

REFERENCES

1. Pang- Ning tan, Michael Steinbach, Vipin kumar, Introduction to Data Mining (2016), Always Learning, Pearson.
2. Gordon S. Linoff, Michael J. A. Berry, Data Mining Techniques for Marketing and customer relationship management (2011), Wiley.
3. Barry Leventhal, An introduction to Data Mining and other techniques for advanced analytics (2010), Journal of Direct, Data and Digital Marketing Practice, Palgrave Macmillan UK, Volume 12, Issue 2, PP 137-153.
4. Mehmed Kantardzic, Data Mining: Concepts, Models, Methods and Algorithms, (2011) John Wiley & Sons.
5. M. A. Deshmukh, R. A. Gulhane, Importance of Clustering in Data Mining (2016), International Journal of Scientific & Engineering Research, Volume 7, Issue 2, PP 247-251.
6. Chandan K. Reddy, Charu C. Aggarwal, Data Clustering Algorithms and Applications (2013), Chapman and Hall/CRC.
7. Mirkin, Boris, Core Data Analysis: Summarization (2019), Correlation and Visualization, Undergraduate Topics in Computer Science Series, Springer.
8. Chengqi Zhang Shichao Zhang Association Rule Mining (2002), Springer.
9. Modi, Krishna, Outlier Analysis Approaches in Data Mining (2016). International Journal of Innovative Research in Technology.
10. Jiawei Han, Michelin Kamber, Jian Pei, Data mining concepts and techniques (2012), Morgan Kaufmann.
11. Parneet KaurManpreet SinghGurpreet Singh Josan, Classification and prediction-based data mining algorithms to predict slow learners in education sector (2015), 3rd International Conference on Recent Trends in Computing, Procedia Computer Science, Volume 57, PP 500-508.
12. K Prasanna Jyothi, Dr R SivaRanjani, Dr Tusar Kanti Mishra, S Ranjan Mishra, A Study of Classification Techniques of Data Mining Techniques in Health-Related Research (2017), International Journal of Innovative Research in Computer and Communication Engineering, Volume 5, Issue 7, PP 13779-13786.

13. Amirah Mohamed Shahiri, Wahidah Husain, Nur'aini Abdul Rashid, A Review on Predicting Student's Performance Using Data Mining Techniques (2015), *Procedia Computer Science*, Volume 72, Special Issue, PP 414-422.
14. Alaa F. Sheta, Sara Elsir M. Ahmed, Hossam Faris, A Comparison between Regression, Artificial Neural Networks and Support Vector Machines for Predicting Stock Market Index (2015), *International Journal of Advanced Research in Artificial Intelligence*, Volume 4, Issue 7, PP 55-63.
15. Festim Halili, Avni Rustemi, Predictive Modeling: Data Mining Regression Technique Applied in a Prototype (2016), *International Journal of computer science and mobile computing*, Volume 5, Issue 8, PP 207-215.
16. John O. Rawlings, Sastry G. Pantula, David A. Dickey, *Applied Regression Analysis: A Research Tool* (1998), Springer.
17. Philippe Esling, Carlos Agon. Time-series data mining. *ACM Computing Surveys*, Association for Computing Machinery, 2012, Volume 45, Issue 1 205.
18. Kenneth D. Lawrence, Stephan Kudyba, Ronald K. Klimberg, *Data Mining Methods and Applications* (2007), Auerbach Publications.
19. Velicer, Wayne & Fava, Joseph. *Time Series Analysis*. (2003).
20. Godwin, Harold & Okafor, Christian. Modified Trend and Seasonal Time Series Analysis for Operations: A Case Study of Soft Drink Production (2012). *International Journal of Engineering Research in Africa*. 7. 63-72.
21. Juang WC, Huang SJ, Huang FD, Cheng PW, Wann SR. Application of time series analysis in modelling and forecasting emergency department visits in a medical centre in Southern Taiwan. *BMJ Open*. 2017 Dec 1;7(11): e018628.
22. Machiwal, D., Jha, M.K. *Methods for Time Series Analysis*. In: *Hydrologic Time Series Analysis: Theory and Practice*. (2012). Springer, Dordrecht.
23. Shaltami, Osama & Bustany, Ilas. *Water quality – A review* (2021).
24. Dwivedi, Anil. *Researches in Water Pollution: A Review*. (2017).
25. Ahmed, Shahid and Ismail, Saba, *Water Pollution and its Sources, Effects & Management: A Case Study of Delhi* (2018), *International Journal of Current Advanced Research*, 07(2), pp. 10436-10442.
26. Yang, Y., Yan, B. & Shen, W. Assessment of point and nonpoint sources pollution in Songhua River Basin (2010), Northeast China by using revised water quality model. *Chin. Geogr. Sci.* 20, 30–36.

27. Dhote, Jayashree & Ingole, Sangita & Chavhan, Dr. Arvind. Review on Waste Water Treatment Technologies (2012). International Journal of Engineering Research and Technology.
28. Zia, Shaikh & Graham, D. & Dolfing, Wastewater Treatment: Biological. (2013).
29. Das Kangabam, R., Bhoominathan, S.D., Kanagaraj, S. et al. Development of a water quality index (WQI) for the Loktak Lake in India (2017). Appl Water Sci 7, 2907–2918
30. Md. Galal Uddin, Stephen Nash, Agnieszka I. Olbert, A review of water quality index models and their use for assessing surface water quality (2021), Ecological Indicators, Volume 122, 107218.
31. Howladar, M.F., Al Numanbakth, M.A. & Faruque, M.O. An application of Water Quality Index (WQI) and multivariate statistics to evaluate the water quality around Maddhapara Granite Mining Industrial Area (2018), Dinajpur, Bangladesh. Environ Syst Res 6, 13.
32. Adelagun, Ruth & Etim, Emmanuel & Godwin, Oko Emmanuel. Application of Water Quality Index for the Assessment of Water from Different Sources in Nigeria (2021).
33. Algalil, Fahd Abd, Aldhyani, Theyazn H. H, Al-Yaari, Mohammed, Alkahtani, Hasan, Maashi, Mashael, Water Quality Prediction Using Artificial Intelligence Algorithms, 2020 6659314.
34. Tiyasha, Tung, T.M. & Yaseen, Z.M. Deep Learning for Prediction of Water Quality Index Classification: Tropical Catchment Environmental Assessment. (2021). Nat Resour Res 30, 4235–4254
35. Sarker, I.H. Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and Research Directions (2021). SN COMPUT. SCI. 2, 420.
36. Bank, Dor & Koenigstein, Noam & Giryes, Raja. Autoencoders. (2020).
37. Zhang, Yue & Liu, Fangai. An Improved Deep Belief Network Prediction Model Based on Knowledge Transfer. (2020). Future Internet. 12. 188. 10.3390/fi12110188.
38. Alzubaidi, L., Zhang, J., Humaidi, A.J. et al. Review of deep learning: concepts, CNN architectures, challenges, applications, future directions (2021). J Big Data 8, 53.
39. Alex Sherstinsky, Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) network, Physica D: Nonlinear Phenomena, 2020, Volume 404, 132306.
40. Benjamin Lindemann, Timo Müller, Hannes Vietz, Nasser Jazdi, Michael Weyrich, A survey on long short-term memory networks for time series prediction, 2021, Procedia CIRP, Volume 99, pp 650-655.

41. Shen, Guizhu & Tan, Qingping & Haoyu, Zhang & Zeng, Ping & Xu, Jianjun. Deep Learning with Gated Recurrent Unit Networks for Financial Sequence Predictions. (2018). *Procedia Computer Science*. 131. 895-903.
42. Fang, Yong & Yang, Shaoshuai & Zhao, Bin & Huang, Cheng. Cyberbullying Detection in Social Networks Using Bi-GRU with Self-Attention Mechanism. (2021). *Information*. 12. 171. 10.3390/info1204017.
43. Bryan Lim, Sercan O. Arik, Nicolas Loeff, Tomas Pfister, Temporal Fusion Transformers for interpretable multi-horizon time series forecasting (2019).
44. Hosna, A., Merry, E., Gyalmo, J. et al. Transfer learning: a friendly introduction (2022). *J Big Data* 9, 102.
45. Day, O., Khoshgoftaar, T.M. A survey on heterogeneous transfer learning (2017). *J Big Data* 4, 29
46. Mignone, Paolo & Pio, Gianvito & Džeroski, Sašo & Ceci, Michelangelo. Multi-task learning for the simultaneous reconstruction of the human and mouse gene regulatory networks. (2020). *Scientific Reports*. 10. 22295. 10.1038/s41598-020-78033-7.
47. Wu, Qingyao, et al. Online Transfer Learning with Multiple Homogeneous or Heterogeneous Source,s 2017. *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 7, pp. 1494–507. *IEEE Xplore*.
48. Gupta, Jaya, et al. Deep Learning (CNN) and Transfer Learning: A Review. *Journal of Physics: Conference Series*, 2022. vol. 2273, no. 1, p. 012029.
49. Shrestha, AK & Basnet, N., The Correlation and regression analysis of physicochemical parameters of River water for the evaluation of percentage contribution to electrical conductivity, 2018, *Journal of Chemistry*, pp. 1-9.
50. Mohammad Mirzavand & Reza Ghazavi, A Stochastic Modelling Technique for Groundwater Level Forecasting in an Arid Environment Using Time Series Methods, 2014, *Water Resour Manage* 29: 1315-1328.
51. Xuedi Zhang, Hui Qian, Jie Chen and Liang Qiao, Assessment of Groundwater Chemistry and Status in a Heavily Used Semi-Arid Region with Multivariate Statistical Analysis, 2014, *Water* 6(8):2212-2232.
52. Semko Rashid, Milad Mohammadan, koorosh Azizi, Prediction of Groundwater level Fluctuation using ANN and ANFIS in Lailakh Plain, 2015, *Journal of renewable Natural Resource Bhutan* 3.1, pp: 77-84.

