

# A Predictive Model of Subcutaneous Glucose Concentration in Type 1 Diabetes Based on Random Forests

Eleni I. Georga, Vasilios C. Protopappas, Demosthenes Polyzos, and Dimitrios I. Fotiadis, *Senior Member, IEEE*

**Abstract**—In this study, an individualized predictive model of the subcutaneous glucose concentration in type 1 diabetes is presented, which relies on the Random Forests regression technique. A multivariate dataset is utilized concerning the s.c. glucose profile, the plasma insulin concentration, the intestinal absorption of meal-derived glucose and the daily energy expenditure. In an attempt to capture daily rhythms in glucose metabolism, we also introduce a time feature in the predictive analysis. The dataset comes from the continuous multi-day recordings of 27 type 1 patients in free-living conditions. Evaluating the performance of the proposed method by 10-fold cross validation, an average RMSE of 6.60, 8.15, 9.25 and 10.83 mg/dl for 15, 30, 60 and 120 min prediction horizons, respectively, was attained.

## I. INTRODUCTION

Achieving tight glycemic control in type 1 diabetes necessitates the proper administration of insulin-based regimen considering the effect of both exogenous and endogenous factors on glucose regulation. The daily management of type 1 diabetes has taken advantage of the recent advances in continuous glucose monitoring (CGM) technologies; however, the wide spectrum of parameters that should be controlled renders it a rather difficult procedure. To this end, the enhancement of CGM devices with predictive models of the subcutaneous (s.c.) glucose profile could help in evaluating an individual's response to therapy, thereby mitigating primarily the incidence of short-term diabetic complications such as hypoglycemia.

The problem of s.c. glucose concentration prediction in type 1 diabetes has been studied in the context of both time-series and machine-learning techniques. The inherent nonlinearity and nonstationarity of the glucose regulatory system limits the predictive capacity (up to 30 min) of the autoregressive models [1, 2] of the CGM time series. This can also be partially attributed to the fact that the auto-correlation function of the s.c. glucose measurements vanishes at about 30 min [3] and, therefore, even a nonlinear technique as the one in [4] cannot show a better behavior. On the other hand, when quantitative information concerning the

carbohydrates intake and the exogenous insulin administration is exploited along with the s.c. glucose signal, predictions for longer horizons are made feasible [5-7]. The input of a Gaussian Process [5] and a Support Vector Regression [6] model was enhanced with physical activity information using real data recorded continuously throughout the observation days. In [7], qualitative descriptors of the lifestyle and the emotional status of the patient were taken into account for the construction of feed-forward neural networks able to provide 75-min predictions.

In this study, the Random Forests (RF) regression technique [8] is employed for the first time in the literature to deal with the problem of s.c. glucose prediction in type 1 diabetes based on a multivariate dataset acquired under free-living conditions. RF is an ensemble of tree predictors that partition the feature space using linear decision boundaries and the final decision is formed by averaging the output of the ensemble. We experiment with 3 different input cases corresponding to combinations of the input variables and we evaluate and compare the predictive accuracy of the RF technique for each one in relation to the prediction horizon.

## II. MATERIALS AND METHODS

### A. Dataset

The dataset includes 27 type 1 diabetic patients following multiple-dose insulin therapy and was collected in the framework of the EU research project METABO [9] from the participating clinical partners. The baseline characteristics of the patients as well as some descriptive statistics of their s.c. glucose profile are given in Table I.

Each patient wore the Guardian Real-Time CGM system (Medtronic Minimed Inc.), which records an average s.c. glucose concentration value with a 5-minute period. The SenseWear Armband body monitoring system (BodyMedia Inc.) was also used for the continuous (every 1 min) recording of the energy expenditure during physical activities. In addition, information regarding the food intake (i.e. type, amount and time) and the insulin injections (i.e. type, dose and time) was recorded on a daily basis using a specially designed paper diary. A dietician analyzed the meal recordings to specify their carbohydrate content.

The 27 patients were placed in 3 groups as follows: group A includes 15 patients for whom we have all the information required, group B includes 5 patients for whom it was not possible to exploit activity data and, group C includes the remaining 7 patients for whom only the CGM signal was available.

