

Properties of Financial Texts

Internet Appendix

Overview

IA1: Trading Volume

Table 1 demonstrates the results of trading volume stationarity testing using the ADF and KPSS procedures. According to the ADF test, both level measure ($\log(Vlm)$) and its first difference ($\Delta\log(Vlm)$) are stationary. First difference of the level ($\Delta\log(Vlm)$) is stationary according to both ADF and KPSS tests. The level measure is not stationary using the KPSS test. **Campbell et al. (1993)** find the level of trading volume to be non-stationary as well. As a result, the ADF test does not have sufficient discriminatory power while the KPSS procedure does. **Figure 1** provides visual evidence; $\log(Vlm)$ is clearly non-stationary, $\Delta\log(Vlm)$ is stationary.

IA2: Business Cycle Prediction Errors

Table 1 shows that the prediction errors are lower in expansions when using stationary sentiment measures. Only days with no prior breaks are used for the evaluation. Prediction errors are statistically significantly lower for both market returns and trading volume. Forecast errors are comparable in magnitude to those obtained using the raw, non-stationary word frequencies (presented in the main paper) further establishing the robustness of residual-based evaluation.

IA3: Quantile Regressions

Table 1 and **Table 2** complement Figure 3 in the main paper. They include estimates for all quantiles, lower and upper bounds, and p-values. The regression specifications are: $R_t = \beta_M M_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$ and $\Delta\log(Vlm)_t = \beta_M M_{t-1} +$

$\beta_{Vlm}\Delta\log(Vlm)_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$. The results are computed using the stationary sentiment measures. Only days with no prior breaks are used for the evaluation. Longest possible time series are used, 1905-2005 for the Dow returns, 1926-2005 for the trading volume. Quantile regression confidence intervals follow [Koenker \(1994\)](#), kernel p-value is based on [Powell \(1991\)](#), bootstrap is pairwise as described in [Koenker \(2005\)](#).

IA4: Prediction Quality

[Table 1](#) includes correlations between prediction residuals from different stationary sentiment measures. It also lists p-values from testing the equality of means of absolute (squared) forecast errors. The correlations are generally very close to 1; all p-values are also near 1. These facts demonstrate that all stationary sentiment measures forecast Dow returns or trading volume on the same days. Some (for example, negative sentiment) predict it a bit more precisely (perhaps due to a higher quality dictionary, negation being more common for the positive words,¹ etc.). Therefore, it makes sense to draw economic conclusions from the estimate with the largest magnitude; none will tangibly increase the breadth of coverage.

¹Intuitively, this makes sense; double negatives are generally frowned upon.

IA1: Trading Volume

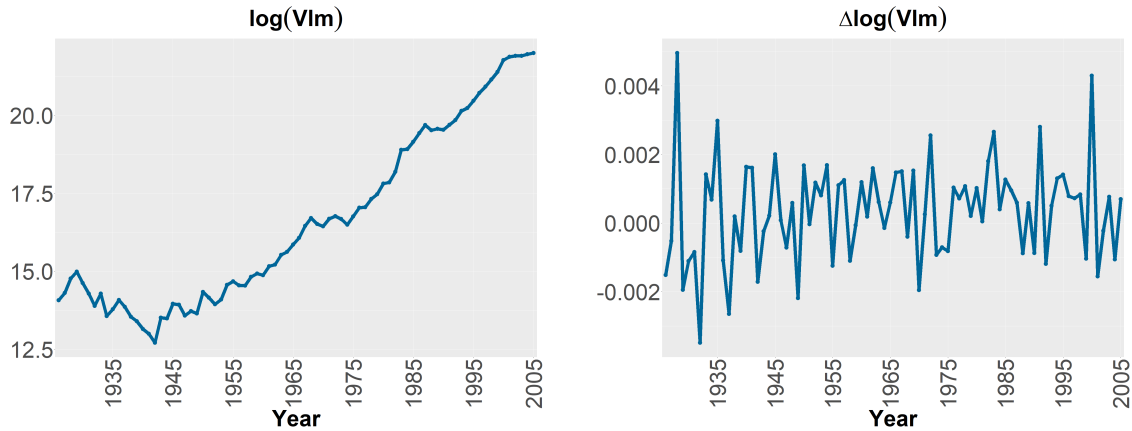
IA1, Table 1: **Trading Volume Stationarity**

Trading volume data covers 1926 - 2005. $\log(Vlm)$, $\Delta\log(Vlm)$ are logarithm of the trading volume aggregated across common shares and its first difference respectively. The ADF test includes intercept and trend, the number of lags is set to 1. The ADF null hypotheses: presence of a unit root (τ_3), unit root without trend (ϕ_3), unit root without trend and without drift (ϕ_2). The KPSS null hypothesis is trend-stationarity. The number of lags (KPSS) is set to $4(T/100)^{0.25}$.

