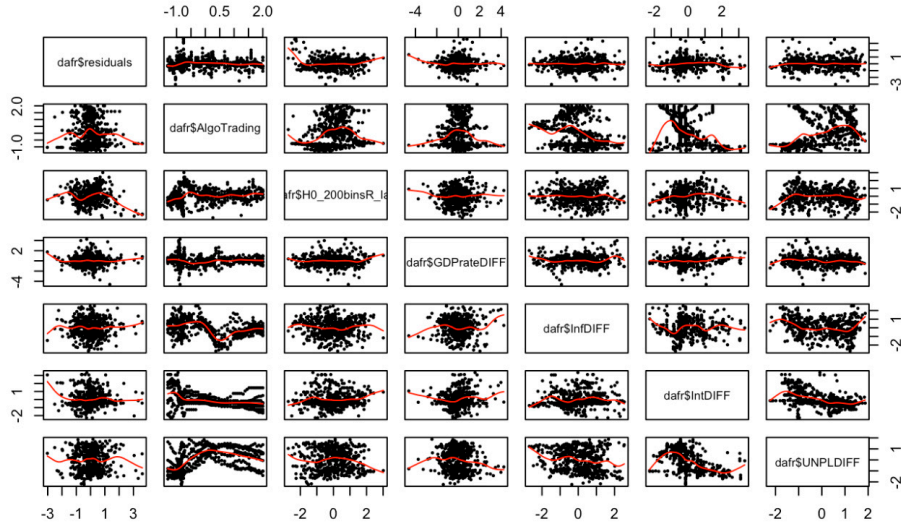


(d)



(e)

Figure S1. Cont.



(f)

Figure SI.1. Residual diagnostics for normal noise model (a-b) $E_{eM,200}$, (c-d) $C_{eM,200}$, (e-f) $h_{eM,200}$.

Fortunately, this also tends to be the least concerning assumption, since a central limit theorem-based results makes the actual distribution of the residuals wash out with enough data. Given the nature of our data from eight different currency pairs, linked in time, any alternative bootstrapping procedure would have to be fairly intricately involved, which would reduce transparency and replicability. We feel that our normal noise assumptions are within what is generally accepted in the community.

SI.2. Comparison without Lagged Term

Social dynamics often contain an important path dependency. Therefore, researchers interested in other influences often include a lagged value of the dependent variable as an independent variable in what is known as dynamic panel data models. Our analysis confirms that the dynamics of the previous bi-monthly period $t-1$ is highly predictive of the next one t . In this sense, our results show the influence of the other tested variables independent from this effect of path dependency. However, there is an ongoing discussion about the practice of including lagged value dependent variable in panel data models [1]. Leaving out the lagged term leads to a lot of autocorrelation in the residuals of the regression, which would violate basic assumptions. Just to make sure, we also ran the exercise without it. The main conclusions drawn from this study are strongly reinforced when running the tests without the lagged term. The influence of algorithmic trading increases. For example, Figure SI.2 shows the case of measure $E_{eM,200}$. $AT_{emp,200}$ is only weakly significant when considering lagged path dependency, but becomes significantly stronger without. The same is shown for measure $h_{eM,200}$, and in general applies to all tests we have seen. We present the version with lagged term in the main article, because it better corresponds to the basic assumptions of linear analysis, and because it presents the more conservative version of our results.

