

Online learning methods for financial time series forecasting

PhD Thesis Defense

Paola Arce

Thesis advisor: Luis Salinas

December 15, 2017



UNIVERSIDAD TÉCNICA
FEDERICO SANTA MARÍA



Departamento de Informática
Universidad Técnica Federico Santa María

- Motivation: High Frequency Trading and Cointegration
- Challenges and Hypothesis
- Contributions
- Limitations and Applicability
- Conclusions and Future Research

Motivation: HFT

High Frequency Trading

High Frequency Trading (HFT) is a technology that implements different trading strategies to buy or sell an stock using sophisticated infrastructure and high speed algorithms.

HFT challenges are related to **time series forecasting** and also **computational** (hardware and software). Today, the volume of data has rapidly increased:

Period	Data precision
1993-2003	seconds
2003-2015	milliseconds
27th July 2015-onwards	microseconds
24th October 2016-onwards	nanoseconds

Table 1: High Frequency data precision

HFT: High Frequency Trading

High Frequency Trading has motivated computer-driven strategies capable of processing large amounts of data in short periods of time.

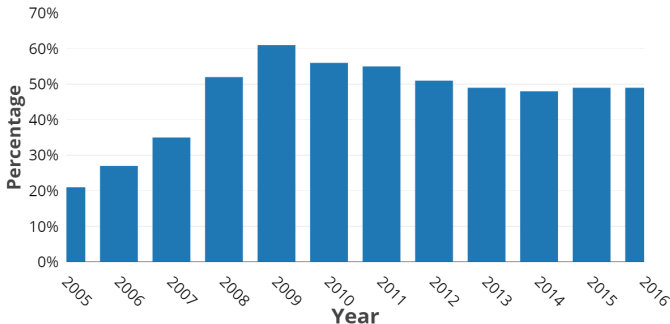


Figure 1: Percentage of HFT trades of total US Equity Trading. Source: TABB group

HFT Revenue

However, the revenues have fallen dramatically mainly because of more regulation, increasing the cost of HFT infrastructure and more competitive algorithms.

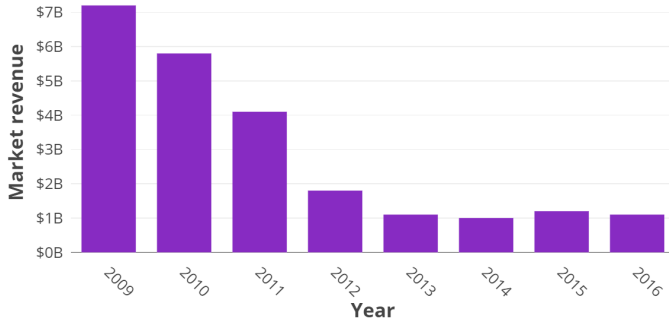


Figure 2: Revenue in the US. Source: TABB group

Motivation: HFT

Arbitrage

The practice of taking advantage of market inefficiency or discrepancies in the market price. Arbitrage contradicts the Efficient Market Hypothesis.

Advantages

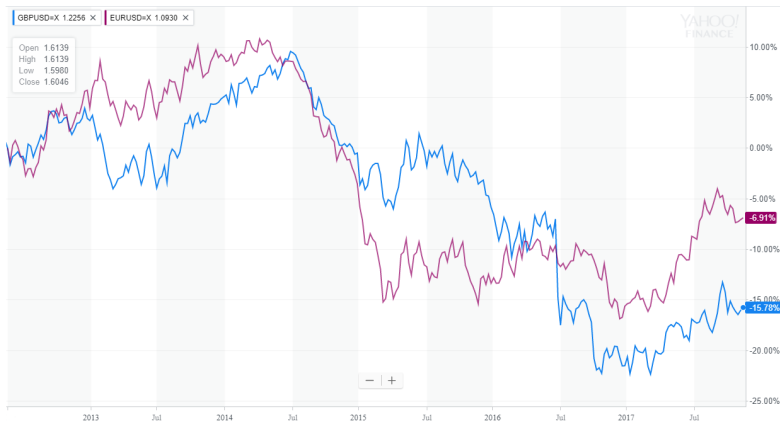
- Increases liquidity
- Allows more efficient price formation
- Reduce the difference between buyers (bid) and sellers (ask) prices.

