

# A Forecasting Procedure involving Bayesian Exponential Smoothing and LASSO

Avner Abrami

November 15, 2014

## Contents

<b>1</b>	<b>Settings and Notations</b>	<b>2</b>
1.1	Assumption 1: High Dimensional Data Analysis . . . . .	2
1.2	Assumption 2: $Y$ Distribution . . . . .	2
<b>2</b>	<b>Model Formulation</b>	<b>2</b>
<b>3</b>	<b>Exponential Smoothing in brief</b>	<b>3</b>
3.1	General Decomposition of Time Series . . . . .	3
3.2	Short history of Exponential Smoothing . . . . .	3
3.3	State Space Models . . . . .	3
<b>4</b>	<b>Holt-Winters Model</b>	<b>4</b>
4.1	Transition Equations . . . . .	4
4.2	Statistical Model . . . . .	4
<b>5</b>	<b>Holt-Winters Model without seasonality : Holt's linear Model</b>	<b>5</b>
5.1	Model Formulation . . . . .	5
5.2	Sum up . . . . .	5
<b>6</b>	<b>Why is our new model compatible with Holt's Linear Model?</b>	<b>6</b>
<b>7</b>	<b>Optimization Program and Algorithm</b>	<b>7</b>
7.1	Optimization Program . . . . .	7
7.2	Variable Projection : General Formulation . . . . .	7
7.2.1	Second step of the algorithm : Determining $\Omega$ . . . . .	7
7.2.2	Sumarizing the Algorithms . . . . .	8
<b>8</b>	<b>Forecasting Procedure</b>	<b>8</b>
	Bibliography	9

# 1 Settings and Notations

One has  $p$  different time series  $Y^1, Y^2, \dots, Y^p$  that we observe for  $t \in [0, n]$ .

They constitute the columns of  $Y$

$$Y = \begin{bmatrix} Y_1^1 & Y_1^2 & \cdots & Y_1^p \\ Y_2^1 & Y_2^2 & \cdots & Y_2^p \\ \vdots & \vdots & \ddots & \vdots \\ Y_n^1 & Y_n^2 & \cdots & Y_n^p \end{bmatrix}$$

## 1.1 Assumption: High Dimensional Data Analysis

$$p \gg n$$

# 2 Model Formulation

$$\forall i \in [1, p], \forall t \in [1, n], Y_{t+1}^i = a_t^i + b_t^i - \sum_{j \neq i} \omega_j^i Y_{t+1}^j + \epsilon_{t+1}^i$$

Where:

$$\Omega = \begin{bmatrix} 1 & \omega_1^1 & \cdots & \omega_1^p \\ \omega_2^1 & 1 & \cdots & \omega_2^p \\ \vdots & \vdots & \ddots & \vdots \\ \omega_p^1 & \omega_p^2 & \cdots & 1 \end{bmatrix}$$

Where the  $(a^i, b^i)^1$  are the smoothing coefficients associated with the  $i^{th}$  time series.  
It is time to see how we construct this model. But first a few words on exponential smoothing.

# 3 Exponential Smoothing in brief

This section has been extracted from *Forecasting with Exponential Smoothing* by Rob J. Hyndlan, Anne B. Koehler, J. Jeith Ord, Ralph D. Snyder.

## 3.1 General Decomposition of Time Series

It is common in business and economics to think of a time series as a combination of various components such as the trend (T), cycle (C), seasonal (S), and irregular or error (E) components. These can be defined as follows:

- **Trend (T)**: The long-term direction of the series
- **Seasonal (S)**: A pattern that repeats with a known periodicity (e.g., 12 months per year, or 7 days per week)
- **Cycle (C)**: A pattern that repeats with some regularity but with unknown and changing periodicity (e.g., a business cycle)
- **Irregular or error (E)**: The unpredictable component of the series

---

<sup>1</sup>to be precisely defined in the next paragraphs

A *purely additive model* can be expressed as :

$$y = T + S + E \quad (1)$$

where the three components are added together to form the observed series.