Research supported by European Commission (Project METABO "Controlling Metabolic Diseases Related to Metabolic Disorders", FP7-ICT-2007-1-216270).

E. I. Georga and D. I. Fotiadis are with the Department of Materials Science and Engineering, University of Ioannina, Ioannina, GR 45110 Greece (e-mail: egeorga@cs.uoi.gr, fotiadis@cc.uoi.gr).

E. I. Georga, V. C. Protopappas and D. Polyzos are with the Department of Mechanical Engineering and Aeronautics, University of Patras, Patras, GR 26500 Greece (e-mail: egeorga@cs.uoi.gr, corresponding author to provide phone: +302651008824; fax: +302651008889; e-mail: vprotop@mech.upatras.gr, polyzos@mech.upatras.gr).

TABLE I. DATASET CHARACTERISTICS

Patient Baseline Characteristics		Descriptive Statistics of the Glucose Dataset	
Gender		Average Hypoglycemic Events Per Day	0.4 (0-2)
No. Female	12		
No. Male	15		
Age (years)		Average Duration of Hypoglycemic Events (min)	43.1 (0-88.33)
Mean ± SD	43.5±13.4		
Range	19-72		
BMI (kg/m <sup>2</sup> )		Average Hyperglycemic Events Per Day	1.64 (0-3.36)
Mean ± SD	25±3.70		
Range	18.75–35.8		
Observation Period (days)		Average Duration of Hyperglycemic Events (min)	129.3 (0-283.75)
Mean ± SD	13.42±3.69		
Range	5-22		

SD=Standard Deviation, The values in parenthesis indicate min and max average values per patient.

B. The Proposed Method

The s.c. glucose concentration at time  $t+l$ , assuming that  $t$  is the time at which the prediction is made and  $l$  is the prediction horizon, is predicted by the RF regression function [8] of the input  $x \in \mathbb{R}^d$  :

$$f(x) = 1/B \sum_{b=1}^B T(x; \Theta_b), \tag{1}$$

where  $B$  is the size of the forest (i.e. number of trees) and  $T(x; \Theta_b)$  denotes the output of the  $b^{th}$  tree characterized by the vectors  $\Theta_b$ . More specifically, each tree is constructed based on an independent set of random vectors,  $\Theta_b$ , generated by a fixed probability distribution. Firstly, the randomness is injected into the building process of each tree by drawing a bootstrap sample from the training dataset, which has the same size as the original dataset. Splitting a node in the tree involves: (i) the selection of  $m$  features randomly from the total set of  $d$  features and then (ii) the determination of the best split among them according to the Gini index. The tree is then grown to its entirety without any pruning. Indeed, this process is repeated until each terminal node is associated with at least  $nodesize$  records.

The input  $x$  in the prediction function  $f$  concerns the following variables: (1) the s.c. glucose profile ( $gl$ ), (2) the plasma insulin concentration ( $I_p$ ), (3) the rate of appearance of exogenous (meal-derived) glucose in plasma ( $Ra$ ), (4) the cumulative amount of exogenous glucose that appeared in plasma ( $SRa$ ), (5) the hour of day (1 - 24) ( $h$ ) and, (6) the cumulative energy expenditure ( $SEE$ ). The time course of  $I_p$  is calculated according to a compartmental approach [10], which describes both the s.c. absorption kinetics for various insulin analogues ranging from rapid- to long-acting ones and the associated plasma dynamics. Regarding  $Ra$ , it is obtained by using the compartmental model of Lehmann *et al.* [11], in which the rate of gastric emptying is a function of the meal carbohydrate content.

In order to model the time delays in the glucose regulation process, we take into account the history of the above input variables. In particular, for each variable  $v$  the successive values within the time window  $[t_v - n_v \Delta t_v, t_v]$  are used for predicting the s.c. glucose concentration at time  $t+l$ :

$$v = [v(t_v - n_v \Delta t_v), \dots, v(t_v - \Delta t_v), v(t_v)], \tag{2}$$

where  $t_v$  is the upper limit of the time window,  $\Delta t_v$  is the sampling period and the parameter  $n_v$  determines the length of the time window. Thus, the total size  $d$  of the input (number of features) is:

$$d = \sum_v (n_v + 1). \tag{3}$$

The value of  $t_v$  is taken equal to  $t$  for the variables  $gl$  and  $SEE$ , whereas it is taken equal to  $t+l$  for the variables  $I_p$ ,  $Ra$  and  $SRa$ . Furthermore, the temporal effect of a variable  $v$  on s.c. glucose, as expressed by the  $n_v \Delta t_v$  quantity, is determined based on our experimental results as well as on theoretical and clinical results found in [2-4, 12, 13].

