

Interpreting Deep Glucose Predictive Models for Diabetic People Using RETAIN*

Maxime De Bois¹[0000-0002-4181-2422], Mounîm A. El Yacoubi²[0000-0002-7383-0588], and Mehdi Ammi³[0000-0003-1763-4045]

¹ CNRS-LIMSI and Université Paris Saclay, Orsay, France
`maxime.debois@limsi.fr`

² Samovar, CNRS, Télécom SudParis, Institut Polytechnique de Paris, Évry, France
`mounim.el_yacoubi@telecom-sudparis.eu`

³ Université Paris 8, Saint-Denis, France
`ammi@ai.univ-paris8.fr`

Abstract. Progress in the biomedical field through the use of deep learning is hindered by the lack of interpretability of the models. In this paper, we study the RETAIN architecture for the forecasting of future glucose values for diabetic people. Thanks to its two-level attention mechanism, the RETAIN model is interpretable while remaining as efficient as standard neural networks.

We evaluate the model on a real-world type-2 diabetic population and we compare it to a random forest model and a LSTM-based recurrent neural network. Our results show that the RETAIN model outperforms the former and equals the latter on common accuracy metrics and clinical acceptability metrics, thereby proving its legitimacy in the context of glucose level forecasting. Furthermore, we propose tools to take advantage of the RETAIN interpretable nature. As informative for the patients as for the practitioners, it can enhance the understanding of the predictions made by the model and improve the design of future glucose predictive models.

Keywords: Deep Learning · Glucose Prediction · Diabetes · Neural Networks · Attention · Interpretability.

1 Introduction

Diabetes is undoubtedly one of the major diseases of the modern world as it has been inputed a total of 1.5 million deaths in 2012 [16]. The every day challenge faced by diabetic people is the regulation of their blood glucose level which is troubled by either the non-production of insulin (type-1 diabetes) or the increasing body resistance to its action (type-2 diabetes). Diabetic people are at risk of facing short terms complications (e.g., coma, death) due to their glycemia falling

* This work is supported by the "IDI 2017" project funded by the IDEX Paris-Saclay, ANR-11-IDEX-0003-02.

too low (hypoglycemia) and also long-term complications (e.g., cardiovascular diseases, blindness) when it gets to high (hyperglycemia).

To help the patients coping with their disease, a lot of technological efforts have been made in the recent years. For instance, by enabling the diabetic patient to forgo the use of lancets to get his or her glucose level, continuous glucose monitoring (CGM) devices (e.g., FreeStyle Libre [10]) are getting more and more common. Besides, we are witnessing the rise of coaching applications specifically made for diabetic people (e.g., mySugr [12]). From a research perspective, current endeavors are focused towards the building of glucose predictive models. Using past glucose values, carbohydrate (CHO) intakes, insulin infusions, and more, the models forecast the future glucose values at horizons varying from 30 minutes (short-term) to 120 minutes (long-term) [11].

Thanks to the increasing availability of data and the access to more computing power, the glucose predictive models are shifting from rather simple models (e.g., autoregressive models [13]), to more complex algorithms from the machine learning and deep learning field. Daskalaki *et al.* have demonstrated the superiority of feed-forward neural networks over the autoregressive models in the context of short-term glucose forecasting [2]. Georga *et al.* explored the usability of extreme learning machines for short-term glucose prediction as well [5]. Recurrent neural networks have recently generated a lot of interest because of their temporal nature, making them particularly suitable for the task of predicting future glucose values [3, 9]. As time-series can be seen as one-dimension images, convolutional neural networks, which are very popular in the image recognition community, have also been tried out for the forecasting of future glucose values [7].

Even though deep models can be effective for the task of glucose prediction, they have a sizable downside: the deeper the model, the more difficult it is to understand its behavior. This is especially an issue for biomedical applications for which it is important to be able to interpret the models in order to understand why a prediction is being made. To address this issue, Georga *et al.* showed that Random Forests (RF), while being highly interpretable, can achieve good performances for the task of glucose prediction [4].

Recently, Choi *et al.* proposed a neural network, called RETAIN, specifically designed for healthcare applications dealing with temporal inputs. Featuring a two-level attention mechanism, the model is meant to be as performant as standard neural networks while being interpretable. This property is highly valuable for the prediction of future glucose values. On one hand, it would help the practitioner design better and safer models by providing a better error analysis tool. On the other hand, for the patient, it would help him or her understand his or her disease better.