A *purely multiplicative model* is written as

$$y = T * S * E \quad (2)$$

where the data are formed as the product of the three components

### 3.2 Short history of Exponential Smoothing

Historically, exponential smoothing describes a class of forecasting methods. In fact, some of the most successful forecasting methods are based on the concept of exponential smoothing. There are a variety of methods that fall into the exponential smoothing family, each having the property that forecasts are weighted combinations of past observations, with recent observations given relatively more weight than older observations. The name 'exponential smoothing' reflects the fact that the weights decrease exponentially as the observations get older.

The idea seems to have originated with Robert G. Brown in about 1944 while he was working for the US Navy as an Operations Research analyst. Independently, Charles Holt was also working on an exponential smoothing method for the US Office of Naval Research (ONR)

### 3.3 State Space Models

State space models allow considerable flexibility in the specification of the parametric structure. In this book, we will use the innovations formulation of the model (e.g., Anderson and Moore 1979; Aoki 1987; Hannan and Deistler 1988). Let  $y_t$  denote the observation at time  $t$ , and let  $x_t$  be a 'state vector' containing unobserved components that describe the level, trend and seasonality of the series. Then a linear innovations state space model can be written as:

$$y_t = w'x_{t-1} + \epsilon_t \quad (3)$$

It is known as the measurement (or observation) equation; it describes the relationship between the unobserved states  $x_{t-1}$  and the observation  $y_t$ .

$$x_t = Fx_{t-1} + g\epsilon_t \quad (4)$$

where  $\epsilon_t$  is a white noise series and  $F$ ,  $g$  and  $w$  are coefficients.

It is known as the transition (or state) equation; it describes the evolution of the states over time.

## 4 Holt-Winters Model

In this section, we will be talking about any of our time series  $Y^i$  but we will omit its index  $i$  for more convenience.

Reminder:  $Y^i$  are the columns of  $Y$ .

$$Y = \begin{bmatrix} Y_1^1 & Y_1^2 & \dots & Y_1^p \\ Y_2^1 & Y_2^2 & \dots & Y_2^p \\ \vdots & \vdots & \ddots & \vdots \\ Y_n^1 & Y_n^2 & \dots & Y_n^p \end{bmatrix}$$

#### 4.1 Transition Equations

The additive Holt-Winters decomposition can be defined through the following transition equations which define the Holt-Winters filter :

$$\begin{array}{ll} \textbf{Mean Equation} & a_t = \alpha(y_t - c_{t-p}) + (1 - \alpha)(a_{t-1} + b_{t-1}) \\ \textbf{Trend equation} & b_t = \beta(a_t - a_{t-1}) + (1 - \beta)b_{t-1} \\ \textbf{Seasonality equation} & c_i = \gamma(y_t - a_{t-1} - b_{t-1}) + (1 - \gamma)c_{t-p} \end{array}$$

where  $\mathbf{y} = (y_1, \dots, y_n)'$  is the vector of observed data,  $\theta = (\alpha, \beta, \gamma)'$  is the smoothing parameters vector, and  $p$  is the seasonal period.

#### 4.2 Statistical Model

A stochastic component is added to the above filter in order to obtain a statistical model :

$$Y_t = (a_{t-1} + b_{t-1} + c_{t-p}) + \epsilon_t$$

$\mu(t) = a_{t-1} + b_{t-1} + c_{t-p}$  is the one-step forecast made at time  $t - 1$  by the Holt-Winters filter, and  $\epsilon_t$  is the forecast error at time  $t$ .

The combination of the above equation and transitions equations give a state representation of the additive Holt-Winters model that will be modified in the next section to remove the seasonality.

### 5 Holt-Winters Model without seasonality : Holt's linear Model

#### 5.1 Model Formulation

We are interested in developing the Holt-Winters method for data without the seasonality component. The corresponding filter is obtained from the initial one through the suppression of the seasonality equation, the  $\gamma$  smoothing parameter, the period  $p$ , and the  $(c_0, \dots, c_{p-1})$  factors.