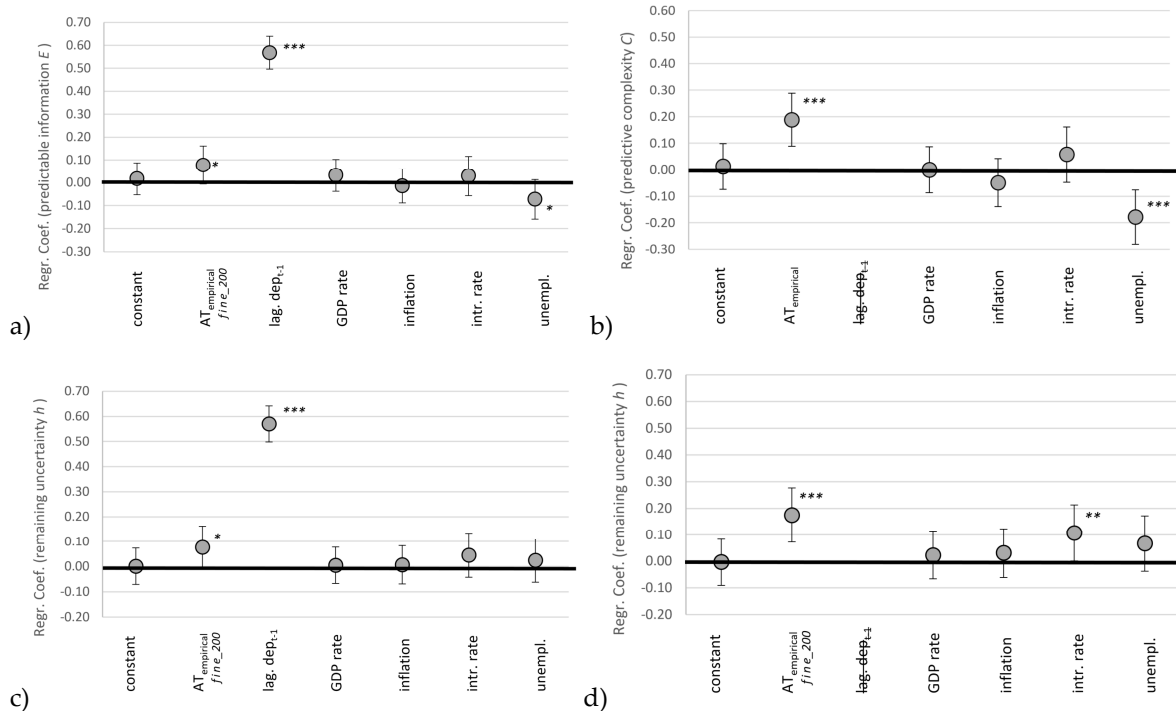


Figure SI.2 Regression coefficients for bi-monthly changes in dynamics measured in 200 fine-grained bins, indicating 95 % confidence intervals with error bars, (a) predictable information E with lagged term; (b) without lagged term; (c) remaining uncertainty h with lagged term; (b) without lagged term. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$ ($N = 520$).

SI.3. Full Array of Models for H1

As discussed in the main article, we also use three complementary estimates for the rise of algorithmic trading (AT), namely empirical, linear and exponential. Given that our information-theoretic dynamical-systems indicators are not as deeply established in the social sciences, and given that there are less agreed upon best practices in their estimation, we test their robustness by using two different methods to calculate them for each bi-monthly interval, namely ϵ -machines (epsilon machines, ϵ_M) [2] and frequency counts (f_q) [3]. Note that predictive complexity is a measure of the associated ϵ -machine [4,5], and we only derive those once, with the Causal State Splitting Reconstruction (CSSR) algorithm.

It shows that our results are quite robust, independent from the derivation method, and independent from the estimate for the rise of algorithmic trading. In Tables SI.1 and SI.2, model (2) and model (7) are marked in bold since these are the ones presented in the main article. As can be seen, these are among the models where algorithmic trading has the least influence, and are therefore a rather conservative estimate of our broader results.

Table SI.1: Tests for bi-monthly changes in coarse-grained (20 bins based) complexity in form predictable information (E) and predictive complexity (C), measured according to frequency counts (f_q) and ϵ -machines (ϵ_M); showing unstandardized beta coefficients, with standard errors in italic parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$ (N = 520).

Dep. Var.	(1) $E_{f_q_20}$	(2) $E_{\epsilon M_20}$	(3) $E_{f_q_20}$	(4) $E_{\epsilon M_20}$	(5) $E_{f_q_20}$	(6) $E_{\epsilon M_20}$	(7) $C_{\epsilon M_20}$	(8) $C_{\epsilon M_20}$	(9) $C_{\epsilon M_20}$
constant	-0.003 (0.035)	0.004 (0.036)	-0.001 (0.035)	0.006 (0.036)	-0.003 (0.035)	0.004 (0.036)	-0.009 (0.041)	-0.008 (0.041)	-0.009 (0.041)
AT _{emp}	-0.149*** (0.042)	-0.118*** (0.043)					-0.232*** (0.049)		
AT _{lin}			-0.199*** (0.045)	-0.173*** (0.046)				-0.231*** (0.051)	
AT _{exp}					-0.150*** (0.042)	-0.117*** (0.043)			-0.228*** (0.048)
dep _{t-1}	0.412*** (0.040)	0.454*** (0.039)	0.395*** (0.041)	0.438*** (0.039)	0.411*** (0.040)	0.454*** (0.039)	0.193*** (0.044)	0.196*** (0.044)	0.195*** (0.044)
GDP _r	0.004 (0.036)	0.013 (0.037)	0.001 (0.036)	0.012 (0.036)	0.004 (0.036)	0.013 (0.037)	-0.016 (0.041)	-0.021 (0.041)	-0.016 (0.041)
infl	0.019 (0.037)	0.007 (0.038)	0.013 (0.037)	-0.002 (0.038)	0.021 (0.037)	0.009 (0.038)	0.006 (0.043)	0.011 (0.043)	0.010 (0.043)
intr	0.064 (0.043)	0.057 (0.044)	0.053 (0.043)	0.044 (0.044)	0.061 (0.043)	0.056 (0.044)	0.012 (0.049)	0.012 (0.050)	0.010 (0.050)
unpl	-0.144*** (0.044)	-0.100** (0.044)	-0.123*** (0.044)	-0.079* (0.045)	-0.146*** (0.044)	-0.102** (0.044)	-0.104** (0.049)	-0.086* (0.050)	-0.107* (0.049)
F(6,513)	48***	42***	50***	43***	48***	41***	16***	16***	16***
adjusted R ²	0.352	0.317	0.361	0.326	0.353	0.318	0.148	0.144	0.148

Table SI.2: Tests for bi-monthly changes in fine-grained (200 bins based) complexity in form predictable information (E) and predictive complexity (C), measured according to frequency counts (f_q) and ϵ -machines (ϵ_M); showing unstandardized beta coefficients, with standard errors in italic parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$ (N = 520).