In this work, we study the use of the RETAIN architecture for the challenging task of the forecast of future glucose values for diabetic people. In particular, we adapt its interpretability feature to regression problems and propose several analysis and visualization tools to interpret the predictions made by the model.

The rest of the paper is structured as follows. First, we describe the RETAIN architecture and how the predictions are interpreted from it. Then, we describe the overall experimental methodology. Finally, we provide the results and analysis of the experiments before concluding.

2 RETAIN

This section presents the RETAIN architecture that has been previously introduced in [1] and its interpretation for time-series forecasting, and in particular for glucose prediction.

2.1 Architecture

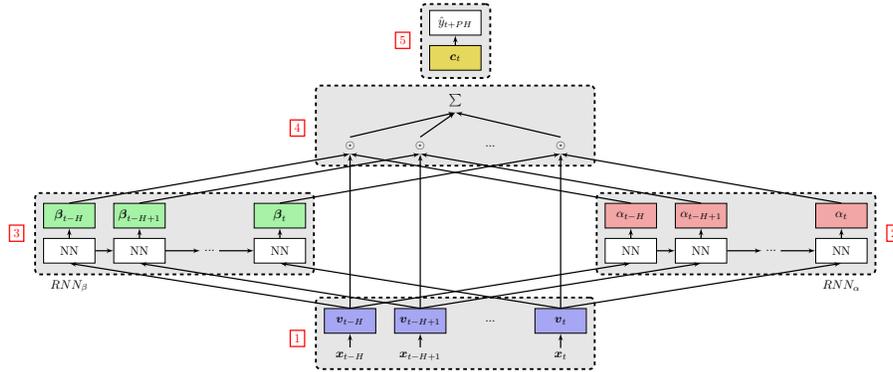


Fig. 1: Graphical overview of the RETAIN model. **Step 1:** The input signals are transformed into embeddings. **Step 2:** Time-level attention weights are computed from the embeddings. **Step 3:** Variable-level attention weights are also computed from the embeddings. **Step 4:** Using the attention weights, the context vector is computed. **Step 5:** The prediction is made from the context vector.

Most of the efficiency of the RETAIN model comes from its two levels of attention: the *time-level* attention (also called *visit-level* attention [1]), and the *variable-level* attention. The general attention mechanism comes from the natural language processing field where it enables the model to understand relationships between words in a sentence [15]. Here, when dealing with temporal inputs (e.g., time-series, events), while the time-level attention makes the network focus on specific time-steps, the variable-level attention emphasizes specific input within the time-steps.

The predictions of the RETAIN model are made following five different steps for which Figure 1 provides a graphical representation. In the following, t refers to the current time-step the prediction is made, r to the number of different

input variables, H to the size of the history (the number of past values for every input variable), PH to the prediction horizon, the subscript $i \in [t - H, t]$ to the i -th time-step, and the subscript $j \in [1, r]$ to the j -th input variable.

- **Step 1:** Each time-step input vector \mathbf{x}_i is linearly transformed into a learnable embedding \mathbf{v}_i following: $\mathbf{v}_i = \mathbf{W}_{emb}\mathbf{x}_i$.
- **Step 2:** These embeddings are given as inputs to a recurrent neural network, RNN_α , which outputs the time-level attention weights α_i (see [1] for more details).
- **Step 3:** Similarly, the embeddings are fed into a second recurrent neural network, RNN_β , which computes the variable-level attention weights β_i (see [1] for more details).
- **Step 4:** Using both attention weights, the context vector \mathbf{c}_t is computed following: $\mathbf{c}_t = \sum_{i=t-H}^t \alpha_i \beta_i \odot \mathbf{v}_i$.
- **Step 5:** The predictions of the model are made by linearly transforming the context vectors: $\hat{y}_{t+PH} = \mathbf{W}\mathbf{c}_t + b$.

The only difference between our architecture and the original one is that we do not compute the attention weights in reverse time-order (Steps 2 and 3) [1], but rather in forward order, as the latter yielded better performances for our application.

2.2 Interpretation

In their original paper, the authors of RETAIN propose a way to interpret the outputs of the RETAIN model in the context of multiclass classification. We propose here an adaptation of the methodology to regression problems.