The combination of the obtained filter with the revised equation,

$$Y_t = a_{t-1} + b_{t-1} + \epsilon_t \tag{5}$$

gives the following state representation for the additive non-seasonal Holt-Winters model :

$$Y = (Y_1, Y_2, \dots, Y_n)^T = M\psi + L\varepsilon \tag{6}$$

where :

$$M = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ \vdots & \vdots \\ 1 & n \end{pmatrix} \in \mathcal{M}_{n \times 2}(\mathbb{R}) \text{ a design matrix}$$

$\psi = (a_0, b_0)' \in \mathbb{R}^2$  the initial conditions vector

$$L = \begin{pmatrix} 1 & 0 & \dots & \dots & 0 \\ & \ddots & \ddots & & \vdots \\ & l_{i,j|i>j} & & \ddots & \vdots \\ = (1 + (i-j)\beta)\alpha & & & \ddots & 0 \\ & & & & 1 \end{pmatrix} \in \mathcal{M}_{n \times n}(\mathbb{R}) \text{ the smoothing matrix}$$

$\varepsilon = (\epsilon_1, \dots, \epsilon_n)' \in \mathbb{R}^n$  the forecast errors vector

## 5.2 Sum up

$$Y_t^i = a_{t-1}^i + b_{t-1}^i + \epsilon_t^i$$

$$Y^i = (Y_1^i, Y_2^i, \dots, Y_n^i)^T = M^i \psi^i + L^i \varepsilon^i \quad (7)$$

## 6 Why is our new model compatible with Holt's Linear Model?

Our model is:

$$\forall i \in [1, p], \forall t \in [1, n], Y_{t+1}^i = a_t^i + b_t^i - \sum_{j \neq i} \omega_j^i Y_{t+1}^j + \epsilon_{t+1}^i$$

Which can easily be rewritten as:

$$Y_{t+1}^i + \sum_{j \neq i} \omega_j^i Y_{t+1}^j = \sum_j \omega_j^i Y_{t+1}^j = (\omega^i)^T Y_{t+1} = a_t^i + b_t^i + \epsilon_{t+1}^i$$

$$\omega_i^i = 1$$

According to section 5:

$$(\omega^i)^T Y_{t+1} = M \psi^i + L^i \epsilon^i$$

But more generally:

$$Y\Omega = \mathcal{M}\Psi + L\epsilon$$

$$Y = \begin{bmatrix} Y_1^1 & Y_1^2 & \dots & Y_1^p \\ Y_2^1 & Y_2^2 & \dots & Y_2^p \\ \vdots & & \ddots & \vdots \\ Y_n^1 & Y_n^2 & \dots & Y_n^p \end{bmatrix}$$

$$\begin{aligned}\Psi &= \begin{bmatrix} \psi^1 & | & \psi^2 & | & \dots & | & \psi^p \end{bmatrix} \\ \mathcal{M} &= \begin{bmatrix} M & 0 & \dots & 0 \\ 0 & M & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & M \end{bmatrix} \\ L &= \begin{bmatrix} L^1 & 0 & \dots & 0 \\ 0 & L^2 & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & L^p \end{bmatrix} \\ \epsilon &= \begin{bmatrix} \epsilon_1^1 & \epsilon_2^1 & \dots & \epsilon_p^1 \\ \epsilon_1^2 & \epsilon_2^2 & \dots & \epsilon_p^2 \\ \vdots & & \ddots & \vdots \\ \epsilon_1^n & \epsilon_2^n & \dots & \epsilon_p^n \end{bmatrix}\end{aligned}$$

## 7 Optimization Program and Algorithm

### 7.1 Optimization Program

In order to assess our statistical model, we would like to solve:

$$\min_{\theta \in [0,1]} \{ \min_{\omega} \{ \|L^{-1}Y\Omega - L^{-1}\mathcal{M}\Psi\|_F + \lambda \|\Omega\| \} \}$$

Where:

$$\theta = \begin{bmatrix} (\alpha^1, \beta^1) & | & (\alpha^2, \beta^2) & | & \dots & | & (\alpha^p, \beta^p) \end{bmatrix}$$

If  $Y_\theta = (L^{-1}Y\Omega)$ ,  $W_\theta = L^{-1}\mathcal{M}\Psi$ , the program becomes:

$$\min_{\theta \in [0,1]} \{ \min_{\Omega} \{ \|Y_\theta\Omega - W_\theta\|_F + \lambda \|\Omega\| \} \}$$

Which is a problem of the form

$$\min_{(\theta \in [0,1]^2, \Omega)} g(\theta, \Omega) \tag{8}$$

It is gonna be a two step algorithm. The first one (apparently the hardest?) is the  $\theta$  optimization. Indeed  $Y_\theta$ ,  $W_\theta$  are polynomial fractions of  $\theta$  (cf  $L^{-1}$ ).

