# Detection of GNSS-TEC Noise Related to the Tonga Volcanic Eruption Using Optimization Machine Learning Techniques and Integrated Data

**Nhung Le** , **Benjamin Männel** , **Luyen K. Bui** , **Mihaela Jarema**, **Thai Chinh Nguyen** , **and Harald Schuh**

**Abstract** Total Electron Content (TEC) is the integral of the electron density along the path between receivers and satellites. TEC measured from Global Navigation Satellite Systems (GNSS) data is valuable to monitor space weather and correct ionospheric models. TEC noise detection is also an essential channel to forecast space weather and research the relationship between the atmosphere and natural phenomena like geomagnetic storms, earthquakes, volcanos, and tsunamis. In this study, we apply optimization machine learning techniques and integrated GNSS and solar activity data to determine GNSS-TEC noise at the International GNSS Service (IGS) stations in the Tonga volcanic region. We investigate 38 indices related to the geomagnetic field and solar wind plasma to select the essential parameters for forecast models. The findings show the best-suited parameters to predict vertical TEC time series: plasma temperature (or Plasma speed), proton density, Lyman alpha, R sunspot, Ap index (or Kp, Dst), and F10.7 index. Applying the Ensemble algorithm to build the TEC forecast models at the investigated IGS stations gets the accuracy from 1.01 to 3.17 TECU. The study also shows that machine learning combined

N. Le (✉) · B. Männel · T. C. Nguyen · H. Schuh
GFZ German Research Centre for Geosciences, Potsdam, Germany
e-mail: nhung@gfz-potsdam.de

T. C. Nguyen
e-mail: nguyenthaichinh@humg.edu.vn

N. Le · H. Schuh
Technische Universität Berlin, Berlin, Germany

N. Le
Hanoi University of Natural Resources and Environment, Hanoi, Vietnam

L. K. Bui
National Center for Airborne Laser Mapping, University of Houston, Houston, USA
e-mail: buikhacluyen@humg.edu.vn

L. K. Bui · T. C. Nguyen
Hanoi University of Mining and Geology, Hanoi, Vietnam

M. Jarema
MathWorks, Munchen, Germany

with integrated data can provide a robust approach to detecting TEC noise caused by seismic activities.

**Keywords** Machine learning · GNSS-TEC forecast · GNSS · Solar activity · Tonga volcanic eruption

## 1 Introduction

Continuous Global Navigation Satellite Systems (GNSS) data can be used for various applications. Satellite signal propagation in space depends uniquely on electron density in the ionosphere [1–3]. Thus, the estimation of Total Electron Contents (TEC) in the ionosphere can provide valuable information to correct errors in GNSS positioning. Furthermore, solar activity is the main factor causing fluctuations in the Earth's electrical and magnetic fields [4]. The hot plasma makes energetic charged particles in space escape from the Sun's gravity and interacts with the Earth's magnetic field. The interaction between the solar wind plasma and Earth's magnetic field leads to some natural phenomena in the atmosphere like auroras, geomagnetic storms, and ionospheric anomalies [5–8]. Monitoring TEC disturbances thereby reflects the solar activity and is an important channel to forecast space weather.

Some other factors also result in ionospheric anomalies in the short term, for example, nuclear explosions [9–12] and rocket launching [13, 14]. Thanks to an increasing number of continuous GNSS stations, the research stream that has been attractive to scientists for almost two decades is TEC anomalies associated with seismic activities [15–19]. Determination of the TEC disturbances related to seismic events applies different methods and monitoring instruments. The French low orbit satellite DEMETER[1] was launched in 2004 to investigate ionospheric disturbances related to earthquakes and volcanos [20]. The multi-purpose GNSS networks like the GEONET in Japan and the SEALION in Southeast Asia have been attached to ionosondes, scintillation monitors, and magnetometers to observe the effect of seismic events on the atmosphere [21]. Some studies revealed the signs of earthquake precursors linked to ionospheric perturbation [22–25]. The Global Ionospheric Maps (GIM) are also a valuable data source for detecting TEC anomalies caused by these seismic activities [26].

So far, there have been two approaches to studying ionospheric fluctuations related to seismic events. The first one is based on the physical mechanism of the seismic wave generation into the atmosphere [27]. The other relies on analyses of statistics and the probability of TEC anomalies in the epicenter regions and time of earthquake occurrences (mainshocks) [28]. However, there is no absolutely certain guarantee about the coincidence between observed ionospheric anomalies in the location and time of the earthquakes with other non-seismic activities. Monitoring TEC noise during periods of low solar activity to study the effect of seismic activities is a proper solution [19]; thus, many remarkable earthquakes resulting in TEC anomalies are

---

[1] https://directory.eoportal.org/web/eoportal.

skipped in investigations. Besides, there is also no common standard to measure ionospheric noise levels. Hence, applying Machine Learning (ML) and integrated data to distinguish TEC noise sources and extract TEC noise caused by seismic events will be carried out in this study.

ML has been a current trend applied to multidisciplinary research fields, especially for space weather forecasts and hazard warnings. The solar wind plasma and geomagnetic data have been used to predict TEC models in a few studies in the literature [29–31]. For example, Claudio Cesaroni et al. [32] used neural networks to predict global TEC at a daily sampling rate with the forecast accuracy of approximately 3 to 5 TECU. Xu Lin et al. [33] implemented the networks of convLSTM (convolutional Long Short-Term Memory) and PredRNN (Predictive Recurrent Neural Network) to correct errors of the delays. However, these criteria have not yet met the requirements of TEC anomaly detection related to seismic activity. Since TEC noise caused by seismic events often remains within a few minutes to a few hours and forecast accuracy of under 3 TECU can overcome TEC disturbances on the equator area or TEC variations at a low active time [34]. Other literary studies used the solar indices in their forecast models, such as Ap and F10.7, to correct ionospheric delays [35] or A.E. and SYM/H indices for TEC nowcasting [36]. Nevertheless, there has been no consistency in the selected indices among the studies.

Therefore, this study combines the optimization ML techniques with statistical hypothesis tests to determine suitable parameters related to the solar and geomagnetic activities for the TEC forecast models. These ML models will be the basis for separate TEC noise sources. Hence, we use the trained ML models to extract vertical TEC (VTEC) noise related to the Tonga volcanic eruption on 15 January 2022. Finally, we apply statistical and spectral techniques to analyze GNSS-VTEC noise at the International GNSS Service (IGS) stations nearby Tonga.

## 2  Study Area, Data, and Methodology

### 2.1  Study Area

The Tonga-Hunga Ha'apai includes small islands along the caldera rim in the Western-South Pacific Ocean. The Tonga volcano has experienced seven Holocene eruptive periods, with the first recorded eruption ~900 years ago [37]. For the latest period, it woke up by 20 December 2021 and ended after a massive explosion with a height of ~30 km at 04:10 UTC on 15 January 2022 [38]. The Tonga volcanic eruption triggered a tsunami with the waves observed thousands of miles away from the Caribbean and Alaska. After four minutes, a shallow earthquake of 5.7 $M_w$ occurred near the epicenter at 20.536°S and 175.382°W [39]. It is considered one of the few volcanoes tracked in detail and with different methods and technologies.
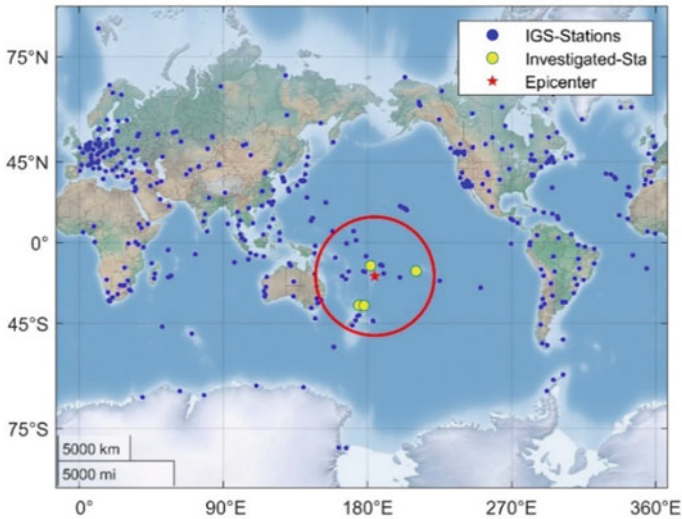
**Fig. 1** Investigated TEC anomalies related to the 2022 January Tonga volcano at the IGS stations

## 2.2 *Data*

We use the GNSS data from four IGS stations surrounding the epicenter of the 2022 January Tonga volcano to study the effects of these seismic events on the ionosphere (Fig. 1). These selected IGS stations must ensure conditions like being located within the radius of perception, the equivalent accuracy, and continuously monitored data. The stations FTNA and THTI are located in French Polynesia, and AUCK and WARK are in New Zealand. The GNSS data are available at the data center of the Crustal Dynamics Data Information System (CDDIS).[2] GNSS observations are the major initial data to compute the TEC time series for building forecast models.