53. Batool A., Samad N., Kazmi S. S., Ghufuran M. A., Imad S., Shafqat M. and Mahmood T. Spring water quality and human health: an assessment of natural springs of margalla hills 2018. Islamabad zone- III. *Int. J. Hydrol.*, 2(1): 41–46.
54. Sakaa, B., Brahmia, N., Chaffai, H. and Hani, A. Assessment of water quality index in an unmonitored river basin using multilayer perceptron neural networks and principal component analysis. *Desalin. 2020. Water Treat.*, 200: 42-54.
55. Setshedi, K. J., Mutingwende, N. and Ngqwala, N. P. The use of artificial neural networks to predict the physicochemical characteristics of water quality in three district municipalities, Eastern Cape Province, 2021. South Africa. *Int. J. Environ. Res. Public Health*, 18(10): 5248-56.
56. Singh, B., Sihag, P., Singh, V. P., Sepahvand, A. and Singh, K. Soft computing technique-based prediction of water quality index. 2021. *Water Supply*, 21(8): 4015-4029.
57. Kulisz, M., Kujawska, J., Przysucha, B. and Cel, W. Forecasting water quality index in groundwater using artificial neural network. 2021. *Energies*, 14(18):58-75.
58. A. Solanki, H. Aggarwal, and K. Khare, Predictive Analysis of Water Quality Parameters using Deep Learning, 2015. *International Journal of Computer Applications*, vol. 125, no. 9, pp. 0975-8887.
59. U. Ahmed, R. Mumtaz, H. Anwar, A.A. Shah, R. Irfan, J. Garca-Nieto, Efficient water quality prediction using supervised machine learning, (2019) *Water* 11, 2210.
60. C.V. Sillberg, P. Kullavanijaya, O. Chavalparit, Water quality classification by integration of attribute realization and support vector machine for the chao phraya river, *Journal of Ecological Engineering* 22 (2021), 70–86.
61. Wang, Y., Zhou, J., Chen, K., Wang, Y., & Liu, L. (2017). Water quality prediction method based on LSTM neural network. In 2017 12th International Conference on Intelligent Systems and Knowledge Engineering (ISKE) (pp. 1-5). IEEE.
62. Zhenbo Li, Fang Peng, Bingshan Niu, Guangyao Li, Jing Wu, Zheng Miao, Water Quality Prediction Model Combining Sparse Auto-encoder and LSTM Network, *IFAC-PapersOnLine*, Volume 51, Issue 17, 2018, Pages 831-836.
63. Xu J, Wang K, Lin C, Xiao L, Huang X, Zhang Y. FM-GRU: A Time Series Prediction Method for Water Quality Based on seq2seq Framework. *Water*. 2021; 13(8):1031.
64. M. Yilma, Z. Kiflie, A. Windsperger, N. Gessese, Application of artificial neural network in water quality index prediction: a case study in little Akaki River, Addis Ababa, Ethiopia, *Modeling Earth Systems and Environment* 4 (2018), 175-187.

65. Srivastava, A., Cano, A. Analysis and forecasting of rivers pH level using deep learning (2022). *Prog Artif Intell* 11, 181–191.
66. G Tan, J Yan, C Gao, and S Yang, Prediction of water quality time series data based on least squares support vector machine, *Procedia Engineering*, Vol. 31, 2012, pp. 1194-1199.
67. WC Leong, A Bahadori, J Zhang, and Z Ahmad, Prediction of water quality index (WQI) using support vector machine (SVM) and least square-support vector machine (LS-SVM), *International Journal of River Basin Management*, Vol. 19, 2021, pp. 149-156.
68. Aldhyani, T. H. H., Al-Yaari, M., Alkahtani, H. & Maashi, M. Water quality prediction using artificial intelligence algorithms. *Applied Bionics and Biomechanics* 2020, 6659314 (2020).
69. Yang, Y. et al. A study on water quality prediction by a hybrid CNN-LSTM model with attention mechanism. *Environmental Science and Pollution Research* (2021) .
70. Kamaraj M, Rangarajan S (2022) Predicting the future land use and land cover changes for Bhavani basin, Tamil Nadu, India, using QGIS MOLUSCE plugin. *Environ Sci Pollut Res* (2022).
71. Arunkumar R, Thambusamy, Velmurugan (2021) An exploratory data analysis process on groundwater quality data. 54:41–48.
72. Haghiabi, Amir Hamzeh, et al. ‘Water Quality Prediction Using Machine Learning Methods’. *Water Quality Research Journal*, vol. 53, no. 1, Feb. 2018, pp. 3–13.
73. Ubah, J.I., Orakwe, L.C., Ogbu, K.N. et al. Forecasting water quality parameters using artificial neural network for irrigation purposes. *Sci Rep* 11, 24438 (2021).
74. Kargar, K. et al. Estimating longitudinal dispersion coefficient in natural streams using empirical models and machine learning algorithms. *Eng. Appl. Comput. Fluid Mech.* 14, 311–322 (2020).
75. Alizadeh, M. J. et al. Effect of river flow on the quality of estuarine and coastal waters using machine learning models. *Eng. Appl. Comput. Fluid Mech.* 12(1), 810–823 (2018).
76. WHO 2019 Drinking-water <https://www.who.int/news-room/factsheets/detail/drinking-water>.
77. Basant, N., Gupta, S., Malik, A., Singh, K.P., 2010. Linear and nonlinear modelling for simultaneous prediction of dissolved oxygen and biochemical oxygen demand of the surface water -A case study. *Chemometr.Intellig.Lab.Syst.*104,172–180.
78. Heddam, S., Kisi, O., 2018. Modelling daily dissolved oxygen concentration using least square support vector machine, multivariate adaptive regression splines and M5 model tree. *J. Hydrol.* 559, 499–509.

79. Li, G., 2006. Stream temperature and dissolved oxygen modelling in the Lower Flint River Basin. PhD Dissertation. University of Georgia, Athens, GA.
80. Wang, Y., Zhou, J., Chen, K., Wang, Y., & Liu, L. (2017). Water quality prediction method based on LSTM neural network. In 2017 12th International Conference on Intelligent Systems and Knowledge Engineering (ISKE) (pp. 1-5). IEEE.
81. Li L, Jiang P, Xu H, Lin G, Guo D, Wu H (2019) Water quality prediction based on recurrent neural network and improved evidence theory: a case study of Qiantang River, China. *Environ Sci Pollut Res* 26(19): 19879–19896.
82. Li L, Jiang P, Xu H, Lin G, Guo D, Wu H (2019) Water quality prediction based on recurrent neural network and improved evidence theory: a case study of Qiantang River, China. *Environ Sci Pollut Res* 26(19): 19879–19896
83. Liu J, Yu C, Hu Z, Zhao Y, Bai Y, Xie M, Luo J (2020) Accurate prediction scheme of water quality in smart mariculture with deep Bi-S-SRU learning network. *IEEE Access* 8:24784–24798
84. Yahya A, Saeed A, Ahmed AN, Binti Othman F, Ibrahim RK, Afan HA, Elshafie A (2019) Water quality prediction model-based support vector machine model for ungauged river catchment under dual scenarios. *Water* 11(6):1231.
85. L. Hu, C. Zhang, C. Hu, and G. Jiang, Use of grey system for assessment of drinking water quality: a case study of Jiaozuo city, China, *Advances in Grey Systems Research*, Springer Berlin Heidelberg, pp. 469-478, 2010
86. M. A. Tirabassi, A statistically based mathematical water quality model for a non-estuarine river system1 *JAWRA Journal of the American Water Resources Association*, Vol. 7, pp. 1221-1237, December 1971.
87. H. Liao and W. Sun. Forecasting and evaluating water quality of Chao Lake based on an improved decision tree method. *Procedia Environmental Sciences* 2 (2010): 970-979
88. M. Valdivia, D.W. Graham, and D. Werner, Climatic, Geographic and Operational Determinants of Trihalomethanes (THMs) in Drinking Water Systems Scientific reports, Vol. 6, 2016.
89. M Tarique, H. Khaleeq, and A. A. ElNour, A Reliable Wireless System for Water Quality Monitoring, Vo. 8, No. 3, 2016.
90. Y. Khan and C. S. See, Predicting and Analyzing Water Quality using Machine Learning: A Comprehensive Model, *IEEE Long Island Systems, Applications and Technology Conference (LISAT)*, 2016.

91. X. Li and J. Song, A New ANN-Markov Chain Methodology for Water Quality Prediction, International Joint Conference on Neural Networks, pp. 12-17 July 2015
92. A. Solanki, H. Aggarwal, & K. Khare, Predictive Analysis of Water Quality Parameters using Deep Learning, International Journal of Computer Applications, vol. 125, no. 9, pp. 0975-8887, Access from Google Scholar, Sept. 2015.
93. Jayti Bhatt, Jignesh Patoliya, Iot Based Water Quality Monitoring System, IRFIC, 21feb,2016.
94. Dong He, Li-Xin Zhang, "The Water Quality Monitoring System Based on WSN,"Institute of Mechanical and electronic information, China University of Geosciences (WuHan), WuHan,China, pp. 3661-3664, 2012.
95. Herzog, S., Worg " otter, " F., Parlitz, U., 2018. Data-Driven Modeling and Prediction of Complex Spatio-Temporal Dynamics in Excitable Media. *Front. Appl. Math. Stat.* 4
96. Tung Tiyasha, T.M., Yaseen, Z.M., 2020. A survey on river water quality modelling using artificial intelligence models: 2000–2020. *J. Hydrol.* 585, 124670.
97. Kannan V, Ramesh R, Sasikumar C. 2005. Study on groundwater characteristics and the effects of discharged effluents from textile units at Karur district. *Journal of Environmental Biology*, 26(2):269-272.
98. Zhang X, Qian H, Chen J, Qiao L. 2014. Assessment of Groundwater Chemistry and Status in a Heavily Used Semi-Arid Region with Multivariate Statistical Analysis. *Water*. 6(8):2212-2232.
99. Ubah, J.I., Orakwe, L.C., Ogbu, K.N. et al. Forecasting water quality parameters using artificial neural network for irrigation purposes. *Sci Rep* 11, 24438 (2021).
100. Haghiabi, Amir Hamzeh, et al. 'Water Quality Prediction Using Machine Learning Methods'. *Water Quality Research Journal*, vol. 53, no. 1, Feb. 2018, pp. 3–13.
101. Krishan G, Singh S, Kumar CP, Gurjar S, Ghosh NC. 2016. Assessment of Water Quality Index (WQI) of Groundwater in Rajkot District, Gujarat, India. *Journal of Earth Science & Climatic Change*. 7(3):1000341.