But then obtaining  $\Omega$  via a Lasso based algorithm would be a quite easy task.

### 7.2 Variable Projection : General Formulation

This section has been extracted from *Estimating Nuisance Parameters in Inverse Problems* by Aleksandr Y. Aravkin and Tristan van Leeuwen.

We consider problems of the form (8). As explained previously, for any given  $\theta$  one can easily find (via a Lasso based Algorithm):

$$\hat{\Omega}(\theta) \in \operatorname{argmin}_{\Omega} g(\theta, \Omega)$$

This condition can be relaxed, and  $\theta$  can be considered a local minimum. Rather than working to solve (8), we can instead focus on the reduced objective:

$$\hat{g}(\theta) = g(\theta, \hat{\Omega}(\theta))$$

This approach is justified by a theorem (Theorem 2.1 page 3 in Prof. Aravkin's paper)

### 7.3 First

We are in a simple case as  $\theta \in [0, 1]^2$  is in a closed and bounded set.

To solve our problem, Prof Aravkin suggests a modified projected gradient method in his paper:

$$\theta^{k+1} = P_{[0,1]^2}[\theta^k - \gamma \nabla_x \hat{g}(\theta^k)] \quad (9)$$

Where P is the projection on  $[0, 1]^2$ .

#### 7.3.1 Second step of the algorithm : Determining $\Omega$

If  $\theta$  is fixed, solving the program for  $\Omega$  is the same as resolving a Lasso Problem.

#### 7.3.2 Sumarizing the Algorithms

**Initialization:**  $(\Omega)_0 = I_p$

Compute  $\theta_0 = \text{ProjectionVariable}$

**while** not converged **do**:

$(\Omega)_{k+1} = \text{Lasso}$

$\theta_{k+1} = \text{ProjectionVariable}$

**end while**

**Output :**  $\Omega, \theta$

## 8 Forecasting Procedure

Let  $Y_{obs}$  be the  $n \times 1$  vector of observations and  $Y_{new}$  the  $h \times 1$  vector of predictions. We suppose that the joint vector  $Y = (Y_{obs}, Y_{new})'$  still follows the model described in 5.1 with the same assumptions :

$$Y\Omega = \begin{pmatrix} Y_{obs}\Omega \\ Y_{new}\Omega \end{pmatrix} = \begin{pmatrix} \mathcal{M}_1 \\ \mathcal{M}_2 \end{pmatrix} \psi + \begin{pmatrix} L_1 & 0 \\ L_{21} & L_2 \end{pmatrix} \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix}$$

As a consequence:

$$\forall i \in [1, n], Y_{new}\Omega = \mathcal{M}_2 \tilde{\psi} + L_{21} L_1^{-1} (Y_{obs} - \mathcal{M}_1 \tilde{\psi}) + L_2(\varepsilon_2)$$

$$\forall i \in [1, n], E[Y_{\text{new}}] = (\mathcal{M}_2 \tilde{\psi} + L_{21} L_1^{-1} (Y_{\text{obs}} - \mathcal{M}_1 \tilde{\psi})) \Omega^{-1}$$

Indeed, by definition  $E[\epsilon] = 0$ .



## References

- [1] HYNDAMN, KOEHLER, ORD, SNYDER (2008) *Forecasting with Exponential Smoothing*
- [2] BERMUDEZ, JD., SEGURA, JV., VERCHER, E. (2010) *Bayesian forecasting with the Holt–Winters model*
- [3] Aleksandr Y. Aravkin and Tristan van Leeuwen (2012) *Estimating Nuisance Parameters in Inverse Problems*