Thirty-eight solar wind plasma and geomagnetic field parameters are analyzed to determine the best-suited predictors for the ML models (Table 4). These data are taken from the world data bank: the space weather prediction center NOAA,[3] USA; the world data center for Geomagnetism,[4] Kyoto, Japan; and the space weather live,[5] Belgium. In addition, the seismic data are collected from the data center GEOFON,[6] GFZ Potsdam, Germany and the U.S. geological survey center USGS.[7] Figure 2 shows eight main parameters of the geomagnetic field and solar wind. The level of solar activity is usually defined via indices such as the sunspot number and the

---

[2] https://cddis.nasa.gov/.

[3] https://www.swpc.noaa.gov/.

[4] http://wdc.kugi.kyoto-u.ac.jp/.

[5] https://www.spaceweatherlive.com/.

[6] https://geofon.gfz-potsdam.de/.
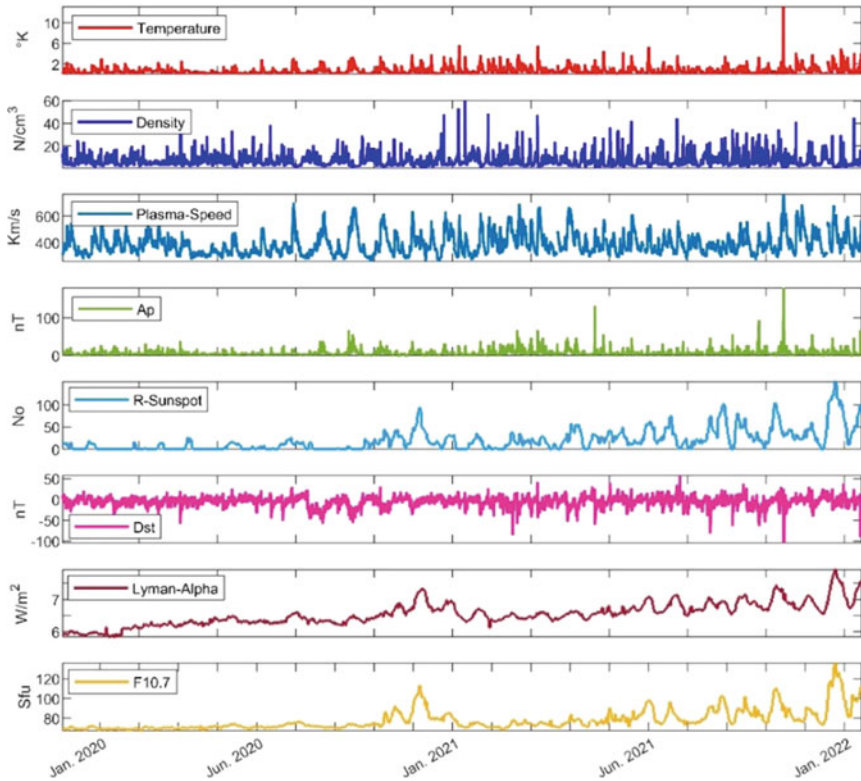
[7] https://www.usgs.gov/.

**Fig. 2** Eight specific features in 38 investigated parameters of the geomagnetic field and solar wind (plasma temperature, plasma density, Plasma speed, Ap index, R sunspot, Dst index, Lyman alpha, and F10.7 index)

solar radio flux at 10.7 cm (F10.7 index). Indices of global geomagnetic activity like Kp, Ap, Cp, and C9 are provided by the German Research Center for Geosciences (GFZ). Dst values of the disturbance storm time index are obtained from the world data center for Geomagnetism in Kyoto, Japan.

## 2.3 Methodology

Regression analysis is a mathematical method that describes the relationship between one or many independent variables and the dependent variable. In machine learning, the independent variables are factors to predict the dependent variable. Therefore, changes in the independent variables will result in changes in the dependent variable. Usually, only main factors should be included in the regression models to optimize

forecast performances and avoid biased conclusions. In this way, we use the parameters related to solar activity as the main factors in ML models to predict the VTEC time series at the IGS stations in Tonga. The differences between forecast models and actual TEC values will be considered anomalies caused by other factors.

ML models used for ionospheric anomaly detection must be sensitive enough to distinguish noise sources while remaining resistant to outliers. Hence, we clean data using filtering algorithms, with the moving window thresholds selected flexibly in the sampling rates and data characteristics.

The study uses integrated data with different units. The solar activity data also vary in an extensive value range from one-thousandth (e.g., Sigma alpha/proton ratio) to thousands (e.g., Plasma temperature), while VTEC time series change from a few (at midday) to hundreds (at midnight) TECU. It might make the initial assumption that higher ranging numbers have superiority of some sort. Furthermore, significant differences in value range among features can decrease convergence progress or saturate too fast for the algorithms based on gradient descent (e.g., Linear and Gaussian) and distance (e.g., SVM). Therefore, these input data should be re-scaled to fit regression models and push up processing speed.

As mentioned, the parameters associated with the solar activity will be predictors in the ML models. These parameters have different characteristics. The regression model's excess or lack of independent variables can decrease prediction performances. Therefore, determining the suitability of predictors in ML models is necessary, known as *"feature selection"*—one of the main hyperparameter tuning techniques in machine learning. The multiple regression analyses combined with statistical tests are applied to select the best-relevant parameters for ML models.

Training the forecast models is based on four ML algorithms, including Linear Regression, Support Vector Machine (SVM), Tree Ensemble, and Regression Trees using the ML toolbox in MATLAB® to select the optimal models. VTEC disturbances caused by seismic activities can last from some minutes to a few hours [34]. To capture the most negligible variations in the VTEC time series, we investigate two cases: (1) using hourly time series of two-year data and (2) using one-minute time series of the 15-day data to predict one day. We extract VTEC noise based on the trained ML models and analyze its physical characteristics by the spectrum method [40]. To this end, we apply Welch's algorithm to estimate the power spectral density [41] and the continuous wavelet transform (CWT) method to compute the spectral magnitude in GNSS-VTEC noise [42, 43].

Figure 3 shows the methodology mentioned above with three main steps. The first one is pre-processing with cleaning raw data, testing characteristics, and re-scaling data. The second step is feature selection using two statistical tests, analysis of variance (ANOVA), and Fisher test. The final step includes preparing input data, splitting data, training forecast models, optimization processing to get the highest performance models, and extracting and analyzing GNSS-TEC noise at the investigated IGS stations.
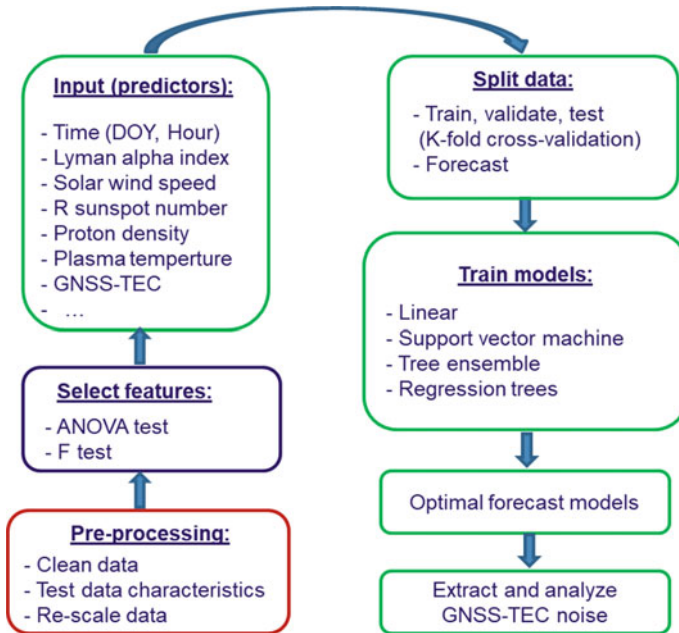
**Fig. 3** Flowchart of GNSS-TEC noise detection based on ML techniques and the integrated data of GNSS and solar activity

## 3    Results and Discussions

### 3.1    Data Pre-Processing

The Moving Median filters outliers at the same thresholds and sliding window size (Fig. 4). The Augmented Dickey-Fuller (ADF) tests the stationarity of the VTEC time series at the IGS stations, and the details are shown in Table 1. The absolute values of the ADF test (in bold italics) are larger than the critical values $t_{critical}$ for all statistical t-test levels (1, 5 and 10%) at the significant statistics (see Table 1). Hence, the VTEC time series appear to be stationary, and they can be used to train the forecast models to detect ionospheric anomalies.