LIST OF PUBLICATIONS

Papers Published in International Conference Proceedings

1. Jitha P Nair, and M. S. Vijaya. Predictive Models for River Water Quality Using Machine Learning and Big Data Techniques - A Survey. 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS), IEEE, 2021, pp. 1747-53. **(SCOPUS INDEXED)**
2. Jitha P Nair, M S Vijaya. Exploratory Data Analysis of Bhavani River Water Quality Index Data. In: Kumar, S., Hiranwal, S., Purohit, S.D., Prasad, M. (eds) Proceedings of International Conference on Communication and Computational Technologies. Algorithms for Intelligent Systems. Springer, Singapore. **(WEB OF SCIENCE INDEXED)**
3. Jitha P Nair, and M. S. Vijaya. River Water Quality Prediction and Index Classification Using Machine Learning. Journal of Physics: Conference Series, vol. 2325, no. 1, Aug. 2022, p. 012011. **(SCOPUS INDEXED)**

Papers Published in International Journals

1. Jitha P Nair, and M. S. Vijaya. Analysing and Modelling Dissolved Oxygen Concentration using Deep Learning Algorithm, International Journal of Mechanical Engineering, Volume No.7, pp. 12-22. **(SCOPUS INDEXED)**
2. Jitha P Nair, and Vijaya, M. S. (2023). Enhanced Water Quality Prediction Model with Seasonal Time Series Data. Journal of Data Acquisition and Processing, 38(1), 1283. **(SCOPUS INDEXED)**
3. Jitha P Nair, and Vijaya, M. S. (2023) Design and development of efficient water quality prediction models using variants of recurrent neural networks. European Chemical Bulletin. Vol.No.12 Special issue 5. **(SCOPUS INDEXED)**
4. Jitha P Nair, and Vijaya, M. S. (2023) Temporal Fusion Transformer: A Deep Learning Approach for Modelling and Forecasting River Water Quality Index. Nair. International Journal of Intelligent Systems and Applications in Engineering, 11(10s), 277-293. **(SCOPUS INDEXED)**

Papers in Review - International Journals

1. TFT Architecture and RNN Variants for Water Quality Prediction of Bharathapuzha River - International Journal of Information Technology
2. Transfer Learning for Improved Deep Neural Network Models in River Water Quality Index Prediction -Journal of Network and Computer Application.
3. Transfer Learning and Temporal Fusion Transformer for Enhanced River Water Quality Index Prediction - Expert System with Application.
4. Heterogeneous Transfer Learning for Modelling and Forecasting River Water Quality Index - Journal of Cloud Computing

APPENDIX A- DATASETS

WQI-PCA Dataset

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF		
1	Temp	pH	Conduct	Turbidity	Phenol	Total	Chloride	CO3	TKM	Ammon	Hardnes	Ca	Hard	Mg	Hard	Sulphate	Sodium	TSS	TDS	FDS	Phospho	Boron	Potassi	BOD	Fluoride	Nitrate	TC	FC	Stalon	Latitude	Longitude	Year	DO	WQI
2	25	7.15	340	2	0	11	21	4	0.1	0.25	188	74	44	12	27.1	300	190	174	0.11	0.1	2.67	0.89	0.12	1.1	88	80	1	76.591	10.931	2016	6.99	76.79		
3	24	7.46	339	2	0	10	21	3.9	0.1	0.25	188	74	44	12.5	27.1	300	190	174	0.11	0.1	2.67	0.89	0.18	1.1	98	80	1	76.591	10.931	2016	6.97	77.1		
4	25	7.5	339	2	0	12	22	4	0.09	0.25	189	74.5	44	12	27.2	300	189	174	0.11	0.1	2.66	0.89	0.18	1.1	118	80	1	76.591	10.931	2016	6.81	77.12		
5	25	7.18	340	2	0	11	21	3.9	0.1	0.25	189	74.5	43.5	12	27.2	300	189	174.5	0.11	0.1	2.66	0.88	0.18	1.1	86	79.5	1	76.591	10.931	2016	7.19	72.45		
6	25	7.45	340	2	0	10	20	4	0.1	0.25	189	74	43.5	12.5	27	300	189	174.5	0.11	0.1	2.67	0.89	0.18	1.2	65	78.5	1	76.591	10.931	2016	7.3	81.9		
7	24	7.05	342	2	0	10	20	4	0.09	0.25	189	73.5	43.5	12.5	27.1	300	190	174	0.11	0.1	2.67	0.87	0.17	1.1	105	79	1	76.591	10.931	2016	7.39	72.47		
8	24	7.4	341	2	0	12	20	4	0.1	0.25	188	73.5	44	12	27.1	300	189	174	0.11	0.1	2.66	0.82	0.17	1.2	113	80	1	76.591	10.931	2016	7.06	81.7		
9	25	7.38	339	2	0	11	21	3.9	0.1	0.25	188	73.5	44	12	27	300	190	174	0.11	0.1	2.66	0.81	0.17	1.2	113	80	1	76.591	10.931	2016	7.02	81.72		
10	25	7.56	340	2	0	12	21	3.9	0.1	0.25	188	74	44	12.5	27.1	300	189	173.5	0.11	0.1	2.66	0.88	0.18	1.2	65	80	1	76.591	10.931	2016	6.97	81.81		
11	25	7.1	340	2	0	11	21	4	0.11	0.25	187.5	74	44	12	27.1	300	188	173	0.11	0.1	2.66	0.82	0.18	1.2	65	80	1	76.591	10.931	2016	7.39	81.86		
12	24	7.27	339	2	0	12	22	4	0.11	0.25	187	74	44.5	12	27.1	300	189	173	0.11	0.1	2.68	0.88	0.18	1.1	88	80	1	76.591	10.931	2016	7.12	77.17		
13	25	7.28	341	2	0	11	22	3.9	0.1	0.25	187.5	74.5	44.5	12	27	300	189	173.5	0.11	0.1	2.68	0.81	0.18	1.1	93	79.5	1	76.591	10.931	2016	7.38	77.18		
14	24	7.3	339	2	0	10	22	3.9	0.1	0.25	187	74.5	44.5	11.5	27.1	300	189	173.5	0.11	0.1	2.67	0.83	0.19	1.2	69	79.5	1	76.591	10.931	2016	7.12	81.81		
15	25	7.56	340	2	0	11	22	4	0.1	0.25	188	74.5	44	11.5	27	300	190	174	0.11	0.1	2.67	0.86	0.19	1.1	108	79.5	1	76.591	10.931	2016	7.27	77.33		
16	25	7.06	341	2	0	12	20	4.1	0.1	0.25	188	75	44	11.5	27.1	300	190	174.5	0.11	0.1	2.68	0.84	0.19	1.2	64	80	1	76.591	10.931	2016	6.86	77.15		
17	24	7.21	341	2	0	10	21	4.1	0.11	0.25	189	74.5	44	11.5	27.1	300	190	174.5	0.11	0.1	2.67	0.87	0.18	1.1	91	80.5	1	76.591	10.931	2016	7.11	72.44		
18	25	7.51	341	2	0	11	21	4.1	0.11	0.25	188	73.5	43.5	12	27	300	190	174	0.11	0.1	2.66	0.86	0.18	1.1	63	80	1	76.591	10.931	2016	6.88	72.37		
19	25	7.35	341	2	0	10	20	4.1	0.1	0.25	188	73.5	43.5	12	27	300	191	174	0.11	0.1	2.66	0.88	0.18	1.1	107	80.5	1	76.591	10.931	2016	6.97	72.44		
20	25	7.37	339	2	0	10	20	4.1	0.1	0.25	187	73	43.5	12	27.1	300	191	174.5	0.11	0.1	2.66	0.88	0.18	1.1	103	80.5	1	76.591	10.931	2016	6.97	76.98		
21	24	7.37	339	2	0	10	20	3.9	0.09	0.25	187	74	44	12	27.2	300	191	174	0.11	0.1	2.68	0.89	0.17	1.2	101	80.5	1	76.591	10.931	2016	7.28	81.86		
22	24	7.2	340	2	0	11	20	3.9	0.09	0.25	187	73	44	12	27.1	300	191	174	0.11	0.1	2.67	0.84	0.17	1.1	103	80.5	1	76.591	10.931	2016	6.97	76.98		
23	25	7.01	339	2	0	10	20	3.9	0.09	0.25	187	74	44	12	27.2	300	191	174	0.11	0.1	2.68	0.89	0.17	1.2	101	80.5	1	76.591	10.931	2016	7.28	81.86		
24	25	7.58	340	2	0	10	21	4.1	0.09	0.25	187.5	74.5	44.5	12.5	27.2	300	191	174	0.11	0.1	2.68	0.88	0.18	1.2	107	80	1	76.591	10.931	2016	6.8	77.02		
25	25	7.14	340	2	0	11	21	4.1	0.1	0.25	188	74.5	44.5	12	27.1	300	190	174.5	0.11	0.1	2.68	0.88	0.18	1.2	107	80	1	76.591	10.931	2016	6.84	81.77		
26	24	7.55	339	2	0	12	21	3.9	0.1	0.25	188	74.5	44.5	12	27.2	300	190	174.5	0.11	0.1	2.67	0.83	0.19	0.9	76	80.5	1	76.591	10.931	2016	7.33	87.95		
27	25	7.06	341	2	0	12	21	4	0.1	0.25	188	73.5	44	12	27.1	300	190	174.5	0.11	0.1	2.68	0.87	0.18	1.1	112	80.5	1	76.591	10.931	2016	6.86	77.06		
28	24	7.53	341	2	0	10	22	4	0.1	0.25	188.5	73.5	44	12	27.1	300	191	174	0.11	0.1	2.68	0.87	0.19	1.1	90	80.5	1	76.591	10.931	2016	7.03	77.2		
29	25	7.49	340	2	0	12	21	4	0.11	0.25	188	73	44	12	27.1	300	191	174	0.11	0.1	2.67	0.9	0.19	0.9	81	80.5	1	76.591	10.931	2016	6.91	87.77		
30	24	7.26	339	2	0	11	22	4.1	0.1	0.25	187	74	44	11.5	27.2	300	191	174	0.11	0.1	2.67	0.83	0.18	1.1	10	80.5	1	76.591	10.931	2016	6.93	72.33		
31	24	7.53	339	2	0	11	21	4	0.11	0.25	187	74	43.5	11.5	27.1	300	191	174	0.11	0.1	2.67	0.83	0.17	1.2	109	80	1	76.591	10.931	2016	6.94	81.73		
32	24	7.39	340	2	0	11	22	4	0.1	0.25	188	74	44	12	27.1	300	190	174	0.11	0.1	2.67	0.83	0.17	1.1	120	80	1	76.591	10.931	2016	7.37	77.19		
33	24	7.23	335	2	0	13	23	4	0.1	0.25	120	72	48	12	28	300	190	170	0.11	0.1	2.68	0.81	0.18	0.9	109	80	1	76.591	10.931	2016	6.83	87.64		
34	24	7.74	335	2	0	13	23	4	0.1	0.25	122	72	48	12	28	300	190	170	0.11	0.1	2.68	0.88	0.18	0.9	15	76	1	76.591	10.931	2016	7.02	87.84		
35	24	7.76	330	2	0	15	23	4	0.1	0.25	124	72	50	12	28	300	185	170	0.11	0.1	2.68	0.84	0.19	0.9	83	75	1	76.591	10.931	2016	6.12	87.45		
36	24	7.47	330	2	0	15	23	4	0.1	0.25	125	70	50	12.5	28.5	300	185	168	0.11	0.1	2.7	0.82	0.19	0.9	92	70	1	76.591	10.931	2016	6.51	87.54		
37	26	7.44	329	2	0	18	25	4	0.1	0.25	125	70	50	12.5	28.5	300	185	168	0.11	0.1	2.7	0.82	0.19	0.9	92	70	1	76.591	10.931	2016	6.17	82.73		
38	25	7.45	325	2	0	18	24	4.2	0.1	0.25	125	70	54	12.5	28.5	300	183	168	0.11	0.1	2.8	0.9	0.18	0.8	93	65	1	76.591	10.931	2016	6.93	82.99		
39	25	7.42	325	2	0	18	25	4.1	0.1	0.25	125	70	54	12	28.5	300	183	168	0.11	0.1	2.8	0.8	0.18	0.8	107	65	1	76.591	10.931	2016	5.94	82.59		
40	24	7.54	323	2	0	18	25	3.9	0.11	0.25	127	68	54	12	29	300	183	165	0.11	0.1	2.9	0.85	0.18	0.8	106	65	1	76.591	10.931	2016	5.99	82.85		
41	24	7.18	320	2	0	19	25	3.9	0.11	0.25	127	69	58	12	29	300	183	165	0.11	0.1	2.9	0.83	0.17	0.8	105	60	1	76.591	10.931	2016	6.29	82.67		
42	26	7.43	318	2	0	19	26	4	0.1	0.25	127	69	58	13	29	300	180	165	0.11	0.1	2.9	0.87	0.17	0.8	105	60	1	76.591	10.931	2016	6.03	87.99		

