

From previous-tick to pre-averaging: Spectra of equidistant transformations for unevenly spaced high-frequency data

Vitali Alexeev ¹ Jun Chen ² Katja Ignatieva ²

¹UTS Business School, University of Technology Sydney, Australia

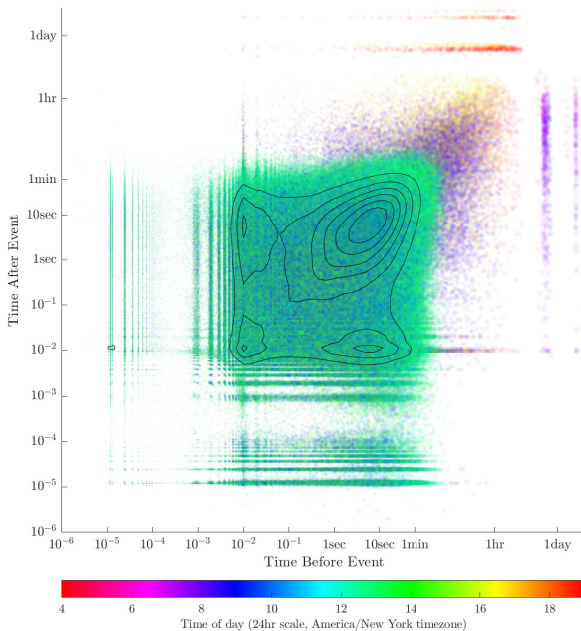
²School of Risk and Actuarial Studies, Business School, UNSW Australia

The 12th Bachelier World Congress
FGV EMap, Rio de Janeiro
July 8-12, 2024

Motivation

High-frequency Data

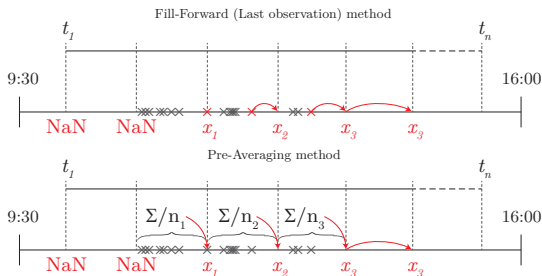
- A vast amount of **data**.
 - More data → More precise estimation;
 - Higher frequency → More persistent noise;
- But... (typically) Irregularly spaced
 - Tick-by-tick trade data;
 - Sentiment from textual data (news and social media);
 - Implications for univariate and, especially, multivariate analyses.



Motivation

Common (calendar) schemes for equidistant data:

- *Previous tick* (also known as *last observation*)
- *Pre-average* (also recently *pre-median*, but mostly for cryptos)



Objective:

- Devise equidistant mesh to control trade-off between the number of discarded observations and the amount of (microstructure) noise.

Microstructure Noise

Let $x(t)$ be the **true log price** of an asset at time t , $x(t)$ satisfies

$$dx(t) = \mu(t, x(t))dt + \sigma(t, x(t))dW(t). \quad (1)$$

$$IV = \int_0^1 \sigma(t, x(t))^2 dt. \quad (2)$$

Let $y(t)$ be the **observed prices**, $\epsilon(t)$ be microstructure noise, it holds

$$y(t) = x(t) + \epsilon(t). \quad (3)$$

Assumption

- ① $\{\epsilon(t)\}$ is stochastically independent of the price process $\{x(t)\}$;
- ② $\{\epsilon(t)\}$ is **white noise** with $E(\epsilon(t)) = 0$ and $\text{Var}(\epsilon(t)) = \eta^2 < \infty$.