### 3.2    Feature Selection

We apply statistical tests and analysis techniques to select the best-suited features from 42 input parameters, in which 38 parameters of the solar wind plasma and geomagnetic field, three parameters of time (Hour, Day of Year, and Year) and one lagged VTEC time series. F-test is used to measure the feature importance via the
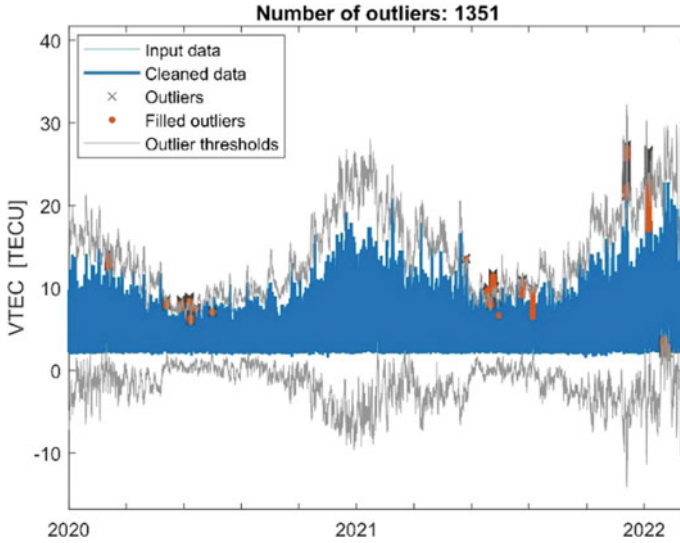
**Fig. 4** Filtering of outliers by the Moving Median algorithm for the VTEC time series (one-minute sampling rate) of the IGS station AUCK (New Zealand) with the confidence interval of 99.7% at the sliding window of 2880 (two-day data)

**Table 1** Test the stationarity of the VTEC time series at the IGS stations AUCK, FTNA, THTI, and WARK using the Augmented Dickey-Fuller algorithm

| Test levels (%) | AUCK | | FTNA | | THTI | | WARK | |
|---|---|---|---|---|---|---|---|---|
| | t | p-value | t | p-value | t | p-value | t | p-value |
| | *− 5.03* | 1.80E-05 | *− 5.63* | 8.95E-07 | *− 5.99* | 1.32E-07 | *− 4.73* | 7.21E-05 |
| 1 | − 3.43 | | − 3.43 | | − 3.43 | | − 3.43 | |
| 5 | − 2.86 | | − 2.86 | | − 2.86 | | − 2.86 | |
| 10 | − 2.57 | | − 2.57 | | − 2.57 | | − 2.57 | |

score ranking. Figure 5 shows the classification results of the univariate features at the IGS stations, with the nine highest-score features displayed on the horizontal axes of the graphs. The scores at the high-latitude stations (WARK and AUCK) reach up to 449.58 and 479.27, while those (THTI and FTNA) are only 217.87 and 268.40, respectively. This finding indicates that solar activity has a more significant impact on TEC observed at high-latitude stations. In addition, the effect of the variables like Kp, Ap, and Lyman alpha remains at the highest level for all the IGS stations, in which Ap and Kp are of identical scores because of their correlation. The detailed results of the feature importance score are presented in Table 4. Regression model predictors should be independent variables to improve processing speed and avoid biased conclusions. Hence, we employ analysis of variances (via ANOVA tests) to detect multicollinearity in the variables under consideration.
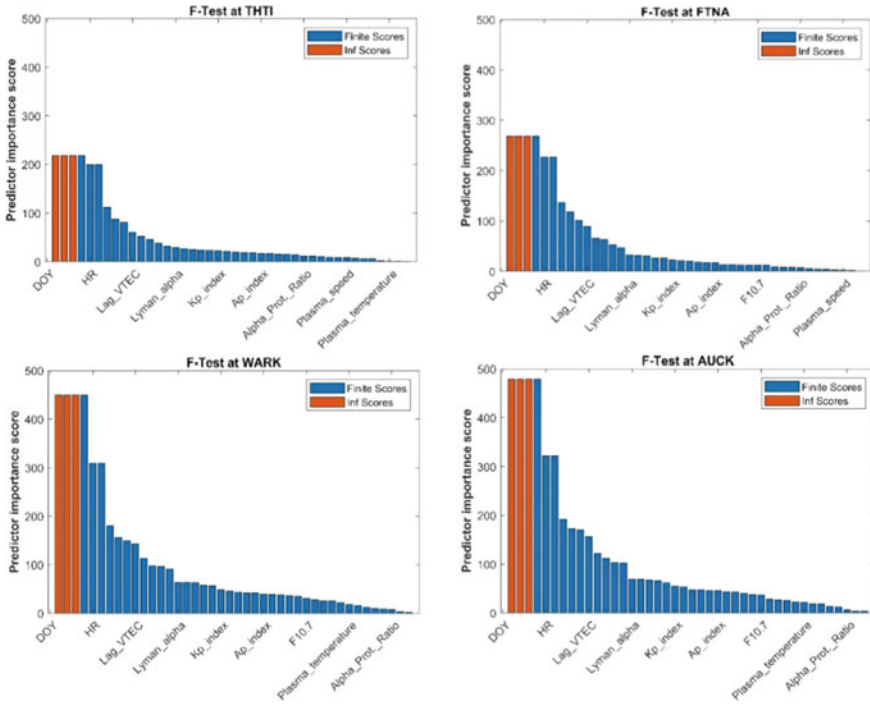
**Fig. 5** Univariate feature ranking for regression models using F-test at four IGS stations AUCK, FTNA, THTI, and WARK. The horizontal axis indicates the specific score features

Table 2 lists 20 features with significant statistics (p-values < 0.05) using the ANOVA tests. The coefficients (Beta) of the t-test show the correlations among the independent variables in the regression models. Tolerance and VIF indicate information of multicollinearity.
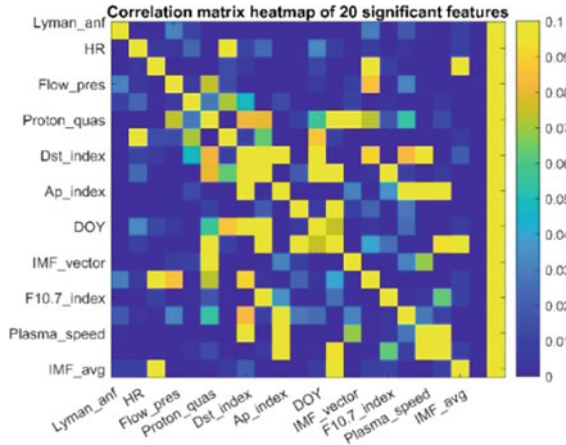
Four features have a high potential for multicollinearity, including plasma speed, flow pressure, IMF magnitude average, and Magnitude IMF vector. Besides, the features with tolerance in the interval from 0.1 to 0.2 should also be checked cross-correlation to enhance forecast performance. The correlation matrix heatmap reveals the pairs of the linear relevant features, such as Proton quasy invariant and plasma temperature, Dst (or Kp) and Ap index (Fig. 6). The detailed information on the correlation matrix of 20 significant features is shown in Table 5. The redundant variables (e.g., IMF magnitude vector, IMF magnitude average, plasma flow pressure) should be rejected before training forecast models.

