

# Online Convex Optimization for Time Series Prediction\*

Research Report: Internship Summer 2018

Saad Christophe      Ackerer Damien<sup>†</sup>

August 31, 2018

## Abstract

This report is a description of the project I realized during my summer internship at Swissquote Group Holding SA in Gland, Switzerland as a quantitative researcher intern. In this project, I worked on a novel approach of time series prediction based on online learning and convex optimization. I am going to show the efficiency and reliability of some algorithms along with their capability of replacing traditional and usually used techniques. In particular, I am going to study a new approach of univariate and multivariate ARMA models based on online convex optimization. [Haz16] [OAS13] [CxF17].

## 1 Motivation

A time series is a sequence of signals taken at successive equally spaced points in time. Given some time series, predicting the next values according to the previous ones can be very useful in the field of data analysis. Many works have been done on this purpose, the **AutoRegressive Moving Average models** (ARMA) are often used for time series prediction. They assume that each new signal is a noisy linear combination of the last  $k$  observations and  $q$  last independent noise terms.

$$X_t = \sum_{i=1}^k \alpha_i X_{t-i} + \sum_{i=1}^q \beta_i \epsilon_{t-i} + \epsilon_t$$

where  $\alpha \in \mathbb{R}^k$  is the autoregressive (AR) coefficient and  $\beta \in \mathbb{R}^q$  is the moving-average (MA) coefficient. There are many techniques for estimating these parameters in order to make the best ARMA prediction. The most commonly used : the **Maximum-Likelihood Estimation**. This estimation technique operates in maximizing the (log)likelihood function. There are many disadvantages of using this method as shown in Table 1.

Maximum-Likelihood Estimation	Online Learning approach
<ul style="list-style-type: none"><li>Huge computational cost</li><li>Strong assumptions on noise terms (independence and identically distributed Gaussian distribution)</li></ul>	<ul style="list-style-type: none"><li>Very low computational cost</li><li>Approach more general, capacity of coping with a wider range of time series</li><li>Noise terms arbitrarily generated</li><li>Handle any type of loss functions</li></ul>

Table 1: Disadvantages of Maximum-Likelihood estimation and advantages of an Online Learning approach

*In which situations is it better to use online convex optimization algorithms and make better predictions than the traditional maximum-likelihood estimation?*

\*Quantitative Asset Management, Swissquote Bank. Email: qam@swissquote.ch

<sup>†</sup>Internship supervisor

## 2 Introduction to online convex optimization

In this section, we give a introduction to the concept of online convex optimization(OCO) along with an example.

### 2.1 Basic Definitions

Here are some basic definitions needed to understand online convex optimization.

**Convex Optimization :** It is a subfield of optimization that studies the problem of minimizing **convex functions** over **convex sets**.

**Online Learning :** It a method of machine learning in which data becomes available in **a sequential order** and is used to **update our best predictor** for future data at each step.

### 2.2 The Online Convex Optimization model

We explain now the idea behind the Online Convex Optimization model along with an [Example](#).

Online convex optimization is defined in a game-theoretic framework, it can be seen as a structured repeated game where:

1. At each iteration  $t$ , an online player chooses an element  $\gamma^t$  from a convex bounded set of decisions  $\kappa \in \mathbb{R}^n$  without knowing the outcome of his choice.
2. After committing his decision, the adversary chooses a convex loss function  $l_t : \kappa \rightarrow \mathbb{R}$  from a bounded family of loss functions.
3. The player then suffer loss  $l_t(\gamma^t)$  and adds this loss to his total sum of losses  $L = \sum_{i=1}^T l_i(\gamma^i)$

Our goal is to minimize the sum of losses over a number  $T$  of iterations. This can be achieved by trying to minimize the convex loss function over the set at each iteration  $t$  in order to make better choices at  $t+1$ .

A **good OCO algorithm** is one with a regret (defined below) that grows sublinearly as a function of  $T$ , this implies that on average the algorithm is performing as the best strategy we can have in hindsight. In other words,

$$R_T = \sup_{\{l_1, \dots, l_T\}} \left( \sum_{t=1}^T l_t(\tilde{\gamma}^t) - \min_{\gamma \in \kappa} \sum_{t=1}^T l_t(\gamma) \right) = o(T)$$

where  $\min_{\gamma \in \kappa} \sum_{t=1}^T l_t(\gamma)$  is the total loss of the best decisions we could have made with the best strategy.

To evaluate the performance of an algorithm, we look at an upper bound of its worst case regret.

**Example 1** (The weighted majority (WM) algorithm). *Let us imagine a game where an online player needs to make one of two choices A and B, he has the assistance of N advisers. This situation can be modeled as a choice from the set  $\kappa$  of advisors with  $|\kappa| = N$ . The player has to choose at iteration t the best advice.*

*One method is to attribute at each iteration t, to each advisor i a weight  $w_t(i)$  which will define the quality of his advice.*

*We then group the advisers in two subsets  $\kappa_A, \kappa_B \subseteq \kappa$ , the advisers that chose respectively A and B.*

*For choices A and B we give weights  $W_t(A) = \sum_{a \in \kappa_A} w_t(a)$  and  $W_t(B) = \sum_{b \in \kappa_B} w_t(b)$*

*See [Algorithm 1](#)*

---

**Algorithm 1** WM algorithm

---

- 1: **Initialization**  $w_1(i) = 1, i = 1, \dots, N$  and Total Loss  $L = 0$
- 2: **for**  $t = 1$  **to**  $T$
- 3:   **Predict**  $\tilde{c} \leftarrow \begin{cases} A & W_t(A) > W_t(B) \\ B & \text{otherwise.} \end{cases}$
- 4:   **Observe**  $c$  and suffer loss  $l_t(\tilde{c})$ ,  $L = L + l_t(\tilde{c})$
- 5:   **Update**  $w_{t+1}(i) \leftarrow \begin{cases} w_i(t) & w_i(t) = c \\ (1-\epsilon)w_i(t) & \text{otherwise for some } \epsilon > 0 \end{cases}$
- 6:   **Update**  $W_{t+1}(A)$  and  $W_{t+1}(B)$
- 7: **end for**
- 8: **return**  $L$

---

### 3 Online convex optimization algorithms

In this section, we present two OCO algorithms, the Online Gradient Descent and the Online Newton step algorithms.

#### 3.1 Settings

We give some setting and notations for the rest of this report.

$\kappa \subseteq \mathbb{R}^n$	convex set
$l_t$	loss function revealed at time $t$
$T$	total number of iterations (run)
$\{\gamma^t\}_{t=1}^T$	time series
$\{\tilde{\gamma}^t\}_{t=1}^T$	predictions series
$\Pi_\kappa(y)$	Euclidean projection of $y$ onto $\kappa$ $\arg \min_{x \in \kappa} \ y - x\ _2$
$\Pi_\kappa^A(y)$	Euclidean projection of $y$ onto $\kappa$ in the norm induced by a matrix $A$ $\arg \min_{x \in \kappa} (y - x)^\top A(y - x)$
$\eta$	learning rate, step sizes we make at each iteration towards the minimum of the convex function
$G$	upper bound of the norm of the gradient $\ \nabla l_t(x)\  \leq G$
$D$	upper bound of the diameter of set $\kappa$ $\sup_{x, y \in \kappa} \ x - y\ _2 \leq D$
$\lambda$	the exp-concavity parameter of the loss functions $\{l_t\}_{t=1}^T$ , this means that $\forall t, e^{-\lambda l_t}$ is concave
<b>Remark.</b>	$\Pi_\kappa(y) = \Pi_\kappa^I(y)$

#### 3.2 Online Gradient Descent (OGD)

Online Gradient Descent is the first OCO algorithm we are going to present.

**Gradient descent**, represented in [Figure 1](#), is an iterative optimization method which consists of iteratively moving the current point in the direction of the function's gradient to reach its extrema. In our case, this function will be the loss function and we will try to reach its minimum.

**Online gradient descent**, [Figure 2](#), is a **first order method**, this means that it only takes into consideration the first order derivative of the convex function (its gradient). [Algorithm 2](#) describes how it operates.