APPENDIX B- CODING

Code to calculate WQI using the WQI-PCA dataset

```
importing pandas as pd
import pandas as pd
Read and store content
of an excel file
read_file = pd.read_excel("WQI-PCA.xlsx")
Write the dataframe object into csv file
read_file.to_csv("WQI-PCA.csv",
index = None,
header=True)
read csv file and convert into a dataframe object
df = pd.DataFrame(pd.read_csv("WQI-PCA.csv"))
show the dataframe
df
Temp = df[df.columns[0]]
Temp
i = 1
Temp=[]
for i in range(len(df)):
s1=28
w1=0.004249273
v1=df[df.columns[0]][i]
q1=(v1/s1)*100
tmp = w1*q1
print(tmp)
Temp.append(tmp)
```

```

pH= df[df.columns[1]]
i=1
pH=[]
for i in range(len(df)):
s2=8.5
w2=0.013997604
v2=df[df.columns[1]][i]
q2=(v2/s2)*100
p2 = w2*q2
print(p2)
pH.append(p2)
"-----"for each attribute similar coding
fc= df[df.columns[26]] # FC
i = 1
fc=[]
for i in range(len(df)):
s27= 60
w27=0.001982994
v27=df[df.columns[26]][i]
q27=(v27/s27)*100
fc27 = w27*q27
print(fc27)
fc.append(fc27)
wqi=[]
total=0
for i in range(len(df)):
total=Temp[i]+pH[i]+cd[i]+tr[i]+pa[i]+ta[i]+ch[i]+cod[i]+tkn[i]+am[i]+ha[i]+cal[i]+mg[i]+sl[i]
+so[i]+tss[i]+tds[i]+fds[i]+ps[i]+b[i]+do[i]+bod[i]+fl[i]+do[i]+ni[i]+tc[i]+fc[i]
wqi.append(total)

```



```

print(wqi)
df = pd.read_csv('WQI-PCA.csv')
df.head()
df["WQI"] = wqi
df.to_csv("WQI-PCA.csv", index=False)
textbf{Code for Exploratory Data Analysis}
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
get_ipython().run_line_magic('matplotlib', 'inline')
sns.set(color_codes=True)
df = pd.read_csv("weather.csv")
# To display the top 5 rows
df.head(5)
df.tail(5)
df.dtypes
df.shape
duplicate_rows_df = df[df.duplicated()]
print("number of duplicate rows: ", duplicate_rows_df.shape)
df.count() # Used to count the number of rows
df.describe()
df = df.drop_duplicates()
df.head(5)
df.count()
print(df.isnull().sum())
df = df.dropna() # Dropping the missing values.
df.count()

```

```

print(df.isnull().sum()) # After dropping the values
plt.figure(figsize=(10,5))
c= df.corr()
sns.heatmap(c,cmap="BrBG",annot=True)
c
# Create the default pairplot
sns.pairplot(df)
df.hist()
sns.boxplot(x=df['Temp'])
sns.boxplot(x=df['pH'])
sns.boxplot(x=df['Conductivity '])
sns.boxplot(x=df['Turbidity'])
sns.boxplot(x=df['Total Alkalinity'])
sns.boxplot(x=df['Cholride'])
sns.boxplot(x=df['COD'])
sns.boxplot(x=df['TKN'])
sns.boxplot(x=df['Ammonia'])
sns.boxplot(x=df['Hardness'])
sns.boxplot(x=df['Sodium'])
sns.boxplot(x=df['TDS'])
sns.boxplot(x=df['FDS'])
sns.boxplot(x=df['Pottassium'])
sns.boxplot(x=df['BOD'])
sns.boxplot(x=df['Fluoride'])
sns.boxplot(x=df['Nitrate-N'])
sns.boxplot(x=df['TC'])
sns.boxplot(x=df['Dew'])
sns.boxplot(x=df['Humidity'])

```

```

sns.boxplot(x=df['Sealevelpressure'])
sns.boxplot(x=df['Precipitation'])
sns.boxplot(x=df['Precipcover'])
sns.boxplot(x=df['Windspeed'])
sns.boxplot(x=df['Cloudcover'])
sns.heatmap(c)
sns.heatmap(c,cmap="BrBG")
d=df.corr(method='pearson')
d

```

Code for Feature Selection

```

from sklearn.datasets import make_regression
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import f_regression
from matplotlib import pyplot

# feature selection
def select_features(X_train, y_train, X_test):
    # configure to select all features
    fs = SelectKBest(score_func=f_regression, k='all')
    # learn relationship from training data
    fs.fit(X_train, y_train)
    # transform train input data
    X_train_fs = fs.transform(X_train)
    # transform test input data
    X_test_fs = fs.transform(X_test)
    return X_train_fs, X_test_fs, fs

# load the dataset
X, y = make_regression(n_samples=1000, n_features=100, n_informative=10, noise=0.1,
random_state=1)

```

```

import pandas as pd
import seaborn as sns
df = pd.read_csv('WQI.csv')
df.head()

#Divide the features into Independent and Dependent Variable
X = df.drop('WQI' , axis =1)
y = df['WQI']

# split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)

# feature selection
X_train_fs, X_test_fs, fs = select_features(X_train, y_train, X_test)

# what are scores for the features
for i in range(len(fs.scores_)):
print('Feature %d: %f % (i, fs.scores_[i]))

# plot the scores
pyplot.bar([i for i in range(len(fs.scores_))], fs.scores_)
pyplot.show()

```

Code for WQI Prediction using GRU using WQI-PCA Dataset

```

import tensorflow as tf
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from tensorflow.keras.layers import Dense, LSTM, Dropout, GRU

# load the dataset
df = pd.read_csv("WQI_PCA.csv")
df

# extract the target variable

```

```

y = df['WQI'].values
# drop the target variable from the dataset
df = df.drop(['WQI'], axis=1)
# normalize the input data using MinMaxScaler
scaler = MinMaxScaler()
X = scaler.fit_transform(df)
X
df
# split the data into training and testing sets
split = int(0.8 * X.shape[0])
X_train, X_test = X[:split, :], X[split:, :]
y_train, y_test = y[:split], y[split:]
# convert the training and testing sets to constants
X_train = tf.constant(X_train, dtype=tf.float32)
y_train = tf.constant(y_train, dtype=tf.float32)
X_test = tf.constant(X_test, dtype=tf.float32)
y_test = tf.constant(y_test, dtype=tf.float32)
# add an additional dimension to the input data
X_train = tf.expand_dims(X_train, axis=-1)
X_test = tf.expand_dims(X_test, axis=-1)
# create the GRU model
model = tf.keras.Sequential()
model.add(tf.keras.layers.GRU(64, input_shape=(X_train.shape[1], 1), return_sequences=True))
model.add(tf.keras.layers.GRU(32))
model.add(tf.keras.layers.Dense(1, activation='linear'))
model.add(Dropout(0.2))
# compile the model
model.compile(optimizer='adam', loss='mean_squared_error')

```

```

# fit the model on the training data
history= model.fit(X _train, y _train, epochs=200,validation _split = 0.2, batch _size=64,
verbose=2)

# make predictions on the test data
y _pred = model.predict(X _test)
y _pred = y _pred.reshape(-1,)

# calculate the evaluation metrics
mae = mean _absolute _error(y _test, y _pred)
mse = mean _squared _error(y _test, y _pred)
rmse = np.sqrt(mse)
r2 = r2 _score(y _test, y _pred)

# print the evaluation metrics
print("MAE:", mae)
print("MSE:", mse)
print("RMSE:", rmse)
print("R2-Score:", r2)

```

Code for WQI Prediction using WQI-SA Dataset and LSTM

```

import pandas
import matplotlib.pyplot as plt
dataset = pandas.read _csv('WQI-SA.csv', usecols=[], engine='python')
# LSTM for Water quality problem with regression framing
import numpy
import matplotlib.pyplot as plt
from pandas import read _csv
import math
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM
from sklearn.preprocessing import MinMaxScaler