Dep. Var.	(1) $E_{f_q_{200}}$	(2) $E_{\epsilon M_{200}}$	(3) $E_{f_q_{200}}$	(4) $E_{\epsilon M_{200}}$	(5) $E_{f_q_{200}}$	(6) $E_{\epsilon M_{200}}$	(7) $C_{\epsilon M_{200}}$	(8) $C_{\epsilon M_{200}}$	(9) $C_{\epsilon M_{200}}$
constant	0.008 <i>(0.036)</i>	0.017 <i>(0.035)</i>	0.008 <i>(0.036)</i>	0.017 <i>(0.035)</i>	0.008 <i>(0.036)</i>	0.017 <i>(0.035)</i>	0.006 <i>(0.037)</i>	0.005 <i>(0.037)</i>	0.006 <i>(0.037)</i>
AT _{emp}	0.094** <i>(0.043)</i>	0.077* <i>(0.042)</i>					0.122*** <i>(0.045)</i>		
AT _{lin}			0.076* <i>(0.045)</i>	0.061 <i>(0.044)</i>				0.163*** <i>(0.047)</i>	
AT _{exp}					0.093** <i>(0.043)</i>	0.078* <i>(0.041)</i>			0.128*** <i>(0.044)</i>
dep _{t-1}	0.510*** <i>(0.038)</i>	0.568*** <i>(0.036)</i>	0.516*** <i>(0.038)</i>	0.572*** <i>(0.036)</i>	0.510*** <i>(0.038)</i>	0.568*** <i>(0.036)</i>	0.507*** <i>(0.038)</i>	0.495*** <i>(0.038)</i>	0.505*** <i>(0.038)</i>
GDP _r	0.040 <i>(0.037)</i>	0.032 <i>(0.035)</i>	0.044 <i>(0.037)</i>	0.035 <i>(0.035)</i>	0.040 <i>(0.037)</i>	0.032 <i>(0.035)</i>	0.014 <i>(0.038)</i>	0.015 <i>(0.037)</i>	0.013 <i>(0.038)</i>
infl	0.001 <i>(0.039)</i>	-0.015 <i>(0.037)</i>	-0.005 <i>(0.038)</i>	-0.020 <i>(0.037)</i>	0.000 <i>(0.038)</i>	-0.016 <i>(0.037)</i>	-0.020 <i>(0.039)</i>	-0.014 <i>(0.039)</i>	-0.020 <i>(0.039)</i>
intr	0.057 <i>(0.045)</i>	0.029 <i>(0.043)</i>	0.052 <i>(0.045)</i>	0.025 <i>(0.043)</i>	0.058 <i>(0.045)</i>	0.030 <i>(0.043)</i>	0.044 <i>(0.045)</i>	0.056 <i>(0.045)</i>	0.048 <i>(0.046)</i>
unpl	-0.109** <i>(0.045)</i>	-0.073* <i>(0.043)</i>	-0.110** <i>(0.046)</i>	-0.074* <i>(0.044)</i>	-0.108** <i>(0.045)</i>	-0.072* <i>(0.042)</i>	-0.016 <i>(0.045)</i>	-0.037 <i>(0.046)</i>	-0.016 <i>(0.045)</i>
F(6,513)	41***	49***	40***	48***	41***	49***	36***	37***	36***
adjusted R ²	0.314	0.356	0.312	0.354	0.314	0.356	0.287	0.293	0.289

SI.4. Full Array of Models for H2

Model (2) in Tables SI.3 and SI.4 is the one presented in the main article. As can be seen, it is among those models in which algorithmic trading has the least influence, and are therefore a rather conservative estimate of our broader results.

Table SI.3: Tests for bi-monthly changes in coarse-grained (20 bins based) remaining uncertainty in form entropy rate (h), measured according to frequency counts (f_q) and ϵ -machines (e_M); showing unstandardized beta coefficients, with standard errors in italic parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$ ($N = 520$).

Dep. Var.	(1) $h_{f_q,20}$	(2) $h_{e_M,20}$	(3) $h_{f_q,20}$	(4) $h_{e_M,20}$	(5) $h_{f_q,20}$	(6) $h_{e_M,20}$
constant	-0.011 <i>(0.040)</i>	-0.014 <i>(0.041)</i>	-0.010 <i>(0.041)</i>	-0.013 <i>(0.041)</i>	-0.011 <i>(0.040)</i>	-0.014 <i>(0.041)</i>
AT _{temp}	-0.262*** <i>(0.048)</i>	-0.209*** <i>(0.049)</i>				
AT _{lin}			-0.248*** <i>(0.051)</i>	-0.187*** <i>(0.052)</i>		
AT _{exp}					-0.262*** <i>(0.048)</i>	-0.206*** <i>(0.049)</i>
dep _{t-1}	0.126*** <i>(0.045)</i>	0.210*** <i>(0.044)</i>	0.134*** <i>(0.045)</i>	0.219*** <i>(0.044)</i>	0.127*** <i>(0.045)</i>	0.212*** <i>(0.044)</i>
GDP _r	-0.051 <i>(0.041)</i>	-0.030 <i>(0.042)</i>	-0.058 <i>(0.041)</i>	-0.035 <i>(0.042)</i>	-0.051 <i>(0.041)</i>	-0.030 <i>(0.042)</i>
infl	0.040 <i>(0.043)</i>	0.006 <i>(0.044)</i>	0.049 <i>(0.043)</i>	0.016 <i>(0.044)</i>	0.045 <i>(0.042)</i>	0.010 <i>(0.043)</i>
intr	0.040 <i>(0.049)</i>	0.001 <i>(0.050)</i>	0.042 <i>(0.050)</i>	0.006 <i>(0.051)</i>	0.037 <i>(0.049)</i>	-0.001 <i>(0.050)</i>
unpl	-0.098** <i>(0.049)</i>	-0.066 <i>(0.050)</i>	-0.081 <i>(0.050)</i>	-0.055 <i>(0.051)</i>	-0.102** <i>(0.049)</i>	-0.070 <i>(0.050)</i>
F(6,513)	17***	13***	16***	12***	17***	13***
adjusted R ²	0.157	0.119	0.148	0.111	0.158	0.119