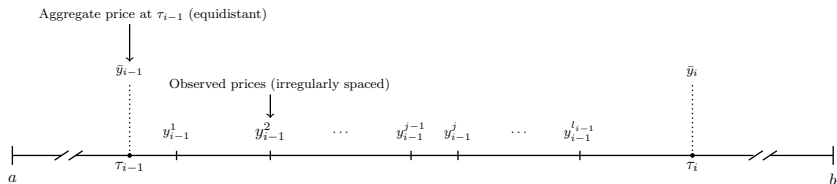
Issue: $RV \xrightarrow{P} IV, \quad n \rightarrow \infty$

Let t_i , $i = 1, \dots, n$, in $a = t_1 < t_2 < \dots < t_n = b$ be observation times on partition $[a, b]$. An equidistant grid τ_i , $i = 1, \dots, m$, spanning the time interval $[a, b]$:

$$a = \tau_1 < \tau_2 < \dots < \tau_m = b, \quad (4)$$

where, generally, $n > m$. The grid partitions $[a, b]$ into $m - 1$ sub-intervals.

Within a sub-interval $(\tau_{i-1}, \tau_i]$, $y_{\tau_{i-1}}^j$ is the j^{th} observation with timestamp t_{i-1}^j , and l_{i-1} is the number of obs located within the sub-interval.



Pre-weighted Sampling Scheme

Let $y_{i-1}^1, \dots, y_{i-1}^{l_{i-1}}$ be the observable prices in a sub-interval $(\tau_{i-1}, \tau_i]$, and $\omega = (\omega_{i-1}^1, \dots, \omega_{i-1}^{l_{i-1}})$ be the vector of weights for these observations such that $\sum_{j=1}^{l_{i-1}} \omega_{i-1}^j = 1$. Then the aggregate price at τ_i :

$$\bar{y}_i = \begin{cases} \sum_{j=1}^{l_{i-1}} \omega_{i-1}^j y_{i-1}^j, & \text{if } l_{i-1} > 0, \\ \bar{y}_{i-1}, & \text{if } l_{i-1} = 0. \end{cases} \quad (5)$$

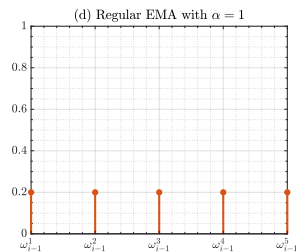
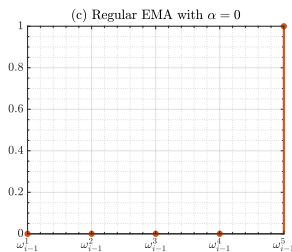
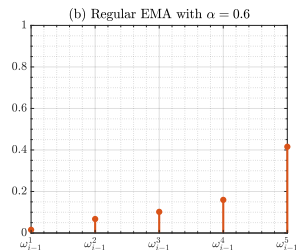
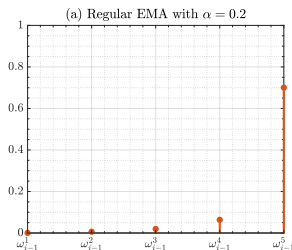
Previous tick: $y_{i-1}^{l_{i-1}}$ and Pre-average: $\frac{1}{l_{i-1}} \sum_{j=1}^{l_{i-1}} y_{i-1}^j$