Overall, there are two outstanding advantages of feature selection based on statistical tests. The first one is a clear classification of the importance of each feature in a regression model. This helps analysts decide which parameters should be used to train forecast models. The second advantage is to detect relevant features without training test models (trial steps) as other feature selection methods (e.g., K nearest neighbor,

**Table 2** ANOVA test of significant features at the IGS station WARK (New Zealand)

| Features | Coefficients (Beta) | t-test | p-values | Tolerance | VIF |
|---|---|---|---|---|---|
| Lag_VTEC | 0.86 | 212.73 | 0.00 | 0.83 | 1.20 |
| HR | − 0.05 | − 11.98 | 0.00 | 0.88 | 1.13 |
| Plasma_speed | − 0.12 | − 9.11 | 0.00 | 0.08 | 12.79 |
| Kp_index | 0.08 | 7.65 | 0.00 | 0.14 | 7.22 |
| Plasma_temperature | 0.07 | 7.58 | 0.00 | 0.14 | 6.99 |
| YEAR | − 0.05 | − 7.42 | 0.00 | 0.34 | 2.91 |
| Proton_quazy_invariant | − 0.04 | − 6.66 | 0.00 | 0.35 | 2.86 |
| BZ_GSE | − 0.07 | − 6.04 | 0.00 | 0.10 | 9.80 |
| DOY | − 0.04 | − 5.76 | 0.00 | 0.29 | 3.41 |
| Flow_pressure | 0.06 | 4.28 | 0.00 | 0.07 | 15.35 |
| Lyman_alpha | 0.04 | 3.94 | 0.00 | 0.13 | 7.91 |
| Sigma_Np | − 0.03 | − 3.90 | 0.00 | 0.33 | 3.07 |
| IMF_magnitude_avg | 0.15 | 3.78 | 0.00 | 0.01 | 118.63 |
| Proton_density | − 0.04 | − 3.28 | 0.00 | 0.08 | 13.27 |
| Dst_index | − 0.02 | − 3.27 | 0.00 | 0.53 | 1.89 |
| Plasma_lat_angle | − 0.01 | − 3.26 | 0.00 | 0.97 | 1.03 |
| F10.7_index | 0.03 | 2.83 | 0.00 | 0.13 | 7.43 |
| Ap_index | 0.02 | 2.30 | 0.02 | 0.17 | 6.05 |
| Magnitude_IMF_vector | − 0.08 | − 2.23 | 0.03 | 0.01 | 85.37 |
| R_sunspot | 0.02 | 2.12 | 0.03 | 0.22 | 4.63 |

**Fig. 6** Correlation matrix heatmap of 20 significant features using the ANOVA test at the IGS station WARK



Gaussian process regression, or neighborhood component analysis). Forecast performance will be significantly improved because performing trial steps via loops on all 42 features is time-consuming. In this study, statistical tests have pointed out the best-suited features to build the regression ML models: Plasma temperature (or plasma speed), proton density, Ap index (or Kp, Dst index), F10.7 index, Lyman alpha, and R sunspot. Together with the infinite variables (H.R., DOY, Lag VTEC), these six features will be used as the predictors (independent variables) in the regression ML models to forecast the GNSS-VTEC time series.

## *3.3 Detection and Analysis of GNSS-TEC Noise*

We investigate 15 mathematical models based on four ML methods (Linear Regression, Regression Trees, SVM, and Tree Ensemble) to predict the VTEC time series for one day. Optimization processing is performed to select the best models for GNSS-VTEC noise detection. Figure 7 presents two forecast models using the one-minute and hourly time series at the THTI station. VTEC prediction using the hourly time series produces the ML models with greater generalizability and robustness to extreme values and outliers. These models will have higher reliability for anomaly detection in deformation analyses and long-term predictions. Nevertheless, based on noise characteristics caused by seismic activities, assessments on both cases (one-minute and hourly time series) should be performed.

Table 3 lists the VTEC forecast accuracy (in root mean square error, RMSE) of the ML models at four IGS stations. The accuracy of the Ensemble algorithm outperforms others. The forecast performance based on the Ensemble reaches the highest, from 1.01 TECU at the WARK station (hourly) to 3.17 TECU at the FTNA station (hourly). In contrast, the linear algorithm shows the lowest, from 1.31 TECU at the AUCK station (hourly) to 5.07 TECU at the FTNA station (one minute). The accuracy of the
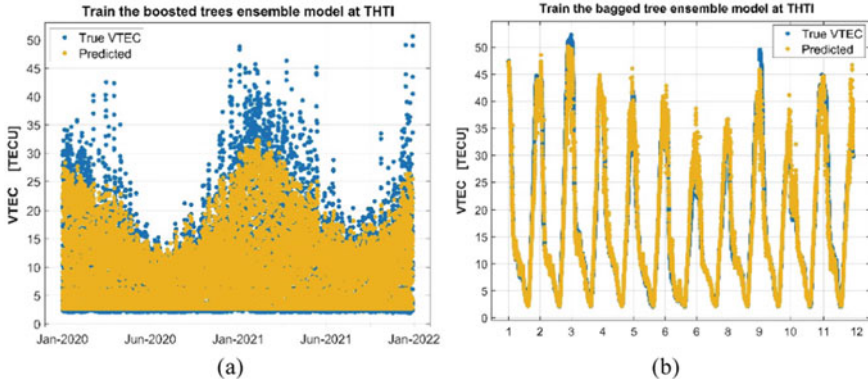
**Fig. 7** Training the GNSS-VTEC forecast model at station THTI: **a** the Boosted Tree Ensemble algorithm using data from 01 January 2020 to 01 January 2021 and **b** the Bagged Tree Ensemble algorithm using data from 01st to 12th January 2022

two offshore stations (FTNA and THTI) is lower than of the inshore stations (WARK and AUCK). It is likely due to poor input data quality, which indicates RMSE of the VTEC time series of 2.86, 2.67, 1.74, and 1.71 TECU for FTNA, THTI, WARK, and AUCK, respectively.

Forecasts of VTEC at the FTNA station (French Polynesia) on 15th January 2022 are shown in Fig. 8. Given the global seismic data (GEOFON and USGS), there were no other remarkable seismic events in the 7000-km radius (from the volcano epicenter in Tonga) within three days 13th, 14th, and 15th of January. The solar activity data have been included to predict the VTEC time series. As a result, the Tonga seismic events on 15th January 2022 are believed to cause the VTEC disturbances at the investigated IGS stations.

To eliminate the effect of systematic errors, we extract the VTEC noise based on the same ML algorithm for all the IGS stations (Fig. 9). The VTEC noise at the THTI

**Table 3** Accuracy of the ML models for one day forecast at four IGS stations

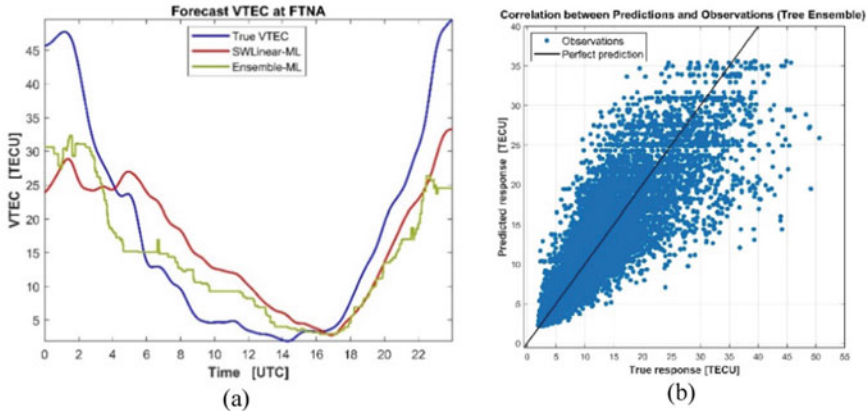| Methods | FTNA RMSE (TECU) | THTI RMSE (TECU) | AUCK RMSE (TECU) | WARK RMSE (TECU) | Input data | Sampling rate |
|---|---|---|---|---|---|---|
| Ensemble | 3.17 | 2.86 | 1.02 | 1.01 | 2 years | Hourly |
| Coarse tree | 3.81 | 3.49 | 1.22 | 1.21 | | |
| SVM | 3.67 | 3.56 | 1.22 | 1.27 | | |
| Linear | 3.83 | 3.78 | 1.31 | 1.32 | | |
| Ensemble | 2.92 | 2.67 | 1.20 | 2.07 | 12 days | One minute |
| Coarse tree | 3.14 | 3.55 | 1.53 | 1.33 | | |
| SVM | 4.54 | 5.28 | 2.38 | 2.65 | | |
| Linear | 5.07 | 5.89 | 2.86 | 2.86 | | |

**Fig. 8** Forecast of VTEC at the station FTNA on 15 January 2022 **(a)** and test the correlation between predictions and observations of the Tree Ensemble model **(b)**

station shows the highest variation, up to 34.47 TECU, followed in turn by FTNA, AUCK, and WARK with corresponding values of 25.09, 21.99, and 17.41 TECU. At the same time, the forecast accuracy (i.e., RMSE) ranges from 1.01 to 3.17 TECU. The TEC fluctuations reach up to 6.5 times the RMSE values of the ML models. This finding shows a correlation between the occurrence of the seismic and the ionosphere anomalies on 15th January. These VTEC anomalies occurred a few hours around the mainshock (at 4.10 UTC, 15 January 2022). However, no positive/negative linear relationship between the TEC fluctuation amplitudes and the distances from the IGS stations to the earthquake epicenter has been seen on the ML models in Fig. 9.