---

**Algorithm 2** Online Gradient Descent

---

- 1: **Input** convex set  $\kappa$ ,  $T$ ,  $\gamma_1 \in \kappa$ , step sizes  $\{\eta_t\}$
- 2: **for**  $t = 1$  **to**  $T$
- 3:   **Predict**  $\tilde{\gamma}^t$
- 4:   **Observe**  $\gamma^t$  and  $l_t(\tilde{\gamma}^t)$
- 5:   **Project and Update:**  $\widetilde{\gamma_{t+1}} = \Pi_\kappa(\gamma^t - \eta_t \nabla l_t)$
- 6: **end for**

---

**OGD** grantees a regret bound of

$$R_t \leq \frac{3}{2} GD\sqrt{T} = \mathcal{O}(GD\sqrt{T})$$

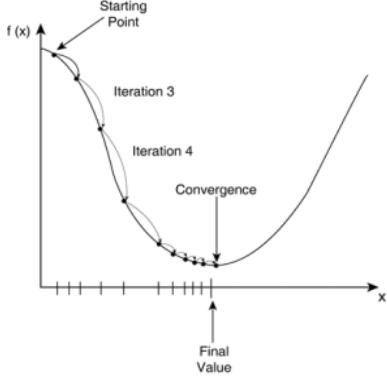


Figure 1: Gradient descent, taking steps in the direction of the gradient to reach the minimum of the convex function

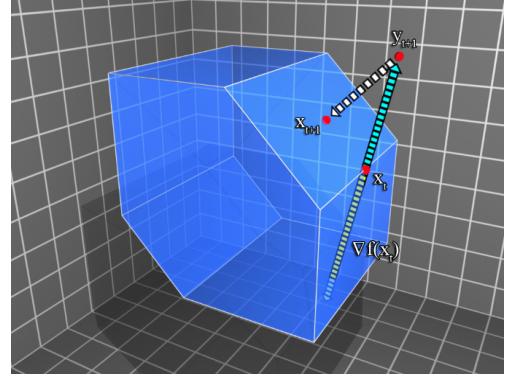


Figure 2: Online gradient descent, moving in the direction of the gradient and then projecting back to the set

### 3.3 Online Newton Step (ONS)

Next, we look at the second algorithm, the Online Newton Step.

**Newton-Raphson** is an iterative method for finding the roots of a differentiable function by moving in a direction defined by the function and its derivative

$$x_{k+1} = x_k - (\nabla f(x_k))^{-1} f(x_k)$$

In our case, we are interested in finding an extrema, so the steps will be taken in the direction defined by the first and second order derivatives (the gradient and the Hessian)

$$x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

Since the Hessian is too expensive to compute, this OCO algorithm will be based on the **quasi-newton method** which uses an approximation of the hessian.

**Online newton step** is a second order method which will indirectly take into consideration the second order derivative of the function (the Hessian) by moving in the direction of the product between the inverse of a matrix  $A_t$  **related to the hessian** and the gradient. [Algorithm 3](#) describes how it operates.

---

#### Algorithm 3 Online Newton Step

---

- 1: **Input** convex set  $\kappa$ ,  $T$ ,  $\gamma_1 \in \kappa$ , step size  $\eta$ , an initial matrix  $A_0 = \frac{1}{\eta^2 D^2} I$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Predict**  $\widetilde{\gamma^t}$
  - 4:   **Observe**  $\gamma^t$  and  $l_t(\widetilde{\gamma^t})$
  - 5:   **Update**  $A_t = A_{t-1} + \nabla l_t \nabla l_t^\top$
  - 6:   **Project and Update:**  $\widetilde{\gamma_{t+1}} = \prod_{\kappa}^{A_t} (\gamma^t - \frac{1}{\eta} A^{-1} \nabla l_t)$
  - 7: **end for**
- 

**ONS** grantees a regret bound of

$$R_t \leq 5(\frac{1}{\lambda} + GD)n \log(T) = \mathcal{O}((GD + \frac{1}{\lambda}) \log(T))$$

**Remark.** There is an efficient way of computing  $A_t^{-1}$  at each iteration using Sherman-Morrison formula

$$(A + uu^\top)^{-1} = A^{-1} - \frac{A^{-1}uu^\top A^{-1}}{1 + u^\top A^{-1}u}$$

where  $A \in \mathbb{R}^{n \times n}$  and  $u \in \mathbb{R}^n$

## 4 Application to univariate ARMA model

These machine learning techniques are applicable to the ARMA model in order to make an estimation of the parameters and make good predictions. But since the sequence of noise terms are never revealed at any time, we cannot use online convex optimization over the space of AR and MA coefficients  $(\alpha, \beta)$ . We use an improper learning approach by approximating the ARMA( $p, q$ ) model by an AR( $M$ ) where  $M = m + k$  for some  $m \in \mathbb{N}$

$$\tilde{X}_t = \sum_{i=1}^M \gamma_i^t X_{t-i}$$

**Remark.** In all this project, only the squared loss function is considered  $l_t(\tilde{x}) = (x - \tilde{x})^2$ , the parameters are calculated according to the squared loss.

### 4.1 Settings

These are the settings of the online learning models:

$$1. \kappa = \{\gamma \in \mathbb{R}^M : |\gamma_i| \leq c \quad i = 1, \dots, M\} \quad \text{such that } |\alpha_i| \leq c \quad i = 1, \dots, k$$

$$2. l_t(\tilde{X}) = (X_t - \tilde{X}_t)^2 = (X_t - \sum_{i=1}^M \gamma_i^t X_{t-i})^2 = l_t(\gamma^t)$$

$$3. \nabla l_t(\gamma^t) = \begin{pmatrix} \frac{\partial}{\partial \gamma_1} ((X_t - \sum_{i=1}^M \gamma_i^t X_{t-i})^2) \\ \frac{\partial}{\partial \gamma_2} ((X_t - \sum_{i=1}^M \gamma_i^t X_{t-i})^2) \\ \vdots \\ \frac{\partial}{\partial \gamma_M} ((X_t - \sum_{i=1}^M \gamma_i^t X_{t-i})^2) \end{pmatrix} = -2(X_t - \tilde{X}_t) \begin{pmatrix} X_{t-1} \\ X_{t-2} \\ \vdots \\ X_{t-M} \end{pmatrix}$$

$$4. G = 2c\sqrt{M} \cdot (X_{max})^2$$

$$5. D = 2c\sqrt{M} = 2\max_{x \in \kappa} \|x\|_2 \geq \sup_{x, y \in \kappa} \|x - y\|_2$$

$$6. \lambda = \frac{1}{M}$$

### 4.2 ARMA OGD

We apply the Online Gradient Descent on the ARMA model as shown in [Algorithm 4](#).

---

**Algorithm 4** ARMA online gradient descent

---

- 1: **Input** Horizon  $M$ , first step size  $\eta_1 = \frac{D}{G}$ , set constraint  $c$ ,  $\gamma^1 \in \kappa$ ,  $T$ , in-sample data  $\{X\}$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Update**  $\eta_t \leftarrow \frac{1}{\sqrt{t}} \eta_1$
  - 4:   **Predict**  $\tilde{X}_t = \sum_{i=1}^M \gamma_i^t X_{t-i}$
  - 5:   **Observe**  $X_t$  and suffer loss  $l_t(\gamma^t)$
  - 6:   **Project and Update:**  $\gamma^{t+1} \leftarrow \Pi_\kappa(\gamma^t - \eta_t \nabla l_t)$
  - 7: **end for**
-

### 4.3 ARMA ONS

We apply the Online Newton Step algorithm on the ARMA model as shown in [Algorithm 5](#).

---

**Algorithm 5** ARMA online newton step

---

- 1: **Input** Horizon  $M$ , step size  $\eta = \frac{1}{2} \min\{\lambda, \frac{1}{4GD}\}$ , set constraint  $c$ ,  $\gamma^1 \in \kappa$ , initial  $(M \times M)$  matrix  $A_0 = \frac{1}{\eta^2 D^2} I_M$ ,  $T$ , in-sample data  $\{X\}$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Predict**  $\tilde{X}_t = \sum_{i=1}^M \gamma_i^t X_{t-i}$
  - 4:   **Observe**  $X_t$  and suffer loss  $l_t(\gamma^t)$
  - 5:   **Update**  $A_t \leftarrow A_{t-1} + \nabla l_t \nabla l_t^\top$
  - 6:   **Project and Update:**  $\gamma^{t+1} \leftarrow \Pi_{\kappa}^{A_t} (\gamma^t - \frac{1}{\eta} A_t^{-1} \nabla l_t)$
  - 7: **end for**
- 

### 4.4 Experiments

We now test these algorithms on synthetic and real data, we compare their predictions with those done by the ARMA MLE by using the root-mean-square-error.

$$RMSE_t = \sqrt{\frac{\sum_{i=1}^t (X_i - \tilde{X}_i)^2}{t}}$$

We set  $M = 10$  and  $c = 1$  for all experiments.

#### a Setting 1: Sanity Check

We generate data using the R function arima.sim with gaussian noise normally distributed  $\mathcal{N}(0, \sigma^2)$ . The parameters used are in [Table 2](#) and the results are shown in [Figure 3](#) and [Figure 4](#).