Recursively for $j = 1, \dots, l_{i-1}$ with initial value for \bar{y}_i

Regular pre-EMA: $(1 - \alpha_k) y_{i-1}^j + \alpha_k \bar{y}_i$

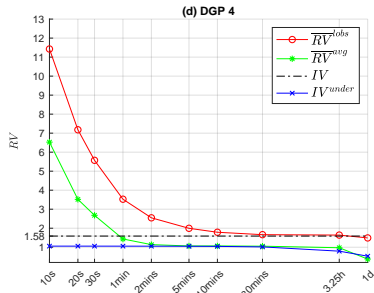
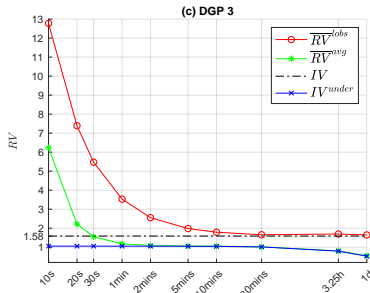
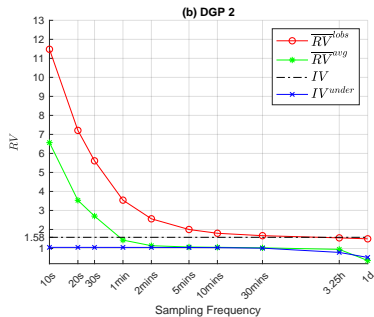
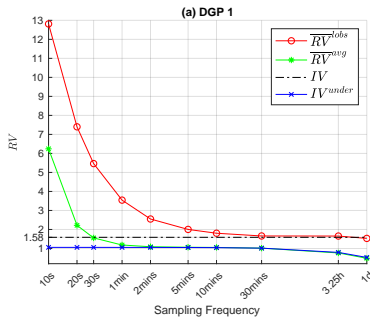
Irregular pre-EMA: $(1 - \exp(-\alpha_k \delta_j)) y_{i-1}^j + \alpha_k \bar{y}_i$

Observation weights in a sub-interval for different α



Sampling Schemes: Pros and Cons

Scheme	Pros	Cons
Previous tick	·easy to implement	·discards most observations
Pre-average	·full use of all data	·equal weights to all obs ·over-emphasizes less recent obs ·underestimates the IV Details
Pre-EMA	·full use of all data ·larger weights to latest obs ·a parametric hybrid	·hyperparameter to determine Details



Simulation Study

True price process	Sample collection method	
	Equidistant collection	Irregular collection
Brownian Motion	DGP 1	DGP 2
Heston Model	DGP 3	DGP 4

- Using Euler discretisation, simulate the continuous GBM and Heston.
- The model parameters are as in [Jacod et al. \(2009\)](#)
- Generate **23,401 obs** each day.
- Employ two approaches to sample observations from the simulated set:
 - **equidistant collection**
 - **irregular sample collection**
- Add **microstructure noise** to the true prices.

Optimal smoothing parameter (varied microstructure noise)

Optimal smoothing parameter (in red) for

Regular Pre-EMA Scheme: $(1 - \alpha_k)y_{i-1}^j$

Irregular Pre-EMA Scheme: $(1 - \exp(-\alpha_k \delta_j))y_{i-1}^j$

Interpretation: 1 \rightarrow *Pre-averaging* and 0 \rightarrow *Previous Tick*

Scheme	Frequency of equidistant mesh (Δ_k)							
	10s	20s	30s	1 min	2 mins	5 mins	10 mins	30 mins
Panel A: Large microstructure noise $\epsilon(t) \sim \mathcal{N}(0, 2 * 0.0005^2)$								
<i>Regular</i>	1	1	1	1	0.8157	0.6750	0.5795	0
<i>Irregular</i>	1	1	1	1	0.8395	0.7013	0.4879	0
Panel B: Medium microstructure noise $\epsilon(t) \sim \mathcal{N}(0, 1 * 0.0005^2)$								
<i>Regular</i>	1	1	1	0.7884	0.5695	0.4655	0.2089	0
<i>Irregular</i>	1	1	1	0.8975	0.6537	0.6004	0.2716	0
Panel C: Small microstructure noise $\epsilon(t) \sim \mathcal{N}(0, 0.5 * 0.0005^2)$								
<i>Regular</i>	1	1	1	0.5931	0.4248	0.3393	0	0
<i>Irregular</i>	1	1	1	0.6948	0.4927	0.3771	0.0461	0

Optimal smoothing parameter (varied liquidity)