Spectral methods are used to analyze the GNSS-VTEC noise at the stations in Tonga's volcanic eruption region and assess VTEC anomalies before and after the mainshock (at 4.10 UTC, 15th January 2022).

Based on Welch's segment averaging estimation at the overlap of 50%, we determine the power spectral density (PSD), in which the power values are computed as follows:

$$y_{dB} = 10 \times \log_{10}(x) \tag{1}$$

where $x$ is the power spectral density computed by the Welch's method.

The frequencies ($f_i$) of the VTEC noise are converted into the normalized frequency ($F_s$) ranging from 0 to 1, to measure the variations of the power spectrum (Figs. 10 and 11) as follow:

$$F_s = \frac{f_i \times \pi}{1440} \tag{2}$$

where $f_i$ is the frequency of the VTEC noise at the IGS stations, and 1440 is the sample number in the VTEC time series.
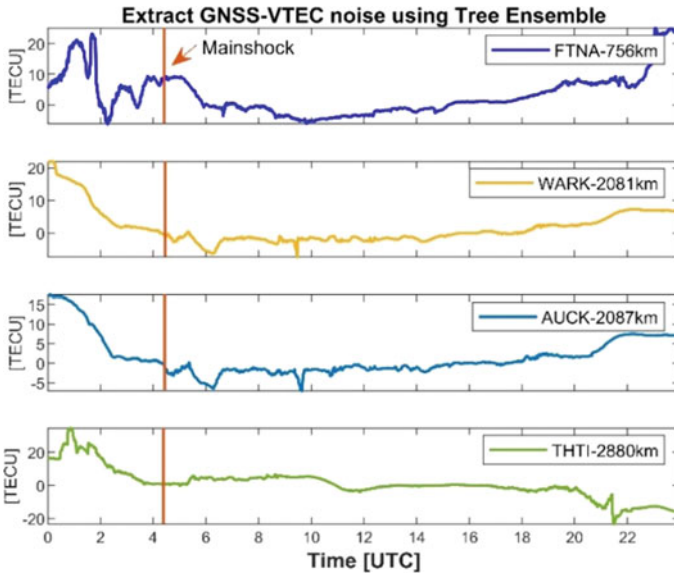
**Fig. 9** Extraction of the GNSS-VTEC noise using the Tree Ensemble models at four stations FTNA, WARK, AUCK, and THTI on 15 January 2022
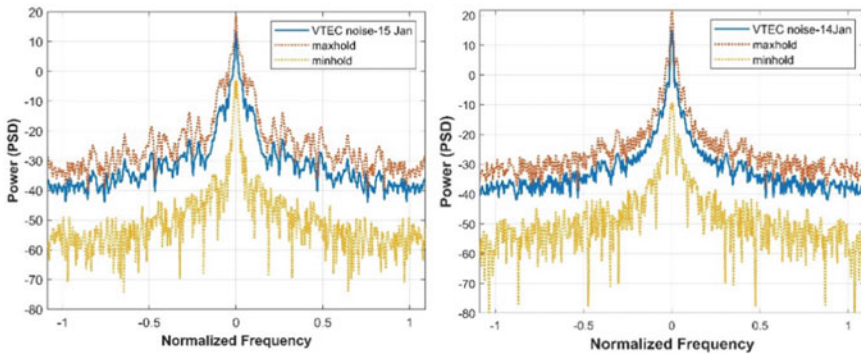


**Fig. 10** Applying the Welch algorithm to estimate the power spectral density of the GNSS-VTEC noise at the station FTNA on 14 and 15 January 2022

At the same normalized frequency, the PSD pattern of VTEC noise on the 15th is rougher than on the 14th of January (Fig. 10).

The fluctuations of PSD at the FTNA station on 15 January 2022 was more significant than others, ranging from $-85.18$ to $31.44$ (15th January 2022); $-77.16$ to $30.00$ (14th January 2022); and from $-78.97$ to $30.96$ (13th January 2022), respectively (Fig. 11). The pattern of the GNSS-VTEC noise variations on the volcanic eruption is denser and more significant compared to other days.
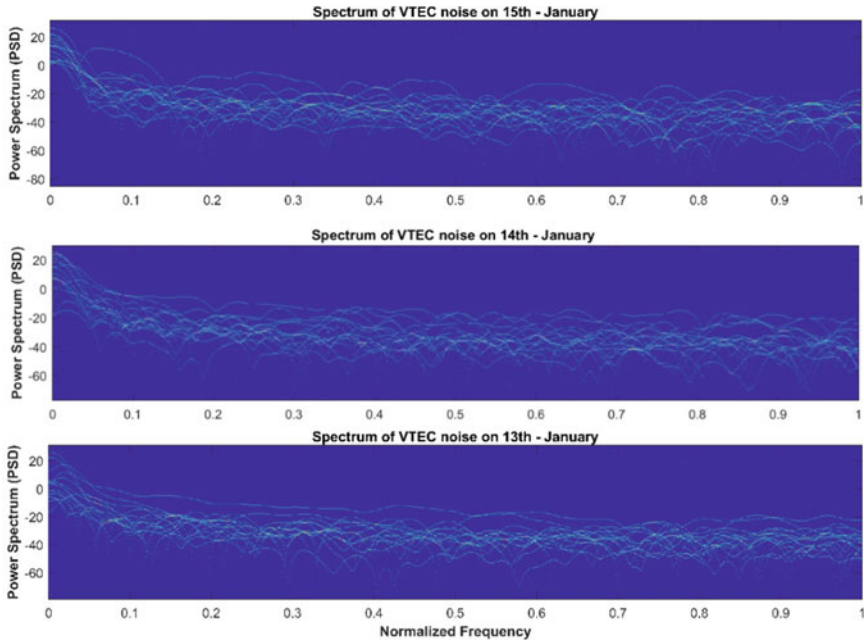
**Fig. 11** Power spectrum density of VTEC at the FTNA station on 15th (top), 14th (middle), and 13 (bottom) January 2022

The continuous wavelet transformation (CWT) on the same sampling frequency band $\left(F_s = \pi / 1440\right)$ is applied to compute the spectral magnitude of the GNSS-VTEC noise at the stations over 24 h. Figure 12 describes the magnitude scalograms of the VTEC noise at four IGS stations on 15th January, in which scalogram is the CWT absolute value. The spectrum magnitude at the FTNA station reaches the highest level (5%), followed by THTI, AUCK, and WARK (spectrum magnitude ranging from 2 to 3%). Besides, the scalogram maps also present earlier and more significant fluctuations for the offshore IGS stations (FTNA and THTI in French Polynesia) compared to the inshore stations (AUCK and WARK in New Zealand). This phenomenon may be the wave consonance of the earthquake and tsunami following the volcanic eruption in the ocean. However, more research is needed to gain a complete picture of the cause-effect relationship between time, space, and noise levels in seismic areas.
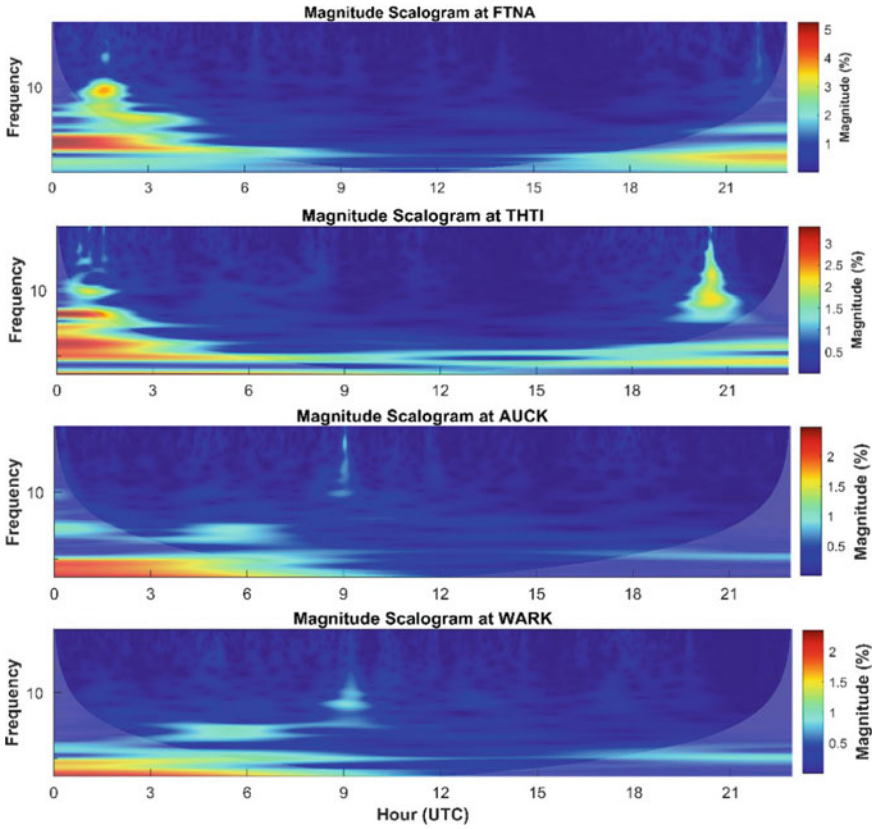
**Fig. 12** Magnitude scalogram maps of the GNSS-VTEC noise at four IGS stations FTNA, THTI, AUCK, and WARK on 15 January 2022