By going through the different operations made in the model, we can express the prediction \hat{y}_{t+PH} in this form:

$$\hat{y}_{t+PH} = \sum_{i=t-H}^t \sum_{j=1}^r x_{i,j} \alpha_i \mathbf{W}(\beta_i \odot \mathbf{W}_{emb}[:, j]) + b \quad (1)$$

We can then express the contribution $\omega(\hat{y}_{t+PH}, x_{i,j})$ of the $x_{i,j}$ input feature on the prediction \hat{y}_{t+PH} as follows:

$$\omega(\hat{y}_{t+PH}, x_{i,j}) = \alpha_i \mathbf{W}(\beta_i \odot \mathbf{W}_{emb}[:, j]) x_{i,j} \quad (2)$$

While the contributions in this form are useful to analyze an individual sample, they are not very practical if we want to perform further analysis and statistics. Instead, we propose to look at the absolute normalized contribution values $\omega_{AN}(\hat{y}_{t+PH}, x_{i,j})$:

$$\omega_{AN}(\hat{y}_{t+PH}, x_{i,j}) = \frac{|\omega(\hat{y}_{t+PH}, x_{i,j})|}{\sum_{i=t-H}^t \sum_{j=1}^r |\omega(\hat{y}_{t+PH}, x_{i,j})|} \quad (3)$$

Taking the absolute values makes the computation of the mean contribution across the samples more representative of the overall contributions, preventing

positive and negative contributions from canceling each other. Normalizing the contributions makes the contributions independent from the prediction value itself, enabling a better comparison between samples.

3 Methods

3.1 Experimental Data

In this study, we use the IDIAB dataset whose collection has been approved by the french ethical committee (ID RCB 2018-A00312-53). It is made of data coming from 5 type-2 diabetic patients (4F/1M, age 58.8 ± 8.28 years old, BMI $30.76 \pm 5.14 \text{ kg/m}^2$, HbA1c $6.8 \pm 0.71 \%$) that have been monitored for 31.8 ± 1.17 days in free-living conditions. Whereas their glucose level (in mg/dL) was recorded through the use of the FreeStyle Libre continuous glucose monitoring device, data related to CHO intakes (in g) and insulin (in units) infusions were manually reported through the mySugr (mySugr GmbH) smartphone coaching app for diabetes.

3.2 Models

In this study, we build global glucose predictive models. Whereas personalized glucose predictive models are often more accurate, global models have the advantage of being easier to train by avoiding overfitting thanks to more training data.

We describe here the preprocessing and training steps of the different models used in this study.

3.3 Data Preprocessing

After splitting the patients into four training patients and one testing patient, we have splitted each training patient’s data into a training set and a validation set following a 75%/25% distribution.

To predict the future glucose values at an horizon of 30 minutes, the models are given as inputs the histories of glucose values, insulin infusions, and CHO intakes of the past 3 hours. For every patient, these inputs are standardized (zero mean and unit variance) w.r.t their respective training set.

3.4 Model Training

The Random Forest (RF) model [14] is one of the two baseline models used in this study. Its main strength is that it provides generally good performances while being easily interpretable. Here, a forest of size 100 is fitted using the mean-squared error (MSE) criterion. The minimum number of samples per leaf has been set to 25 to reduce the overfitting of the model to the training set.

Our second baseline, the LSTM model, has been implemented with an architecture that matches the computational complexity of the RETAIN model described below. In particular, every time-step input variables are embedded into a learnable vector of size 64. These embeddings are then given to a 2-layer LSTM model with 128 units per layer. The latter has been trained to minimize the MSE loss function with the Adam optimizer (learning rate of 10^{-3} , mini-batch size of 50). To prevent the overfitting of the network to the training set, the early stopping methodology (patience of 25) has been used.

As for the LSTM model, the RETAIN model has an embedding size of 64. Both RNN_{α} and RNN_{β} are made of one layer of 128 LSTM units. Similarly, the Adam optimizer (same learning rate and mini-batch size) with the early stopping methodology was used to fit the model.

All the hyperparameters have been optimized by grid search on the validation set on a subspace delimited by manual search.

3.5 Evaluation

The models have been evaluated with a 4-fold cross-validation on the training patients followed by a leave-one-(patient)-out cross-validation.

Four different metrics have been used: the Root-Mean-Squared Error (RMSE), the Mean Percentage Absolute Error (MAPE), the Time Lag (TL), and the Continuous Glucose-Error Grid Analysis (CG-EGA).