Optimal smoothing parameter (in red) for

Regular Pre-EMA Scheme: $(1 - \alpha_k)y_{i-1}^j$

Irregular Pre-EMA Scheme: $(1 - \exp(-\alpha_k \delta_j))y_{i-1}^j$

Interpretation: 1 \rightarrow *Pre-averaging* and 0 \rightarrow *Previous Tick*

Scheme	Frequency of equidistant mesh (Δ_k)							
	10s	20s	30s	1 min	2 mins	5 mins	10 mins	30 mins
Panel A: High liquidity $\bar{\lambda}_h = (60, 60, 20, 14, 6, 4, 4, 4, 20, 30, 32, 36, 50)$								
<i>Regular</i>	1	1	1	0.7380	0.6473	0.6385	0.5477	0.5023
<i>Irregular</i>	1	1	1	0.8302	0.7453	0.7272	0.5017	0
Panel B: Intermediate $\bar{\lambda} = (40, 30, 10, 7, 3, 2, 2, 2, 10, 15, 16, 18, 25)$								
<i>Regular</i>	1	1	1	0.7884	0.5695	0.4655	0.2089	0
<i>Irregular</i>	1	1	1	0.8975	0.6537	0.6004	0.2716	0
Panel C: Low $\bar{\lambda}_l = (20, 15, 5, 3, 1, 1, 1, 1, 5, 7, 8, 9, 12)$								
<i>Regular</i>	1	1	1	1	0.6114	0.3729	0.2218	0
<i>Irregular</i>	1	1	1	1	0.7716	0.4334	0.2282	0

Overall in-sample estimation accuracy

Using the DMW test, the number of estimators that a particular estimator can beat in terms of its estimation accuracy. MPF = marginal performance w.r.t. sampling frequencies; MPS = marginal performance w.r.t. sampling schemes.

Δ_k	RV_k^{lobs}	RV_k^{exp}	RV_k^{iexp}	RV_k^{ravg}	MPF
10s	0	4	3	1	0.0175
20s	2	10	9	6	0.0591
30s	5	12	11	7	0.0766
1min	8	28	28	18	0.1794
2min	13	28	28	25	0.2057
5min	19	24	24	24	0.1991
10min	19	19	19	18	0.1641
30min	12	11	11	11	0.0985
MPS	0.1707	0.2976	0.2910	0.2407	1.0000

Overall out-of-sample estimation accuracy

Using the DMW test, the number of estimators that a particular estimator can beat in terms of its estimation accuracy. MPF = marginal performance w.r.t. sampling frequencies; MPS = marginal performance w.r.t. sampling schemes.

Δ_k	RV_k^{lobs}	RV_k^{exp}	RV_k^{iexp}	RV_k^{ravg}	MPF
10s	0	4	3	1	0.0175
20s	2	10	9	6	0.0590
30s	5	12	11	7	0.0764
1min	8	28	28	18	0.1790
2min	13	28	28	24	0.2031
5min	19	24	24	24	0.1987
10min	19	19	19	18	0.1638
30min	13	11	12	11	0.1026
MPS	0.1725	0.2969	0.29265	0.2380	1.0000

Data Summary and Optimal smoothing parameter

- **Tick-by-tick trades** for NASDAQ:GOOG.
- January 2005 to December 2014.
- 45,951,783 observations across 2,513 days (after data cleaning).
- i.e., 18,285 obs per day.
- Two **sub-samples**, allowing performance comparison by means of in- and out-of-sample tests.

Optimal smoothing parameter (in **red**) for

Regular Pre-EMA Scheme: $(1 - \alpha_k)y_{i-1}^j$

Irregular Pre-EMA Scheme: $(1 - \exp(-\alpha_k \delta_j)) y_{i-1}^j$

Interpretation: 1 \rightarrow *Pre-averaging* and 0 \rightarrow *Previous Tick*

Scheme	Frequency of equidistant mesh (Δ_k)							
	10s	20s	30s	1min	2min	5min	10min	30min
<i>Regular</i>	0.7411	0.4847	0.3740	0.1552	0	0	0	0
<i>Irregular</i>	0.7015	0.4239	0.3124	0.0776	0	0	0	0

Overall in-sample estimation accuracy

Using the DMW test, the number of estimators that a particular estimator can beat in terms of its estimation accuracy. MPF = marginal performance w.r.t. sampling frequencies; MPS = marginal performance w.r.t. sampling schemes.