First, based on previous studies showing the existence of a strong dependency between glucose samples which are 30 or fewer minutes apart [2-4], we utilize the measurements of the  $gl$  variable in the last 30 min (i.e.  $n_{gl} = 6$ ,  $\Delta t_{gl} = 5$  min) with respect to the time  $t$ . Furthermore, we exploit the values of the  $I_p$  and the  $Ra$  variables within the last 30 min (i.e.  $n_{Ra, I_p} = 6$ ,  $\Delta t_{Ra, I_p} = 5$  min) with respect to the time  $t+l$  aiming at capturing both their magnitude and trend. To clarify that the upcoming values of these variables within the time interval  $[t, t+l]$  are computed by the compartmental models based on the insulin and meal recordings until the current time  $t$ . In addition, we have exploited the area under the  $Ra$  curve over the last 90 min (i.e.  $n_{SRa} = 9$ ,  $\Delta t_{SRa} = 10$  min) with respect to the time of prediction  $t+l$ , in accordance with studies on type 1 diabetes concerning the absorption of meal-derived glucose into the systemic circulation [12]. In particular, we introduce the variable  $SRa$ , which represents the cumulative amount of exogenous glucose inserted in the plasma over time (calculated every 10 min). Similarly, the short-term effects of physical activities and exercise on glucose variability [13] were treated by introducing the variable  $SEE$ , which expresses the energy expenditure over the last 3 hrs (i.e.  $n_{SEE} = 18$ ,  $\Delta t_{SEE} = 10$  min) in the form of a vector calculated cumulatively every 10 min. Furthermore, an attempt was made to capture the circadian rhythms of glucose concentration [14] by simply using as input the hour of day at which the prediction is made (i.e.  $n_h = \Delta t_h = 0$ ).

In order to elucidate the predictive capability of the input variables, we investigate 3 different input cases. In the first case, denoted herein as Case 1, the prediction of s.c. glucose is made based only upon the past s.c. glucose profile ( $gl$ ). In the second case, denoted as Case 2, the  $I_p$ ,  $Ra$ ,  $SRa$  and  $h$  input variables are also added in the input of the predictive function. The last case, namely Case 3, results from the addition of the  $SEE$  variable to the input of Case 2. Obviously Case 1 can be applied to all patient groups, Case 2 can be applied to groups A and B, and Case 3 only to group A. Finally, predictions are performed for four values of the prediction horizon  $l$ , i.e.  $l = 15, 30, 60$  and  $120$  min.

### C. Model Training and Evaluation

The predictive performance of the proposed method is evaluated individually for each patient by employing a 10-fold cross validation procedure. The number of trees in the forest,  $B$ , that minimizes the 10-fold cross validation RMSE ( $RMSE_{10\text{-fold}}$ ) is found by testing exhaustively all the values in the range  $[1, 100]$ . The  $RMSE_{10\text{-fold}}$  is defined as follows:

$$RMSE_{10\text{-fold}} = \frac{1}{10} \sum_{k=1}^{10} \sqrt{\frac{1}{N_k} \sum_{j=1}^{N_k} (y_j - f(x^j))^2}, \quad (4)$$

where  $N_k$  represents the size of the  $k^{\text{th}}$  fold,  $y_j$  is the actual value of glucose associated with the input  $x^j$  and  $f(x^j)$  is the glucose value computed by the RF. The generalization error of the RF is also affected by the values of the parameters  $m$  and  $nodesize$  which are used during the tree induction process. However, their determination is usually made *a priori* and, the values commonly used in the literature for regression are  $d/3$  (where  $d$  is the number of input features) and 5, respectively [8].