Disadvantages

- Removal of human decision making could lead to failures in untested or incorrect strategies
- fairness

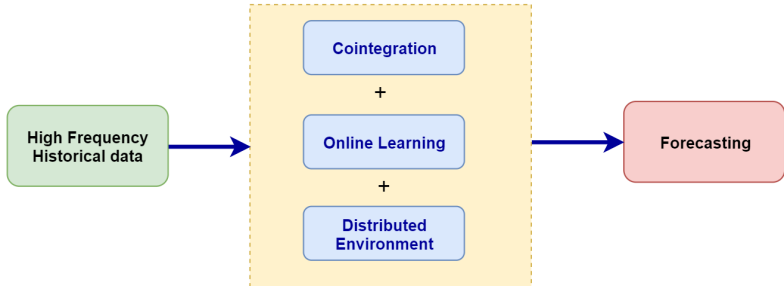
Motivation: Cointegration

Cointegration reflects the idea of that some set of non-stationary variables cannot wander too far from each other.



Contribution

This thesis contribution has several components:



High Frequency data

High frequency time series are commonly:

- Non-stationary
- Irregularly spaced over time.
- Heteroscedastic
- Potentially cointegrated



This implies that classical statistical models as such can't be used. Therefore, financial time series forecasting is still a challenge.

Strong definition

A **strictly stationary time series** $y_t : r.v$ with values in \mathbb{R} defined for $t \in \mathbb{Z}$, is one for which the probabilistic behaviour of every collection of values $\{y_{t_1}, y_{t_2}, \dots, y_{t_L}\}$ is identical to that of the time shifted set, more precisely:

$$P\{y_{t_1} \leq c_1, \dots, y_{t_L} \leq c_L\} = P\{y_{t_1+h} \leq c_1, \dots, y_{t_L+h} \leq c_L\}$$

$\forall L \in \mathbb{N}, \forall h \in \mathbb{Z}$, where c_1, \dots, c_L are constants.

Weak definition

A weakly stationary time series is a process where the mean, variance and auto covariance do not change over time:

$$E(y_t) = \mu \quad \forall t \in \mathbb{N}$$

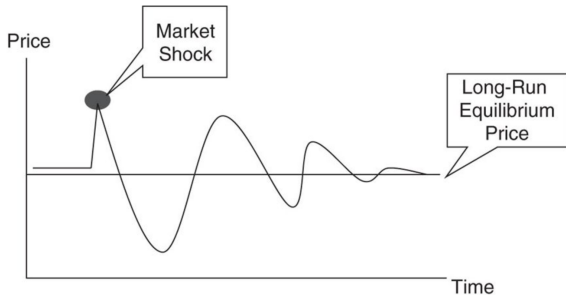
$$E(y_t^2) = \sigma^2 \quad \forall t \in \mathbb{N}$$

$$\lambda(s, t) = \lambda(s + h, t + h) \quad \forall s, t \in \mathbb{N}, \forall h \in \mathbb{Z}$$

with $\lambda(s, t) = E[(y_s - \mu)(y_t - \mu)]$

Stationarity: consequences

One of the consequences of stationary processes is how they recover from a shock. A shock represents an unexpected change in a variable in a particular period of time. For stationary time series, shocks to the system will gradually die away.



Non-stationary processes

There are different types of non-stationary time series models often found in economics. Unit root process is one of them commonly used to model prices:

Unit root process

Unit root or $I(1)$ processes are also called stochastic trend if they have the following form:

$$y_t = \mu + y_{t-1} + \epsilon_t$$

where ϵ_t is a stationary process. When $\mu = 0$ the process is called **pure random walk** and when $\mu \neq 0$ the process is called **random walk with drift**.