Δ_k	RV_k^{lobs}	RV_k^{exp}	RV_k^{iexp}	RV_k^{ravg}	MPF
10s	4	16	16	9	0.1119
20s	11	28	25	21	0.2114
30s	16	24	24	20	0.2090
1min	20	20	20	15	0.1866
2min	16	15	15	13	0.1468
5min	9	8	8	8	0.0821
10min	5	4	4	4	0.0423
30min	1	0	0	3	0.0100
MPS	0.2040	0.2861	0.2786	0.2313	1.0000

Overall out-of-sample estimation accuracy

Using the DMW test, the number of estimators that a particular estimator can beat in terms of its estimation accuracy. MPF = marginal performance w.r.t. sampling frequencies; MPS = marginal performance w.r.t. sampling schemes.

Δ_k	RV_k^{lobs}	RV_k^{exp}	RV_k^{iexp}	RV_k^{ravg}	MPF
10s	8	17	17	14	0.1425
20s	13	26	22	28	0.2265
30s	17	20	20	21	0.1985
1min	17	17	17	16	0.1705
2min	13	12	12	12	0.1247
5min	9	8	8	8	0.0840
10min	5	4	4	4	0.0434
30min	1	0	0	3	0.0102
MPS	0.2112	0.2646	0.2545	0.2697	1.0000

Summary

- Connect the *previous tick sampling scheme* and the *pre-averaging sampling scheme* by the *pre-EMA sampling scheme*;
- Apply *pre-EMA sampling scheme* to irregular high frequency data;
- Develop the method to find the optimal α in the *regular pre-EMA sampling scheme* and the method can also be applied to the *irregular pre-EMA sampling scheme*;
- Conduct comparison between RVs at various sampling frequencies with different sampling schemes.

Thanks!

Pre-averaging Sampling Scheme

[Back to Main](#)

Assumption 3

There are L observations in each sub-interval $(\tau_{i-1}, \tau_i]$ for $i = 1, \dots, m$.

Lemma

Under stated assumptions, it holds

$$E(RV^{avg}) = \frac{1 + 2L^2}{3L^2} IV + 2\frac{m}{L}\eta^2 \quad (6)$$

where, RV^{avg} is the RV based on the pre-averaging sampling scheme. Thus

$$\lim_{L \rightarrow \infty} E(RV^{avg}) = \frac{2}{3} IV \quad (7)$$

Optimal α for Regular Pre-EMA Sampling Scheme

[Back to Main](#)

Let θ_n be the true value of IV, $\tilde{\theta}_n$ be the noisy but unbiased estimator of θ_n , we assume

- ① $\tilde{\theta}_n = \theta_n + \nu_n$, with $\mathbb{E}[\nu_n | \mathcal{F}_{n-1}, \theta_n] = 0$;
- ② $\theta_n = \theta_{n-1} + \vartheta_n$, with $\mathbb{E}[\vartheta_n | \mathcal{F}_{n-1}] = 0$;
- ③ $\Upsilon_n = \sum_{j=1}^J \lambda_j \tilde{\theta}_{n+j}$, where $1 \leq J < \infty$, $\lambda_j \geq 0 \quad \forall j$ and $\sum_{j=1}^J \lambda_j = 1$.