Besides the  $RMSE_{10\text{-fold}}$ , the assessment of the predictive accuracy of the proposed method is also based on the average correlation coefficient ( $r_{10\text{-fold}}$ ) resulting from the 10-fold cross validation:

$$r_{10\text{-fold}} = \frac{1}{10} \sum_{k=1}^{10} r_k, \quad (5)$$

where  $r_k$  is the correlation coefficient regarding the  $k^{\text{th}}$  fold.

In addition, the obtained predictions are evaluated with the aid of the Clarke's Error Grid Analysis (EGA) [15], which addresses the errors from a clinical point of view. This analysis defines 5 zones (i.e. A-E) so that the predicted-observed s.c. glucose concentration points within zones A and B are considered to be clinically acceptable, whereas the ones associated with zones C-E are likely to result in an unnecessary or erroneous treatment.

### III. RESULTS

Table II reports the average value and the corresponding standard deviation of  $RMSE_{10\text{-fold}}$  and  $r_{10\text{-fold}}$  for all input cases and for different patient groups. Regarding Group A, we observe that Case 1 results in a sufficient low error and a high degree of correlation concerning short-term (i.e. for 15 min and 30 min) predictions of the s.c. glucose concentration. However, as indicated by both measures, the predicted glucose profile deviates significantly from the real one for medium- to long-term prediction horizons (i.e. for 60 min and 120 min). The introduction of the  $I_p$ ,  $Ra$ ,  $SRa$  and  $h$  input variables in Case 2 reduces the average  $RMSE_{10\text{-fold}}$  associated with the 15-min and 30-min predictions by 29% and 42%, respectively, in comparison with Case 1. Regarding the 60-min and 120-min predictions, they are improved by 53% and 60%, respectively, compared to Case 1. In Case 3, in which the  $SEE$  variable is additionally used, the average  $RMSE_{10\text{-fold}}$  does not exceed 9 mg/dl for 15-min and 30-min predictions and 11 mg/dl for 60-min and 120-min predictions being further improved compared to Case 2 (i.e.

TABLE II. AVERAGE ERROR OF THE RF PREDICTIVE MODEL

Case#	Prediction Horizon							
	15 min		30 min		60 min		120 min	
	RMSE	$r$	RMSE	$r$	RMSE	$r$	RMSE	$r$
<b>Group A</b>								
Case 1	9.84 (2.07)	0.98 (0.01)	15.37 (2.47)	0.95 (0.02)	23.43 (3.64)	0.88 (0.05)	31.04 (6.09)	0.75 (0.11)
Case 2	6.99 (1.43)	0.99 (0.00)	8.98 (1.59)	0.98 (0.01)	11.00 (1.58)	0.97 (0.01)	12.38 (2.10)	0.97 (0.01)
Case 3	6.60 (1.32)	0.99 (0.00)	8.15 (1.65)	0.99 (0.01)	9.25 (1.39)	0.98 (0.01)	10.83 (2.76)	0.97 (0.01)
<b>Group B</b>								
Case 1	10.22 (0.68)	0.98 (0.00)	16.82 (2.42)	0.95 (0.02)	26.01 (5.04)	0.87 (0.04)	35.03 (7.10)	0.75 (0.08)
Case 2	7.49 (0.62)	0.99 (0.00)	10.02 (1.51)	0.98 (0.01)	12.71 (3.00)	0.97 (0.01)	13.97 (3.65)	0.97 (0.01)
<b>Group C</b>								
Case 1	11.33 (2.17)	0.97 (0.02)	17.64 (2.77)	0.92 (0.05)	26.05 (3.94)	0.82 (0.07)	35.37 (6.41)	0.61 (0.10)

by 6%, 9%, 16% and 13%). Furthermore, the average correlation coefficient  $r_{10\text{-fold}}$  is almost 1 in Cases 2, 3 for all different prediction horizons. The RF glucose predictive model exhibits similar behavior for groups B and C.

The average percentages of the predicted-observed points falling into the different zones of the Clarke's EGA, as applied to group A, are presented in Table III. It can be seen that practically all points lie in zones A and B even for higher prediction horizons. This analysis further supports that the predictions, and particularly the medium- to long-term ones, become more safe in Cases 2 and 3 where more than 95% of points lie within zone A for all horizons, as is also shown by the average  $r_{10\text{-fold}}$  obtained for these cases. Finally, similar were the results of this analysis for groups B and C.