Table SI.4: Tests for bi-monthly changes in fine-grained (200 bins based) remaining uncertainty in form entropy rate (h), measured according to frequency counts (h_{fq}) and ϵ -machines (h_{eM}); showing unstandardized beta coefficients, with standard errors in italic parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$ ($N = 520$).

Dep. Var.	(1) $h_{fq_{200}}$	(2) $h_{eM_{200}}$	(3) $h_{fq_{200}}$	(4) $h_{eM_{200}}$	(5) $h_{fq_{200}}$	(6) $h_{eM_{200}}$
constant	0.002 <i>(0.036)</i>	0.001 <i>(0.036)</i>	0.000 <i>(0.036)</i>	-0.001 <i>(0.036)</i>	0.002 <i>(0.036)</i>	0.001 <i>(0.036)</i>
AT _{emp}	0.088** <i>(0.043)</i>	0.076* <i>(0.043)</i>				
AT _{lin}			0.152*** <i>(0.045)</i>	0.119*** <i>(0.045)</i>		
AT _{exp}					0.096** <i>(0.042)</i>	0.082* <i>(0.042)</i>
dep _{t-1}	0.566*** <i>(0.037)</i>	0.570*** <i>(0.036)</i>	0.547*** <i>(0.037)</i>	0.559*** <i>(0.037)</i>	0.563*** <i>(0.037)</i>	0.568*** <i>(0.036)</i>
GDP _r	0.010 <i>(0.036)</i>	0.005 <i>(0.036)</i>	0.010 <i>(0.036)</i>	0.005 <i>(0.036)</i>	0.010 <i>(0.036)</i>	0.005 <i>(0.036)</i>
infl	0.008 <i>(0.038)</i>	0.007 <i>(0.038)</i>	0.020 <i>(0.038)</i>	0.015 <i>(0.038)</i>	0.009 <i>(0.038)</i>	0.007 <i>(0.038)</i>
intr	0.017 <i>(0.044)</i>	0.045 <i>(0.044)</i>	0.034 <i>(0.044)</i>	0.057 <i>(0.044)</i>	0.020 <i>(0.044)</i>	0.047 <i>(0.044)</i>
unpl	0.016 <i>(0.043)</i>	0.024 <i>(0.043)</i>	-0.007 <i>(0.044)</i>	0.007 <i>(0.044)</i>	0.016 <i>(0.043)</i>	0.024 <i>(0.043)</i>
F(6,513)	46***	45***	47***	46***	46***	45***
adjusted R ²	0.340	0.337	0.349	0.342	0.341	0.338

SI.5. Histograms of Bid-Ask Spreads

Figure SI.3 shows the raw (non-standardized) bid-ask spreads for our eight currency pairs. The mean and standard deviation naturally varies, depending on the relative value of the currency.