Although there are a few signs of seismic precursors on the scalogram maps at the offshore stations, the potential for earthquake prediction using GNSS-TEC data has remained low in terms of probability and statistics thus far. Besides, the seasonal characteristics of the TEC time series change in diurnal, annual, and 11-year solar cycles [44, 45]. Therefore, using the stationary time series of a few years to predict one day is an optimal solution for forecast performances to balance training time (or computation speed) and forecast accuracy. Nonetheless, further investigations on a time series longer than 11 years should be conducted to comprehensively assess the accuracy of TEC noise detection based on ML and integrated data of GNSS and solar activity.

# 4 Conclusion

This study has provided the results of the GNSS-TEC noise detection at four IGS stations associated with the Tonga volcano. Overall, the combination of regression ML techniques with integrated data of GNSS and the solar activity for TEC anomaly detection is a robust statistical solution. Depending on the input data quality, the accuracy of TEC noise detection over four investigated IGS stations ranges from ~ 1.0 to ~ 5.9 TECU. The Ensemble algorithm gets the highest performance (from 1.01 to 3.17 TECU), while Linear Regression is the least effective (1.32 to 5.89 TECU). Statistical tests play a crucial role in the hyperparameter tuning step to select the most relevant predictors. The ML-based forecast models using integrated data are potential applications for near real-time TEC anomaly warning and for adjustments of global ionospheric models in GNSS positioning. Extending investigations on ionospheric anomalies associated with seismic activities should be conducted for a better view of the cause-effect relationship between seismic events and other natural phenomena in the Earth's climate system.

**Conflict of Interests**
The authors declare that they have no competing interests.

**Authors' Contributions**
All of the authors have fair contributions.

# Appendix

**Table 4** F-test of the feature importance to select the best-fitted predictors for the regression ML models corresponding to four IGS stations AUCK, FTNA, THTI, and WARK

| Features (Predictors) | AUCK Scores | FTNA Scores | THTI Scores | WARK Scores |
|---|---|---|---|---|
| DOY (Day of year) | *Inf* | *Inf* | *Inf* | *Inf* |
| H.R. (Hour) | *Inf* | *Inf* | *Inf* | *Inf* |
| Lag_VTEC | *Inf* | *Inf* | *Inf* | *Inf* |
| Lyman_alpha | 479.27 | 268.40 | 217.87 | 449.58 |
| Kp_index | 321.47 | 227.27 | 199.24 | 309.34 |
| Ap_index | 321.47 | 227.27 | 199.24 | 309.34 |
| F10.7_index | 192.34 | 136.94 | 52.56 | 180.41 |
| Plasma_temperature | 173.27 | 63.10 | 80.56 | 156.56 |
| Alpha/Proton Density Ratio | 170.07 | 65.89 | 59.91 | 98.17 |
| Plasma_speed | 156.17 | 101.20 | 87.59 | 142.80 |
| R_sunspot | 122.22 | 89.63 | 24.98 | 113.19 |
| Sigma-Alpha/Proton_ratio | 112.66 | 65.89 | 59.91 | 98.17 |
| IMF_magnitude_avg | 103.27 | 26.44 | 26.96 | 96.15 |
| Sigma_T | 102.00 | 32.50 | 38.01 | 91.54 |
| Sigma_IMF_vector | 69.08 | 17.10 | 17.37 | 63.60 |
| Magnitude_IMF_vector | 68.92 | 18.21 | 17.17 | 63.07 |
| Flow_pressure | 67.52 | 22.68 | 14.91 | 63.83 |
| Sigma_V | 66.09 | 20.90 | 23.20 | 57.63 |
| RMS_BZ_GSE | 61.25 | 12.31 | 10.81 | 56.62 |
| Plasma_beta | 54.64 | 30.49 | 19.17 | 48.88 |
| Magnetosonic_mach_num | 53.18 | 31.62 | 22.57 | 46.36 |
| RMS_BY_GSE | 46.98 | 12.72 | 15.71 | 43.73 |
| Sigma_flow_latitude | 46.93 | 6.03 | 6.46 | 41.63 |
| Proton_quasy_invariant | 45.98 | 13.04 | 13.97 | 42.10 |
| Dst_index | 45.42 | 52.58 | 17.37 | 63.60 |
| Elecrtric_field | 43.30 | 26.85 | 31.70 | 39.02 |
| Alfven_mach_num | 42.33 | 12.19 | 11.78 | 38.27 |
| RMS_BX_GSE | 39.55 | 7.71 | 11.50 | 35.00 |
| Sigma_flow_longitude | 37.27 | 8.10 | 7.14 | 31.32 |
| Plasma_flow_latitude | 36.72 | 17.13 | 28.88 | 39.77 |
| BY_GSE | 28.07 | 2.61 | 9.29 | 27.89 |
| Bx_GSE/GSM | 26.79 | 12.99 | 24.15 | 25.42 |
| BY_GSM | 25.46 | 2.78 | 9.27 | 25.11 |
| YEAR | 22.07 | 47.07 | 19.02 | 10.06 |

(continued)

**Table 4** (continued)

| Features (Predictors) | AUCK Scores | FTNA Scores | THTI Scores | WARK Scores |
|---|---|---|---|---|
| BZ_GSM | 21.95 | 12.08 | 18.38 | 21.92 |
| RMS_magnitude | 18.38 | 4.97 | 0.89 | 16.45 |
| BZ_GSE | 18.06 | 8.62 | 2.40 | 18.14 |
| Plasma_flow_longitude | 13.03 | 0.43 | 1.24 | 8.58 |
| Proton_density | 12.23 | 20.16 | 9.04 | 12.57 |
| Long_Avg_IMF | 6.47 | 4.59 | 1.37 | 4.16 |
| Sigma_Np | 4.05 | 9.18 | 6.38 | 2.55 |
| Lat_Avg_IMF | 3.64 | 1.68 | 21.60 | 8.11 |

