

Implementation of probabilistic forecasts in a GPU for big data problems.

Juan R. Trapero^{a,*}, Enrique Holgado de Frutos^a, Francisco Ramos^b

^a*University of Castilla-La Mancha, Department of Business Administration, 13071 Ciudad Real, Spain*

^b*University of Castilla-La Mancha, Department of Electrical, Electronics, Control and Communications Engineering, 13071 Ciudad Real, Spain*

Abstract

High Performance Computing via General Purpose Graphical Processing Unit (GPU) is a potential instrument to speed up computational times. In a world where big data is becoming a revolution, GPU could play an important role. This work intends to analyze the performance of GPU by implementing the calculation of probabilistic forecasts based on single exponential smoothing in conjunction with simulated predictive distributions. Essentially, supply chain companies must deal with a high number of forecasts at SKU level. In this context, reducing the computational times can be a source of a competitive advantage. Since the forecasts are usually made independently between SKUs, this problem can be easily parallelized and GPU computing can exploit such parallelization. To the best of authors knowledge, this is the first time GPU is applied to a supply chain demand forecasting context. Firstly, we will show how to adapt the programming of probabilistic forecasts in a parallel fashion. Then, real data coming from a manufacturer company will be used to illustrate the differences between GPU and traditional CPU computing. The results show that GPU can significantly increase the computational speedup ratio more than 30 times with respect to traditional CPU computing.

Keywords: forecasting, GPU, big data, supply chain management

*Corresponding author.

Email addresses: juanramon.trapero@uclm.es (Juan R. Trapero),
Enrique.Holgado@alu.uclm.es (Enrique Holgado de Frutos),
francisco.ramos@uclm.es (Francisco Ramos)

1. Introduction

The world digitalization is generating a massive amount of data that can be obtained from different sources [1]. In a supply chain context those data sources are generated from enterprise resource planning (ERP) systems, distributed manufacturing environments, orders and shipment logistics, social media feeds, customer buying patterns, product lifecycle operations, and technology-driven data sources such as global positioning systems (GPS), radio frequency based identification (RFID) tracking, mobile devices and surveillance videos, among others [2, 3]. Dealing with such an amount and variety of information has challenged the traditional data analysis methods, [1, 4]. Therefore, organisations have to cope with big datasets characterized by 4Vs: large volume, velocity, variety, and value. Currently, the “Big Data” concept is starting to be defined in terms of the 5Vs model, which added the “Veracity” dimension (related to proper data governance and privacy concerns), [5, 6]. A recent International Data Corporation (IDC) forecast suggests that the Big Data technology will grow at a rate of 26.4% compound annual growth, reaching 41.5 billion through 2018, [2]. Essentially, the bigger the data, the more challenging it becomes to manage and analyse in order to deliver useful business insights. Manyika et al. [7] indicate that 300 billion dollars of potential annual value can be generated in US healthcare if organisations or governments can capture big data’s value.

One of the main improvements that big data may bring is more accurate forecasts, which implies improving customer service levels, while lowering inventory costs, waste, and working capital [8]. Focusing on how the forecasts can be improved based on such big data, several articles indicate that predictive analytics is the key to reap the benefits of big data. It seems that predictive analytics comprises all the new techniques able to deal with big data that traditional forecasting methods cannot, [4]. Regarding the differences between predictive analytics and traditional forecasting methods, predictive analytics neglects statistical significance (given that we move from a small sample framework typical from the traditional statistical analysis to a big sample, even the entire population). Predictive analytics also acknowledges the importance of computational efficiency, [4]. Waller and Fawcett [9] classifies predictive analytics as a subset of data science and indicate that whereas forecasting is about predicting the future, predictive analytics adds

36 questions regarding what would have happened in the past, given different
37 conditions. Although most of the research is focused on how Big Data Pre-
38 dictive Analytics (BDPA) improves corporate financial performance, it is also
39 interesting to note that the impact of BDPA on other aspects of the triple
40 botton-line (environmental and social sustainability) remains still as an open
41 question [10].

42 Despite the acclaimed benefits of BDPA, there is a lack of empirical works
43 that define at operational level how the required forecasts can be improved,
44 or at least, how the traditional procedures should be changed to adapt them
45 to the predictive analytics context. Nguyen et al. [11] carried out a literature
46 review about big data analytics in supply chain management and concluded
47 that only 3 papers out of 88 are centered on demand forecasting.

48 Nikolopoulos and Petropoulos [12] intend to shed some light into the
49 problem of forecasting big data focusing on the optimization stage of the
50 forecasting process. In that reference, the authors analyze the influence of
51 the complex optimization routines on the forecasting performance. Their
52 results show that less complex routines more oriented to face big data prob-
53 lems do not reduce significantly the forecasting accuracy, although it does
54 improve the computation speed. Interestingly, they also note that the fore-
55 casting discipline has not developed too much progress in the incorporation of
56 the Information and Communication Technologies for improving forecasting
57 performance. The necessity of a compromise between computational speed
58 and forecasting accuracy is also pointed out in a retail industry context in
59 [13]. For instance, Seaman in [13] quantifies the forecasting problem for Wal-
60 mart with upwards of a trillion forecasts being needed, calculating over 10
61 millions forecasts a second. In this sense, Seaman [13] indicates the “par-
62 allelizability” of the forecasting models used as a key factor to improve the
63 speed of forecasting, as well as, the underlying computational infrastructure
64 to support the parallelization. In the same sense, Zhong et al. [14] indicate
65 that traditional serial algorithms, models, and mechanisms are inefficient for
66 Big Data and new data parallelism and advanced approaches are a current
67 challenge.