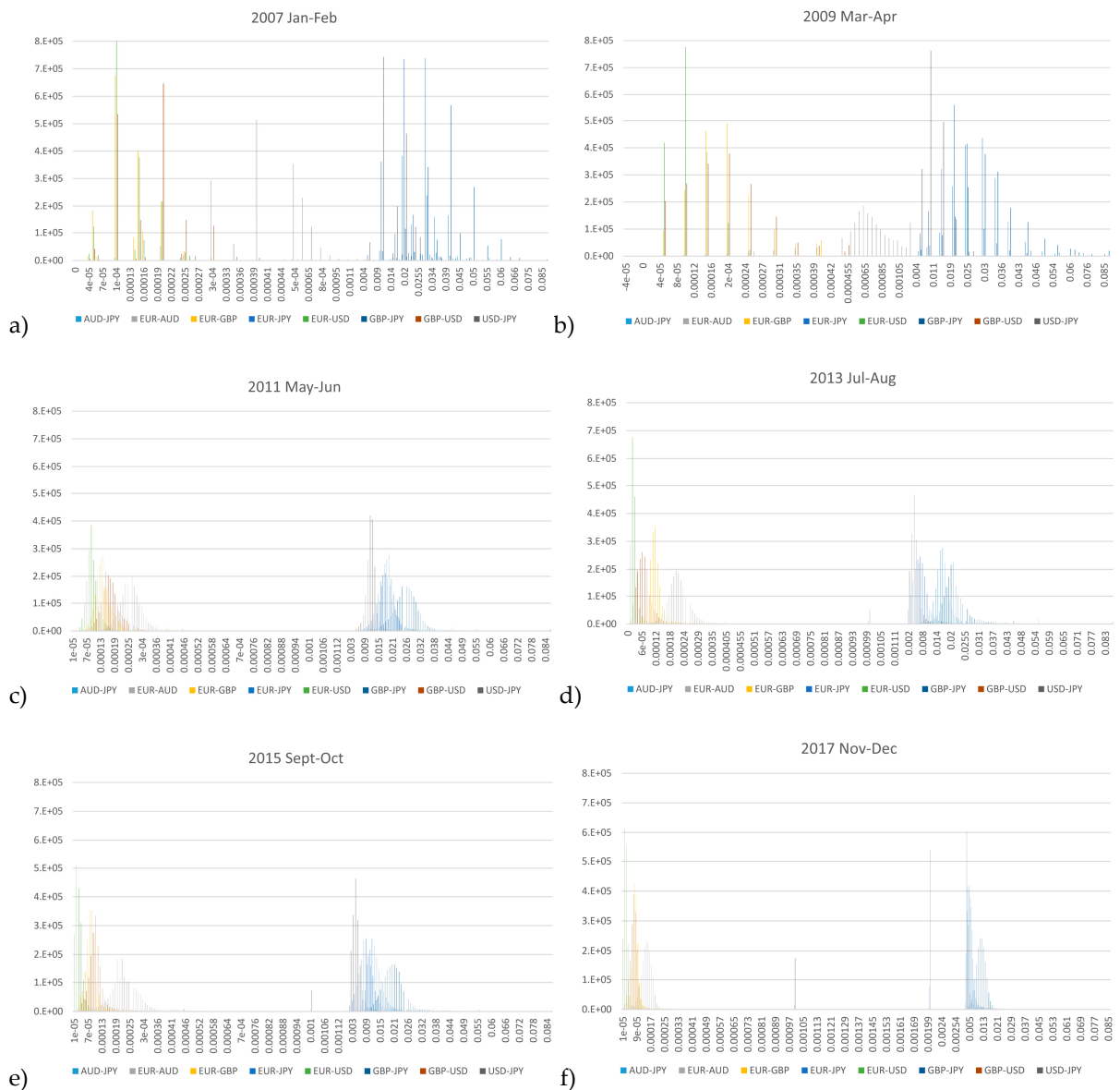


Figure SI.3 Histograms of unbidden bid-ask spreads for (a) Jan-Feb 2007; (b) Mar-Apr 2009; (c) May-Jun 2011; (d) Jul-Aug 2013; (e) Sep-Oct 2015; (f) Nov-Dec 2017.

Figure SI.4 shows the standardized bid-ask spreads of all eight currency pairs together, standardized by subtracting the respective mean and dividing it by its standard deviation. It is noticeable that the unbidden distribution became more uniform (less extreme events), which is then expressed in increased entropy (uncertainty) in our analysis.