	Statistic		Critical Values		
	$\log(Vlm)$	$\Delta\log(Vlm)$	10%	5%	1%
ADF, τ_3	-15.63	-154.40	-3.12	-3.41	-3.96
ADF, ϕ_2	81.42	7,946.11	4.03	4.68	6.09
ADF, ϕ_3	122.10	11,919.16	5.34	6.25	8.27
KPSS	26.09	0.002	0.12	0.15	0.22

IA1, Figure 1: **Trading Volume**

Values are aggregated (yearly means). Left panel depicts $\log(Vlm)$, right panel $\Delta\log(Vlm)$. $\log(Vlm)$, $\Delta\log(Vlm)$ are logarithm of the trading volume aggregated across common shares and its first difference respectively.



IA2: Forecast Errors in Expansions and Recessions

IA2, Table 1: Prediction Errors and Business Cycle

All sentiment measures are computed using a two-step procedure and are stationary. MAE and MSE are compared using two-sample Welch's t-test; RMSE (square root of MSE) is reported instead to match units. Return prediction errors are in basis points (bps), trading volume in percentage points (%pp). All p-values less than 0.001 are reported as 0. Only days with no prior breaks are included. Recessions are from NBER (USRECD).

Dow Returns, 1905-2005						
$R_t = \beta_M M_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$						
	MAE (bps)			RMSE (bps)		
	Exp	Rec	p-value	Exp	Rec	p-value
ΔPos	63.847	91.614	0	90.340	140.822	0
ΔNeg	63.814	91.533	0	90.321	140.775	0
$\Delta Pess$	63.805	91.532	0	90.302	140.777	0
Trading Volume, 1926-2005						
$\Delta \log(Vlm)_t = \beta_M M_{t-1} + \beta_{Vlm} \Delta \log(Vlm)_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$						
	MAE (%pp)			RMSE (%pp)		
	Exp	Rec	p-value	Exp	Rec	p-value
ΔPos	14.048	17.183	0	19.841	23.389	0
ΔNeg	13.998	17.149	0	19.789	23.345	0
$\Delta Pess$	14.017	17.177	0	19.812	23.366	0

IA3: Quantile Regressions

IA3, Table 1: Quantile Regressions, Daily Dow, Stationary Sentiment

Panel A: ΔPos					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	-2.11	-4.36	0.82	0.15	0.12
$\tau = 0.2$	0.88	-0.52	2.33	0.30	0.35
$\tau = 0.3$	2.58	1.18	3.82	0	0
$\tau = 0.4$	2.42	1.37	3.72	0	0
$\tau = 0.5$	2.59	1.51	3.57	0	0
$\tau = 0.6$	2.78	1.83	4.01	0	0
$\tau = 0.7$	3.91	2.65	5.00	0	0
$\tau = 0.8$	4.95	3.25	6.17	0	0
$\tau = 0.9$	4.10	2.47	6.08	0	0
Panel B: ΔNeg					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	-1.45	-3.97	0.40	0.27	0.27
$\tau = 0.2$	-2.11	-3.44	-0.47	0.02	0.02
$\tau = 0.3$	-3.18	-4.23	-1.83	0	0
$\tau = 0.4$	-3.82	-4.81	-2.50	0	0
$\tau = 0.5$	-3.72	-5.04	-2.87	0	0
$\tau = 0.6$	-3.51	-4.63	-2.49	0	0
$\tau = 0.7$	-3.15	-4.31	-2.10	0	0
$\tau = 0.8$	-3.47	-5.02	-2.06	0	0
$\tau = 0.9$	-3.89	-5.82	-1.98	0	0
Panel C: $\Delta Pess$					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	-0.05	-2.09	2.69	0.98	0.97
$\tau = 0.2$	-1.96	-3.64	-0.33	0.03	0.04
$\tau = 0.3$	-3.59	-4.69	-2.43	0	0
$\tau = 0.4$	-3.88	-4.88	-2.89	0	0
$\tau = 0.5$	-3.98	-4.96	-2.95	0	0
$\tau = 0.6$	-4.27	-5.33	-3.24	0	0
$\tau = 0.7$	-5.27	-6.47	-3.85	0	0
$\tau = 0.8$	-5.48	-6.71	-4.20	0	0
$\tau = 0.9$	-5.47	-7.34	-3.34	0	0