Both the RMSE and MAPE metrics give a measure of the accuracy of the prediction. The TL metric provides an estimate of the time gained by doing the prediction and is computed as the time-shift (in minutes) that maximizes the correlation between the true and the predicted glucose values. Finally, the CG-EGA measures the clinical acceptability of the predictions [6]. By analyzing both the prediction accuracy and the accuracy of the variation between two consecutive predictions, the CG-EGA classifies the prediction either as an accurate prediction (AP), a benign error (BE), or an erroneous prediction (EP). For a model to be clinically acceptable, it needs to have high AP and low EP rates.

4 Results & Discussion

4.1 Experimental Results

The performances of the three models are shown in Table 1 and Table 2. With an average deterioration of 1.4% in RMSE/MAPE/TL when compared to the LSTM model, the RETAIN model displays a comparable prediction accuracy. Its clinical acceptability is also very similar to the LSTM model.

When compared to the RF model, the RETAIN model shows an improvement of 8.5%, 8.4%, 20.7% in the RMSE, MAPE, and TL metrics respectively. It also has a better clinical acceptability with a lower EP rate (-9.7% which comes at the cost of a slightly lower AP rate (-3.3% of the remaining room for improvement)).

Overall, these results are showing that the RETAIN model is a legitimate model for the task of glucose prediction.

Table 1: Performances of the models with mean \pm standard deviation, averaged on the population.

Model	RMSE	MAPE	TL
RF	19.23 \pm 6.73	9.37 \pm 1.58	15.31 \pm 3.38
LSTM	17.52 \pm 5.52	8.35 \pm 1.30	12.01 \pm 2.36
RETAIN	17.60 \pm 4.90	8.58 \pm 0.84	12.14 \pm 2.53

Table 2: Clinical acceptability of the models with mean \pm standard deviation, averaged on the population.

Model	CG-EGA		
	AP	BE	EP
RF	86.00 \pm 4.37	10.79 \pm 3.59	3.21 \pm 0.84
LSTM	85.67 \pm 3.28	11.46 \pm 2.47	2.87 \pm 0.95
RETAIN	85.54 \pm 5.41	11.56 \pm 4.50	2.90 \pm 0.95

4.2 Interpreting the RETAIN Model

The real strength of the RETAIN model, however, lies in its interpretability. We propose here several different visualization tools for the analysis of the behavior of the RETAIN model. To ease the reading, we will refer to the contribution as the absolute normalized contribution, presented in Section 2.2.

First, by looking at the individual maximum contribution of the input variables, we can see if each of them has ever contributed significantly to the prediction. Figure 2 plots the maximum contribution of the model inputs related to the 3-hour histories of glucose values, insulin infusions, and CHO intakes. We can see that the the older an input value is, the less contribution it has. The decrease in the contribution is faster for the insulin and CHO signals (close to zero after 30 minutes) than for the glucose signal (close to zero after 60 minutes). This suggests that it is not usefull in this context to use histories that are longer than one hour. Reducing the number of past values inputed to the model should increase the performances by making it harder to overfit and should reduce the training time. Such an analysis is not possible with a standard LSTM model.

From a different perspective, we can look at the behavior of the model when an event occurs. Figure 3 depicts the behavior of the model following the occurrence of two different events: insulin infusions and CHO intakes. We can compare these plots to the mean contribution when no event has occurred in the last hour with Figure 4.

When either one of the events occurs, we can see that the glucose value that has the most importance is not the current glucose value, but the previous one (which is the value 5 minutes before the event). This specific value keeps a relative high importance as the time moves on. This shows that, when an event occurs, the model uses the last glucose value before the event as a value

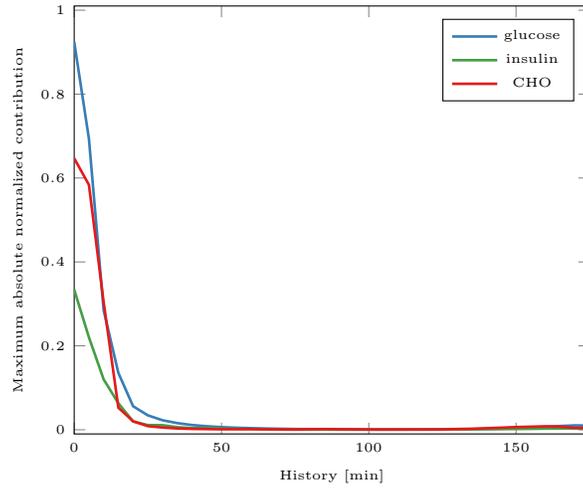


Fig. 2: Maximum absolute normalized contribution of the input signals (history of glucose, insulin, and CHO).