Integration of order d

A time series \mathbf{y} is said to be $I(d)$, or integrated of order d , if after differentiating the variable d times, we get an stationary process, more precisely:

$$(1 - L)^d \mathbf{y} \sim I(0),$$

where $I(0)$ is a stationary time series and L is the lag operator:

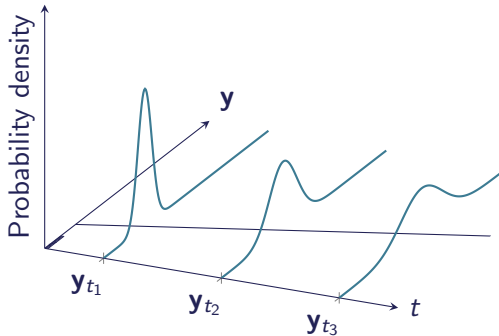
$$(1 - L)\mathbf{y} = \Delta \mathbf{y} = \mathbf{y}_t - \mathbf{y}_{t-1} \quad \forall t$$

The screenshot displays a trading platform interface for the EUR/USD currency pair. The main window shows a 'TICK' chart with price movement over time. A circular inset provides a magnified view of a specific price fluctuation, showing a sharp drop followed by a recovery. The interface includes a top toolbar with various charting tools and a bottom volume bar.

Heteroscedastic

Heteroscedasticity

Heteroscedasticity identifies **non-constant volatility**. Volatility gives a measure of uncertainty about future returns. It can either be obtained from historical returns or implied from derivatives such as options.



Cointegration definition

Cointegration

The relationship between non-stationary time series can be modelled if a **stationary linear combination** is shown to exist. When this happens it is said they have a **long-run equilibrium** relationship and the variables are **cointegrated**.



Cointegration definition

Let $\mathbf{y} = \{\mathbf{y}^1, \dots, \mathbf{y}^I\}$ be a set of I time series of order $I(1)$ which are said to be cointegrated if a vector $\boldsymbol{\beta} = [\beta(1), \dots, \beta(I)]^\top \in \mathbb{R}^I$ exists such that the time series,

$$\mathbf{Z}_t := \boldsymbol{\beta}^\top \mathbf{y} = \beta(1)\mathbf{y}^1 + \dots + \beta(I)\mathbf{y}^I \sim I(0).$$

The Vector Error correction

The Vector Error correction model (VECM), considering p lags, is used to model cointegrated time series:

$$\Delta \mathbf{y}_t = \underbrace{\Omega \mathbf{y}_{t-1}}_{\text{Error correction term}} + \underbrace{\sum_{i=1}^{p-1} \phi_i^* \Delta \mathbf{y}_{t-i}}_{\text{Autoregressive term}} + \underbrace{\mathbf{c}}_{\text{Intercept}} + \underbrace{\epsilon_t}_{\text{i.i.d}}$$

where

- $\mathbf{y}_t = [y_t^1, \dots, y_t^l]^\top$,
- the errors \mathbf{e}_t are assumed to follow i.i.d l -dimensional multivariate normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$,

The Vector Error correction

The Vector Error correction model (VECM), considering p lags, is used to model cointegrated time series:

$$\Delta \mathbf{y}_t = \underbrace{\Omega \mathbf{y}_{t-1}}_{\text{Error correction term}} + \underbrace{\sum_{i=1}^{p-1} \phi_i^* \Delta \mathbf{y}_{t-i}}_{\text{Autorregressive term}} + \underbrace{\mathbf{c}}_{\text{Intercept}} + \underbrace{\epsilon_t}_{\text{i.i.d}}$$

where

- $\Omega = \alpha \beta^\top$. The columns of β contains the cointegration vectors and the rows of α correspond with the adjusted vectors.

VECM limitations

The Vector Error correction model is used to model cointegrated time series but only with **batch data and rarely used with high frequency data** mainly due to **computational limitations**:

- Calculation of the cointegration vectors β is obtained by the Johansen procedure which is of order $O(n^3)$.
- VECM parameters are obtained using the ordinary least squares (**OLS**) method.