68 Demand planning and inventory management requires probabilistic fore-
69 casts to optimize supply chain decisions. For instance, safety stocks and
70 reorder points in order up to level replenishment policies are based on those
71 probabilistic forecasts. Here, probabilistic forecasts refers to forecasts of a
72 certain quantile, i.e., a forecast of the demand mean plus a forecast of the
73 demand variability. In general terms, the variability is a measure of demand

74 uncertainty. Traditionally, these estimations are done via serial computing.
75 Although probabilistic forecasts can be based on theoretical models [15],
76 where it is assumed that the underlying demand process is correctly captured
77 by the forecasting model, the uncertainty in the parameter estimates are usu-
78 ally neglected [16]. A possible solution to overcome that limitation is to use
79 simulated prediction distributions [17, p. 77]. In fact, for some complex
80 models, simulation is the only method available to compute predictive dis-
81 tributions [17, p. 79]. However, the main inconvenient is the computational
82 burden.

83 This work aims at exploring the use of GPU to improve the computa-
84 tional speed of calculating probabilistic forecasts via simulated prediction
85 distributions. Essentially, we show how the forecasting method should be
86 implemented in a parallel fashion to take advantage of the potential benefits
87 of GPU. We particularize our results for a well-known method as single expo-
88 nential smoothing for point forecasting and simulated prediction distributions
89 through Monte Carlo experiments for variability forecasts. These techniques
90 will be implemented in MATLAB. In this work, we intend to implement
91 the algorithms with minimum changes with respect to traditional MATLAB
92 programming, i.e., without programming in CUDA (Compute Unified De-
93 vice Architecture). Following that philosophy, we will use some “built-in”
94 functions already available in MATLAB that supports GPU computing, min-
95 imizing the programming effort. Real data coming from a manufacturer will
96 be used to illustrate the methodology.

97 This article is organised as follows: Section 2 explains the main differences
98 of a GPU with respect to traditional CPU. Section 3 describes the exponen-
99 tial smoothing forecasting method expressed as a transfer function and how
100 the smoothing parameter is determined via the grid and search optimization
101 routine. Additionally, in the same section, the simulated probabilistic fore-
102 cast methodology will be explained. Section 4 shows the evaluation metrics
103 in terms of forecasting accuracy and computational speed. Section 5 is de-
104 voted to test those algorithms with real data coming from a manufacturer
105 shipments data. Finally, section 6 summarizes the main conclusions.

106 **2. What can GPUs do for probabilistic demand forecasting?**

107 GPUs were designed for graphical representation purposes characterized
108 by tedious and repetitive, although not complicated, calculations with thou-
109 sands of millions of them per second. The idea was easy: we divide the

110 desired image processing in independent units of work (e.g. dividing the
111 whole image in small chunks), separately calculate each of these units of
112 work, and then obtain the final output by composing the results of all units
113 of work.

114 The introduction of CUDA, which is a general purpose parallel com-
115 puting architecture that is more user-friendly and mature, has made GPU
116 computing to evolve towards the General Purpose GPU. In general terms,
117 GPUs are recommended when the following criteria are met: i) computa-
118 tionally extensive: the time spent in calculations should highly exceed the
119 time spent on transferring data to and from GPU memory, as memory access
120 becomes a bottleneck when using GPUs, i.e., data must be sent from CPU
121 to GPU before calculations and returned afterwards to main memory. As
122 the communication between CPU and GPU is typically performed through
123 an AGP or PCI Express bus, memory access is slower than with the usual
124 CPU architecture; ii) massively parallel: the calculations must be divisible in
125 hundreds/thousands of independent units of work. These calculations must
126 be the same for every chunk of data that is to be processed in a GPU core
127 at the same time.

128 Demand forecasting fits both criteria perfectly and, hence, can greatly
129 benefit from the use of GPUS. On the one hand, forecasting algorithms can
130 be individually applied to each SKU, as they are (mainly) independent to
131 each other, what makes them easily parallelizable. On the other hand, each
132 algorithm relies on a Monte Carlo experiments with (many) replications,
133 being computationally demanding. Each of these replications could as well
134 be performed in a different core. Therefore, parallel computing is expected to
135 significantly reduce the computation time, what might become a competitive
136 advantage in a big data driven world.

137 However, parallel programming is a different programming paradigm from
138 traditional structured programming or object-oriented programming used for
139 most applications. Some libraries, such as CUDA¹ (Computer Unified De-
140 vice Architecture), were developed to help in the process of shifting this
141 paradigm for a variety of applications, not only graphics processing. Thank-
142 fully, some mathematical frameworks, such as MATLAB, provide their own
143 built-in functions to leverage the change of paradigm with slight changes
144 respect to traditional programming.