**Table 5** The correlation matrix of 20 significant features at station WARK

| Features | Feature codes | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Lyman_alpha | 1 | 1.000 | -0.022 | -0.004 | 0.031 | 0.010 | -0.070 | -0.050 | -0.132 | -0.680 | -0.028 | -0.005 | -0.643 | -0.097 | -0.005 | 0.030 | -0.492 | 0.018 | 0.001 | -0.021 | 0.009 |
| HR | 2 | -0.022 | 1.000 | 0.004 | 0.003 | 0.022 | -0.004 | 0.339 | 0.011 | 0.024 | -0.007 | -0.003 | 0.033 | 0.008 | -0.010 | 0.006 | -0.032 | 0.002 | -0.011 | -0.006 | 0.010 |
| Lat_Avg_IMF | 3 | -0.004 | 0.004 | 1.000 | -0.023 | -0.016 | -0.081 | 0.014 | 0.006 | -0.012 | -0.126 | -0.002 | 0.012 | -0.012 | -0.080 | 0.115 | 0.009 | -0.587 | -0.001 | -0.013 | 0.097 |
| Flow_pressure | 4 | 0.031 | 0.003 | -0.023 | 1.000 | 0.006 | 0.074 | 0.013 | -0.023 | -0.030 | -0.203 | -0.051 | 0.005 | -0.267 | 0.007 | 0.086 | -0.071 | 0.032 | -0.115 | -0.797 | -0.117 |
| Plasma_lat_angle | 5 | 0.010 | 0.022 | -0.016 | 0.006 | 1.000 | 0.028 | 0.071 | 0.048 | -0.037 | 0.013 | -0.009 | -0.038 | -0.027 | 0.010 | -0.008 | 0.002 | -0.015 | -0.002 | 0.004 | -0.018 |
| Proton_quazy_invariant | 6 | -0.070 | -0.004 | -0.081 | 0.074 | 0.028 | 1.000 | -0.015 | 0.082 | 0.080 | -0.055 | 0.010 | 0.056 | 0.281 | 0.120 | 0.073 | 0.017 | 0.055 | -0.318 | 0.013 | -0.173 |
| Lag_VTEC | 7 | -0.050 | 0.339 | 0.014 | 0.013 | 0.071 | -0.015 | 1.000 | 0.009 | 0.063 | -0.015 | -0.001 | 0.087 | 0.009 | -0.001 | -0.016 | -0.097 | -0.003 | -0.034 | -0.020 | 0.005 |
| Dst_index | 8 | -0.132 | 0.011 | 0.006 | -0.023 | 0.048 | 0.082 | 0.009 | 1.000 | 0.103 | 0.127 | -0.031 | 0.194 | -0.019 | -0.021 | 0.090 | -0.002 | 0.084 | 0.107 | -0.067 | 0.020 |
| YEAR | 9 | -0.680 | 0.024 | -0.012 | -0.030 | -0.037 | 0.080 | 0.063 | 0.103 | 1.000 | -0.010 | 0.019 | 0.672 | 0.101 | -0.008 | 0.003 | 0.201 | 0.009 | -0.014 | 0.005 | -0.010 |
| Ap_index | 10 | -0.028 | -0.007 | -0.126 | -0.203 | 0.013 | -0.055 | -0.015 | 0.127 | -0.010 | 1.000 | -0.066 | -0.005 | -0.054 | 0.031 | -0.733 | 0.036 | 0.102 | 0.128 | 0.211 | -0.032 |
| Sigma_Np | 11 | -0.005 | -0.003 | -0.002 | -0.051 | -0.009 | 0.010 | -0.001 | -0.031 | 0.019 | -0.066 | 1.000 | 0.005 | 0.235 | 0.008 | 0.022 | -0.002 | 0.028 | -0.022 | -0.157 | -0.017 |
| DOY | 12 | -0.643 | 0.033 | 0.012 | 0.005 | -0.038 | 0.056 | 0.087 | 0.194 | 0.672 | -0.005 | 0.005 | 1.000 | 0.074 | -0.009 | 0.003 | -0.026 | 0.025 | -0.014 | -0.045 | 0.006 |
| Plasma_temperature | 13 | -0.097 | 0.008 | -0.012 | -0.267 | -0.027 | 0.281 | 0.009 | -0.019 | 0.101 | -0.054 | 0.235 | 0.074 | 1.000 | -0.095 | 0.041 | 0.026 | -0.037 | -0.597 | 0.094 | 0.167 |
| Magnitude_IMF_vector | 14 | -0.005 | -0.010 | -0.080 | 0.007 | 0.010 | 0.120 | -0.001 | -0.021 | -0.008 | 0.031 | 0.008 | -0.009 | -0.095 | 1.000 | -0.042 | 0.000 | 0.031 | 0.069 | -0.008 | -0.960 |
| Kp_index | 15 | 0.030 | 0.006 | 0.115 | 0.086 | -0.008 | 0.073 | -0.016 | 0.090 | 0.003 | -0.733 | 0.022 | 0.003 | 0.041 | -0.042 | 1.000 | -0.009 | -0.120 | -0.206 | -0.156 | 0.009 |
| F10.7_index | 16 | -0.492 | -0.032 | 0.009 | -0.071 | 0.002 | 0.017 | -0.097 | -0.002 | 0.201 | 0.036 | -0.002 | -0.026 | 0.026 | 0.000 | -0.009 | 1.000 | 0.003 | 0.004 | 0.063 | 0.010 |
| BZ_GSE | 17 | 0.018 | 0.002 | -0.587 | 0.032 | -0.015 | 0.055 | -0.003 | 0.084 | 0.009 | 0.102 | 0.028 | 0.025 | -0.037 | 0.031 | -0.120 | 0.003 | 1.000 | 0.028 | 0.010 | -0.053 |
| Plasma_speed | 18 | 0.001 | -0.011 | -0.001 | -0.115 | -0.002 | -0.318 | -0.034 | 0.107 | -0.014 | 0.128 | -0.022 | -0.014 | -0.597 | 0.069 | -0.206 | 0.004 | 0.028 | 1.000 | 0.325 | -0.161 |
| Proton_density | 19 | -0.021 | -0.006 | -0.013 | -0.797 | 0.004 | 0.013 | -0.020 | -0.067 | 0.005 | 0.211 | -0.157 | -0.045 | 0.094 | -0.008 | -0.156 | 0.063 | 0.010 | 0.325 | 1.000 | 0.009 |
| IMF_magnitude_avg | 20 | 0.009 | 0.010 | 0.097 | -0.117 | -0.018 | -0.173 | 0.005 | 0.020 | -0.010 | -0.032 | -0.017 | 0.006 | 0.167 | -0.960 | 0.009 | 0.010 | -0.053 | -0.161 | 0.009 | 1.000 |

# References

1. Nathan, B., Christos, C.: Radio propagation and adaptive antennas for wireless communication links: terrestrial, atmospheric and ionospheric. Wiley, Hoboken, New Jersey (2007)
2. Robert, D.H., Jonh, K.H.: The high-latitude ionosphere and its effects on radio propagation. Cambridge University Press, New York (2003)
3. Hornbostel, A.: Propagation problems in satellite navigation. Proc. WFMN07 Chemnitz Ger., pp. 42–49, (2007). Retrieved from http://archiv.tu-chemnitz.de/pub/2007/0210/
4. Medžida, M., Natraš, R., Džana, H., Dževad, K.: Investigation of ionospheric variations and sudden disturbances as a source of GNSS errors and earthquake precursor. Sci. J. Civ. Eng. (2017). Retrieved from https://publik.tuwien.ac.at/files/publik_270748.pdf
5. Meyer-Vernet, N.: Basics of the Solar Wind. The United States of America by Cambridge University Press, New York (2007)