Recently, online learning algorithms have been proposed to solve problems with large data sets mainly because of:

- Their simplicity
- They process one instance at a time
- The hypothesis is updated every time new data arrives

Algorithm 1 Structure of an Online Learning System

- 1: Receives input \mathbf{x}_t
 - 2: Makes prediction $\hat{\mathbf{y}}_t$
 - 3: Receives response \mathbf{y}_t
 - 4: Incurs loss $l_t(\mathbf{y}_t, \hat{\mathbf{y}}_t)$
-

Performance is later measured after T trials as:

$$L_T = \sum_{t=1}^T l_t(\mathbf{y}_t, \hat{\mathbf{y}}_t)$$

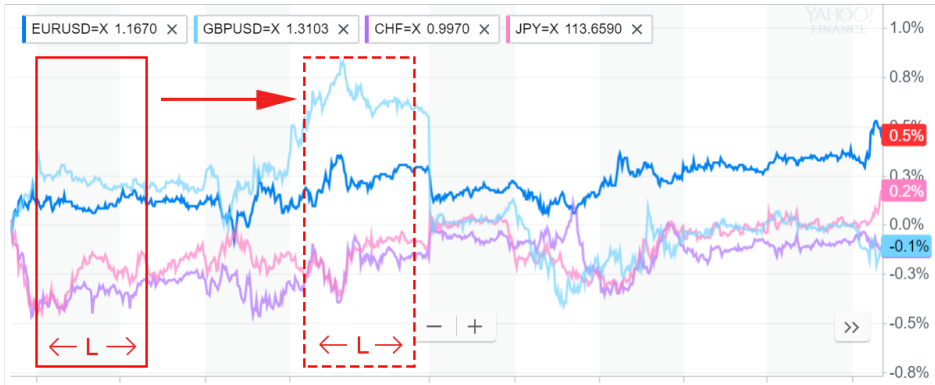
Proposal

To adapt VECM to be used with high frequency data, where the best parameters are updated when new data arrives.

Hypothesis

An algorithm based on cointegration and high performance computing will allow faster forecasting algorithms for financial time series to be obtained while maintaining good accuracy levels.

Adaptive VECM



VECM(p) matrix form

$$\Delta \mathbf{y}_t = \underbrace{\Omega \mathbf{y}_{t-1}}_{\text{Error correction term}} + \underbrace{\sum_{i=1}^{p-1} \phi_i^* \Delta \mathbf{y}_{t-i}}_{\text{Autoregressive term}} + \underbrace{\mathbf{c}}_{\text{Intercept}} + \underbrace{\epsilon_t}_{\text{i.i.d}}$$

If we have N data points, with $N > p$, the VECM matrix form is as follows:

$$\underbrace{\begin{bmatrix} \Delta \mathbf{y}_{p+1}^T \\ \Delta \mathbf{y}_{p+2}^T \\ \vdots \\ \Delta \mathbf{y}_N^T \end{bmatrix}}_{\mathbf{Y}} = \underbrace{\begin{bmatrix} \mathbf{y}_p^T \beta & \Delta \mathbf{y}_p^T & \Delta \mathbf{y}_{p-1}^T & \cdots & \Delta \mathbf{y}_2^T & 1 \\ \mathbf{y}_{p+1}^T \beta & \Delta \mathbf{y}_{p+1}^T & \Delta \mathbf{y}_p^T & \cdots & \Delta \mathbf{y}_3^T & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{y}_{N-1}^T \beta & \Delta \mathbf{y}_{N-1}^T & \Delta \mathbf{y}_{N-2}^T & \cdots & \Delta \mathbf{y}_{N-p-1}^T & 1 \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} \alpha^T \\ \phi_1^{*T} \\ \phi_2^{*T} \\ \vdots \\ \phi_{p-1}^{*T} \\ \mathbf{c}^T \end{bmatrix}}_{\mathbf{W}} + \underbrace{\begin{bmatrix} \epsilon_{p+1}^T \\ \epsilon_{p+2}^T \\ \vdots \\ \epsilon_N^T \end{bmatrix}}_{\mathbf{E}}$$

VECM can be solved using a sliding windows of data, the problem can be expressed as follows:

$$\mathbf{X}(t) = \begin{bmatrix} \mathbf{x}_{t-L}^\top \\ \vdots \\ \mathbf{x}_t^\top \end{bmatrix}, \mathbf{Y}(t) = \begin{bmatrix} \mathbf{y}_{t-L}^\top \\ \vdots \\ \mathbf{y}_t^\top \end{bmatrix}.$$