$\alpha$	(0.6, -0.5, 0.4, -0.4, 0.3)
$\beta$	(0.3, -0.2)
$\sigma$	0.3
run	2000
training run (ONS and OGD)	2000

Table 2: Settings 1

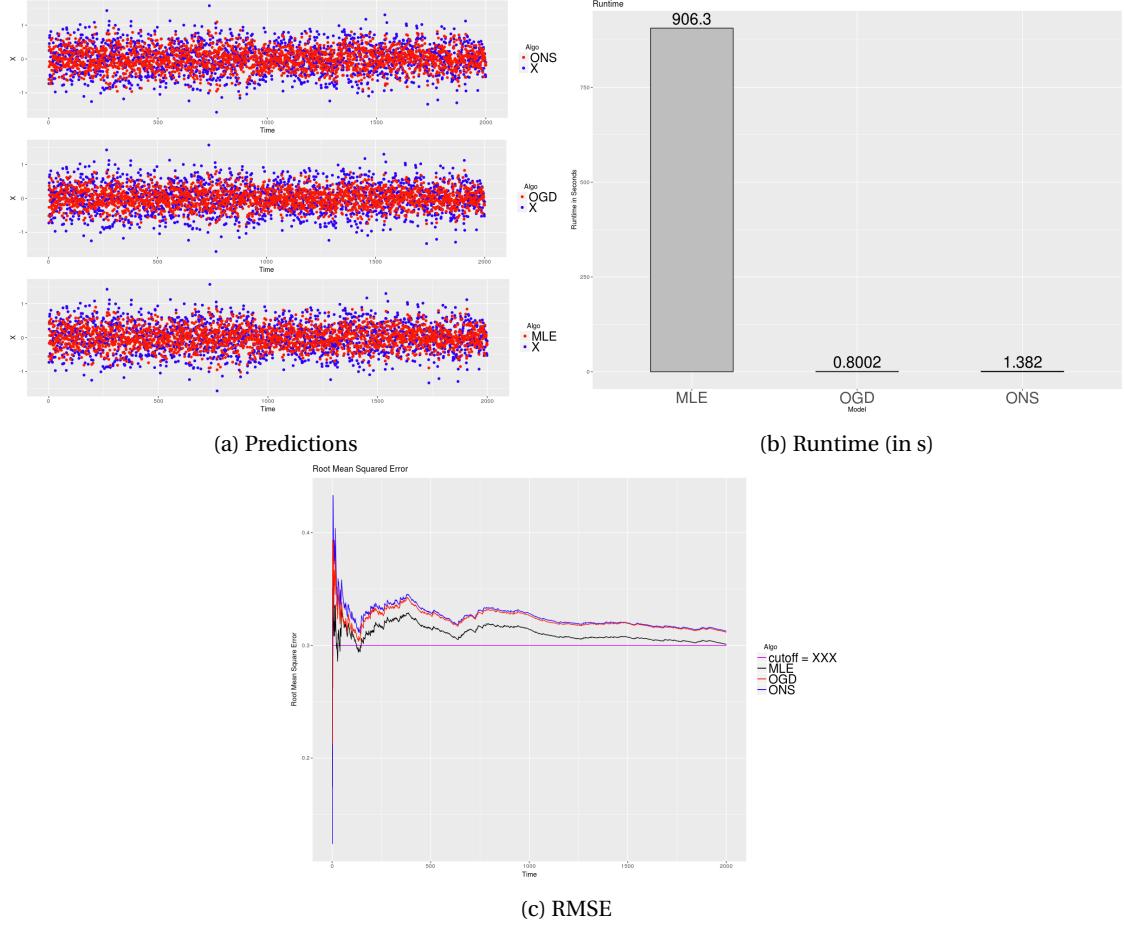


Figure 3: Setting 1 results with MLE having a growing window for predictions

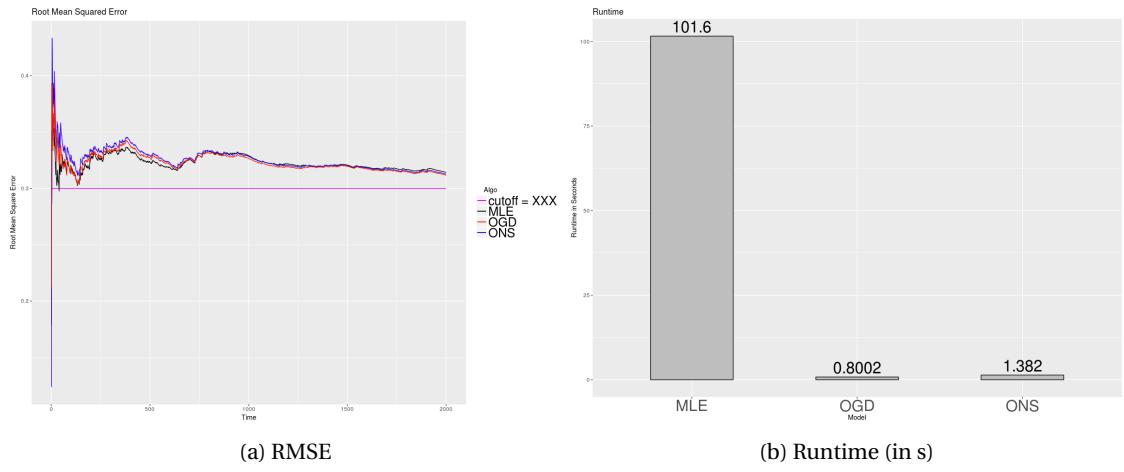


Figure 4: Setting 1 results with MLE having a fixed window of 252 for predictions

## b Setting 2: Abrupt change

We combine two series of synthetic data using arima.sim with different parameters to check how the algorithms react to an abrupt change. The noise are normally distributed  $\mathcal{N}(0, \sigma^2)$ . The parameters used are in [Table 3](#) and the results are shown in [Figure 5](#) and [Figure 6](#).

$\alpha_1$	(0.6, -0.5, 0.4, -0.4, 0.3)
$\alpha_2$	(-0.4, -0.1, 0.65, -0.1, 0.1)
$\beta_1$	(0.3, -0.2)
$\beta_2$	(-0.7, 0.4)
$\sigma_1$	0.3
$\sigma_2$	0.1
run	1000
training run (ONS and OGD)	1000