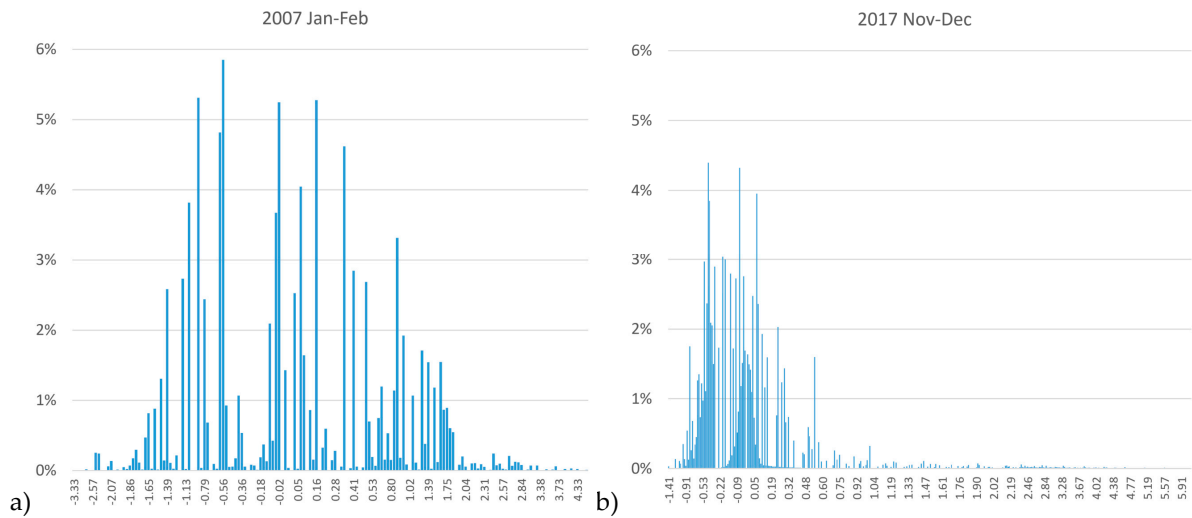
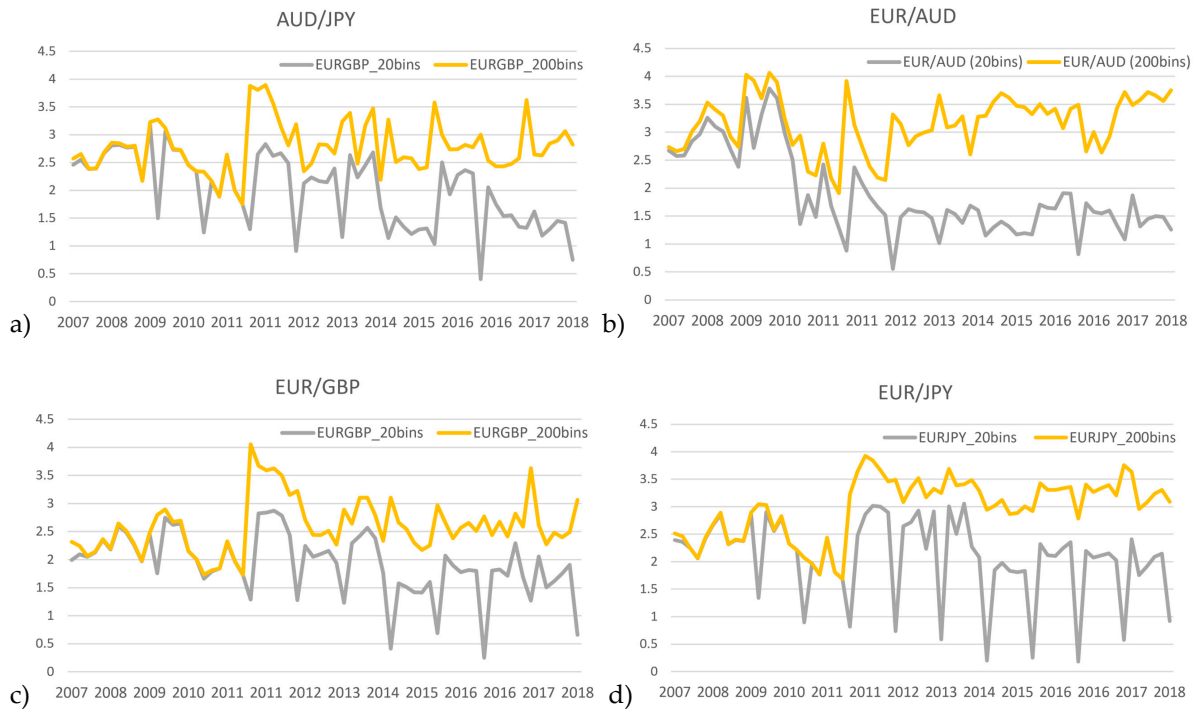


Figure SI.4 Histograms of standardized bid-ask spreads for (a) Jan-Feb 2007; (b) Nov-Dec 2017.

SI.6. Conditional Entropy Plots

In line with the chain rule of entropy, Figure SI.5 shows the diverging uncertainty between the more coarse-grained and the more fine-grained perspective. In 2007, both were still similar. The area between both levels of uncertainty is the conditional uncertainty (conditioned on the more coarse-grained resolution level).