The optimal solution using a OLS is then:

$$\mathbf{W}(t)_* = (\mathbf{X}(t)^\top \mathbf{X}(t))^{-1} \mathbf{X}(t)^\top \mathbf{Y}(t) \quad (1)$$

$$= \left(\sum_{i=0}^L \mathbf{x}_{t-i} \mathbf{x}_{t-i}^\top \right)^{-1} \sum_{i=0}^L \mathbf{x}_{t-i} \mathbf{y}_{t-i}^\top \quad (2)$$

Algorithm 2 Online Ordinary Least Squares

Input:

$\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$: N input vectors

$\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$: N targets

L : sliding window size ($L < N$)

Output:

$\{f(\mathbf{x}_{L+1}), \dots, f(\mathbf{x}_N)\}$: model predictions

- 1: Initialize $\mathbf{A} = \sum_{t=1}^L \mathbf{x}_t \mathbf{x}_t^\top$ and $\mathbf{b} = \sum_{t=1}^L \mathbf{x}_t \mathbf{y}_t^\top$
 - 2: **for** $t = L + 1$ to N **do**
 - 3: read new \mathbf{x}_t
 - 4: $\mathbf{A} = \mathbf{A} + \mathbf{x}_t \mathbf{x}_t^\top - \mathbf{x}_{t-L-1} \mathbf{x}_{t-L-1}^\top$
 - 5: output prediction $f(\mathbf{x}_t) = \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{x}_t$
 - 6: Read new \mathbf{y}_t
 - 7: $\mathbf{b} = \mathbf{b} + \mathbf{x}_t \mathbf{y}_t^\top$
 - 8: **end for**
-

How to choose L ?: Ω matrix properties

Since $\Omega = \alpha\beta^\top$ we can also say that:

$$\Omega \mathbf{y} = \alpha\beta^\top \mathbf{y} = \alpha \mathbf{Z}_t \sim I(0)$$

The matrix Ω has the following properties:

- If $\text{rank}(\Omega) = 0$ there is no cointegration
- If $\text{rank}(\Omega) = I$ i.e full rank, Ω is non-singular or invertible, therefore we can express \mathbf{y} as follows:

$$\mathbf{y} = \Omega^{-1} \alpha \mathbf{Z}_t \sim I(0),$$

then the time series are not $I(1)$ but stationary.

- If $\text{rank}(\Omega) = r$, $0 < r < I$ then, r cointegration vectors exist such that:

$$\Omega \mathbf{y} = \alpha \mathbf{Z}_t \sim I(0).$$

How to choose L?

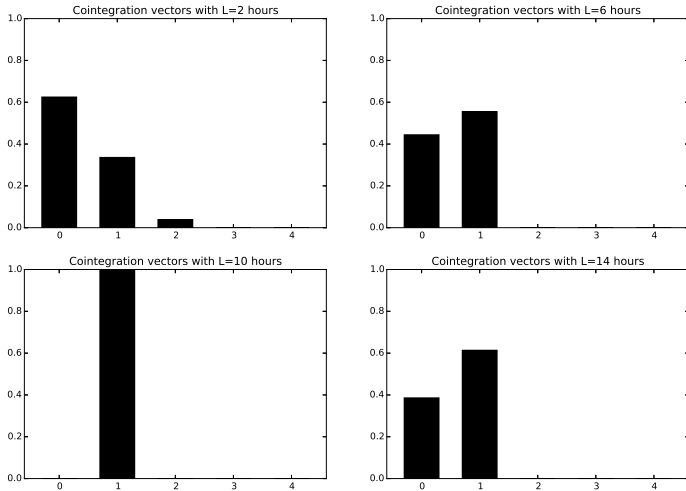


Figure 3: Distribution of the number of cointegration vectors using $p = 1$ lags.

Since $r = 0$ means no cointegration and $r = I = 4$ reveals that no process is $I(1)$ but stationary. The interesting cases of cointegration are those where r lies strictly between 0 and I , i.e. $0 < r < I$.