Optimal α minimizes the objective function

$$\alpha^* \equiv \arg \min_{\alpha} \mathbb{E}[Q(RV_{\alpha}^{\exp}(n), \theta_n)] \quad (8)$$

or, equivalently:

$$\tilde{\alpha}^* \equiv \arg \min_{\tilde{\alpha}} \mathbb{E}[Q(RV_{\alpha}^{\exp}(n), \tilde{\theta}_n)]. \quad (9)$$

Difference in Size of Microstructure Noise

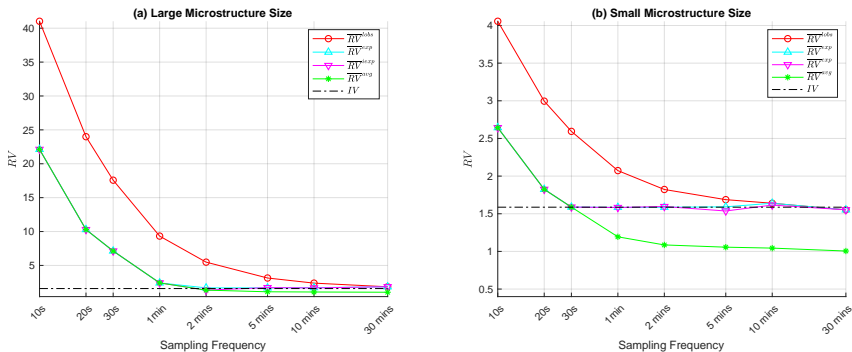
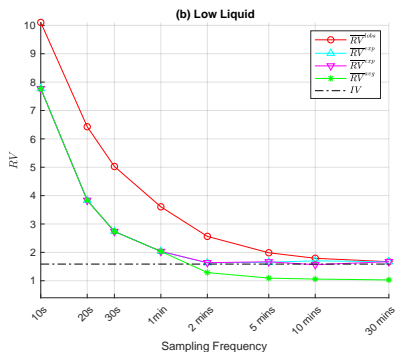
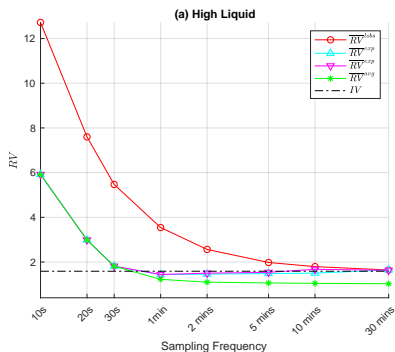


Figure: $\epsilon(t) \sim \mathcal{N}(0, 2 * 0.0005^2)$ v.s. $\epsilon(t) \sim \mathcal{N}(0, 0.5 * 0.0005^2)$

Difference in Liquidity



Diebold-Mariano-West tests results (5%)

“✓” = RV from the row has greater estimation accuracy than RV from the column.
 “X” = reject the null hypothesis in favour of the RV from the corresponding column.
 “-” = both estimators are comparable in terms of estimation accuracy.

(A) RV_k^{lobs} estimator.

RV_k^{lobs} (...)	10s	20s	30s	1min	2min	5min	10min	30min	1d
10s	-	X	X	X	X	X	X	X	✓
20s	✓	-	X	X	X	X	X	X	✓
30s	✓	✓	-	X	X	X	X	X	✓
1min	✓	✓	-	-	X	X	X	X	✓
2min	✓	✓	✓	✓	-	X	X	-	✓
5min	✓	✓	✓	✓	✓	-	-	✓	✓
10min	✓	✓	✓	✓	✓	-	-	✓	✓
30min	✓	✓	✓	✓	-	X	X	-	✓
1d	X	X	X	X	X	X	X	X	-

(B) RV_k^{exp} estimator.

RV_k^{exp} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	X	X	X	X	X	X	X
20s	✓	-	X	X	X	X	X	X
30s	✓	✓	-	X	X	X	X	-
1min	✓	✓	✓	-	-	✓	✓	✓
2min	✓	✓	✓	-	-	✓	✓	✓
5min	✓	✓	✓	X	X	-	✓	✓
10min	✓	✓	✓	X	X	X	-	✓
30min	✓	✓	-	X	X	X	X	-

(C) RV_k^{iexp} estimator.