```

```

from sklearn.metrics import mean_squared_error
# convert an array of values into a dataset matrix
def create_dataset(dataset, look_back=1):
    dataX, dataY = [], []
    for i in range(len(dataset)-look_back-1):
        a = dataset[i:(i+look_back), 0]
        dataX.append(a)
        dataY.append(dataset[i + look_back, 0])
    return numpy.array(dataX), numpy.array(dataY)

# fix random seed for reproducibility
numpy.random.seed(7)

# load the dataset
dataframe = read_csv('weather.csv', usecols=[38], engine='python')
dataset = dataframe.values
dataset = dataset.astype('float32')

# normalize the dataset
scaler = MinMaxScaler(feature_range=(0, 1))
dataset = scaler.fit_transform(dataset)

# normalize the dataset
scaler = MinMaxScaler(feature_range=(0, 1))
dataset = scaler.fit_transform(dataset)

# split into train and test sets
train_size = int(len(dataset) * 0.67)
test_size = len(dataset) - train_size
train, test = dataset[0:train_size,:], dataset[train_size:len(dataset),:]
train, test = dataset[0:train_size,:], dataset[train_size:len(dataset),:]

# reshape into X=t and Y=t+1
look_back = 1

```

```

trainX, trainY = create_dataset(train, look_back)
testX, testY = create_dataset(test, look_back)

# reshape input to be [samples, time steps, features]
trainX = numpy.reshape(trainX, (trainX.shape[0], 1, trainX.shape[1]))
testX = numpy.reshape(testX, (testX.shape[0], 1, testX.shape[1]))

# create and fit the LSTM network
model = Sequential()
model.add(LSTM(4, input_shape=(1, look_back)))
model.add(Dense(1))
model.compile(loss='mean_squared_error', optimizer='adam')
model.fit(trainX, trainY, epochs=10, batch_size=1, verbose=2)

# make predictions
trainPredict = model.predict(trainX)
testPredict = model.predict(testX)

print(trainPredict.shape)
print(testPredict.shape)

# invert predictions
trainPredict = scaler.inverse_transform(trainPredict)
trainY = scaler.inverse_transform([trainY])
testPredict = scaler.inverse_transform(testPredict)
testY = scaler.inverse_transform([testY])

y_pred=testPredict
y_test=testX

print( (y_test),(y_pred))

from sklearn import metrics

testScore = metrics.mean_absolute_error(testY[0], testPredict[:,0])

print('Test Score: %.4f MAE' % (testScore))

from sklearn import metrics

```



```

# calculate mean squared error
testScore = metrics.mean_squared_error(testY[0], testPredict[:,0])
print('Test Score: %.5f MSE' % (testScore))

# calculate root mean squared error
testScore = math.sqrt(mean_squared_error(testY[0], testPredict[:,0]))
print('Test Score: %.4f RMSE' % (testScore))

from sklearn import metrics

# calculate r squared
testScore = metrics.r2_score(testY[0], testPredict[:,0])
print('Test Score: %.2f R2 score' % (testScore))

from sklearn import metrics

# calculate Explained Variance
testScore = metrics.explained_variance_score(testY[0], testPredict[:,0])
print('Test Score: %.2f Explained Variance' % (testScore))

```

Code for WQI Prediction using TFT

```

EPOCHS = 200
INLEN = 32
HIDDEN = 64
LSTMLAYERS = 2
ATTHEADS = 1
DROPOUT = 0.1
BATCH = 32
N_FC = 12      # default forecast horizon
RAND = 42      # set random state
N_SAMPLES = 100 # number of times a prediction is sampled from a probabilistic model
N_JOBS = 3     # parallel processors to use; -1 = all processors
# default quantiles for QuantileRegression
QUANTILES = [0.01, 0.05, 0.1, 0.2, 0.25, 0.5, 0.75, 0.8, 0.9, 0.95, 0.99]

```

```

TRAIN = "20201201" # train/test split
MSEAS = 12      # max seasonality to check: months
ALPHA = 0.05    # significance level for seasonality test
FIGSIZE = (9, 6)
qL1, qL2, qL3 = 0.01, 0.05, 0.10
# percentiles of predictions: lower bounds
qU1, qU2, qU3 = 1-qL1, 1-qL2, 1-qL3
# upper bounds derived from lower bounds
label_q1 = f'{int(qU1 * 100)} / {int(qL1 * 100)} percentile band'
label_q2 = f'{int(qU2 * 100)} / {int(qL2 * 100)} percentile band'
label_q3 = f'{int(qU3 * 100)} / {int(qL3 * 100)} percentile band'
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from darts import TimeSeries, concatenate
from darts.dataprocessing.transformers import Scaler
from darts.models import TFTModel, NaiveSeasonal, NaiveDrift, ExponentialSmoothing
from darts.utils.statistics import check_seasonality, extract_trend_and_seasonality
from darts.metrics import mape
from darts.utils.timeseries_generation import datetime_attribute_timeseries
from darts.utils.likelihood_models import QuantileRegression
from darts.utils.utils import ModelMode, SeasonalityMode, TrendMode
pd.set_option("display.precision",2)
np.set_printoptions(precision=2, suppress=True)
pd.options.display.float_format = '{:,.2f}'.format
df = pd.read_csv('WQI-SA.csv')
series = TimeSeries.from_dataframe(df, 'Month','WQI')
ts = series / TimeSeries.from_series(series.time_index.days_in_month)

```

```

ts = ts.astype(np.float32)
print(type(ts))
df = ts.pd_dataframe()
print(df.shape[0])
print(df.shape[1])
ts = TimeSeries.from_series(df["WQI"])
plt.figure(100, figsize=(18, 5))
ts.plot();
mseas=MSEAS
for m in range(2, MSEAS):
is_seasonal, mseas = check_seasonality(ts, m=m, alpha=ALPHA, max_lag=MSEAS)
if is_seasonal:
break
print("seasonal? " + str(is_seasonal))
if is_seasonal:
print("There is seasonality of order {} months'.format(mseas))
ts_trend, ts_seas = extract_trend_and_seasonality(ts=ts, freq=mseas)
# train/test
if isinstance(TRAIN, str):
split = pd.Timestamp(TRAIN)
else:
split = TRAIN
ts_train, ts_test = ts.split_after(split)
#print(ts_train.count())
#scale the time series on the training settransformer = Scaler()
transformer = Scaler()
ts_ttrain = transformer.fit_transform(ts_train)
print(len(ts_ttrain))

```

```

ts_ttest = transformer.transform(ts_test)
ts_t = transformer.transform(ts)
# create covariates: year, month, and integer index series
cov = datetime_attribute_timeseries(ts, attribute="year", one_hot=False)
cov = cov.stack(datetime_attribute_timeseries(ts, attribute="month", one_hot=False))
cov = cov.stack(TimeSeries.from_times_and_values(
times=ts.time_index,
values=np.arange(len(ts)),
columns=["linear_increase"]))
cov = cov.astype(np.float32)
# train/test
train_cov, test_cov = cov.split_after(split)
# rescale the covariates: fit on the training set
scaler = Scaler()
scaler.fit(train_cov)
tcov = scaler.transform(cov)
# naive seasonal forecast
modelNs = NaiveSeasonal(K=mseas)
modelNs.fit(ts_train)
ts_predNs = modelNs.predict(len(ts_test))
# naive drift (trend) forecast
modelNd = NaiveDrift()
modelNd.fit(ts_train)
ts_predNd = modelNd.predict(len(ts_test))
ts_predN = ts_predNd + ts_predNs - ts_train.last_value()
plt.figure(100, figsize=(18, 5))
ts.plot(label="actual")
ts_predN.plot(label="naive forecast")

```

```

plt.title("Naive Forecast (MAPE: {:.2f} %)".format(mape(ts _test, ts _predN)))
plt.legend();
# search space for best theta value: check 100 alternatives
modelX = ExponentialSmoothing(
seasonal _periods=mseas,
seasonal=ModelMode.MULTIPLICATIVE)
modelX.fit(ts _train)
ts _predX = modelX.predict( n=len(ts _test),
num _samples=N _SAMPLES)
print(ts _test)
print(ts _predX)
plt.figure(100, figsize=(18, 5))
ts.plot(label="actual")
ts _predX.plot(label="Exponential Smoothing")
plt.title("Exponential Smoothing (MAPE: {:.2f} %)".format(mape(ts _test, ts _predX)))
plt.legend();
#temporal fusion transformer
model = TFTModel( input _chunk _length=INLEN,
output _chunk _length=N _FC,
hidden _size=HIDDEN,
lstm _layers=LSTMLAYERS,
num _attention _heads=ATTHEADS,
dropout=DROPOUT,
batch _size=BATCH,
n _epochs=EPOCHS,
likelihood=QuantileRegression(quantiles=QUANTILES),
# loss _fn=MSELoss(),
random _state=RAND,

```

```

force_reset=True)
# training
model.fit( ts_ttrain,
future_covariates=tcov,
verbose=True)
# testing: generate predictions
ts_tpred = model.predict( n=len(ts_test),num_samples=N_SAMPLES,n_jobs=N_JOBS)
print("Do the predictions constitute a probabilistic time series?", ts_tpred.is_probabilistic)
# testing: helper function: plot predictions
def plot_predict(ts_actual, ts_test, ts_pred):
# plot time series, limited to forecast horizon
plt.figure(figsize=FIGSIZE)
ts_actual.plot(label="actual")
# plot actual
ts_pred.plot(low_quantile=qL1, high_quantile=qU1, label=label_q1) # plot U1 quantile band
ts_pred.plot(low_quantile=qL2, high_quantile=qU2, label=label_q2) # plot U2 quantile band
ts_pred.plot(low_quantile=qL3, high_quantile=qU3, label=label_q3) # plot U3 quantile band
ts_pred.plot(central_quantile="mean", label="expected")
# plot "mean" or median=0.5
plt.title("TFT: test set (MAPE: {:.2f} %)".format(mape(ts_test, ts_pred)))
plt.legend();
# testing: call helper function: plot predictions
ts_pred = transformer.inverse_transform(ts_tpred)
plot_predict(ts, ts_test, ts_pred)
# testing: call helper function: plot predictions, focus on test set
ts_pred = transformer.inverse_transform(ts_tpred)
ts_actual = ts[ ts_tpred.start_time(): ts_tpred.end_time() ]
# actual values in forecast horizon