Percentage of cointegration (PC)

$$PC = \frac{\#\{it \mid it \text{ has } r \text{ c.v. with } 0 < r < I\}}{\#it} \times 100$$

- This data was collected from **Dukascopy**, a free database which gives access to the **Swiss Foreign Exchange marketplace**.
- Tests were carried out using four foreign exchange rates all related to USD: **EURUSD, GBPUSD, USDCHF and USDJPY**.
- The tests were done using **10-seconds frequency** from ask prices from the 11th to the 15th of August 2014.



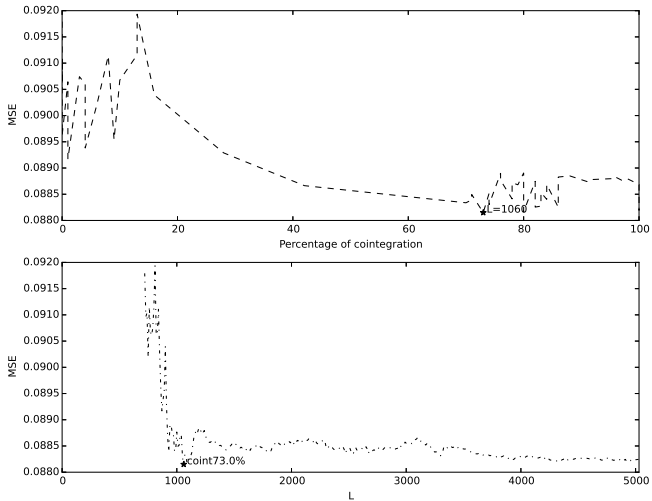


Figure 4: MSE versus the percentage of cointegration considering 1000 iterations.

Adaptive VECM algorithm (AVECM)

1. It uses VECM only considering a sliding windows of size L of data.
2. We proposed to choose L and the number of lags p in order to maximise the percentage of cointegration PC in the near past. This process was done every time that new data was processed.
3. And to use a distributed environment to make this parameters search, since this is the most expensive routine.

Unit root tests

Augmented Dickey Fuller (ADF) test with lags $p = 1, 2, 3, 4, 5$.

Variable	ADF(1)	ADF(2)	ADF(3)	ADF(4)	ADF(5)
EURUSD	-0.052	-0.054	-0.054	-0.054	-0.054
GBPUSD	-0.744	-0.784	-0.805	-0.837	-0.846
USDCHF	-0.476	-0.493	-0.493	-0.495	-0.502
USDJPY	0.357	0.360	0.360	0.367	0.367
Δ EURUSD	-128.4*	-128.4*	-96.85*	-89.12*	-89.12*
Δ GBPUSD	-131.4*	-112.7*	-102.5*	-92.86*	-88.29*
Δ USDCHF	-127.8*	-127.8*	-96.94*	-88.82*	-80.79*
Δ USDJPY	-135.1*	-135.1*	-101.2*	-101.2*	-101.2*

Table 2: Unit roots tests for EURUSD, GBPUSD, USDCHF and USDJPY at 10-second frequency.

Performance accuracy

	MSE			U-statistic	
	AVECM	ARIMA	p-value	AVECM	ARIMA
EURUSD	1.0702 e-09	1.1481 e-09	9.250 e-12	0.6863	0.7108
GBPUSD	1.6630 e-09	1.7408 e-09	6.951 e-02	0.6866	0.7025
USDCHF	5.8503 e-10	6.3545 e-10	2.899 e-14	0.6803	0.7091
USDJPY	6.3483 e-06	6.5194 e-06	6.853 e-05	0.6964	0.7057

The function that get optimal parameters L and p was implemented using MPI in Python. Tests ran in our cluster **hpc.utfsm.cl** considering 2 servers Xeon E5-2667 (2.90GHz) of 24 cores each (48 cores in total) and 24GB RAM.

Parameter settings:

- The L parameter was always chosen between 100 and 4000 (1 hour=360 data points) and p always took values between 1 and 5.
- Parameter *nparams* represents the number of pairs (L,p) used to maximise the percentage of cointegration.