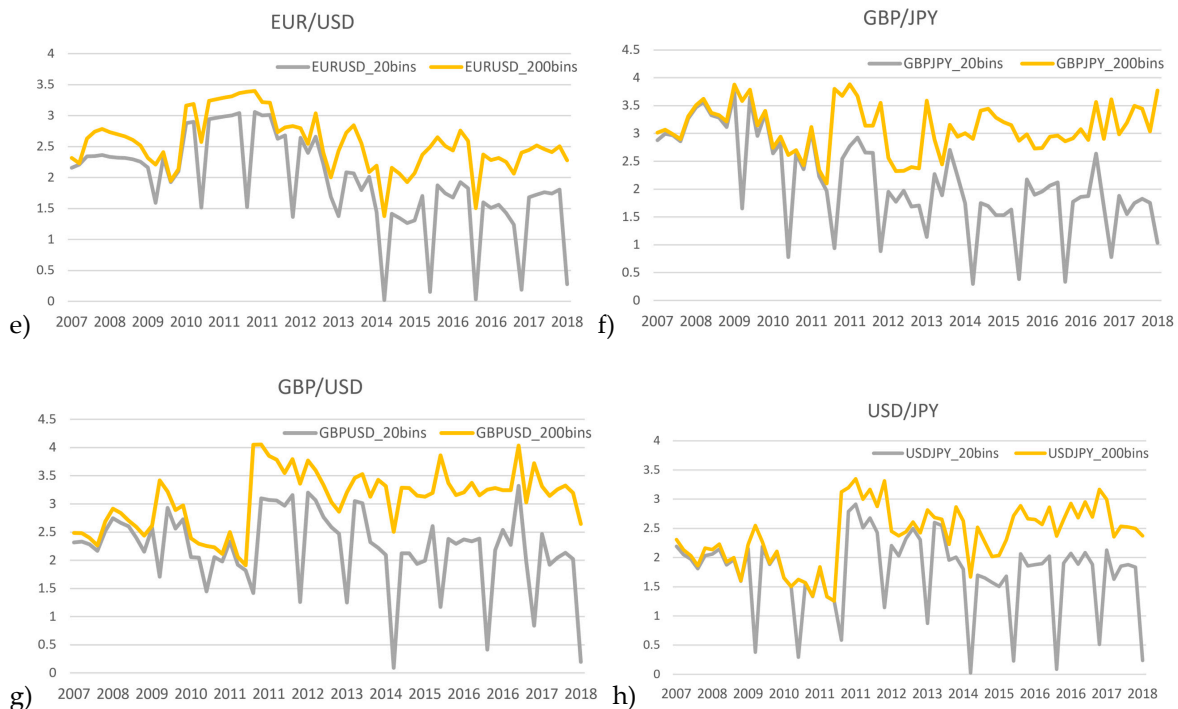


Figure SI.5 Conditional entropy plots for all currencies.

SI.7. Decreasing Bid-Ask Spreads

Gaining some intuition about the evolution of bid-ask spreads, Figure SI.6 shows visual evidence for a changing bid-ask spreads over the decade. In agreement with Hendershott and Moulto [6] the descriptive data shows an increase of the spread during the 2008 financial crisis (Lehman Brothers filed bankruptcy in September 2008) (in 2011, Hendershott et al. [7] speculated that there could be increases of temporary nature regarding realized spreads, related to asymmetries exploited by liquidity suppliers during early phases of algorithmic trading. From a decade long perspective, this might be the case, but it also seems to be the case that the increase in spread rather is linked to the financial crisis and its accompanying turmoil per se). Over the decade, however, bi-monthly bid-ask spreads decreased by half, sometimes more. In Jan-Feb 2007, the average bid-ask spread of EUR/USD was precisely one-hundredth of a cent higher than in Nov-Dec 2017 (0.00013 vs. 0.00003). Figure SI.6b shows that a similar decreasing tendency also applies to the bi-monthly standard deviation of the bid-ask spread over the same period.

It is important to point out that our relatively short time window cannot eliminate the possibility that the rather large and volatile bid-ask spread around 2008-2009 is rather the result of the global financial crisis.

Our analysis of predicting the bi-monthly means and standard deviations with our six IVs (Table SI.5), confirms that both decreases are linked to a strongly and monotonically increasing tendency that is in line with the rise of algorithmic trading (our independent variable AT, in its three different versions, namely empirical (AT_{emp}), linear (AT_{lin}) and exponential (AT_{exp})). The strongest predictor is the lagged path dependency term, closely followed by our AT variable. Interest rate and unemployment rate are also significant predictors, but less important in terms of effect size. GDP growth rate and inflation, which have seen important variances over the decade, do not play a significant role in predicting changing bid-ask tendencies. Additionally, to the negative association between algorithmic trading and bid-ask spread, we

can also add that the standard deviation decreased in association with our increased algorithmic trading tendency, which gives us first indications in terms of temporal predictability.