Table 3: Settings 2

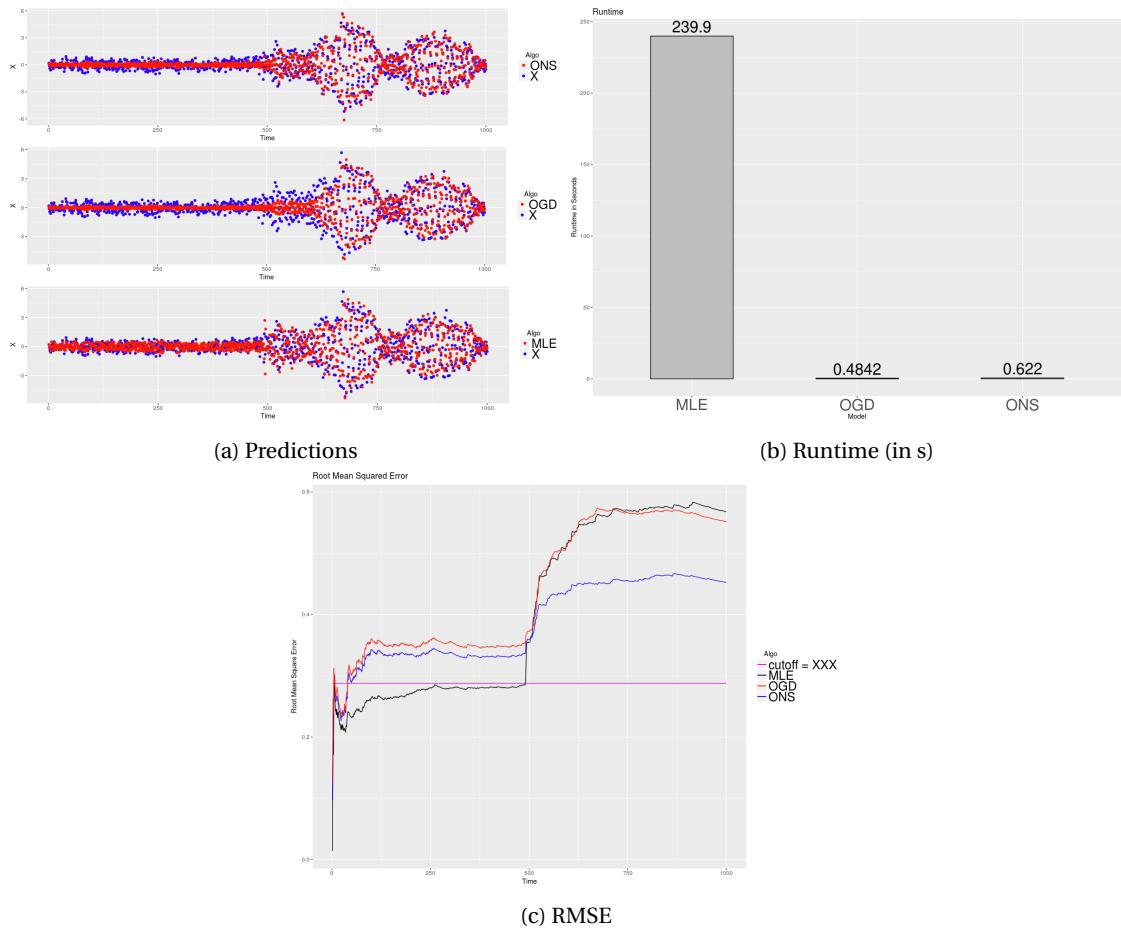


Figure 5: Setting 2 results with MLE having a growing window for predictions

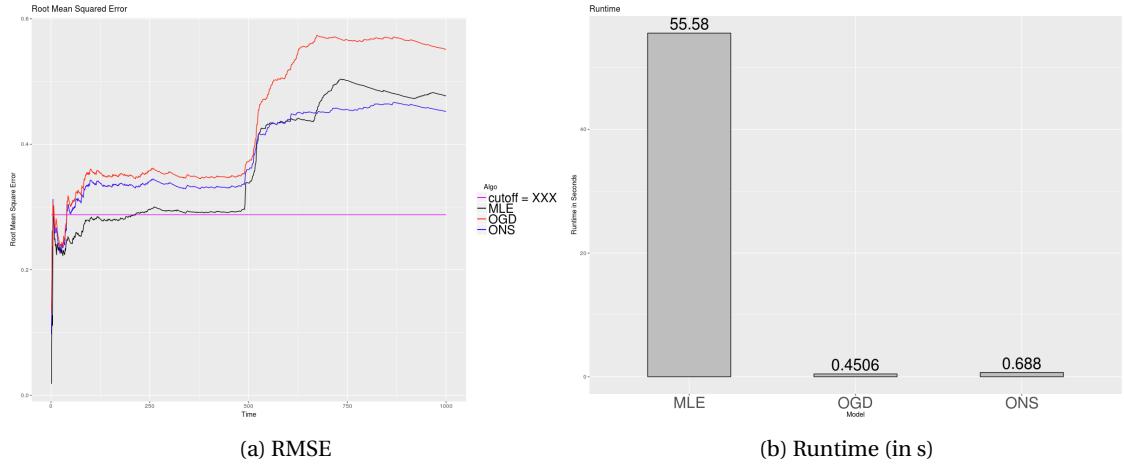


Figure 6: Setting 2 results with MLE having a fixed window of 252 for predictions

### c Setting 3: Real Data, Apple Daily Volumes

We are now interested and testing these algorithm on real data, apple daily volumes ([Figure 7](#)). The parameters used are in [Table 4](#) and the results are shown in [Figure 8](#).

run	1000
training run (ONS and OGD)	1200
MLE window size	252

Table 4: Settings 3

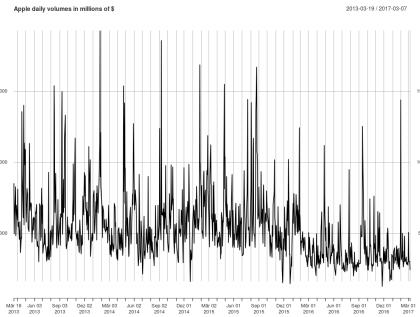


Figure 7: Apple Daily Volumes in millions of \$

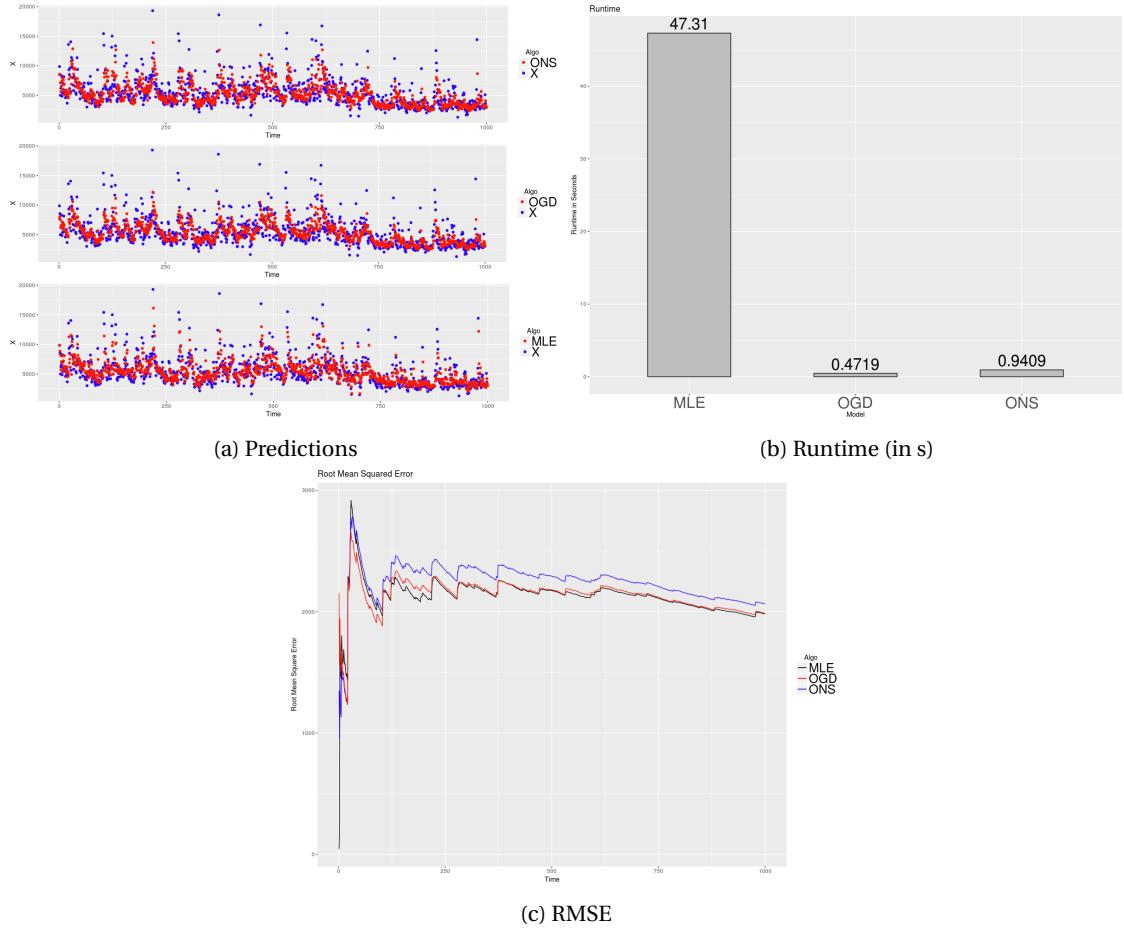


Figure 8: Setting 3 results

#### d Setting 4: Real Data, half-hourly electricity demand in England and Wales

Another example of real data time series prediction, half-hourly electricity demand in England and Wales (Figure 9). The parameters used are in Table 5 and the results are shown in Figure 10 and in Figure 11.