Execution times

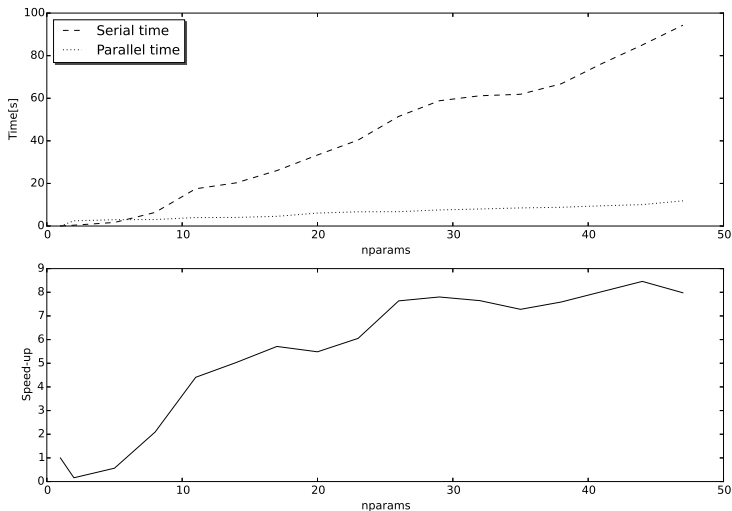


Figure 5: Computing time of sequential and parallel algorithm is shown in the upper figure. Speed-up is shown below.




Cointegration Cointegration relations change with time. We empirically showed that the number of cointegration vectors is sensitive to the number of lags but also to the amount of data considered. The percentage of cointegration is used to find VECM parameters.

Online learning The adaptation of the VECM to be used in an online context expressing the OLS method with matrix actualization.

Distributed environment The use of a distributed environment to calculate L and p and obtain a response below 10 seconds.

- Cointegration information can be used as an integration tool to detect arbitrage opportunities.
- Use with stocks, forex rates and crypto-currency?
- Applicable to other non-stationary but cointegrated time series

- AVECM improves performance measures by finding parameters of L and p maximising the percentage of cointegration.
- AVECM performance is superior to ARIMA and the naive random walk model in terms of MSE and U -statistic.
- We showed that VECM computational limitations can be tackled using a an online approach in a distributed environment.

-  Paola Arce, Jonathan Antognini, Werner Kristjanpoller, and Luis Salinas, *An online vector error correction model for exchange rates forecasting*, Proceedings of the International Conference on Pattern Recognition Applications and Methods, 2015, pp. 193–200.
-  Paola Arce, Jonathan Antognini, Werner Kristjanpoller, and Luis Salinas, *Fast and adaptive cointegration based model for forecasting high frequency financial time series*, Computational Economics (2017), 1–14.
-  Paola Arce and Luis Salinas, *Online ridge regression method using sliding windows*, XXXI International Conference of the Chilean Computer Science Society (2012).

- To add **matrix optimizations** to obtain VECM parameters.
- To add a regularization parameter to get better generalisation capabilities using **Ridge Regression** and the **Aggregating Algorithm for Regression**.
- To include **more explicative variables** such as bid-ask spread and change in volume.
- The online approach for **other econometrics models** it is also worthy of study.

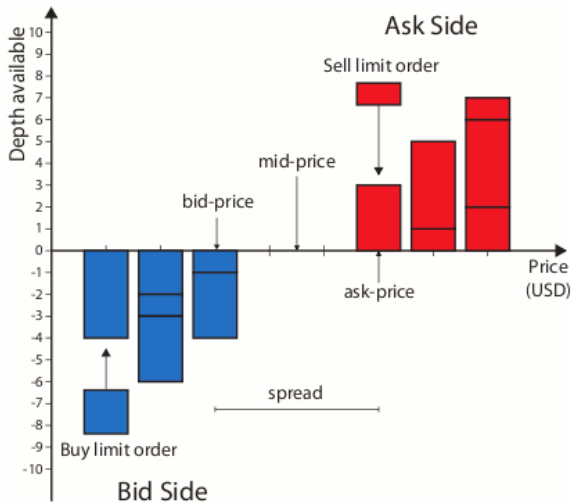
Thank you for your attention
Questions?