of reference. On the other hand, for both insulin and CHO signals, when their respective event occurs, the contribution of the value of the event is relatively high for the next 20 minutes. However, after this time, the contribution of the event is close to zero and the mean contribution profile becomes similar to the one for which no event has occurred in the past hour, depicted by Figure 4.

5 Conclusion

In this study, we have presented the application of the RETAIN model proposed by Choi *et al.* [1] to the challenging task of 30-minutes ahead-of-time glucose prediction for diabetic people. Using a two-level attention mechanism, the RETAIN model is able to produce interpretable predictions, which is highly valuable in the context of a biomedical field.

We have evaluated the model on a type-2 diabetic population of 5 patients and compared it against a Random Forest and a LSTM-based recurrent neural network. By being interpretable while respectively equalling and outperforming the LSTM and the RF models, we show that the RETAIN model is very promising.

In the future, we plan to extend the study to another dataset, namely the Ohio T1DM dataset [8]. In particular, this dataset comprises 6 type-1 diabetic patients with similar data. Also, thanks to the interpretability of the RETAIN model, we plan to explore variants of its architecture and input data (e.g., physical activity measures).

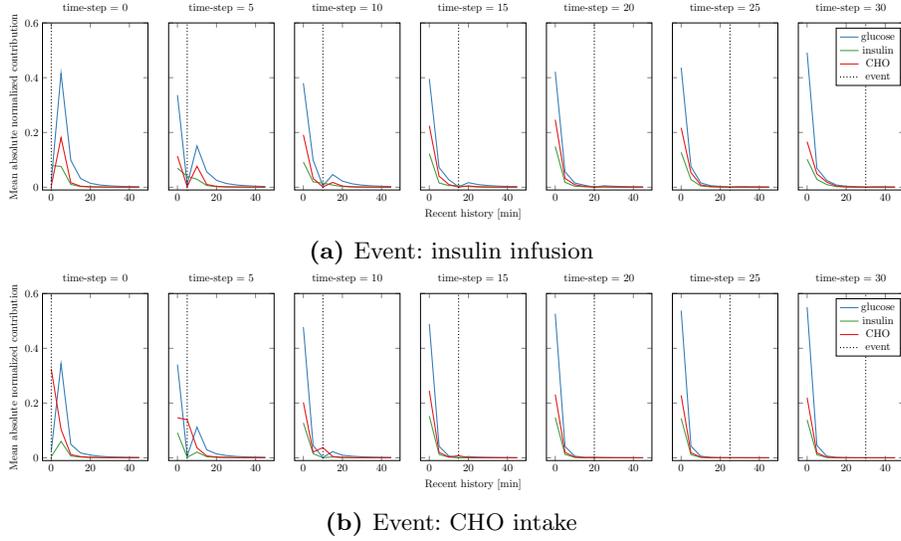


Fig. 3: Mean evolution through time of the absolute normalized contribution of the input signals (history of glucose, insulin, and CHO) after the occurrence of an event: Figure 3a, insulin infusion; and Figure 3b, CHO intake.

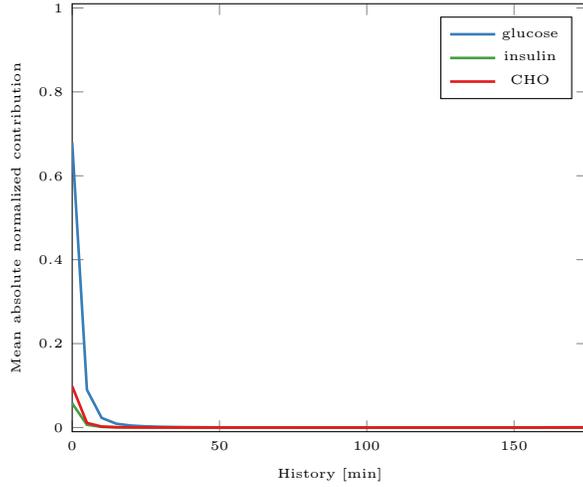


Fig. 4: Mean absolute normalized contribution of the input signals when no event (CHO intake or insulin infusion) occurred in the last hour.