run	1500
training run (ONS and OGD)	2000

Table 5: Settings 4

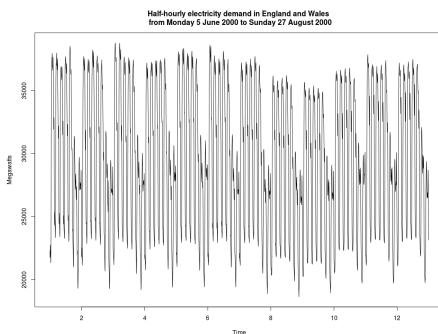


Figure 9: Half-hourly electricity demand in England and Wales in Megawatts

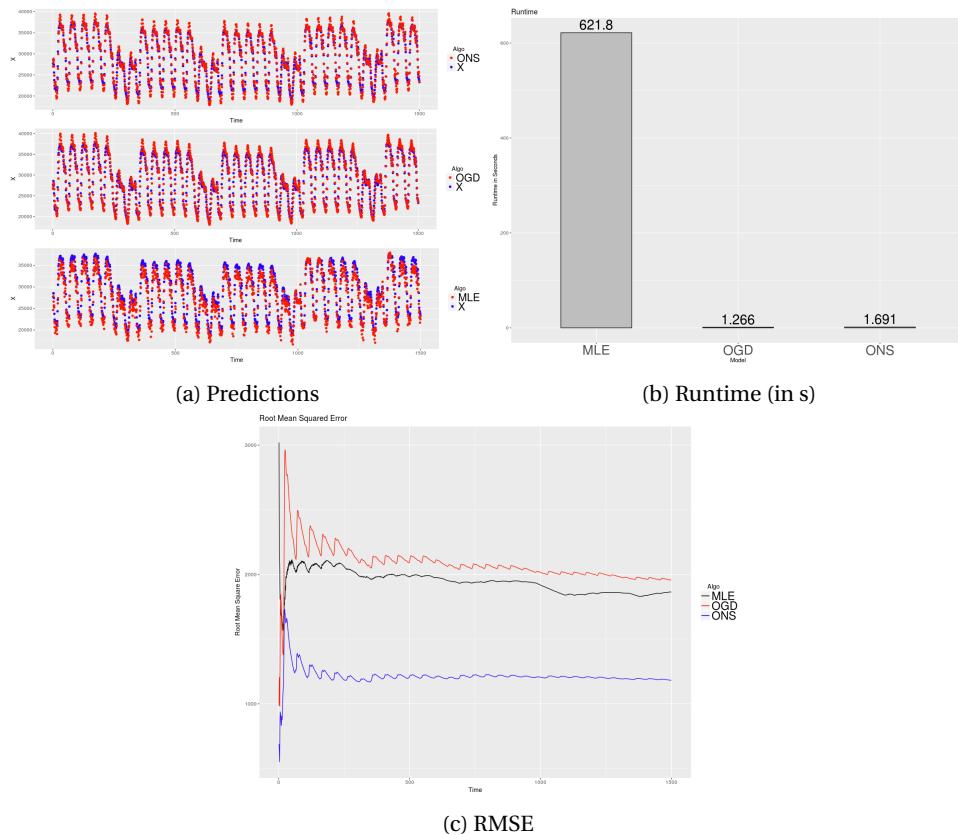


Figure 10: Setting 4 results with MLE having a growing window for predictions

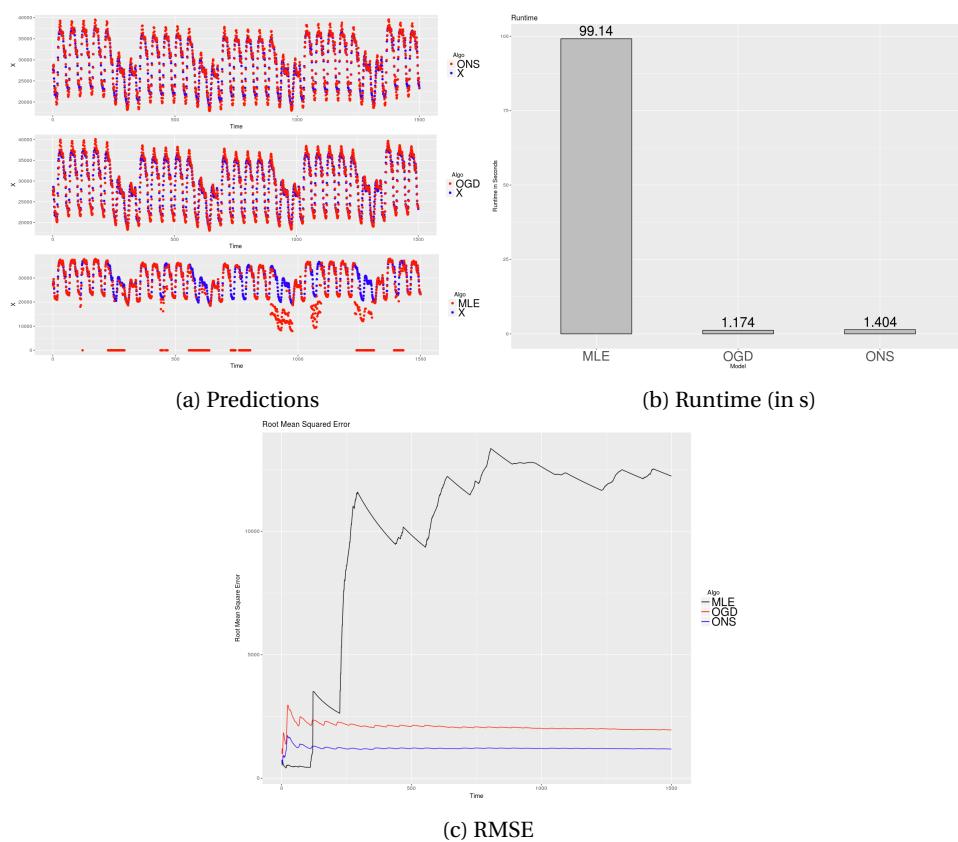


Figure 11: Setting 4 results with MLE having a fixed window of 252 for predictions

## 5 Application to multivariate ARMA process

In this section, we apply Online Convex Optimization to multivariate ARMA models. We test our algorithms on both synthetic and real data.

### 5.1 Some notations

Here are some notations and definitions that we will use in the rest of the report.

$A[i, j]$	is the $i^{th}$ row and $j^{th}$ column element of an $(m \times n)$ matrix $A$
$\ A\ _2$	Frobenius norm of a matrix $A$ $\ A\ _2 = \sqrt{\sum_i \sum_j A[i, j]^2}$
$\ A\ _W$	Frobenius norm of a matrix $A$ induced by a matrix $W$ $\ A\ _W = \ W^{\frac{1}{2}}(A - X)W^{-\frac{1}{2}}\ _2$
$\ A\ _*$	nuclear norm of a matrix $A$ $\ A\ _* = \sum_{i=1}^{\min(m, n)} \sigma_i(A)$
$\ A\ _{\max}$	norm max of a matrix $A$ $\ A\ _{\max} = \max_{i,j}  A[i, j] $
$\Pi_{\kappa}^W(A)$	Euclidean projection of $A$ onto $\kappa$ in the norm induced by a matrix $W$ $\arg \min_{X \in \kappa} \ A - X\ _W$
$\Pi_{\kappa}(A)$	Euclidean projection of $A$ onto $\kappa$ $\arg \min_{X \in \kappa} \ A - X\ _2$
$\Pi_{\mathcal{B}(*, \rho)}(A)$	Projection of $A$ onto the nuclear norm ball of radius $\rho$
	<i>Procedure:</i>
	1-Make a singular values decomposition of $A = U\Sigma V^T$
	2-Project greatest $\min(m, n)$ singular values onto an $l1$ -ball of radius $\rho$ ( <a href="#">Algorithm 10</a> )
	3-Make a new diagonal matrix $\Sigma^*$ with diagonal the projection of the singular values
	4-Return the matrix $U\Sigma^*V^T$
	The returned matrix may have a lower rank, this is what we call low rank approximation.

**Low rank approximation:** It is a technique of machine learning where a matrix of rank  $r$  is approximated to another matrix of lower rank. The best approximation is measured by minimizing a predefined cost function. This technique is used to simplify the model.

**Definition 1** (Cointegration). *Let  $x$  be a collection of  $n$  time series. We say that  $x_t$  is cointegrated of order  $(d, p)$  if each of its components are integrated of order  $d$  and there exists a linear combination  $\mu_1 x_{1,t} + \mu_2 x_{2,t} + \dots + \mu_n x_{n,t}$  which is of order  $(d - p)$ , the vector  $\mu$  is called the cointegration vector.*

$$x_t \sim CI(d, p) \iff x_{1,t}, x_{2,t}, \dots, x_{n,t} \sim I(d) \wedge \exists_{\mu \neq 0} \mu^\top x_t \sim I(d-p)$$

In general, we have  $x_{i,t} \sim I(1)$  and  $\Delta x_t = x_t - x_{t-1} \sim I(0)$  (stationary), in this case, the variables follow a common trend.

### 5.2 The VARMA process (Stationary Time Series)

If we wish to forecast a **stationary** series not only based upon its own past realizations, but additionally taking realizations of  $n - 1$  other stationary series into account, then we can model the series as **Vector AutoRegressive Moving Average model** VARMA( $k, q$ ). It can be represented by the following equation:

$$\tilde{x}_t = \sum_{i=1}^k \Phi_i x_{t-i} + \sum_{i=1}^q \Theta_i \epsilon_{t-i} + \epsilon_t$$

where  $x_t, \epsilon_t \in \mathbb{R}^n$  and  $\Phi_i, \Theta_i \in \mathbb{R}^{n \times n}$  are the coefficients matrices which introduce cross-dependencies between the series.

Each prediction will be an  $n$ -dimensional vector containing the forecast of each time series.

As realized for the univariate ARMA process, we approximate the VARMA( $k, q$ ) to a VAR( $M$ ) process:

$$\tilde{x}_t = \sum_{i=1}^M \Gamma_i x_{t-i}$$

## a Settings

1.  $\kappa = \{\Gamma_i \in \mathbb{R}^{n \times n} : \|\Gamma_i\|_{\max} \leq 1 \quad i = 1, \dots, M\} \quad \|A\|_{\max} \leq 1 \iff \max_{i,j} |A[i, j]| \leq 1$
2. Let  $\psi = \{\Gamma_1, \Gamma_2, \dots, \Gamma_M\}$
3.  $l_t(\tilde{x}_t) = \|x_t - \tilde{x}_t\|_2^2 = (x_{1,t} - \tilde{x}_{1,t})^2 + \dots + (x_{n,t} - \tilde{x}_{n,t})^2 = (x_{1,t} - \sum_{i=1}^M \Gamma_i[1,]x_{t-i})^2 + \dots + (x_{n,t} - \sum_{i=1}^M \Gamma_i[n,]x_{t-i})^2$
4.  $\nabla_{\Gamma_i} l_t = \frac{\partial l_t}{\partial \Gamma_i} = \frac{\partial(x_{1,t} - \tilde{x}_{1,t})^2}{\partial \Gamma_i} + \dots + \frac{\partial(x_{n,t} - \tilde{x}_{n,t})^2}{\partial \Gamma_i} = -2(x_{1,t} - \tilde{x}_{1,t}) \frac{\partial \tilde{x}_{1,t}}{\partial \Gamma_i} - \dots - 2(x_{n,t} - \tilde{x}_{n,t}) \frac{\partial \tilde{x}_{n,t}}{\partial \Gamma_i}$   
 $= -2(x_{1,t} - \tilde{x}_{1,t}) \begin{pmatrix} \frac{\partial \Gamma_i[1,]x_{t-i}}{\partial \Gamma_i[1,1]} & \frac{\partial \Gamma_i[1,]x_{t-i}}{\partial \Gamma_i[1,2]} & \dots \\ \frac{\partial \Gamma_i[1,]x_{t-i}}{\partial \Gamma_i[2,1]} & \dots & \\ \vdots & & \end{pmatrix} - \dots - 2(x_{n,t} - \tilde{x}_{n,t}) \begin{pmatrix} \frac{\partial \Gamma_i[n,]x_{t-i}}{\partial \Gamma_i[n,1]} & \frac{\partial \Gamma_i[n,]x_{t-i}}{\partial \Gamma_i[n,2]} & \dots \\ \frac{\partial \Gamma_i[n,]x_{t-i}}{\partial \Gamma_i[2,1]} & \dots & \\ \vdots & & \end{pmatrix}$   