```

```

plot _predict(ts _actual, ts _test, ts _pred)
# testing: collect the prediction percentiles in a dataframe dfY
dfY = pd.DataFrame()
# helper method: calculate percentiles of predictions
def predQ(ts _tpred, q):
ts _t = ts _tpred.quantile _timeseries(q)
# percentile of predictions
ts = transformer.inverse _transform(ts _t)
# inverse-transform the percentile
s = TimeSeries.pd _series(ts)
header = "Q" + format(int(q*100), "02d")
dfY[header] = s
# fill Actual column into the new dataframe
dfY["Actual"] = TimeSeries.pd _series(ts _actual)
# call helper function: percentiles of predictions
quantiles = [0.5, qU1, qU2, qU3, qL3, qL2, qL1]
_ = [predQ(ts _tpred, q) for q in quantiles]
dfY.iloc[np.r _[0:2, -2:0]]
# backtest the model: predictions vs actual values
# start: at 10 % of time series
STARTBACKT = ts.get _timestamp _at _point(point=max(0.2, INLEN))
#pd.Timestamp("19500301")
print("start:",STARTBACKT)
# forecast horizon: until end of test period
tdiff = ( pd.to _datetime(ts.end _time()).to _period("M") -
pd.to _datetime(STARTBACKT).to _period("M")).n
print("months:",tdiff)
ts _tbacktest = model.historical _forecasts( series=ts _t,
past _covariates=None,

```

```

future_covariates=tcov,
start=STARTBACKT,
num_samples=N_SAMPLES,
forecast_horizon=tdiff,
stride=12,
last_points_only=False,
retrain=False,
verbose=True)

# backtesting: helper function: plot backtest predictions
def plot_backtest(ts_backtest, ts_actual, transformer):
plt.figure(figsize=FIGSIZE)
ts_actual.plot(label="actual")
ts_backtest.plot(low_quantile=qL1, high_quantile=qU1, label=label_q1)
ts_backtest.plot(low_quantile=qL3, high_quantile=qU3, label=label_q3)
plt.legend()
MAPE = "MAPE: {:.2f} %".format(mape(ts_actual, ts_backtest))
plt.title( "TFT: backtest from " + format(ts_backtest.start_time(), "%Y. %m") +
" to " + format(ts_backtest.end_time(), "%Y. %m") + " (" + MAPE + ")");

# backtesting: call helper function: plot backtest
ts_backtest = transformer.inverse_transform(ts_tbacktest)
plot_backtest( ts_backtest=concatenate(ts_backtest),
ts_actual=ts,
transformer=transformer)

# create future covariates: year, month, integer index
FC_HORIZ = 12

# months after end of training set
start = ts.start_time()

n_per = len(ts_train) + FC_HORIZ

# set a maximum horizon to create covariates for

```



```

# create covariates from beginning of training data to end of forecast horizon
ts_year = datetime_attribute_timeseries(
pd.date_range(start=start, periods=n_per, freq="MS"), #, closed="right"),
attribute="year",
one_hot=False)
ts_month = datetime_attribute_timeseries(
pd.date_range(start=start, periods=n_per, freq="MS"),
attribute="month",
one_hot=False)
cov = ts_year.stack(ts_month)
# combine year and month with integer index as third covariate
cov = cov.stack(TimeSeries.from_times_and_values(
times=cov.time_index,
values=np.arange(n_per),
columns=['linear_increase']))
cov = cov.astype(np.float32)
tcov = scaler.transform(cov)
print("start:", cov.start_time(), "; end:", cov.end_time())
# generate future, out-of-sample predictions
ts_tpred = model.predict(n=FC_HORIZ, future_covariates=tcov, num_samples=N
_SAMPLES)
print("start:", ts_tpred.start_time(), "; end:", ts_tpred.end_time())
ts_pred = transformer.inverse_transform(ts_tpred)
plt.figure(figsize=FIGSIZE)
ts_actual = ts.slice_intersect(other=ts_pred)
ts_actual.plot(label="actual")
ts_pred.plot(low_quantile=qL1, high_quantile=qU1, label=label_q1)
# plot U1 quantile band
ts_pred.plot(low_quantile=qL3, high_quantile=qU3, label=label_q3) # plot U3 quantile band

```

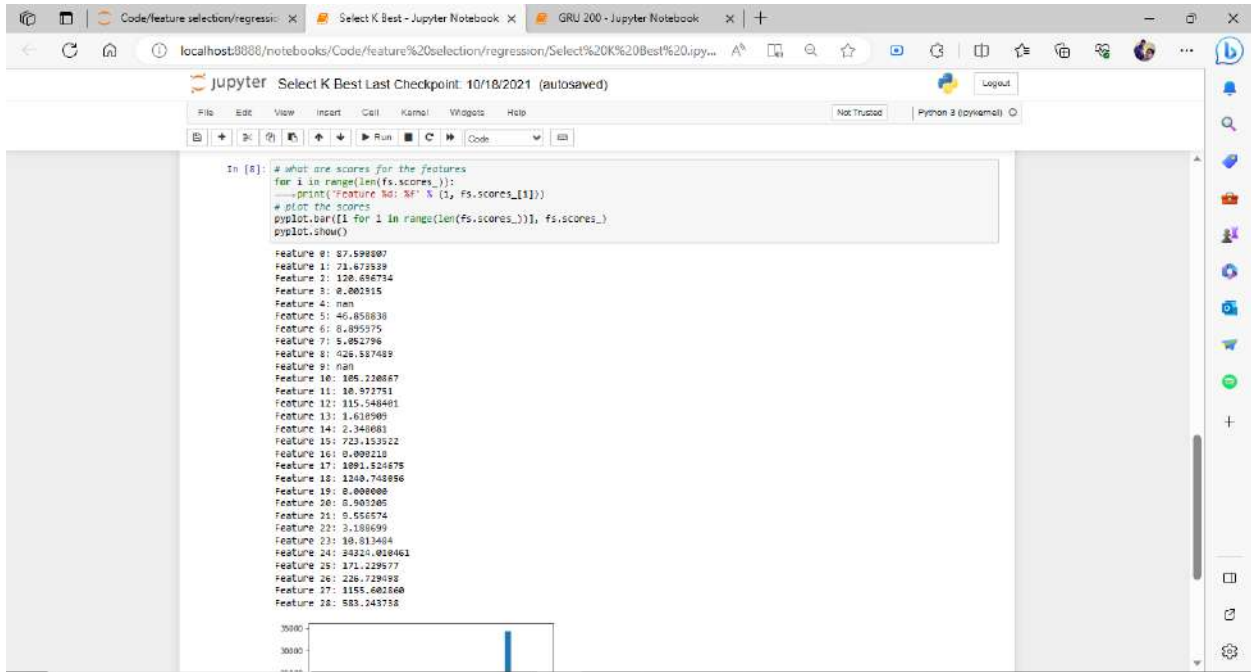
```

ts_pred.plot(central_quantile="mean", label="expected")           # plot "mean" or median=0.5
plt.title( "TFT: forecast" + " from " + format(ts_tpred.start_time(), "%Y. %m") + "
to " + format(ts_tpred.end_time(), "%Y. %m"));
plt.legend();
# forecast: collect the prediction percentiles in a dataframe
dfY = pd.DataFrame()
# call helper function predQ: percentiles of predictions
_ = [predQ(ts_tpred, q) for q in quantiles]
dfY["Actual"] = TimeSeries.pd_series(ts_actual)
dfY.iloc[np.r_[0:2, -2:0]]
ts_tpred.plot(central_quantile="mean", label="expected")       # plot "mean" or median=0.5
plt.legend();
print(ts_actual)
print(ts_pred)

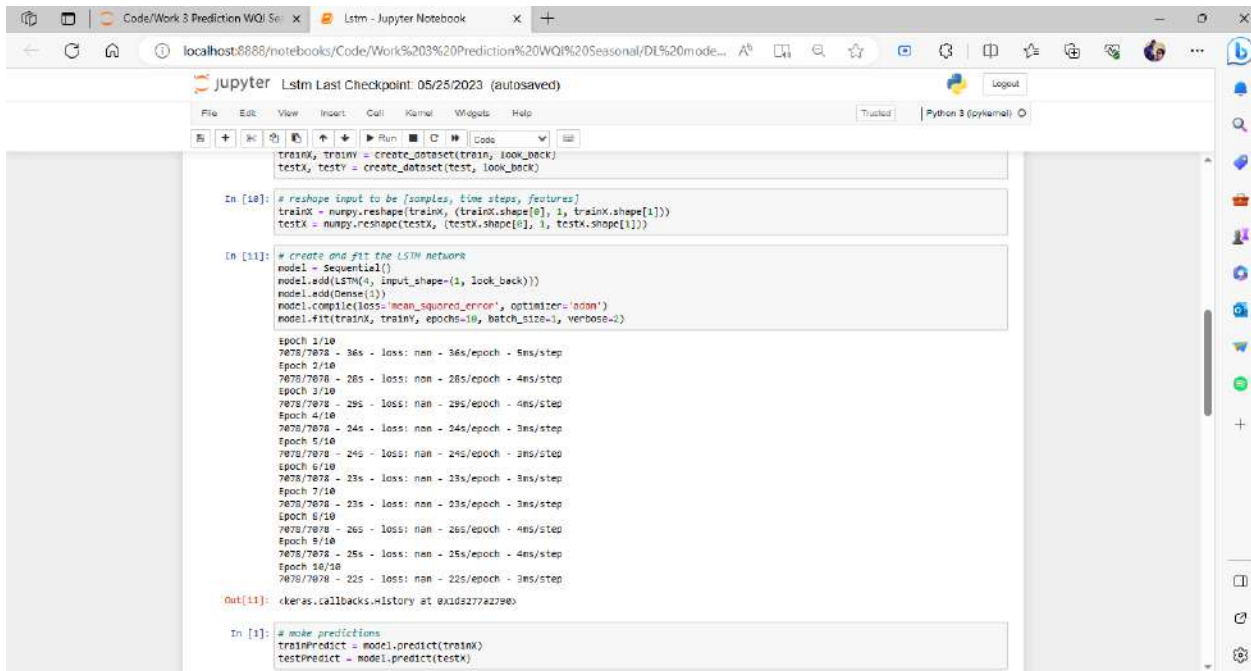
```