### IV. DISCUSSION

In this work, a combination of compartmental models and RF for regression was proposed for the prediction of the s.c. glucose concentration in type 1 diabetes. A multivariate dataset covering the most prevalent regulators of glucose levels was utilized for this purpose. Comparisons of the numerical accuracy and the clinical significance of the generated predictions were made for 3 different input cases.

For the first time RF regression is being used in a glucose prediction scheme. Despite being a highly linear technique, RF is able to produce generalization errors that compare favorably to those of non-linear ones allowing potentially medical interpretation of its results. An additional innovative feature of this work is the manipulation of the model's input: (i) by introducing input variables that were eventually proved to contribute to significantly better predictions and (ii) by quantifying the effect of these inputs to the prediction accuracy through the experimentation with different input cases. More specifically, the output of the compartmental models was exploited the most by using not only the time history of the  $Ra$  and the  $I_p$  signals up to the time at which the prediction is made (i.e.  $t$ ) but also their future values at the time interval  $[t, t+I]$ . In addition, the area under the  $Ra$  curve was first exploited herein through the  $SRa$  variable. Regarding physical activity, we treated it differently from [5,

TABLE III. CLARKE'S EGA FOR GROUP A

Case #	Prediction Horizon											
	15 min			30 min			60 min			120 min		
	Clarke's EGA Zones											
	A	B	C-E	A	B	C-E	A	B	C-E	A	B	C-E
Case 1	98.14	1.63	0.23	92.36	6.53	1.11	79.95	17.16	2.90	69.43	26.35	4.22
Case 2	99.22	0.68	0.10	97.92	1.68	0.40	96.31	3.02	0.67	95.53	3.57	0.90
Case 3	99.26	0.62	0.13	98.23	1.39	0.38	97.59	1.90	0.51	96.43	2.75	0.82

6] by taking into account the time history of energy expenditure through the *SEE* variable. Finally, although comparative results were not provided herein, our experiments revealed the relationship between the hour of prediction and the glucose profile as in [7].

The results derived from our dataset demonstrate that the short-term predictions are feasible even based solely on the recent profile of glucose (Case 1). The effect of the prediction horizon on the obtained accuracy was visible in all input cases in all patient groups. However, when exploiting the full set of input variables (Case 3), the average  $RMSE_{10-fold}$  associated with the 60-min and 120-min predictions for group A was greatly improved by 61% and 65%, respectively, compared to Case 1, reaching 9.25 mg/dl and 10.83 mg/dl. The inclusion of all input variables amplified the linear correlation between the predicted and the observed glucose signals as well, which was also inferred by the Clarke's EGA. The results for Case 3 also show the importance of information about daily activities. Having compared the proposed method with those reported in the literature, we found that only the autoregressive model proposed by Gani *et al.* [2] gives more accurate 30-min predictions (average  $RMSE=3.42$  mg/dl). However, those results concern stationary segments of the CGM dataset. Moreover, a direct comparison with [7] is not feasible, since in that study a generic model is built for all patients.

Part of the limitations of the proposed method comes from the compartmental models which are used. The fact that these models do not consider the effect of some important factors on the absorption of s.c. insulin (e.g. site of injection, skin temperature) and carbohydrates (e.g. fats, fibers, glycemic index), and the fact that they are applied using population parameters inevitably introduce some error in the prediction process. In addition, there are other technical issues that need to be examined such as the optimal observation period required for data collection and the rate of model updating. Besides these, further clinical validation of the proposed method is required.

## V. CONCLUSIONS

A systematic study based on RF for regression and a multivariate dataset was presented for the prediction of the s.c. glucose concentration in type 1 diabetes. High-accuracy short- and long-term predictions were derived in the case where all the available information is used. A challenging prospective could be to exploit the level of transparency the RF technique provides through the interpretation of the rules extracted from the ensemble of trees, as well as through the evaluation of the variables importance based on the measures yielded by RF.