Since  $\Gamma_i[j,]x_{t-i} = \Gamma_i[j,1]x_{1,t-i} + \Gamma_i[j,2]x_{2,t-i} + \dots + \Gamma_i[j,n]x_{n,t-i}$

 $= -2(x_{1,t} - \tilde{x}_{1,t}) \begin{pmatrix} x_{1,t-i} & x_{2,t-i} & \dots & x_{n,t-i} \\ 0 & 0 & \dots & \\ 0 & \dots & & \\ \vdots & & & \end{pmatrix} - \dots - 2(x_{1,t} - \tilde{x}_{1,t}) \begin{pmatrix} 0 & 0 & \dots & \\ 0 & \dots & & \\ \vdots & & & \\ x_{1,t-i} & x_{2,t-i} & \dots & x_{n,t-i} \end{pmatrix}$   
 $= -2 \begin{pmatrix} (x_{1,t} - \tilde{x}_{1,t})x_{t-i}^\top \\ (x_{2,t} - \tilde{x}_{2,t})x_{t-i}^\top \\ \vdots \\ (x_{n,t} - \tilde{x}_{n,t})x_{t-i}^\top \end{pmatrix}$
5.  $G = 2n.(X_{\max})^2$  where  $X_{\max} = \max_i x_{i,t}$
6.  $D = 2n = 2\max_{X \in \kappa} \|X\|_2 \geq \sup_{X, Y \in \kappa} \|X - Y\|_2$
7.  $\lambda = \frac{1}{n}$

## b VARMA OGD

Online Gradient Descent for multivariate ARMA is described in [Algorithm 6](#).

---

### Algorithm 6 VARMA online gradient descent

---

- 1: **Input** Horizon  $M$ , first step size  $\eta_1 = \frac{D}{G}$ ,  $\psi^1 \in \kappa$ ,  $T$ , in-sample data  $\{x\}$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Update**  $\eta_t \leftarrow \frac{1}{\sqrt{t}}\eta_1$
  - 4:   **Predict**  $\tilde{x}_t = \sum_{i=1}^M \Gamma_i^t x_{t-i}$
  - 5:   **Observe**  $x_t$  and suffer loss  $l_t(\psi^t)$
  - 6:   **Project and Update:**  $\Gamma_i^{t+1} \leftarrow \Pi_\kappa(\Gamma_i^t - \eta_t \nabla_{\Gamma_i} l_t)$
  - 7: **end for**
- 

## c VARMA ONS

Online Newton Step for multivariate ARMA is described in [Algorithm 7](#)

---

**Algorithm 7** VARMA online newton step

---

- 1: **Input** Horizon  $M$ , step size  $\eta = \frac{1}{2} \min \{\lambda, \frac{1}{4GD}\}$ ,  $\psi^1 \in \kappa$ , initial list of  $M$  ( $n \times n$ ) matrices  $A_i = \frac{1}{\eta^2 D^2} I_n$ ,  $T$ , in-sample data  $\{x\}$
- 2: **for**  $t = 1$  **to**  $T$
- 3:   **Predict**  $\tilde{x}_t = \sum_{i=1}^M \Gamma_i^t x_{t-i}$
- 4:   **Observe**  $x_t$  and suffer loss  $l_t(\psi^t)$
- 5:   **Update**  $A_i^t \leftarrow A_i^{t-1} + \nabla_{\Gamma_i} l_t \nabla_{\Gamma_i} l_t^\top$
- 6:   **Project and Update:**  $\Gamma_i^{t+1} \leftarrow \Pi_{\kappa}^{A_i^t} (\Gamma_i^t - \frac{1}{\eta} (A_i^t)^{-1} \nabla_{\Gamma_i} l_t)$
- 7: **end for**

---

### 5.3 The EC-VARMA process (Nonstationary Cointegrated Time Series)

The EC-VARMA model can be defined by the following equation (error-correction VARMA):

$$\Delta x_t = \Pi x_{t-1} + \sum_{i=1}^{k-1} \Gamma_i \Delta x_{t-i} + \sum_{i=1}^q \Theta_i \epsilon_{t-i} + \epsilon_t$$

where  $\Pi = -I + \sum_{i=1}^k \Phi_i$  is low rank and  $\Gamma_j = -(\Phi_{j+1} + \dots + \Phi_k)$  for  $j = 1, \dots, k-1$  and with  $x_t$  cointegrated ([Definition 1](#))  $x_t \in \mathbb{R}^n \sim I(1)$  and  $\Delta x_t \sim I(0)$

The error correction term  $\Pi x_{t-1}$  is used to keep and control the long-run dynamic of the time series. An example of how error-correction models operate in [Example 2](#).

We can also apply online convex optimization on EC-VARMA:

$$\begin{aligned} \Delta \tilde{x}_t &= \widehat{\Pi} x_{t-1} + \sum_{i=1}^M \Gamma_i \Delta x_{t-i} \\ \tilde{x}_t &= x_{t-1} + \widehat{\Pi} x_{t-1} + \sum_{i=1}^M \Gamma_i \Delta x_{t-i} \end{aligned}$$

**Example 2** (How do error-correction models work). *Imagine two  $I(1)$  time series  $x_t$  and  $y_t$  with a long-run equilibrium  $x_t = \beta y_t$ ,  $x_t - \beta y_t = 0 \sim I(0)$ ,  $x_t, y_t \sim CI(1, 1)$  with cointegration vector  $(1, -\beta)$ . Now if at time  $t$ ,  $x_t > \beta y_t$ , there will be a disturbance in the equilibrium. We face 3 choices to recover:*

1. a decrease in  $x_t$  and/or an increase in  $y_t$
2. an increase in  $x_t$  and a faster increase in  $y_t$
3. a decrease in  $x_t$  and a slower decrease in  $y_t$

This can be modeled as:

$$\begin{aligned} \Delta x_t &= -\alpha_1(x_{t-1} - \beta y_{t-1}) + \epsilon_1 \\ \Delta y_t &= \alpha_2(x_{t-1} - \beta y_{t-1}) + \epsilon_2 \end{aligned}$$

Where  $\alpha_i > 0$  represents the speed of adjustment to the long run equilibrium. In this case we will have  $\begin{pmatrix} \Delta x_t \\ \Delta y_t \end{pmatrix} = \Pi \begin{pmatrix} x_{t-1} \\ y_{t-1} \end{pmatrix} + \epsilon_t$  with  $\Pi = \begin{pmatrix} -\alpha_1 & \beta\alpha_1 \\ \alpha_2 & -\beta\alpha_2 \end{pmatrix}$  which will bear the cointegrating vector and the speed of adjustment information.

#### a Settings

1.  $\kappa = \{\Gamma_i, \widehat{\Pi} \in \mathbb{R}^{n \times n} : \|\widehat{\Pi}\|_* \leq \rho, \|\Gamma_i\|_{\max} \leq 1 \quad i = 1, \dots, M\}$
2. Let  $\psi = \{\widehat{\Pi}, \Gamma_1, \Gamma_2, \dots, \Gamma_M\}$
3.  $l_t(\tilde{x}_t) = \|x_t - \tilde{x}_t\|_2^2 = (x_{1,t} - \tilde{x}_{1,t})^2 + \dots + (x_{n,t} - \tilde{x}_{n,t})^2$   
 $= (x_{1,t} - x_{t-1} - \widehat{\Pi}[1,] x_{t-1} - \sum_{i=1}^M \Gamma_i[1,] \Delta x_{t-i})^2 + \dots + (x_{n,t} - x_{t-1} - \widehat{\Pi}[n,] x_{t-1} - \sum_{i=1}^M \Gamma_i[n,] \Delta x_{t-i})^2$

4.  $\nabla_{\Gamma_i} l_t = \frac{\partial l_t}{\partial \Gamma_i} = -2 \begin{pmatrix} (x_{1,t} - \tilde{x}_{1,t}) \Delta x_{t-i}^\top \\ (x_{2,t} - \tilde{x}_{2,t}) \Delta x_{t-i}^\top \\ \vdots \\ (x_{n,t} - \tilde{x}_{n,t}) \Delta x_{t-i}^\top \end{pmatrix}$
5.  $G = 2n.(X_{max})^2$  where  $X_{max} = \max_i x_{i,t}$
6.  $D = 2n = 2\max_{X \in \kappa} \|X\|_2 \geq \sup_{X, Y \in \kappa} \|X - Y\|_2$
7.  $\lambda = \frac{1}{n}$

## b EC-VARMA OGD

[Algorithm 8](#) describes the Online Gradient Descent technique applied to the ec-varma model.