Stationarity

A time series $y_t : \mathbb{Z} \rightarrow \mathbb{R}$ is said to be strictly stationary if, for a finite set of L indices $\{t_1, \dots, t_L\} \forall L \in \mathbb{Z}_{\geq 0}$, the joint probability distribution of the random variables $\{y_{t_1}, y_{t_2}, \dots, y_{t_L}\}$ remains unchanged for any time shifted h . The probability distribution for the random variable y_{t_1} is:

$$P\{y_t \leq c_1\} = P\{t \in \mathbb{Z} : y_t \in]-\infty, c_1]\} \quad (3)$$

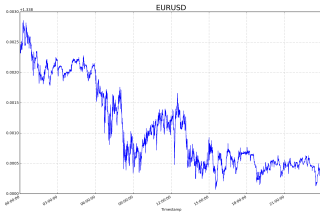
The order book



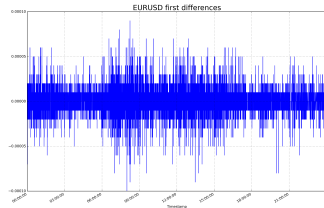
Unit root process

Integration example

High frequency financial time series are commonly an $I(1)$ process

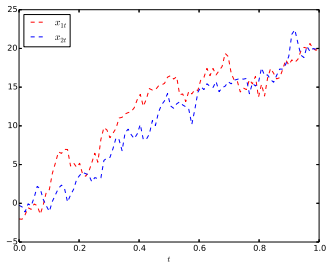


Since their differences are stationary ($I(0)$ process).



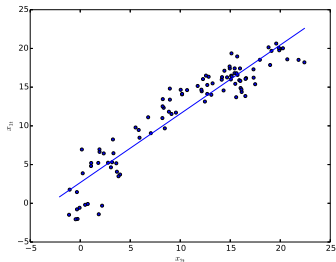
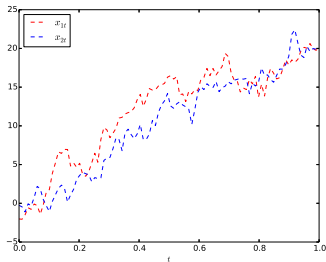
Spurious regression

Standard regression techniques are applied to non-stationary data. The use of non-stationary data can lead to spurious regressions.



Spurious regression

Standard regression techniques are applied to non-stationary data. The use of non-stationary data can lead to spurious regressions.



Random walk

The random walk model is defined as:

$$\mathbf{y}_t = \mathbf{y}_{t-1} + \epsilon_t \quad (4)$$

The naive forecast of the time series difference $\hat{\mathbf{y}}_{t+1}$ for the random walk model is defined as:

$$\hat{\mathbf{y}}_{t+1} = \mathbf{y}_t + \hat{\epsilon}_{t+1} \quad (5)$$

where $\hat{\epsilon}_{t+1} = \epsilon_t$.

A process can be modelled as an ARIMA(p, d, q) if $\mathbf{x}_t = \Delta^d \mathbf{y}_t$, is an ARMA(p, q). An ARMA(p, q) model is the following:

$$\mathbf{x}_t = \sum_{i=1}^p \phi_i \mathbf{x}_{t-i} + \sum_{j=1}^q \theta_j \epsilon_{t-j} \quad (6)$$

with coefficients $\phi_p \neq 0$, $\theta_q \neq 0$ and $\sigma_\epsilon^2 > 0$.

AVECM Algorithm

Input:

\mathbf{y} : matrix with N input vectors and I time series

j : Starting point of testing

it : Ending point of testing

ps : list of p values

Ls : list of L values ($L < N$)

m : Iterations to determine parameters ($m < N - L$)

Output:

$\{\hat{\mathbf{y}}[1], \dots, \hat{\mathbf{y}}[it]\}$: prediction vectors

1: **for** $i = j$ to it **do**

2: $\mathbf{Y} \leftarrow \mathbf{y}[:, i - 1]$

3: $L, p \leftarrow \text{get_best_params}(Ls, ps, m, \mathbf{Y})$

4: $model = \text{VECM}(\mathbf{Y}, L, p)$

5: $\hat{\mathbf{y}}[i - j] = model.predict()$

6: **end for**