IA3, Table 2: Quantile Regressions, Trading Volume, Stationary Sentiment

Panel A: ΔPos					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	0.01	-0.72	0.58	0.98	0.98
$\tau = 0.2$	0.47	0.14	0.80	0.01	0.03
$\tau = 0.3$	0.34	0.08	0.67	0.06	0.06
$\tau = 0.4$	0.38	0.11	0.57	0.01	0.02
$\tau = 0.5$	0.34	0.04	0.62	0.05	0.04
$\tau = 0.6$	0.42	0.11	0.70	0.02	0.01
$\tau = 0.7$	0.45	0.17	0.68	0.01	0.01
$\tau = 0.8$	0.55	0.21	0.86	0	0.01
$\tau = 0.9$	0.92	0.25	1.51	0.01	0
Panel B: ΔNeg					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	-1.71	-2.26	-1.07	0	0
$\tau = 0.2$	-1.33	-1.70	-1.10	0	0
$\tau = 0.3$	-1.42	-1.73	-1.11	0	0
$\tau = 0.4$	-1.37	-1.61	-1.17	0	0
$\tau = 0.5$	-1.37	-1.61	-1.19	0	0
$\tau = 0.6$	-1.31	-1.49	-1.11	0	0
$\tau = 0.7$	-1.24	-1.49	-0.95	0	0
$\tau = 0.8$	-1.41	-1.77	-1.03	0	0
$\tau = 0.9$	-1.83	-2.30	-1.15	0	0
Panel C: $\Delta Pess$					
	β_M	LB	UB	p_{boot}	p_{ker}
$\tau = 0.1$	-0.99	-1.56	-0.49	0	0
$\tau = 0.2$	-1.09	-1.46	-0.77	0	0
$\tau = 0.3$	-1.15	-1.47	-0.81	0	0
$\tau = 0.4$	-1.02	-1.24	-0.77	0	0
$\tau = 0.5$	-1.13	-1.42	-0.83	0	0
$\tau = 0.6$	-1.15	-1.37	-0.89	0	0
$\tau = 0.7$	-1.11	-1.42	-0.76	0	0
$\tau = 0.8$	-1.22	-1.49	-0.84	0	0
$\tau = 0.9$	-1.64	-2.24	-1.04	0	0

IA4: Prediction Quality

IA4, Table 1: Prediction Quality by Measure

Residuals are from: $R_t = \beta_M M_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$ and $\Delta \log(Vlm)_t = \beta_M M_{t-1} + \beta_{Vlm} \Delta \log(Vlm)_{t-1} + \beta_R R_{t-1} + \beta_{RSq} R_{t-1}^2 + C + \epsilon$. Dow spans 1905-2005, trading volume data covers 1926 - 2005, only days with no prior breaks are included. Two-sided Welch's t-test is used to compare MAE (MSE).

Dow, Welch's t-test				
	ΔPos	ΔNeg	$\Delta Pess$	
ΔPos	1	0.955	0.947	} MAE
ΔNeg	0.989	1	0.993	
$\Delta Pess$	0.984	0.995	1	
} MSE				
Dow, Correlations				
	$\epsilon_{\Delta Pos}$	$\epsilon_{\Delta Neg}$	$\epsilon_{\Delta Pess}$	
$\epsilon_{\Delta Pos}$	1	0.999	1.000	} Pearson
$\epsilon_{\Delta Neg}$	0.969	1	1.000	
$\epsilon_{\Delta Pess}$	0.982	0.978	1	
} Kendall				
Volume, Welch's t-test				
	ΔPos	ΔNeg	$\Delta Pess$	
ΔPos	1	0.776	0.872	} MAE
ΔNeg	0.868	1	0.901	
$\Delta Pess$	0.927	0.940	1	
} MSE				
Volume, Correlations				
	$\epsilon_{\Delta Pos}$	$\epsilon_{\Delta Neg}$	$\epsilon_{\Delta Pess}$	
$\epsilon_{\Delta Pos}$	1	0.997	0.999	} Pearson
$\epsilon_{\Delta Neg}$	0.945	1	0.999	
$\epsilon_{\Delta Pess}$	0.968	0.961	1	
} Kendall				

References

- Campbell, J. Y., Grossman, S. J., & Wang, J. (1993, 11). Trading volume and serial correlation in stock returns. *The Quarterly Journal of Economics*, 108, 905–939. doi: 10.2307/2118454
- Koenker, R. (1994). Confidence intervals for regression quantiles. *Asymptotic Statistics*, 349-359. doi: 10.1007/978-3-642-57984-4_29
- Koenker, R. (2005). *Quantile regression*. Cambridge University Press.
- Powell, J. (1991). *Nonparametric and semiparametric methods in econometrics and statistics : proceedings of the fifth international symposium in economic theory and econometrics* (W. A. Barnett, J. Powell, & G. E. Tauchen, Eds.). Cambridge University Press.