---

### Algorithm 8 EC-VARMA online gradient descent

---

- 1: **Input** Horizon  $M$ , first step size  $\eta_1 = \frac{D}{G}$ , nuclear norm ball radius  $\rho$ ,  $\psi^1 \in \kappa$ ,  $T$ , in-sample data  $\{x\}$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Update**  $\eta_t \leftarrow \frac{1}{\sqrt{t}} \eta_1$
  - 4:   **Predict**  $\tilde{x}_t = x_{t-1} + \hat{\Pi} x_{t-1} + \sum_{i=1}^M \Gamma_i \Delta x_{t-i}$
  - 5:   **Observe**  $x_t$  and suffer loss  $l_t(\psi^t)$
  - 6:   **Project and Update:**  $\Gamma_i^{t+1} \leftarrow \Pi_\kappa(\Gamma_i^t - \eta_t \nabla_{\Gamma_i} l_t)$
  - 7:   **Project and Update:**  $\hat{\Pi}^{t+1} \leftarrow \Pi_{\mathcal{B}(*, \rho)}(\hat{\Pi}^t - \eta_t \nabla_{\hat{\Pi}} l_t)$
  - 8: **end for**
- 

## c EC-VARMA ONS

[Algorithm 9](#) describes the Online Newton Step technique applied to the ec-varma model.

---

### Algorithm 9 EC-VARMA online newton step

---

- 1: **Input** Horizon  $M$ , step size  $\eta = \frac{1}{2} \min \{\lambda, \frac{1}{4GD}\}$ ,  $\psi^1 \in \kappa$ , initial list of  $M$  ( $n \times n$ ) matrices  $A_i = \frac{1}{\eta^2 D^2} I_n$ ,  $T$ , in-sample data  $\{x\}$
  - 2: **for**  $t = 1$  **to**  $T$
  - 3:   **Predict**  $\tilde{x}_t = x_{t-1} + \hat{\Pi} x_{t-1} + \sum_{i=1}^M \Gamma_i \Delta x_{t-i}$
  - 4:   **Observe**  $x_t$  and suffer loss  $l_t(\psi^t)$
  - 5:   **Update**  $A_{i+1}^t \leftarrow A_{i+1}^{t-1} + \nabla_{\Gamma_i} l_t \nabla_{\Gamma_i} l_t^\top$
  - 6:   **Update**  $A_1^t \leftarrow A_1^{t-1} + \nabla_{\hat{\Pi}} l_t \nabla_{\hat{\Pi}} l_t^\top$
  - 7:   **Project and Update:**  $\Gamma_i^{t+1} \leftarrow \Pi_\kappa^A(\Gamma_i^t - \frac{1}{\eta} (A_{i+1}^t)^{-1} \nabla_{\Gamma_i} l_t)$
  - 8:   **Project and Update:**  $\hat{\Pi}^{t+1} \leftarrow \Pi_{\mathcal{B}(*, \rho)}(\hat{\Pi}^t - \frac{1}{\eta} (A_1^t)^{-1} \nabla_{\hat{\Pi}} l_t)$
  - 9: **end for**
- 

## 5.4 Experiments

We test these multivariate algorithms on synthetic and real data. In all the experiments, we set  $M = 10$  and  $\rho = 1$ .

**Remark.** *The MLE fails to predict very often, here we choose an appropriate part of the data in which it can make predictions without throwing errors.*

### a Setting 5: Sanity Check 1

Synthetic data generated by the R function arima.sim with correlated noise: noise terms are normally distributed with mean the value of the previous noise term ([Figure 12](#)). The parameters used are in [Table 6](#) and the results are shown in [Figure 13](#).

$\alpha_1$	(0.11, -0.5)
$\alpha_2$	(0.31, -0.2)
$\alpha_3$	(0.6, -0.4)
$\beta_1$	(0.41, -0.39, -0.685, 0.1)
$\beta_2$	(0.1, -0.1, -0.5, 0.2)
$\beta_3$	(0.41, -0.1, -0.5, 0.2)
$\sigma_i$	0.05
$\sigma_2$	0.01
$\sigma_3$	0.15
run	500
training run	5000
MLE window size	252

Table 6: Settings 5

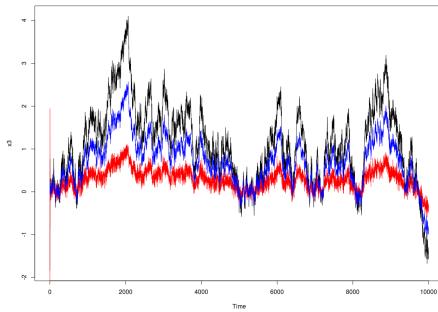


Figure 12: Sythethic data

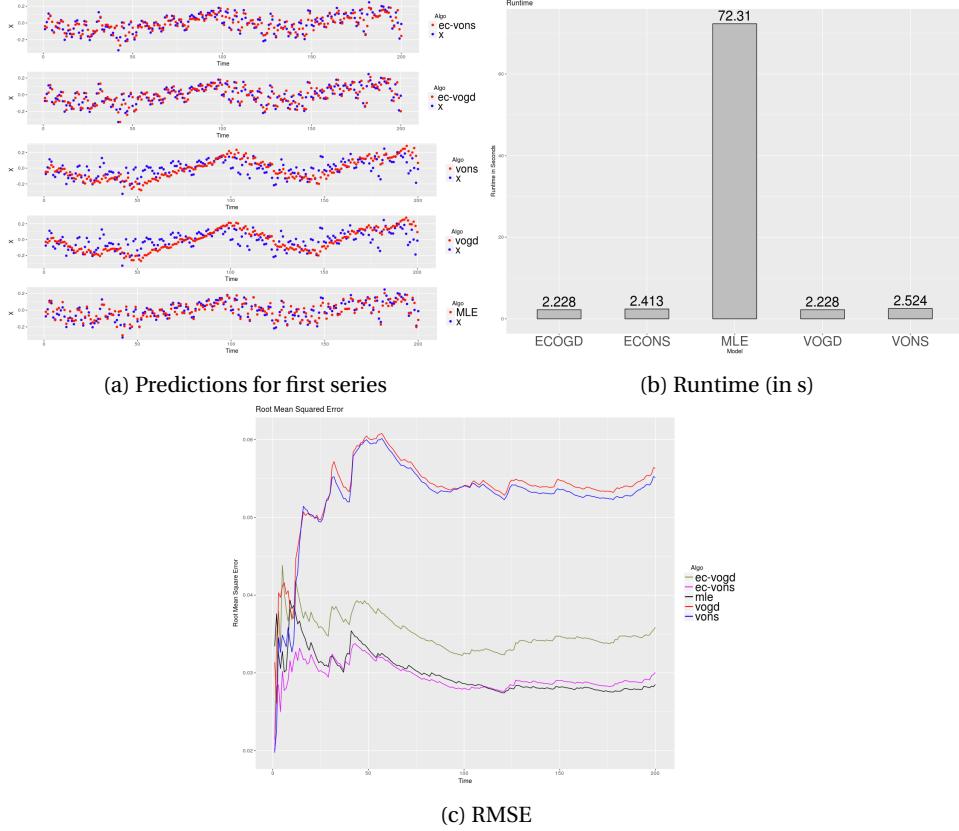


Figure 13: Setting 5 results

### b Setting 6: Sanity Check 2

Synthetic data generated by the R function arima.sim with gaussian noise([Figure 14](#)). The parameters used are in [Table 7](#) and the results are shown in [Figure 15](#).

$\alpha_1$	(0.11, -0.5)
$\alpha_2$	(0.31, -0.2)
$\alpha_3$	(0.6, -0.4)
$\beta_1$	(0.41, -0.39, -0.685, 0.1)
$\beta_2$	(0.1, -0.1, -0.5, 0.2)
$\beta_3$	(0.41, -0.1, -0.5, 0.2)
$\sigma_i$	0.05
run	200
training run	5000
MLE window size	252