¹For the interested reader see <https://developer.nvidia.com/cuda-toolkit>

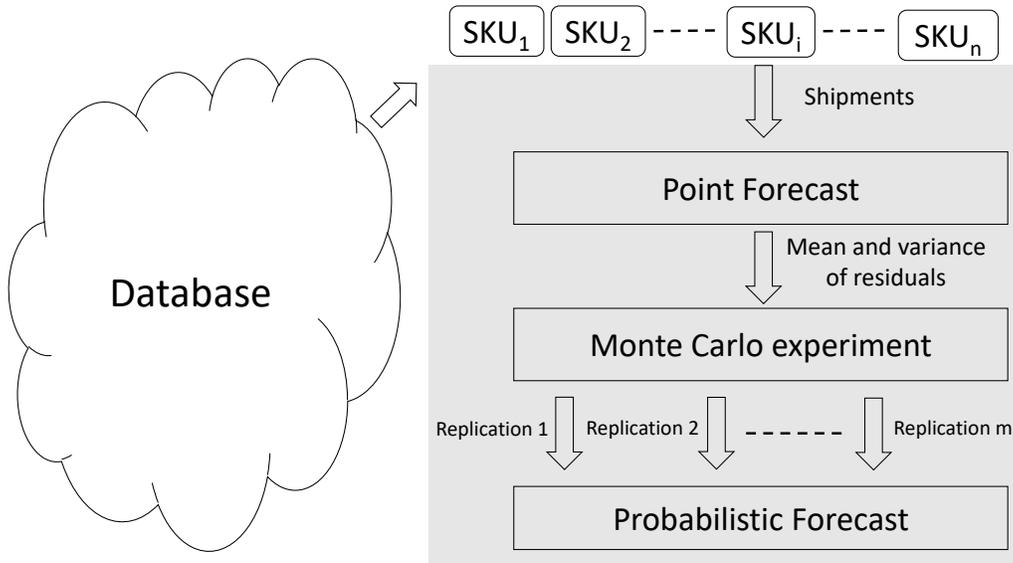


Figure 1: Scheme of probabilistic forecasts calculation.

145 Currently, GPU is becoming more popular in large-scale, data mining
 146 [18], dynamic programming [19, 20], graphs problems [21], discrete location
 147 theory [22], genetic algorithms [23], although its application to predictive
 148 analytics problems is still scarce [24].

149 3. Probabilistic Forecasts

150 In this section, we will explain how to compute the probabilistic forecasts
 151 following the scheme showed in Figure 1, which is divided in three stages.
 152 First, the shipments information for n SKUs is obtained from the company
 153 database, and for each SKU, the point forecast is computed. Second, with
 154 the point forecasts, the mean and variance of residuals are computed and
 155 fed into the Monte Carlo experiment which generates m replications. Third,
 156 the calculation of the probability forecasts is obtained based on those repli-
 157 cations. Below, we describe in detail each stage.

158 *3.1. Point forecasting*

159 *3.1.1. Single exponential smoothing (SES)*

160 First, we need to calculate the point forecast using SKU shipments as an
 161 input, where the point forecasting algorithm yields a measure of the central
 162 tendency of the forecast density function. Among the possible forecasting
 163 methods, Single Exponential Smoothing (SES) has been chosen for two reasons.
 164 Firstly, SES is widely used in business applications [25, 26]. Secondly,
 165 it is consistent with common practice, in which a company may not be using
 166 the optimal forecasting model due to either a lack of expertise or the availability
 167 of only a limited set of forecasting models in their software, [27]. Such
 168 a forecasting method can be formulated as, [28]

$$F_{t+h}^h = \alpha y_t + (1 - \alpha)F_t^h, \quad (1)$$

169 where $0 \leq \alpha \leq 1$ is the smoothing constant, y_t is the actual value and F_t^h
 170 is the h-step ahead forecast, both at time t . h is the forecasting horizon. Given
 171 the recursive nature of exponential smoothing, it is necessary to initialize the
 172 algorithm. We initialize with $F_1^h = y_1$, which is a simple approach because
 173 when the data sample is big, as we expect in a big data framework, the
 174 influence of the initialization is minor [28].

175 In this work, we will use the transfer function form of SES since it is
 176 easier to allow its parallelization. Recall that SES expressed as a transfer
 177 function [29] is:

$$TF_{SES} = Z \left(\frac{F_{t+h}^h}{y_t} \right) = \frac{\alpha}{1 - (1 - \alpha)Z^{-1}} \quad (2)$$

178 where Z is the Z-transform, such as $Z^{-1}Y_t = Y_{t-1}$.

179 *3.1.2. Optimization. Grid and search vs. fmincon*

180 The α value in (1) is optimized by using a grid and search algorithm that
 181 is also known as parameter sweep or exhaustive search [12]. This algorithm
 182 starts a search for values of α ranging from 0 to 1. The parameter that
 183 controls the routine is the number of steps (n_a) that we divide the whole
 184 interval, where the updating step is $s = \frac{1-0}{n_a}$. For each value of α , the Mean
 185 Squared Error (MSE) is calculated as follows:

$$MSE = \frac{1}{N} \sum_{t=1}^N (y_t - F_t^1)^2 \quad (3)$$

186 where N is the sample size used for training. The chosen α is the one that
 187 minimizes one-step ahead MSE over the training set. Authors in [12] show
 188 that this simple method is able to provide a similar forecast accuracy than
 189 other sophisticated alternatives with a considerable reduction of computa-
 190 tional time.

191 The GPU optimizes the exponential smoothing parameter by applying
 192 grid and search routine previously introduced. To compute the MSE, we
 193 need to apply SES for each value of α . This can be done with the “built-
 194 in” MATLAB command **filter**, [30], using the SES on the hold-in sample
 195 expressed as a transfer function, see (2). As a result of this stage, the optimal
 196 α value is obtained.

197 Note that this optimization can be done either serially or in parallel.
 198 Procedure 1 describes the algorithm to obtain the exponential smoothing
 199 parameter in serial computing. In turn, Procedure 2 shows how to obtain
 200 the same value in a parallel manner. Essentially, by computing in parallel we
 201 have removed the outer loop with regards to the serial version. Note that,
 202 such an outer loop is bigger than the inner loop because $n \gg n_a$

Procedure 1 Grid and search optimization (serial)

- 1: **for** $i = 1$ **to** n **do**
 - 2: Select sales of SKU i ($y_{i,t}$)
 - 3: **for** $j = 1$ **to** n_a **do**
 - 4: Compute point forecasts of $y_{i,t}$ with SES for a determined α_{ij}
 - 5: Compute MSE_j over the training set
 - 6: **end for**
 - 7: Determine $\alpha_i = \alpha_j$ such as α_j minimizes MSE_j over the train set
 - 8: **end for**
-

Procedure 2 Grid and search optimization (parallel)

- 1: **for** $j = 1$ **to** n_a **do**
 - 2: Compute point forecasts with SES for a determined α_j for all SKUs in parallel
 - 3: Compute MSE_j over the training set for all SKUs in parallel
 - 4: **end for**
 - 5: Determine $\alpha_i = \alpha_j$ such as α_j minimizes MSE_j over the train set for all SKUs.
-

203 *3.2. Monte Carlo experiment*

204 Apart from the point forecasts, another output available is the residuals
205 in the test set, i.e., the difference between forecasts and actual values for
206 a determined SKU. Note that we use the residuals in the test set instead
207 of the training set to avoid underestimating the magnitude of the residuals
208 variance.