APPENDIX C- SCREENSHOTS

Feature Selection -Select K Best



Implementation of LSTM



Implementation of GRU

```

# create the gru model
model = tf.keras.Sequential()
model.add(tf.keras.layers.GRU(64, input_shape=(X_train.shape[1], 1), return_sequences=True))
model.add(tf.keras.layers.GRU(32))
model.add(tf.keras.layers.Dense(1, activation="linear"))
model.add(Cropout(0.2))

In [12]: # compile the model
model.compile(optimizer="adam", loss="mean_squared_error")

In [13]: # fit the model on the training data
history = model.fit(X_train, y_train, epochs=100, validation_split = 0.3, batch_size=64, verbose=2)

Epoch 1/200
106/106 - 23s - loss: 7182.5903 - val_loss: 7513.1860 - 23s/epoch - 218ms/step
Epoch 2/200
106/106 - 6s - loss: 6819.4688 - val_loss: 6984.8735 - 6s/epoch - 66ms/step
Epoch 3/200
106/106 - 6s - loss: 6325.4604 - val_loss: 6516.7266 - 6s/epoch - 59ms/step
Epoch 4/200
106/106 - 6s - loss: 5763.6600 - val_loss: 6084.0688 - 6s/epoch - 58ms/step
Epoch 5/200
106/106 - 7s - loss: 5497.7173 - val_loss: 5716.2432 - 7s/epoch - 66ms/step
Epoch 6/200
106/106 - 7s - loss: 5222.7930 - val_loss: 5387.4741 - 7s/epoch - 65ms/step
Epoch 7/200
106/106 - 6s - loss: 5031.6566 - val_loss: 5059.6816 - 6s/epoch - 55ms/step
Epoch 8/200
106/106 - 6s - loss: 4877.2943 - val_loss: 4779.5208 - 6s/epoch - 58ms/step
Epoch 9/200
106/106 - 6s - loss: 4693.2188 - val_loss: 4526.9693 - 6s/epoch - 61ms/step
Epoch 10/200
...

In [14]: # make predictions on the test data
y_pred = model.predict(X_test)

```

Implementation of Temporal Fusion Transformer

```

model = TFTModel(
    input_chunk_length=EM_FU,
    output_chunk_length=N_FC,
    hidden_size=HIDDEN,
    lstm_layers=LSTM_LAYERS,
    num_attention_heads=ATTN_HEADS,
    dropout=DROP_OUT,
    batch_size=BATCH,
    n_epochs=EPOCHS,
    likelihood_quantile_regression(quantiles=QUANTILES),
    # loss_fn=MSELoss(),
    random_state=RAND,
    force_reset=TRUE)

In [16]: # training
model.fit(
    ts_train,
    future_covariates=cov,
    verbose=True)

10 | lstm_encoder | LSTM | 86.6 K
11 | lstm_decoder | LSTM | 86.6 K
12 | post_lstm_gsn | _GateAddNorm | 8.4 K
13 | static_enrichment_gsn | _GatedResidualNetwork | 20.9 K
14 | multihead_attn | _InterpretableMultiheadAttention | 16.6 K
15 | post_attn_gsn | _GateAddNorm | 8.4 K
16 | feed_forward_block | _GatedResidualNetwork | 16.8 K
17 | pre_output_gsn | _GateAddNorm | 8.4 K
18 | output_layer | Linear | 715

-----
292 K Trainable params
9 Non-trainable params
292 K Total params
1.170 Total estimated model param size (MB)

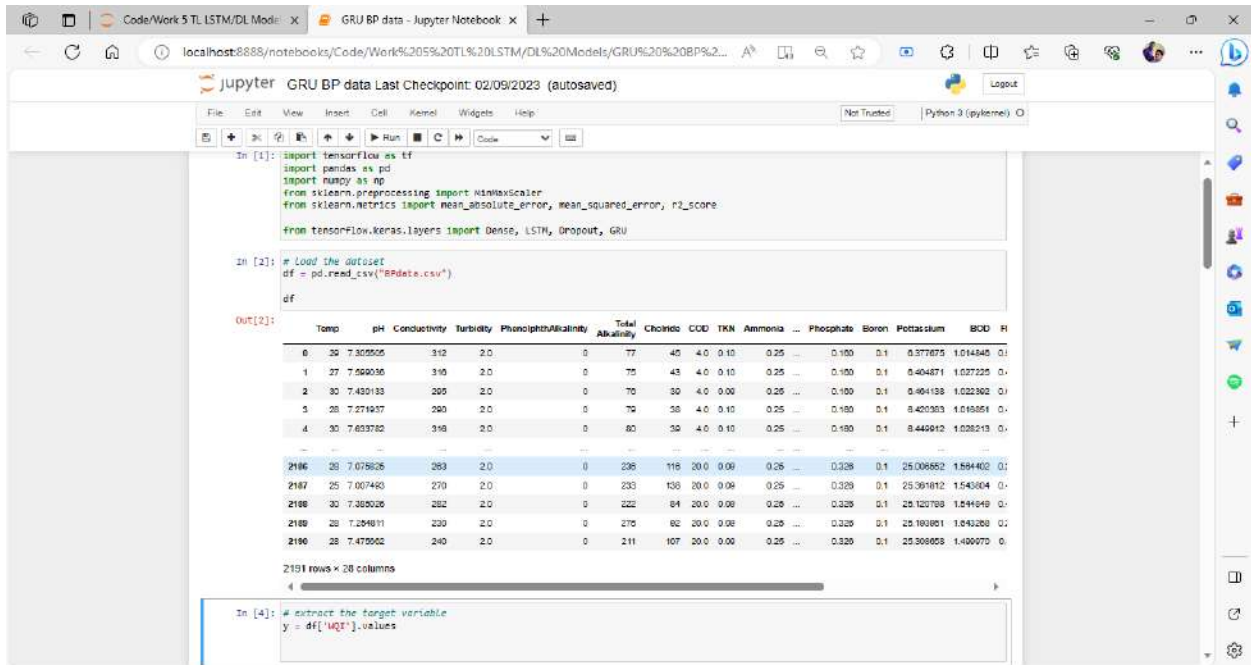
Training: 0it [00:00, ?it/s]
2022-11-24 11:25:45 pytorch_lightning.utilities.rank_zero INFO: "trainer.fit" stopped: "max_epochs=200" reached.

Out[16]: <dartts.models.forecasting_tft_model.TFTModel at 0x2d315e839c0>

In [17]: # testing: generate predictions
ts_tested = model.predict(
    n_hlnts_test=N_SAMPLES_TEST,
    num_samples=N_SAMPLES,
    n_100s=N_100S)

```

Implementation of GRU with Bharathapuzha Data



```
In [1]: import tensorflow as tf
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from tensorflow.keras.layers import Dense, LSTM, Dropout, GRU

In [2]: # Load the dataset
df = pd.read_csv("BPdata.csv")
df

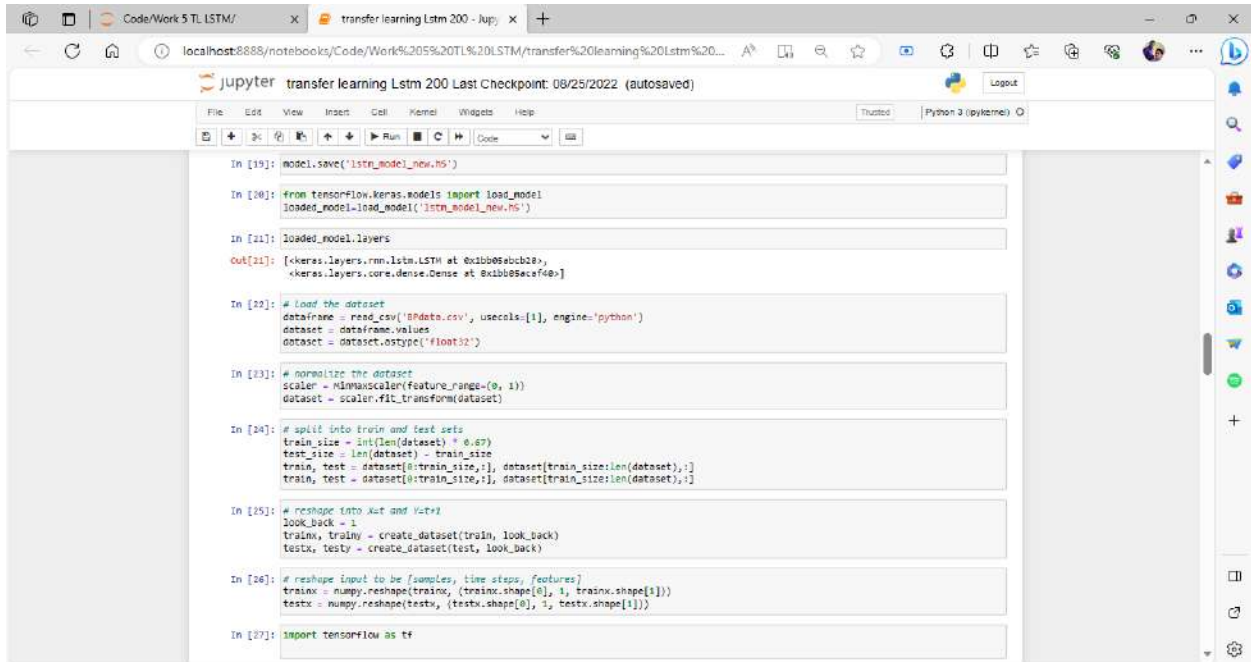
Out[2]:
```

	Temp	pH	Conductivity	Turbidity	Phosphite/Alkalinity	Total Alkalinity	Chloride	COD	TRN	Ammonia	Phosphate	Boron	Potassium	BOD	H
0	29	7.305505	312	2.0	0	77	40	4.0	0.10	0.25	0.100	0.1	0.377875	1.014845	0.1
1	27	7.502036	310	2.0	0	75	43	4.0	0.10	0.25	0.100	0.1	0.404871	1.027225	0.1
2	30	7.432133	295	2.0	0	70	30	4.0	0.00	0.25	0.100	0.1	0.404139	1.023302	0.1
3	28	7.271937	280	2.0	0	70	38	4.0	0.10	0.25	0.180	0.1	0.402383	1.019851	0.1
4	30	7.633782	318	2.0	0	80	30	4.0	0.10	0.25	0.180	0.1	0.446912	1.028213	0.1
...
2190	28	7.075625	283	2.0	0	238	118	20.0	0.00	0.25	0.328	0.1	25.008882	1.594402	0.1
2187	25	7.007493	270	2.0	0	233	138	20.0	0.00	0.25	0.328	0.1	25.391812	1.543404	0.1
2188	30	7.305526	282	2.0	0	222	84	20.0	0.00	0.25	0.328	0.1	25.120780	1.544949	0.1
2189	28	7.284011	230	2.0	0	270	90	20.0	0.00	0.25	0.328	0.1	25.193901	1.649288	0.1
2190	28	7.475002	240	2.0	0	211	107	20.0	0.00	0.25	0.320	0.1	25.308058	1.400970	0.1

```
2191 rows x 28 columns
```

```
In [4]: # extract the target variable
y = df["BOD"].values
```