Table 7: Settings 6

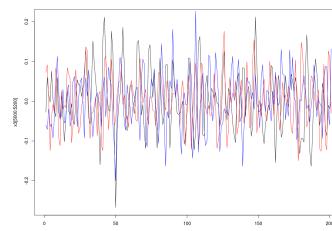


Figure 14: Sythethic data

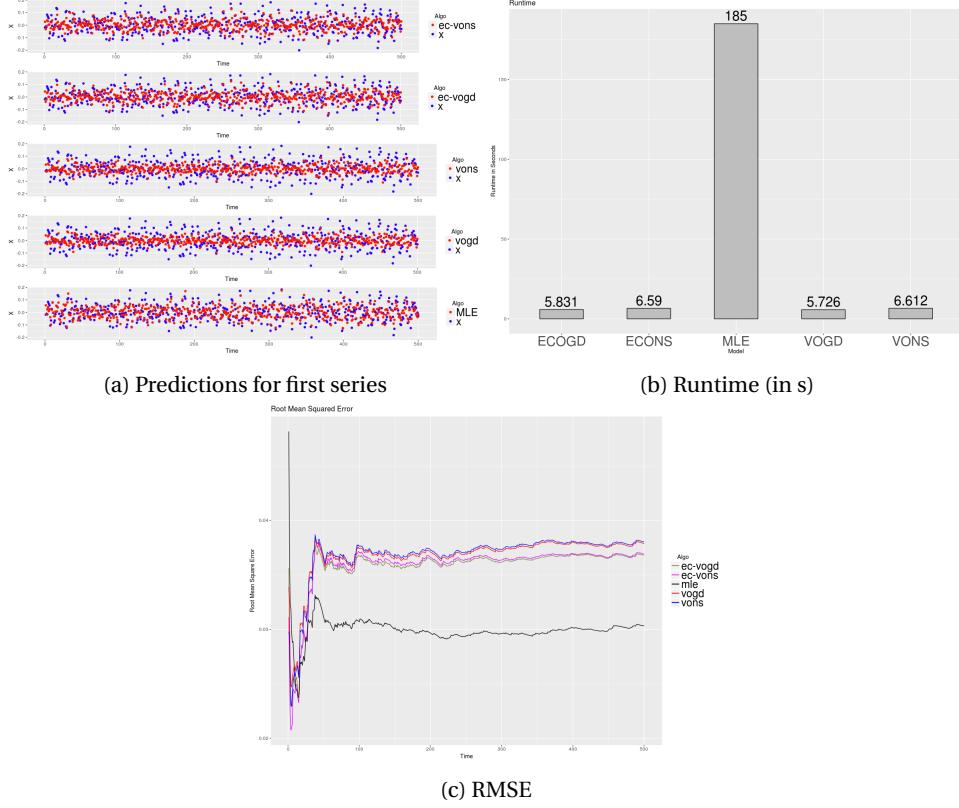


Figure 15: Setting 6 results

### c Setting 7: Real Data, bitcoin volume traded in EUR per 15 min on different markets

We test the univariate and multivariate algorithms on real data: Bitcoin volume traded in EUR per 15 min on different markets (Coinbase, Bitstamp, Btcde), data represented in Figure 16. The parameters used are in Table 8 and the results are shown in Figure 17.

run	500
training run	3000
MLE window size	252

Table 8: Setting 7

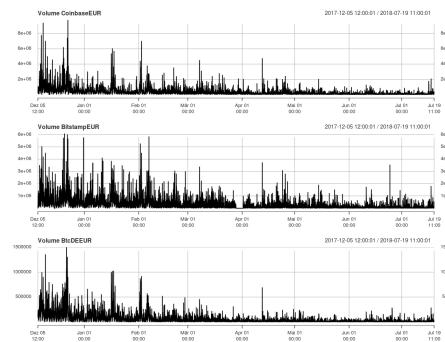


Figure 16: Bitcoin volume traded in EUR per 15 min in \$

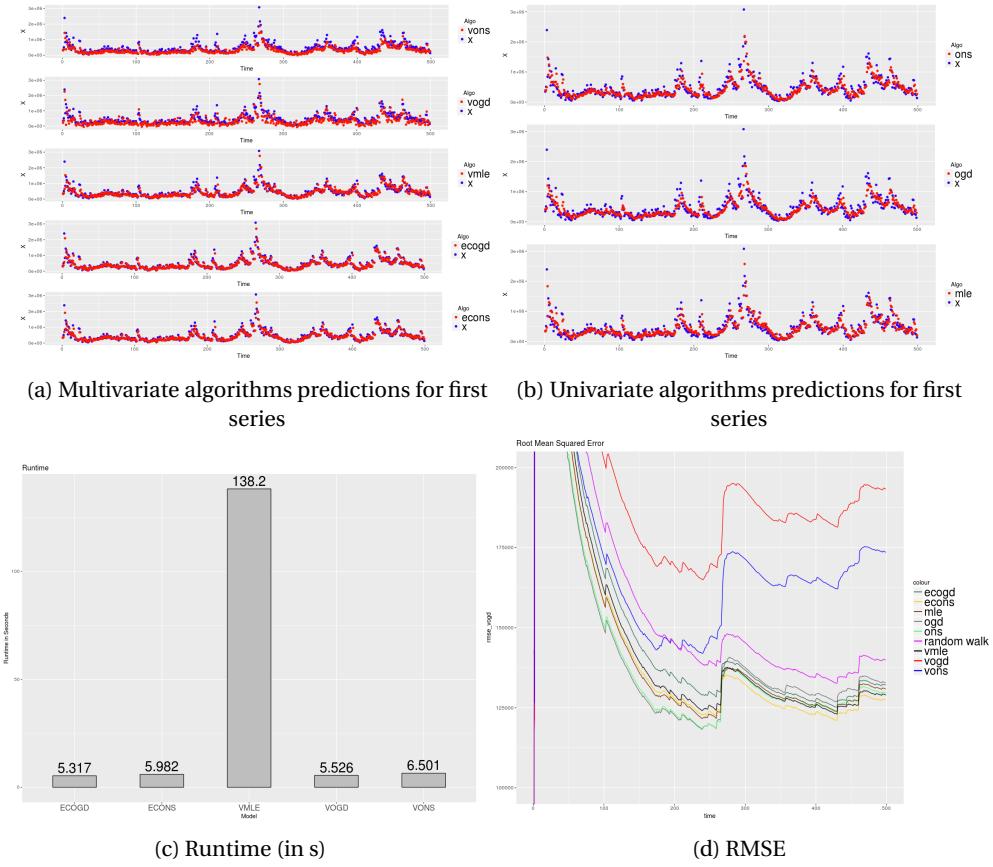


Figure 17: Setting 7 results

## 6 Package presentation and content (more details for usage below)

### 6.1 Utility functions

squared_loss	Computes the squared loss function of two vectors
inv_sm	Computes efficiently $(A + uu^\top)^{-1}$ using the Sherman-Morrison formula
sqrtm_diag	Computes the square root of a matrix using diagonalization
sqrtm_dbit	Computes the square root of a matrix using Denman-Beavers iteration
projection	Computes $\prod_{\kappa}^A(v)$ with $\kappa$ a set with constraint $c$ , $\kappa = \{v \in \mathbb{R}^n :  v_i  \leq c \quad i = 1, \dots, n\}$
mprojection	Computes $\prod_{\kappa}^W(A)$ with $\kappa$ a set with constraint $c$ , $\kappa = \{A \in \mathbb{R}^{n \times n} : \ A\ _{max} \leq c\}$
nnpbprojection	Computes $\prod_{\mathcal{B}(*, \rho)}(A)$ with $\mathcal{B}(*, \rho)$ a nuclear norm ball of radius $\rho$
l1projection	Computes the projection of a vector into an $l_1$ -ball of radius $\rho$ $\{v \in \mathbb{R}^n : \sum_i  v_i  \leq \rho\}$ <b>(details of this function below in Algorithm 10)</b>

### 6.2 OCO functions

arma_ogd	ARMA Online Gradient Descent
arma_ons	ARMA Online Newton Step
varma_ogd	VARMA Online Gradient Descent
varma_ons	VARMA Online Newton Step
ecvarma_ogd	EC-VARMA Online Gradient Descent
ecvarma_ons	EC-VARMA Online Newton Step

---

**Algorithm 10** Projection of a vector onto an l1-ball  $\{v \in \mathbb{R}^n : \sum_i |v_i| \leq \rho\}$ 

---

```
1: Input Vector  $v$ , radius  $\rho$ 
2: Settings  $y = v$ ,  $p = 0$ ,  $z = -\rho$ 
3: While  $y \neq \emptyset$ 
4:   Choose pivot  $k \leftarrow$  Random element from  $y$ 
5:   Set  $\mathcal{H} \leftarrow \{t \in v : t > k\}$ ,  $\mathcal{L} \leftarrow \{t \in v : t < k\}$ ,  $\mathcal{M} \leftarrow |\{t \in v : t = k\}|$ 
6:   if  $(s + Mk + \sum_{t \in \mathcal{H}} t) / (p + M + |\mathcal{H}|) < k$ 
7:      $s \leftarrow s + Mk + \sum_{t \in \mathcal{H}} t$ 
8:      $p \leftarrow p + M + |\mathcal{H}|$ 
9:      $y \leftarrow \mathcal{L}$ 
10:    else  $y \leftarrow \mathcal{H}$ 
11:   end while
12:    $\theta = s/p$ 
13:   for  $i = 1$  to  $n$ 
14:     Set  $\text{proj}_i = \max(v_i - \theta, 0)$ 
15:   end for
16:   Return  $\text{proj}$ 
```

---

## 7 Conclusion

We considered a new approach of time series prediction based on online learning, we presented the online convex optimization model along with two algorithms and their applications on univariate and multivariate ARMA models. We realized some experiences to test the reliability of these algorithms. The main results of this report can be resumed as follows:

1. There is an incomparable difference of running time between the Maximum-Likelihood estimation and the OCO algorithms.
2. The Maximum-likelihood estimation makes better predictions on synthetic data (provided a growing window)
3. On real data, the OCO algorithms make better performance since this implies fast adaptation from noisy data and regime shifts.

### 7.1 Future research

Online convex optimization can be studied for many other applications such as the ARCH model and universal portfolio selection.

## References

- [CXF17] Aveleen Bijral Christopher Xie and Juan Lavista Ferres. Online Prediction Methods for Nonstationary Time Series. 2017.
- [Haz16] Elad Hazan. *Introduction to Online Convex Optimization*. 2016.
- [OAS13] Shie Mannor Oren Anava, Elad Hazan and Ohad Shamir. Online Learning for Time Series Prediction. *JMLR: Workshop and Conference Proceeding*, 2013, 2013.