209 The Monte Carlo replications are $x_{ij,t}^h$, with $i = 1, 2, \dots, n$ and $j =$
210 $1, 2, \dots, m$, where n and m are the number of SKUs and replications, re-
211 spectively. To compute such replications mean and variance of residuals on
212 the test set are calculated and used for generating Gaussian random numbers
213 ($\epsilon_{ij,t}^h$) that are added to the h -step ahead point forecast ($F_{i,t}^h$) for a determined
214 SKU i , such as:

$$x_{ij,t}^h = F_{i,t}^h + \epsilon_{ij,t}^h \quad (4)$$

215 where m is the number of replications in the Monte Carlo experiment and n
216 is the total number of SKUs.

217 According to MATLAB 2018b release, the fastest algorithm to generate
218 random normal $\epsilon_{ij,t}^h$ is Philox 4x32 generator with 10 rounds, and for that
219 reason, it is the one chosen in this study. This generator was introduced in
220 2011, and it was specifically designed for high performance in highly parallel
221 systems such as GPUs [30]. In this study we have used a normal random
222 number generator because other alternatives as using either different sta-
223 tistical distributions or resampling methods as bootstrapping do not have
224 available a built-in function in MATLAB, and thus, they cannot be used
225 directly on a GPU.

226 *3.3. Probabilistic forecast*

227 Finally, the m forecast replications ($x_{ij,t}^h$) for a determined forecasting
228 horizon h can be used to build a predictive distribution of a certain SKU i .
229 Considering that distribution and a target quantile, the MATLAB command
230 **prctile** provides the empirical quantile of interest. In a supply chain context,
231 the quantile of interest is given by the Cycle Service Level (CSL) ([31]), which
232 represents the asymmetry between stock out and holding costs.

233 In summary, the pseudocode of the methodology presented here is shown
234 in procedures 3 and 4, where procedure 3 explains the calculation of prob-
235 abilistic forecasts from a serial computing point of view typically used in
236 CPU, whereas procedure 4 shows its parallel implementation thought for a
237 GPU. Note that the parallel version has removed the inner loop required to
238 generate m Monte Carlo replications.

Procedure 3 Calculation of probabilistic forecasts (CPU serial computing)

- 1: **for** $i = 1$ **to** n **do**
 - 2: Optimize in serial SES parameter α for the i SKU on the training set (see Procedure 1)
 - 3: Compute h-step ahead point forecasts $F_{i,t}^h$ with the parameter optimized in 2
 - 4: Compute residuals on the test set: $e_{i,t}^h = F_{i,t}^h - y_{i,t}$
 - 5: **for** $j = 1$ **to** m **do**
 - 6: Generate Gaussian random numbers $\epsilon_{ij,t}^h$ from mean and variance of residuals $e_{i,t}^h$
 - 7: Build j replication: $x_{ij,t}^h = F_{i,t}^h + \epsilon_{ij,t}^h$
 - 8: **end for**
 - 9: Compute h-step ahead probabilistic forecast based on m replications for a given CSL and for a certain SKU i
 - 10: **end for**
-

239 4. Evaluation metrics

240 In this work, we evaluate two key aspects in forecasting: accuracy and
241 computational speed. The accuracy is measured by an overall metric as the
242 Average Relative Mean Absolute Error (AvgRelMAE) [32]. The AvgRelMAE
243 can be computed as:

$$AvgRelMAE = \left(\prod_{i=1}^n r_i \right)^{1/n} \quad (5)$$

244 where $r_i = \frac{MAE_i^{SES}}{MAE_i^N}$. MAE_i^{SES} is the Mean Absolute Error (MAE) produced
245 by SES on the test set for SKU i ; and MAE_i^N is the MAE of the Naïve
246 method that is used as a benchmark for the same SKU i also in the same
247 test set. n is the total number of SKUs. Note that, the Naïve forecasting
248 method uses the last observation as a forecast for the next period. Thus, if
249 AvgRelMAE is lower than 1, it means that SES is more accurate than Naïve.

250 The differences in computational speed can be easily measured with the
251 speedup ratio, [22]. Speedups are computed as CPU time required to do
252 the calculations divided by GPU time counterpart, such as: $speedup =$
253 CPU_{time}/GPU_{time} .