6. Gurnett, D.A., Amitava, B.: Introduction to Plasma Physics with Space and Laboratory Applications, vol. 17997, no. 383 (2017)
7. Hoffert, M.I.: The effects of solar variability on climate, vol. 19, no. C. The National Academies Press, Washington, D.C. (2012)
8. Ninla Elmawati Falabiba: The sun solar wind heliosphere. Springer, Dordrecht Heidelberg London New York (2019)
9. Jin, S., Jin, R., Liu, X.: GNSS Atmospheric Seismology. Springer Nature Singapore Pte Ltd. (2019)
10. Huang, C.Y., Helmboldt, J.F., Park, J., Pedersen, T.R., Willemann, R.: Ionospheric detection of explosive events. Rev. Geophys. **57**(1), 78–105 (2019). https://doi.org/10.1029/2017RG000594
11. Obayashi, T.: Upper atmospheric disturbances due to high altitude nuclear explosions. Planet. Space Sci. **10**, 47–63 (1963). https://doi.org/10.1016/0032-0633(63)90006-0
12. Park, J., Grejner-Brzezinska, D.A., Von Frese, R.R.B., Morton, Y., Gaya-Pique, L.R.: On using traveling ionospheric disturbances to detect underground nuclear tests. Inst. Navig. Int. Tech. Meet. (ITM) **2**, 1581–1589 (2012)
13. Mabie, J., Bullett, T., Moore, P., Vieira, G.: Identification of rocket-induced acoustic waves in the ionosphere. Geophys. Res. Lett. **43**(20), 11024–11029 (2016). https://doi.org/10.1002/2016GL070820
14. Lin, C.H., et al.: Ionospheric shock waves triggered by rockets. Ann. Geophys. **32**(9), 1145–1152 (2014). https://doi.org/10.5194/angeo-32-1145-2014
15. Heki, K., Fujimoto, T.: Atmospheric modes excited by the 2021 August eruption of the Fukutoku-Okanoba volcano, Izu–Bonin Arc, observed as harmonic TEC oscillations by QZSS. Earth Planets Sp., **74**(1) (2022). https://doi.org/10.1186/s40623-022-01587-5
16. Hasbi, A.M., et al.: Ionospheric and geomagnetic disturbances during the 2005 Sumatran earthquakes. J. Atmos. Solar Terr. Phys. **71**(17–18), 1992–2005 (2009). https://doi.org/10.1016/j.jastp.2009.09.004
17. Pulinets, S.A., Legen'ka, A.D., Hegai, V.V., Kim, V.P., Korsunova, L.P.: Ionosphere disturbances preceding earthquakes according to the data of ground based station of the vertical ionospheric sounding wakkanai. Geomagn. Aeron. **58**(5), 686–692 (2018). https://doi.org/10.1134/S0016793218050110
18. Korsunova, L.P., Khegai, V.V.: Possible short-term precursors of strong crustal earthquakes in japan based on data from the ground stations of vertical ionospheric sounding. Geomagn. Aeron. (2018). https://doi.org/10.1134/S0016793218010085
19. Shi, K., Liu, X., Guo, J., Liu, L., You, X., Wang, F.: Pre-earthquake and coseismic ionosphere disturbances of the Mw 6.6 Lushan earthquake on 20 April 2013 monitored by CMONOC. Atmos. (Basel) **10**(4), 1–21 (2019). https://doi.org/10.3390/ATMOS10040216
20. Zlotnicki, J., Li, F., Parrot, M.: Ionospheric disturbances recorded by DEMETER satellite over active volcanoes: from august 2004 to december 2010. Int. J. Geophys. **2013** (2013). https://doi.org/10.1155/2013/530865
21. Ishii, M.: Extreme Space Weather Research in Japan, vol. 1957. Elsevier Inc. (2018)
22. Akyol, A.A., Arikan, O., Arikan, F.: A machine learning-based detection of earthquake precursors using ionospheric data. Radio Sci. **55**(11), 1–21 (2020). https://doi.org/10.1029/2019RS006931
23. Sharma, G., Champati ray, P.K., Mohanty, S., Kannaujiya, S.: Ionospheric TEC modelling for earthquakes precursors from GNSS data. Quat. Int. **462**, 65–74 (2017). https://doi.org/10.1016/j.quaint.2017.05.007
24. Ulukavak, M., Yalcinkaya, M.: Precursor analysis of ionospheric GPS-TEC variations before the 2010 M7.2 Baja California earthquake. Geomatics Nat. Hazards Risk, **8**(2), 295–308 (2017). https://doi.org/10.1080/19475705.2016.1208684
25. Goto, S.I., Uchida, R., Igarashi, K., Chen, C.H., Kao, M., Umeno, K.: Preseismic ionospheric anomalies detected before the 2016 Taiwan earthquake. J. Geophys. Res. Sp. Phys. **124**(11), 9239–9252 (2019). https://doi.org/10.1029/2019JA026640
26. Tariq, M.A., Shah, M., Hernández-Pajares, M., Iqbal, T.: Pre-earthquake ionospheric anomalies before three major earthquakes by GPS-TEC and GIM-TEC data during 2015–2017. Adv. Sp. Res. **63**(7), 2088–2099 (2019). https://doi.org/10.1016/j.asr.2018.12.028

27. Nina, A. et al.: Variation in natural short-period ionospheric noise, and acoustic and gravity waves revealed by the amplitude analysis of a VLF radio signal on the occasion of the Kraljevo earthquake (Mw = 5.4). Sci. Total Environ. **710**, 136406 (2020). https://doi.org/10.1016/j.sci totenv.2019.136406

28. Zhao, S., Shen, X.H., Zhima, Z., Zhou, C.: The very low-frequency transmitter radio wave anomalies related to the 2010 Ms 7.1 Yushu earthquake observed by the DEMETER satellite and the possible mechanism. Ann. Geophys. **38**(5), 969–981 (2020). https://doi.org/10.5194/angeo-38-969-2020

29. Sun, W. et al.: Forecasting of ionospheric vertical total electron content (TEC) using LSTM networks. Proc. 2017 Int. Conf. Mach. Learn. Cybern. (ICMLC) **2**, 340–344 (2017). https://doi.org/10.1109/ICMLC.2017.8108945

30. Liu, L., Zou, S., Yao, Y., Wang, Z.: Forecasting global ionospheric TEC using deep learning approach. Sp. Weather **18**(11), 1–12 (2020). https://doi.org/10.1029/2020SW002501

31. Ruwali, A., Kumar, A.J.S., Prakash, K.B., Sivavaraprasad, G., Ratnam, D.V.: Implementation of hybrid deep learning model (LSTM-CNN) for ionospheric TEC forecasting using GPS data. IEEE Geosci. Remote Sens. Lett. **18**(6), 1004–1008 (2021). https://doi.org/10.1109/LGRS.2020.2992633

32. Cesaroni, C. et al.: Neural network based model for global total electron content forecasting. J. Sp. Weather Sp. Clim. **10** (2020). https://doi.org/10.1051/swsc/2020013

33. Lin, X. et al.: A Spatiotemporal Network Model for Global Ionospheric TEC Forecasting (2022)

34. Heki, K.: Advances in Ionospheric Research: Current Understanding and Challenges - Ionospheric Disturbances Related to Earthquakes. Wiley/AGU Online Library (2021)

35. Mallika, L.I., Ratnam, D.V., Raman, S., Sivavaraprasad, G.: Machine learning algorithm to forecast ionospheric time delays using Global Navigation satellite system observations. Acta Astronaut. **173**, 221–231 (2020). https://doi.org/10.1016/j.actaastro.2020.04.048

36. Zhukov, A., Sidorov, D., Mylnikova, A., Yasyukevich, Y.: Machine learning methodology for ionosphere total electron content nowcasting. Int. J. Artif. Intell. **16**(1), 144–157 (2018). https://doi.org/10.13140/rg.2.2.19349.83685

37. Global Volcanism Program | Raikoke. Retrieved from https://volcano.si.edu/volcano.cfm?vn=290250

38. Ripple effect_ What the Tonga eruption could mean for tsunami research _ National Oceanic and Atmospheric Administration. Retrieved from https://volcano.si.edu/volcano.cfm?vn=243040

39. GEOFON Program GFZ Potsdam: 9C Seismic Network. Retrieved from http://geofon.gfz-pot sdam.de/eqinfo/list.php?datemin=2022-01-15&datemax=2022-01-15&latmax=&lonmin=&lonmax=&latmin=&magmin=&fmt=html&nmax=

40. Stoica, P., Moses, R.L.: Spectral Analysis of Signals, vol. 4, no. 1. Pearson Prentice Hall (2005)

41. Welch, P.: The use of fast Fourier transform for the estimation of power spectra: a method based on time averaging over short, modified periodograms. IEEE Trans. Audio Electroacoust. **15**(2), 70–73 (1967). https://doi.org/10.1109/TAU.1967.1161901

42. Lilly, J.M.: Element analysis: a wavelet-based method for analyzing time-localized events in noisy time series. Proc. R. Soc. A Math. Phys. Eng. Sci. (2017). https://doi.org/10.1098/rspa.2016.0776

43. Lilly, J.M., Olhede, S.C.: Generalized morse wavelets as a superfamily of analytic wavelets. IEEE Trans. Signal Process. **60**(11), 6036–6041 (2012). https://doi.org/10.1109/TSP.2012.2210890

44. Hudson, T.S., Horseman, A., Sugier, J.: Diurnal, seasonal, and 11-yr solar cycle variation effects on the virtual ionosphere reflection height and implications for the Met Office's lightning detection system, ATDnet. J. Atmos. Ocean. Technol. **33**(7), 1429–1441 (2016). https://doi.org/10.1175/JTECH-D-15-0133.1

45. Zheng, W., et al.: Diurnal, seasonal, annual, and semi-annual variations of ionospheric parameters at different latitudes in East Asian sector during ascending phase of solar activity. Solar Terr. Phys. **3**(2), 45–53 (2017). https://doi.org/10.12737/22594