## REFERENCES

- [1] G. Sparacino, F. Zanderigo, S. Corazza, A. Maran, A. Facchinetti, and C. Cobelli, "Glucose concentration can be predicted ahead in time from continuous glucose monitoring sensor time-series", *IEEE Trans. Biomed. Eng.*, vol. 54, no. 5, pp. 931-937, May 2007.
- [2] A. Gani, A. V. Gribok, Y. Lu, W. K. Ward, R. A. Vigersky, and J. Reifman, "Universal glucose models for predicting subcutaneous glucose concentration in humans", *IEEE Trans. Inf. Technol. Biomed.*, vol. 14, no. 1, pp. 157-165, Jan. 2010.
- [3] B. P. Kovatchev and W. L. Clarke, "Peculiarities of the continuous glucose monitoring data stream and their impact on developing closed-loop control technology", *J. Diabetes. Sci. Technol.*, vol. 2, no. 1, pp. 158-163, Jan. 2008.
- [4] C. Perez-Gandia, A. Facchinetti, G. Sparacino, C. Cobelli, E. J. Gomez, M. Rigla, A. de Leiva, and M. E. Hernando, "Artificial neural network algorithm for online glucose prediction from continuous glucose monitoring", *Diabetes Technol. Ther.*, vol. 12, no. 1, pp. 81-88, Jan. 2010.
- [5] J. J. Valletta, A. J. Chipperfield, and C. D. Byrne, "Gaussian process modelling of blood glucose response to free-living physical activity data in people with type 1 diabetes" in *Proc. 31st Annu. Int. Conf. IEEE Eng. Med. Biol. Soc.*, Minneapolis, 2009, pp. 4913-4916.
- [6] E. I. Georga, V. C. Protopappas, and D. Polyzos, "Prediction of glucose concentration in type 1 diabetic patients using support vector regression" in *Proc. 10th Int. Conf. Inf. Technol. Appl. Biomed.*, Corfu, 2010.
- [7] S. M. Pappada, B. D. Cameron, P. M. Rosman, A. E. Bourey, and T. J. Papadimos, "Neural network-based real-time prediction of glucose in patients with insulin-dependent diabetes", *Diabetes Technol. Ther.*, vol. 13, no. 2, pp. 135-141, Feb. 2011.
- [8] L. Breiman, "Random forests", *Machine Learning*, vol. 45, no. 1, pp. 5-32, Oct. 2001.
- [9] E. I. Georga, V. C. Protopappas, A. Guillen, G. Fico, D. Ardigo, M. T. Arredondo, T. P. Exarchos, D. Polyzos, and D. I. Fotiadis, "Data mining for blood glucose prediction and knowledge discovery in diabetic patients: the METABO diabetes modeling and management system." in *Proc. 31st Annu. Int. Conf. IEEE Eng. Med. Biol. Soc.*, Minneapolis, 2009.
- [10] C. Tarin, E. Teufel, J. Pico, and J. Bondia, and H. J. Pfliederer, "A comprehensive pharmacokinetic model of insulin glargine and other insulin formulations", *IEEE Trans. Biomed. Eng.*, vol. 52, no. 12, pp. 1994-2005, Dec. 2005.
- [11] E. D. Lehmann and T. Deutsch, "A physiological model of glucose-insulin interaction in type 1 diabetes mellitus", *J. Biomed Eng.*, vol. 14, no. 3, pp. 235-242, May 1992.
- [12] M. E. Pennant, L. C. Bluck, M. L. Marcovecchio, B. Salgin, R. Hovorka, and D. B. Dunger, "Insulin administration and rate of glucose appearance in people with type 1 diabetes", *Diabetes Care*, vol. 31, no. 11, pp. 2183-2187, Nov. 2008.
- [13] American Diabetes Association, "Physical activity / exercise and diabetes", *Diabetes Care*, vol. 27, no. 1, pp. 58-62, Jan. 2004.
- [14] W. Huang, K. M. Ramsey, B. Marcheva, and J. Bass, "Circadian rhythms, sleep and metabolism", *J. Clin. Invest.*, vol. 121, no. 6, pp. 2133-2141, June 2011.
- [15] B. P. Kovatchev, L. A. Gonder-Frederick, D. J. Cox, and W. L. Clarke, "Evaluating the accuracy of continuous glucose monitoring sensors: Continuous glucose-error grid analysis illustrated by TheraSense FreeStyle Navigator data", *Diabetes Care*, vol. 27, no. 8, pp. 1922-1928, Aug. 2004.