Procedure 4 Calculation of probabilistic forecasts (GPU parallel computing)

- 1: Optimize in parallel SES parameter α for the all SKUs (see Procedure 2)
 - 2: **for** $i = 1$ **to** n **do**
 - 3: Compute h-step ahead point forecasts $F_{i,t}^h$ with the parameter optimized in 2
 - 4: Compute residuals on the test set: $e_{i,t}^h = F_{i,t}^h - y_{i,t}$
 - 5: Generate m Gaussian random numbers $\epsilon_{ij,t}^h$ from mean and variance of residuals $e_{i,t}^h$, where $j = 1, \dots, m$
 - 6: Build m replications: $x_{ij,t} = F_{i,t}^h + \epsilon_{ij,t}^h$
 - 7: **end for**
 - 8: Compute h-step ahead probabilistic forecast based on m replications for a given CSL and for all SKUs
-

254 5. Experimental results

255 In this section we compare the computational time taken to compute
256 probabilistic forecasts for a real dataset that belongs to an international com-
257 pany based on UK devoted to provide personal care products. The dataset
258 has been studied in detail by Barrow and Kourentzes [33]. It comprises 173
259 SKUs, with 229 observations per SKU. Manufacturer shipments are used to
260 estimate the demand, which is mainly continuous, that is, there are not many
261 zeros in the data. The dataset has been divided in two parts. The hold-in
262 sample or training set that consists of 80% of the dataset (183 observations)
263 and the hold-out sample or test set based on the last part of the data (46
264 observations). The optimization is carried out in the training set and the
265 AvgRelMAE is assessed in the test set for all the SKUs considered. The
266 forecasting horizon is one step ahead ($h = 1$).

267 These experiments will be run in a machine with two Intel (R) Xeon
268 processors at 3.00 GHz with 128 GB of RAM. The GPU card is a Quadro
269 GP100, 16 GB of Memory and 3584 CUDA Parallel-Processing Cores.

270 5.1. Analysis of accuracy

271 The calculation of the probabilistic forecasts requires the use of the grid
272 and search optimization routine, whose performance depends on the grid size
273 parameter (n_a), which should be determined.

274 This section analyzes the influence of n_a on the forecasting accuracy
 275 perspective. Figure 2 shows the AvgRelMAE provided by the forecast-
 276 ing method SES for all the 173 SKUs and different forecasting horizons
 277 $h = 1, 3, 5$. In that graph, it can be observed the forecasting accuracy
 278 obtained by the considered optimization methods. On the one hand, the
 279 typical optimization function of MATLAB (**fmincon** that finds minimum of
 280 constrained nonlinear multivariable function, [30]) in a dashed line and, on
 281 the other hand, the grid and search routine, which was previously introduced,
 282 in a solid line. That figure shows that the required n_a that achieves the same
 283 forecasting accuracy than **fmincon** depends on the forecasting horizon. For
 284 instance, when $h = 1$ less than 5 steps are needed to obtain the same forecast-
 285 ing accuracy than **fmincon**. For horizons greater than 1, a grid size (n_a) of
 286 20 offers a similar accuracy to the one obtained by the **fmincon** function.
 287 These results agree with those obtained by Nikolopoulos and Petropoulos
 288 [12] for a different dataset, where small grid sizes yield similar forecasting
 289 accuracies with respect to more complex optimization routines. Therefore,
 290 for the experiments in this study, the grid size will be set to a conservative
 291 $n_a = 20$.

292 5.2. Analysis of computational times

293 Figure 3 depicts an example of shipments and forecasts for a certain SKU.
 294 The upper panel of that figure shows the shipments data for that SKU. The
 295 lower panel shows the actual shipments and forecasts for the test set. The
 296 forecast errors in the test set are used to obtain the mean and standard
 297 deviation that will be fed to the random normal number generator to build
 298 the m Monte Carlo replications.

299 In our dataset we have a total of 173 SKUs sales data. Our interest is to
 300 analyze the computational times for a bigger number of SKUs. To do that,
 301 we have resampled the original dataset to obtain as many SKUs sales data
 302 as we need on the basis of the original dataset.

303 Table 1 shows the speedup ratios calculated for different volumes of SKUs
 304 from 100 up to 5000. The first row shows the speed up ratio of the MATLAB
 305 function **fmincon** run on a CPU with respect to the grid and search ($n_a = 20$)
 306 run in a GPU. Such a ratio varies between 32 and 42, that is, the GPU option
 307 was more than 30x times faster than the CPU alternative. The second and
 308 third row computes the speed up ratio using the same optimization technique
 309 (grid and search) with different values of n_a . For values of $n > 100$ a lower
 310 value of n_a provides higher speed up ratios.