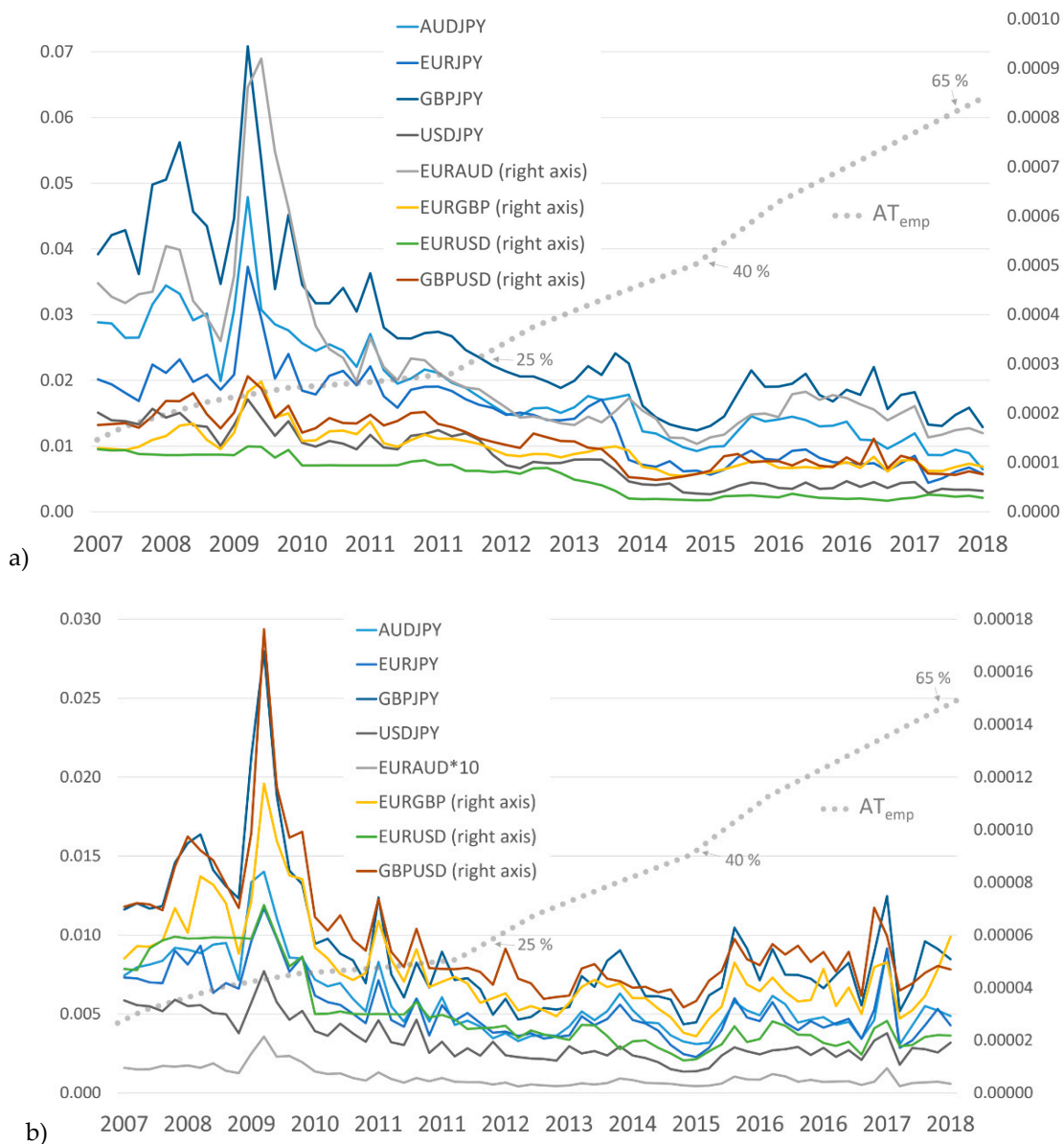


Figure SI.6 Bi-monthly bid-ask spreads Jan 2007 - Dec 2017: (a) mean; (b) standard deviation. Dotted line shows algorithmic trading empirical average with linear extrapolation from Figure 3 (scale indicated by inserted % labels).

Table SI.5: Tests for bi-monthly changes in mean (M) and standard deviation (St) of bid-ask spread; showing unstandardized beta coefficients, with standard errors in italic parentheses. *** p<0.01, ** p<0.05, * p<0.1 (N = 520).

	(1)	(2)	(3)	(4)	(5)	(6)
Dep. Var.	M	St	M	St	M	St
constant	-0.022 (0.018)	-0.013 (0.026)	-0.018 (0.017)	-0.010 (0.025)	-0.022 (0.028)	-0.013 (0.026)
AT _{emp}	-0.194*** (0.029)	-0.099*** (0.033)				
AT _{lin}			-0.276*** (0.032)	-0.189*** (0.036)		
AT _{exp}					-0.197*** (0.029)	-0.108*** (0.033)
dep _{t-1}	0.690*** (0.032)	0.686*** (0.032)	0.622*** (0.034)	0.632*** (0.034)	0.684*** (0.032)	0.680*** (0.033)
GDP _t	0.001 (0.018)	0.035 (0.026)	-0.004 (0.017)	0.039 (0.026)	0.001 (0.018)	0.036 (0.026)
infl	-0.003 (0.019)	-0.023 (0.027)	-0.003 (0.018)	-0.036 (0.027)	0.001 (0.019)	-0.023 (0.027)
intr	0.061*** (0.022)	0.070** (0.031)	0.058*** (0.021)	0.056* (0.031)	0.058*** (0.022)	0.067** (0.031)
unpl	-0.070*** (0.022)	-0.068** (0.032)	-0.056*** (0.022)	-0.055* (0.032)	-0.074*** (0.022)	-0.070** (0.032)
F(6,513)	453***	167***	481***	176***	455***	168***
adjusted R ²	0.839	0.658	0.847	0.669	0.84	0.659

References

- Allison, P. *Statistical Horizons*. June 2, 2015.
- Darmon, D. *Statistical Methods for Analyzing Time Series Data Drawn from Complex Social Systems*. PhD Thesis, University of Maryland 2015, Supervised by Michelle Girvan and William Rand, doi:10.13016/M2V93N.
- James, R.G.; Ellison, C.J.; Crutchfield, J.P. dit: a Python package for discrete information theory. *Journal of Open Source Software* 2018, 3, 738.
- Crutchfield, J.P. Between order and chaos. *Nat Phys* 2012, 8, 17–24, doi:10.1038/nphys2190.
- Shalizi, C.R.; Crutchfield, J.P. Computational Mechanics: Pattern and Prediction, Structure and Simplicity. *Journal of Statistical Physics* 2001, 104, 817–879, doi:10.1023/A:1010388907793.
- Hendershott, T.; Moulton, P.C. Automation, speed, and stock market quality: The NYSE's Hybrid. *Journal of Financial Markets* 2011, 14, 568–604, doi:10.1016/j.finmar.2011.02.003.
- Hendershott, T.; Jones, C.M.; Menkveld, A.J. Does Algorithmic Trading Improve Liquidity? *The Journal of Finance* 2011, 66, 1–33, doi:10.1111/j.1540-6261.2010.01624.x.