Acknowledgment

We would like to thank the diabetes health network Revesdiab and Dr. Sylvie JOANNIDIS for their help in building the IDIAB dataset used in this study.

References

1. Choi, E., Bahadori, M.T., Sun, J., Kulas, J., Schuetz, A., Stewart, W.: Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. In: *Advances in Neural Information Processing Systems*. pp. 3504–3512 (2016)
2. Daskalaki, E., Prountzou, A., Diem, P., Mougiakakou, S.G.: Real-time adaptive models for the personalized prediction of glycemic profile in type 1 diabetes patients. *Diabetes technology & therapeutics* **14**(2), 168–174 (2012)
3. De Bois, M., El Yacoubi, M., Ammi, M.: Prediction-coherent lstm-based recurrent neural network for safer glucose predictions in diabetic people (accepted at ICONIP 2019)
4. Geoga, E.I., Protopappas, V.C., Polyzos, D., Fotiadis, D.I.: A predictive model of subcutaneous glucose concentration in type 1 diabetes based on random forests. In: *2012 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*. pp. 2889–2892. IEEE (2012)
5. Geoga, E.I., Protopappas, V.C., Polyzos, D., Fotiadis, D.I.: Online prediction of glucose concentration in type 1 diabetes using extreme learning machines. In: *Engineering in Medicine and Biology Society (EMBC), 2015 37th Annual International Conference of the IEEE*. pp. 3262–3265. IEEE (2015)
6. Kovatchev, B.P., Gonder-Frederick, L.A., Cox, D.J., Clarke, W.L.: Evaluating the accuracy of continuous glucose-monitoring sensors: continuous glucose-error grid analysis illustrated by the sense freestyle navigator data. *Diabetes Care* **27**(8), 1922–1928 (2004)
7. Li, K., Daniels, J., Liu, C., Herrero-Vinas, P., Georgiou, P.: Convolutional recurrent neural networks for glucose prediction. *IEEE Journal of Biomedical and Health Informatics* (2019)
8. Marling, C., Bunescu, R.: The ohio1dm dataset for blood glucose level prediction. In: *The 3rd International Workshop on Knowledge Discovery in Healthcare Data, Stockholm, Sweden* (2018)
9. Mirshekarian, S., Bunescu, R., Marling, C., Schwartz, F.: Using lstms to learn physiological models of blood glucose behavior. In: *Engineering in Medicine and Biology Society (EMBC), 2017 39th Annual International Conference of the IEEE*. pp. 2887–2891. IEEE (2017)
10. Ólafsdóttir, A.F., Attvall, S., Sandgren, U., Dahlqvist, S., Pivodic, A., Skrtic, S., Theodorsson, E., Lind, M.: A clinical trial of the accuracy and treatment experience of the flash glucose monitor freestyle libre in adults with type 1 diabetes. *Diabetes technology & therapeutics* **19**(3), 164–172 (2017)
11. Oviedo, S., Vehí, J., Calm, R., Armengol, J.: A review of personalized blood glucose prediction strategies for t1dm patients. *International journal for numerical methods in biomedical engineering* **33**(6), e2833 (2017)
12. Rose, K., Koenig, M., Wiesbauer, F.: Evaluating success for behavioral change in diabetes via mhealth and gamification: Mysugrs keys to retention and patient engagement. *Diabetes Technology & Therapeutics* **15**, A114 (2013)

13. Sparacino, G., Zanderigo, F., Corazza, S., Maran, A., Facchinetti, A., Cobelli, C.: Glucose concentration can be predicted ahead in time from continuous glucose monitoring sensor time-series. *IEEE Transactions on biomedical engineering* **54**(5), 931–937 (2007)
14. Svetnik, V., Liaw, A., Tong, C., Culberson, J.C., Sheridan, R.P., Feuston, B.P.: Random forest: a classification and regression tool for compound classification and qsar modeling. *Journal of chemical information and computer sciences* **43**(6), 1947–1958 (2003)
15. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: *Advances in neural information processing systems*. pp. 5998–6008 (2017)
16. World Health Organization, et al.: *Global report on diabetes*. World Health Organization (2016)