Implementation of Transfer Learning with LSTM based Pre-trained Model



```
In [19]: model.save('lstm_model_new.h5')

In [20]: from tensorflow.keras.models import load_model
loaded_model=load_model('lstm_model_new.h5')

In [21]: loaded_model.layers

Out[21]: [<keras.layers.rnn.lstm.LSTM at 0x1b0b0ebcb2>,
<keras.layers.core.dense.Dense at 0x1b0b09caf4e>]

In [22]: # load the dataset
dataframe = read_csv('BPdata.csv', usecols=[1], engine='python')
dataset = dataframe.values
dataset = dataset.astype('float32')

In [23]: # normalize the dataset
scaler = MinMaxScaler(feature_range=(0, 1))
dataset = scaler.fit_transform(dataset)

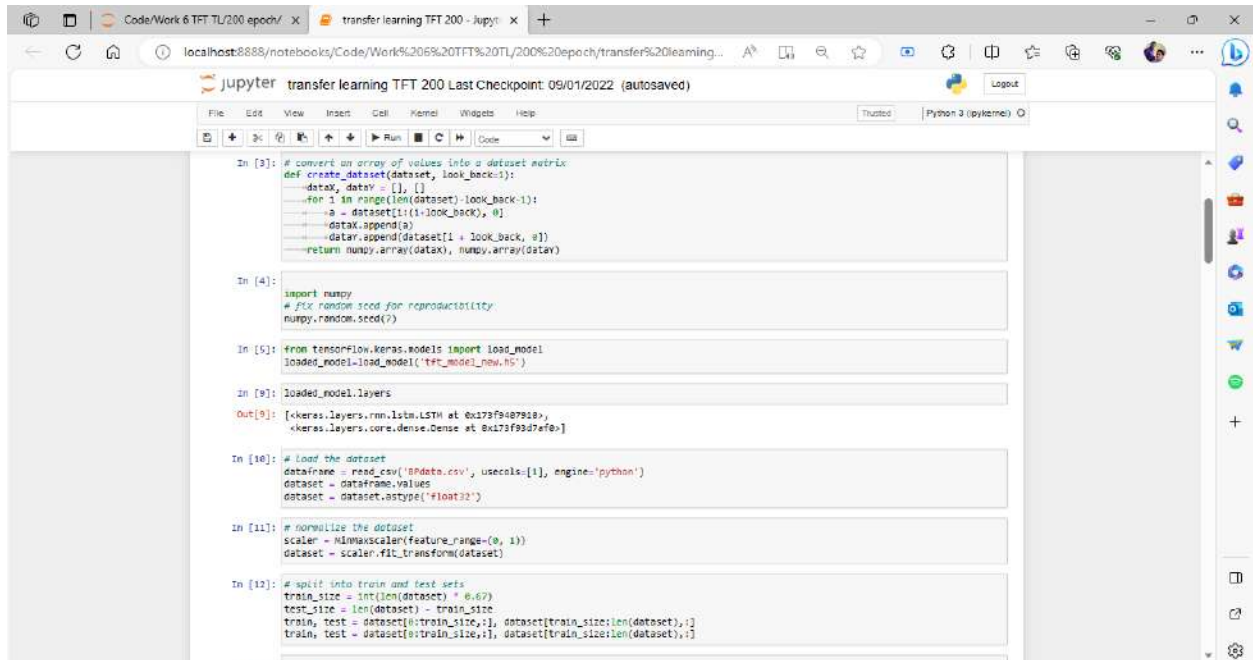
In [24]: # split into train and test sets
train_size = int(len(dataset) * 0.87)
test_size = len(dataset) - train_size
train, test = dataset[0:train_size,:], dataset[train_size:len(dataset),:]
train, test = dataset[0:train_size,1], dataset[train_size:len(dataset),1]

In [25]: # reshape into x and y+1
look_back = 1
trainx, trainy = create_dataset(train, look_back)
testx, testy = create_dataset(test, look_back)

In [26]: # reshape input to be [samples, time steps, features]
trainx = numpy.reshape(trainx, (trainx.shape[0], 1, trainx.shape[1]))
testx = numpy.reshape(testx, (testx.shape[0], 1, testx.shape[1]))

In [27]: import tensorflow as tf
```

Implementation of Transfer Learning with TFT based Pre-trained Model



```
In [3]: # convert an array of values into a dataset matrix
def create_dataset(dataset, look_back=1):
    datax, datay = [], []
    for i in range(len(dataset)-look_back-1):
        a = dataset[i:(i+look_back), 0]
        datax.append(a)
        datay.append(dataset[i + look_back, 1])
    return numpy.array(datax), numpy.array(datay)

In [4]:
import numpy
# fix random seed for reproducibility
numpy.random.seed(7)

In [5]: from tensorflow.keras.models import load_model
loaded_model=load_model('tft_model_h5')

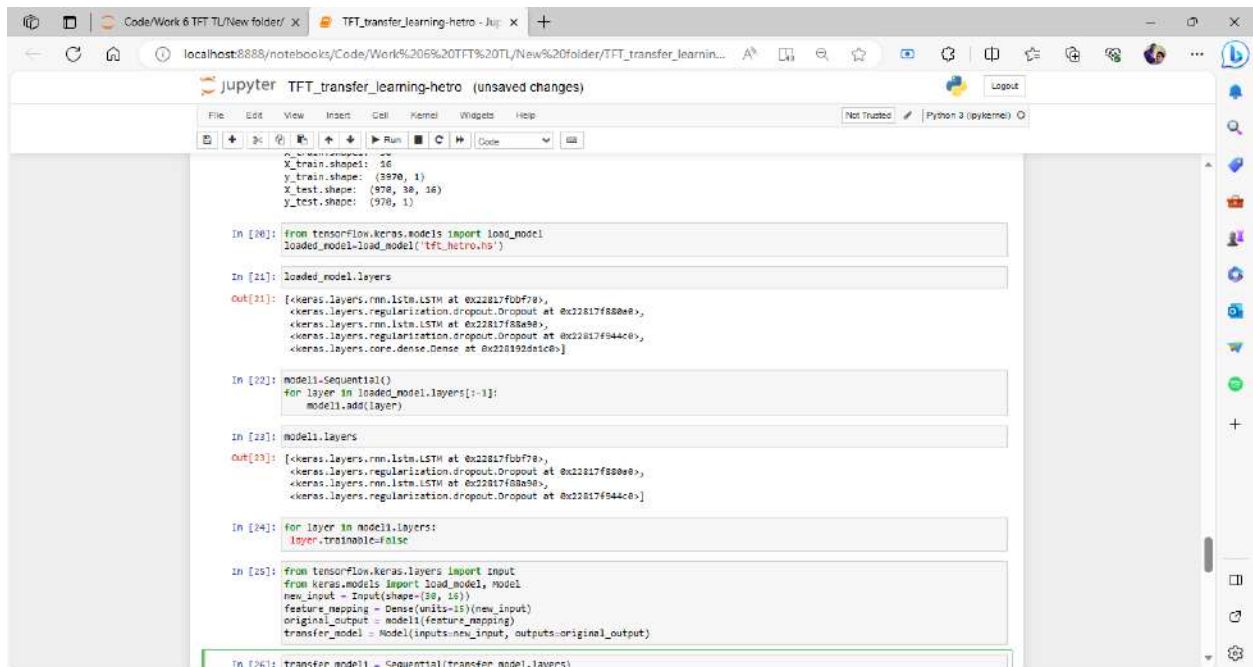
In [9]: loaded_model.layers
Out[9]: [keras.layers.rnn.lstm.LSTM at 0x173f9487918,
keras.layers.core.dense.Dense at 0x173f93d7ef0]

In [10]: # load the dataset
dataframe = read_csv('Bdata.csv', usecols=[1], engine='python')
dataset = dataframe.values
dataset = dataset.astype('float32')

In [11]: # normalize the dataset
scaler =MinMaxScaler(feature_range=(0, 1))
dataset = scaler.fit_transform(dataset)

In [12]: # split into train and test sets
train_size = int(len(dataset) * 0.67)
test_size = len(dataset) - train_size
train, test = dataset[0:train_size,:], dataset[train_size:len(dataset),:]
train, test = dataset[0:train_size,:], dataset[train_size:len(dataset),:]
```

Implementation of Heterogeneous Transfer Learning with TFT Pre-trained Model



```
X_train.shape: 16
y_train.shape: (3970, 1)
X_test.shape: (970, 16)
y_test.shape: (970, 1)

In [20]: from tensorflow.keras.models import load_model
loaded_model=load_model('tft_hetro.h5')

In [21]: loaded_model.layers
Out[21]: [keras.layers.rnn.lstm.LSTM at 0x22817f0b7f0,
keras.layers.regularization.dropout.Dropout at 0x22817f820e0,
keras.layers.rnn.lstm.LSTM at 0x22817f82a00,
keras.layers.regularization.dropout.Dropout at 0x22817f944c0,
keras.layers.core.dense.Dense at 0x228192551c0]

In [22]: model=Sequential()
for layer in loaded_model.layers[:-1]:
    model.add(layer)

In [23]: model.layers
Out[23]: [keras.layers.rnn.lstm.LSTM at 0x22817f0b7f0,
keras.layers.regularization.dropout.Dropout at 0x22817f820e0,
keras.layers.rnn.lstm.LSTM at 0x22817f82a00,
keras.layers.regularization.dropout.Dropout at 0x22817f944c0]

In [24]: for layer in model.layers:
    layer.trainable=False

In [25]: from tensorflow.keras.layers import Input
from keras.models import load_model, Model
new_input = Input(shape=(30, 16))
feature_mapping = Dense(units=15)(new_input)
original_output = model(feature_mapping)
transfer_model = Model(inputs=new_input, outputs=original_output)

In [26]: transfer_model = Sequential(transfer_model.layers)
```