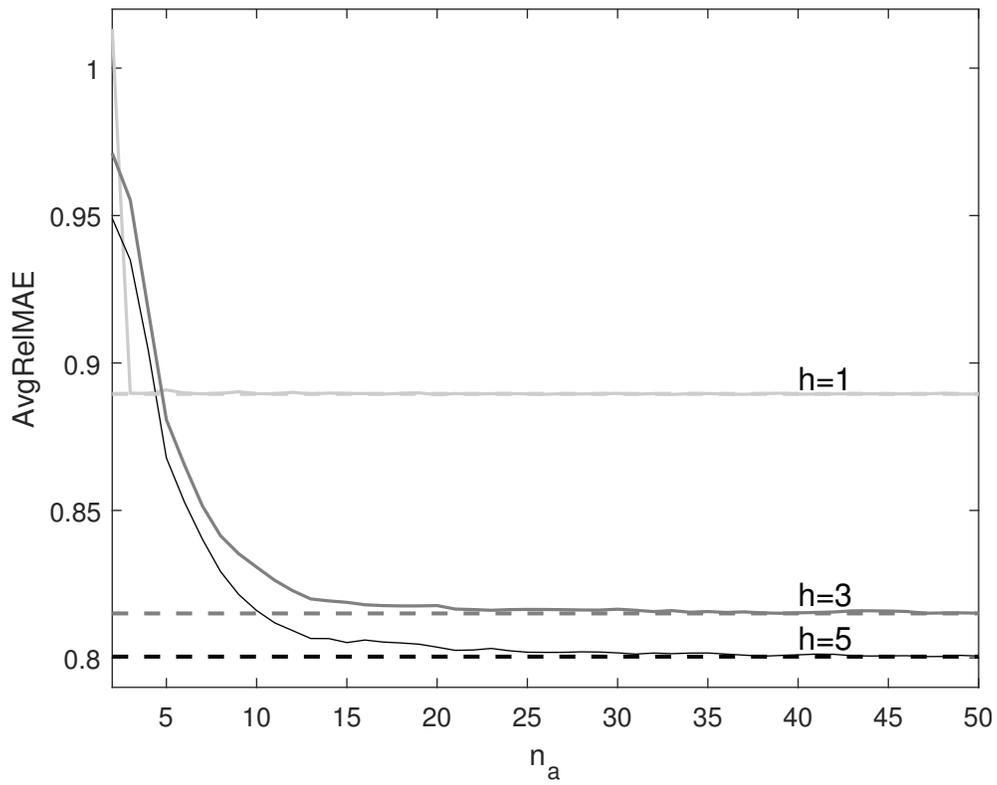


Figure 2: AvgRelMAE versus grid size (n_a) obtained for both fmincon (dashed line) and grid and search (solid line) optimization for all SKUs.

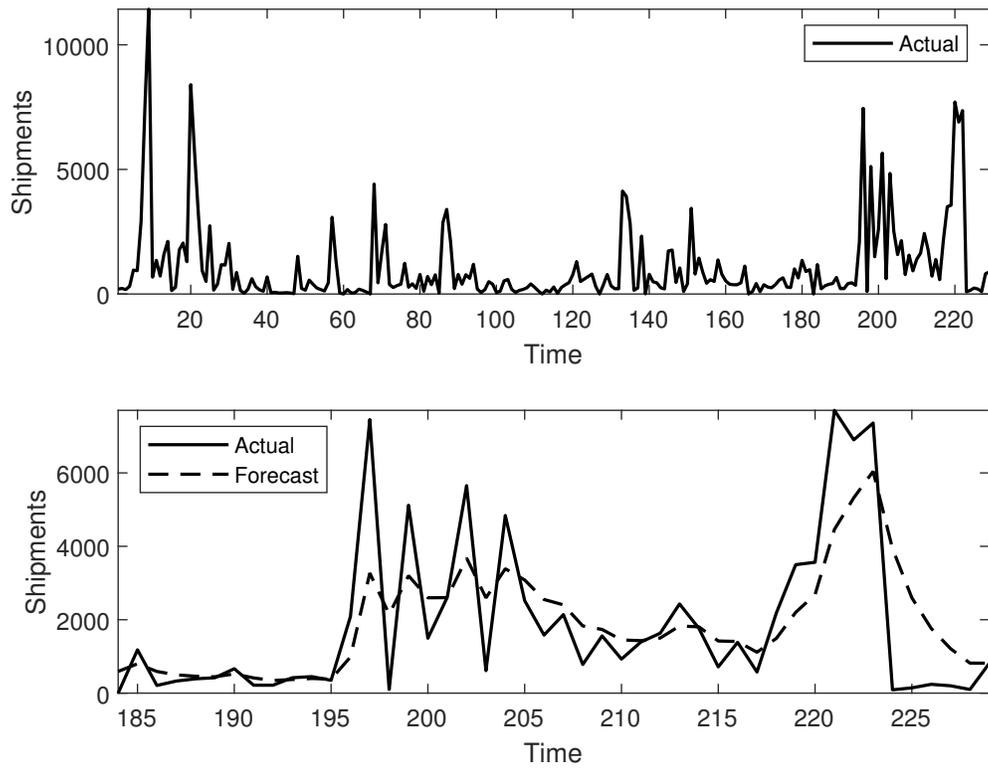


Figure 3: Example of the shipments for a certain SKU. Upper panel shows all the observations. Lower panel shows the forecasts (dashed line) and actual observations (solid line) for the test set.

311 It is interesting to note that we there is a value of n that maximizes the
 312 speed up ratio and it is $n = 1000$. That is an interesting fact, since, a priori,
 313 one could wrongly think that the higher the value of n , the higher the speed
 314 up ratio of the GPU.

315 This is an important conclusion that can be utilized when scaling this
 316 problem. Essentially, for a big data framework where the number of SKUs
 317 may compromise the memory size, MATLAB allows to work with different
 318 GPU cards (in parallel²) and this experiment shed some light about how
 319 many SKUS should be distributed in a GPUs cluster to obtain the maximum
 320 speed up ratio, i.e., the lower computational times.

Speed up	n			
	100	1000	2500	5000
fmincon	32	42	40	38
grid and search ($n_a=40$)	28	39	37	34
grid and search ($n_a=20$)	29	41	40	36

Table 1: Speed up ratios calculated for different values of n and rounded to the nearest integer.

321 Figure 4 shows the computational time required to calculate the proba-
 322 bilistic forecasts when using the CPU and GPU. The results show that the
 323 GPU with a grid and search optimization is faster than the CPU, regardless
 324 of the optimization routine chosen. Moreover, that difference is bigger when
 325 the number of SKUs (n) increases. In GPU the optimization routine was
 326 grid and search.