RV_k^{iexp} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	X	X	X	X	X	X	X
20s	✓	-	X	X	X	X	X	X
30s	✓	✓	-	X	X	X	X	-
1min	✓	✓	✓	-	-	✓	✓	✓
2min	✓	✓	✓	-	-	✓	✓	✓
5min	✓	✓	✓	X	X	-	✓	✓
10min	✓	✓	✓	X	X	X	-	✓
30min	✓	✓	-	X	X	X	X	-

(D) RV_k^{ravg} estimator.

RV_k^{ravg} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	X	X	X	X	X	X	X
20s	✓	-	X	X	X	X	X	X
30s	✓	✓	-	X	X	X	X	X
1min	✓	✓	✓	-	X	X	-	✓
2min	✓	✓	✓	✓	-	✓	✓	✓
5min	✓	✓	✓	✓	X	-	✓	✓
10min	✓	✓	✓	-	X	X	-	✓
30min	✓	✓	✓	X	X	X	X	-

Diebold-Mariano-West tests results (5%) - GOOG.OQ

“✓” = RV from the row has greater estimation accuracy than RV from the column.
 “✗” = reject the null hypothesis in favour of the RV from the corresponding column.
 “-” = both estimators are comparable in terms of estimation accuracy.

(E) RV_k^{lobs} estimator.

RV_k^{lobs} (...)	10s	20s	30s	1min	2min	5min	10min	30min	1d
10s	-	✗	✗	✗	✗	-	-	✓	-
20s	✓	-	✗	✗	-	-	✓	✓	-
30s	✓	✓	-	-	-	✓	✓	✓	-
1min	✓	✓	-	-	✓	✓	✓	✓	✓
2min	✓	-	-	✗	-	✓	✓	✓	✓
5min	-	-	✗	✗	✗	-	✓	✓	✓
10min	-	✗	✗	✗	✗	✗	-	✓	✗
30min	✗	✗	✗	✗	✗	✗	✗	-	✗
1d	-	-	✗	✗	✗	✗	✓	✓	-

(F) RV_k^{exp} estimator.

RV_k^{exp} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	✗	-	-	-	✓	✓	✓
20s	✓	-	-	✓	✓	✓	✓	✓
30s	-	-	-	✓	✓	✓	✓	✓
1min	-	✗	✗	-	✓	✓	✓	✓
2min	-	✗	✗	✗	-	✓	✓	✓
5min	✗	✗	✗	✗	✗	-	✓	✓
10min	✗	✗	✗	✗	✗	✗	-	✓
30min	✗	✗	✗	✗	✗	✗	✗	-

(G) RV_k^{iexp} estimator.

RV_k^{iexp} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	✗	-	-	-	✓	✓	✓
20s	✓	-	-	-	✓	✓	✓	✓
30s	-	-	-	✓	✓	✓	✓	✓
1min	-	-	✗	-	✓	✓	✓	✓
2min	-	✗	✗	✗	-	✓	✓	✓
5min	✗	✗	✗	✗	✗	-	✓	✓
10min	✗	✗	✗	✗	✗	✗	-	✓
30min	✗	✗	✗	✗	✗	✗	✗	-

(H) RV_k^{ravg} estimator.

RV_k^{ravg} (...)	10s	20s	30s	1min	2min	5min	10min	30min
10s	-	✗	✗	✗	-	-	✓	✓
20s	✓	-	-	✓	✓	✓	✓	✓
30s	✓	-	-	✓	✓	✓	✓	✓
1min	✓	✗	✗	-	✓	✓	✓	✓
2min	-	✗	✗	✗	-	✓	✓	✓
5min	-	✗	✗	✗	✗	-	✓	✓
10min	✗	✗	✗	✗	✗	✗	-	✓
30min	✗	✗	✗	✗	✗	✗	✗	-