327 6. Conclusions

328 This article presents a methodology to implement in a GPU the calcula-
 329 tion of probabilistic forecasts in a parallel manner, where the point forecast-
 330 ing algorithm is a widely used forecasting technique as SES and the quantile
 331 forecasts can be obtained empirically via simulated predictive distributions.
 332 This simulation approach to compute qunatile forecasts has two main advan-
 333 tages: i) it incorporates the uncertainty of the parameter estimations that
 334 is usually neglected and; ii) for some complex forecasting models is the only

²Recent versions of MATLAB permit to use the **parfor** command on GPU cards

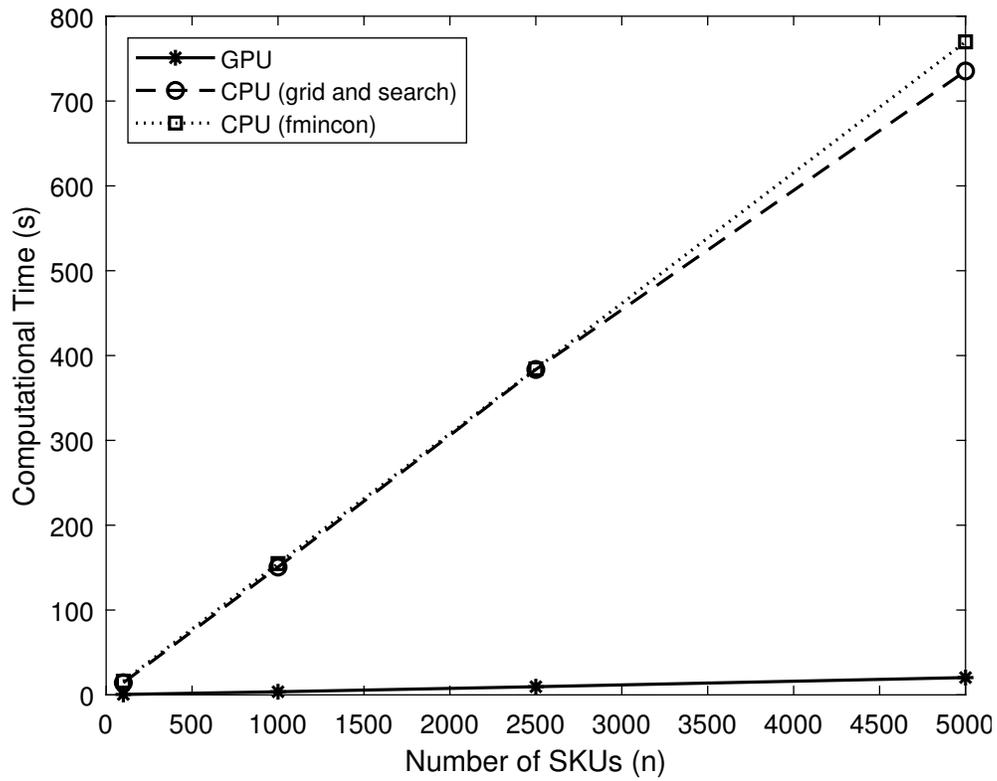


Figure 4: Computational time required for calculating the probabilistic forecast for different number of SKUs (n).

335 approach to compute predictive distributions. However, its main limitation
336 is the computational burden associated.

337 In this work, such a limitation is alleviated through the use of GPU com-
338 puting in parallel implementations of the algorithms involved. This method-
339 ology is implemented in MATLAB using “built-in” functions to circumvent
340 the use of CUDA language. Thus, simply usual MATLAB coding is necessary
341 to apply the proposed approach.

342 To the best of authors’ knowledge, this is the first time that GPU com-
343 puting is used for the problem of supply chain forecasting. The main result
344 show that GPU computing obtains similar forecasting accuracy in a reduced
345 computational time.

346 The improvements of GPU over CPU is limited for the GPU card memory.
347 Apart from expecting future GPU cards with bigger memories, a possibility
348 to relax that limitation in big data applications is to use a cluster of GPU
349 cards. Recent MATLAB 2018b version is able to use several GPUs in parallel
350 too. Additionally, in this work, we have shown that there is an optimal data
351 size to potentially distributing among GPU cards in a cluster, if it is necessary
352 to scale the problem.

353 Although SES is a widely used forecasting method, when the data ex-
354 hibits trend and/or seasonality, more advances exponential smoothing models
355 should be used. Further research should address the application of the pro-
356 posed procedure with other exponential smoothing forecasting models. The
357 only limitation to apply the procedure proposed here is that the forecast-
358 ing model must be able to be represented as a filter. For instance, ARIMA
359 models fit perfectly and can also be used.

360 **Acknowledgment**

361 This work was supported by the European Regional Development Fund
362 and Spanish Government (MINECO/FEDER, UE) under the project with
363 reference DPI2015-64133-R. Juan R. Trapero also thanks Raul Fernandez
364 Rodriguez for his help with the GPU installation.

365 **Bibliography**

- 366 [1] D. Blazquez, J. Domenech, Big data sources and methods for social and
367 economic analyses, *Technological Forecasting and Social Change* 130
368 (2018) 99 – 113.

- 369 [2] K. Govindan, T. Cheng, N. Mishra, N. Shukla, Big data analytics and
370 application for logistics and supply chain management, *Transportation*
371 *Research Part E: Logistics and Transportation Review* (2018).
- 372 [3] S. Tiwari, H. Wee, Y. Daryanto, Big data analytics in supply chain
373 management between 2010 and 2016: Insights to industries, *Computers*
374 *& Industrial Engineering* 115 (2018) 319 – 330.
- 375 [4] A. Gandomi, M. Haider, Beyond the hype: Big data concepts, methods,
376 and analytics, *International Journal of Information Management* 35
377 (2015) 137 – 144.
- 378 [5] G. Bello-Orgaz, J. J. Jung, D. Camacho, Social big data: Recent
379 achievements and new challenges, *Information Fusion* 28 (2016) 45 –
380 59.
- 381 [6] R. Addo-Tenkorang, P. T. Helo, Big data applications in
382 operations/supply-chain management: A literature review, *Computers*
383 *& Industrial Engineering* 101 (2016) 528 – 543.
- 384 [7] J. Manyika, M. Chui, P. Groves, D. Farrell, S. Kuiken, E. Doshi, Open
385 data: unlocking innovation and performance with liquid information,
386 McKinsey Global Institute. Available on-line., 2013.
- 387 [8] C. W. Chase, Using big data to enhance demand-driven forecasting and
388 planning, *The Journal of Business Forecasting* 32 (2013) 27–32.
- 389 [9] A. M. Waller, E. S. Fawcett, Data science, predictive analytics, and big
390 data: A revolution that will transform supply chain design and manage-
391 ment, *Journal of Business Logistics* 34 (2013) 77–84.
- 392 [10] B. T. Hazen, J. B. Skipper, J. D. Ezell, C. A. Boone, Big data and
393 predictive analytics for supply chain sustainability: A theory-driven re-
394 search agenda, *Computers & Industrial Engineering* 101 (2016) 592 –
395 598.
- 396 [11] T. Nguyen, L. Zhou, V. Spiegler, P. Ieromonachou, Y. Lin, Big data ana-
397 lytics in supply chain management: A state-of-the-art literature review,
398 *Computers & Operations Research* 98 (2018) 254 – 264.

- 399 [12] K. Nikolopoulos, F. Petropoulos, Forecasting for big data: Does subop-
400 timality matter?, *Computers & Operations Research* (2017).
- 401 [13] B. Seaman, Considerations of a retail forecasting practitioner, *International Journal of Forecasting* (2018).
402
- 403 [14] R. Y. Zhong, S. T. Newman, G. Q. Huang, S. Lan, Big data for supply
404 chain management in the service and manufacturing sectors: Challenges,
405 opportunities, and future perspectives, *Computers & Industrial Engineering* 101 (2016) 572 – 591.
406
- 407 [15] J. R. Trapero, M. Cards, N. Kourentzes, Empirical safety stock esti-
408 mation based on kernel and GARCH models, *Omega* 84 (2019) 199 –
409 211.
- 410 [16] D. Prak, R. Teunter, A general method for addressing forecasting un-
411 certainty in inventory models, *International Journal of Forecasting* 35
412 (2019) 224 – 238. Special Section: Supply Chain Forecasting.
- 413 [17] R. J. Hyndman, A. B. Koehler, J. K. Ord, R. D. Snyder, *Forecasting with*
414 *Exponential Smoothing: The State Space Approach*, Springer-Verlag,
415 Berlin, 2008.
- 416 [18] C. H. Nadungodage, Y. Xia, J. J. Lee, M. Lee, C. S. Park, GPU ac-
417 celerated item-based collaborative filtering for big-data applications, in:
418 2013 IEEE International Conference on Big Data, pp. 175–180.
- 419 [19] V. Boyer, D. E. Baz, M. Elkihel, Solving knapsack problems on GPU,
420 *Computers & Operations Research* 39 (2012) 42 – 47. Special Issue on
421 Knapsack Problems and Applications.
- 422 [20] V. Curtis, C. Sanches, A low-space algorithm for the subset-sum prob-
423 lem on GPU, *Computers & Operations Research* 83 (2017) 120 – 124.
- 424 [21] B. Nogueira, R. G. Pinheiro, A CPU-GPU local search heuristic for
425 the maximum weight clique problem on massive graphs, *Computers &*
426 *Operations Research* 90 (2018) 232 – 248.
- 427 [22] G. J. Lim, L. Ma, GPU-based parallel vertex substitution algorithm for
428 the p-median problem, *Computers & Industrial Engineering* 64 (2013)
429 381 – 388.

- 430 [23] J. R. Cheng, M. Gen, Accelerating genetic algorithms with GPU com-
431 puting: A selective overview, *Computers & Industrial Engineering* 128
432 (2019) 514 – 525.
- 433 [24] I. M. Coelho, V. N. Coelho, E. J. da S. Luz, L. S. Ochi, F. G. Guimares,
434 E. Rios, A GPU deep learning metaheuristic based model for time series
435 forecasting, *Applied Energy* 201 (2017) 412 – 418.
- 436 [25] E. S. Gardner, Exponential smoothing: The state of the art, *Journal*
437 *of Forecasting* 4 (1985) 1–28.
- 438 [26] E. S. Gardner, Exponential smoothing: The state of the art, Part II,
439 *International Journal of Forecasting* 22 (2006) 637–666.
- 440 [27] R. Fildes, Research into Forecasting Practice, *Foresight: The Interna-*
441 *tional Journal of Applied Forecasting* (2017) 39–46.
- 442 [28] J. K. Ord, R. Fildes, N. Kourentzes, *Principles of Business Forecasting*,
443 *Wessex Press Publishing Co.*, 2nd edition, 2017.
- 444 [29] S. M. Disney, D. R. Towill, On the bullwhip and inventory variance
445 produced by an ordering policy, *Omega* 31 (2003) 157 – 167.
- 446 [30] MATLAB, Release 2018a, *The MathWorks Inc.*, Natick, Massachusetts,
447 2018.
- 448 [31] E. Silver, D. Pyke, D. Thomas, *Inventory and Production Manage-*
449 *ment in Supply Chains. Fourth Edition*, *CRC Press. Taylor and Francis*
450 *Group*, 2017.
- 451 [32] A. Davydenko, R. Fildes, Measuring forecasting accuracy: The case
452 of judgmental adjustments to sku-level demand forecasts, *International*
453 *Journal of Forecasting* 29 (2013) 510 – 522.
- 454 [33] D. Barrow, N. Kourentzes, Distributions of forecasting errors of forecast
455 combinations: implications for inventory management, *International*
456 *Journal of Production Economics* 177 (2016) 24–33.