*Article*

# Spatiotemporal Variation, Meteorological Driving Factors, and Statistical Models Study of Lake Surface Area in the Yellow River Basin

Li Tang and Xiaohui Sun *

College of Mining Engineering, Taiyuan University of Technology, Taiyuan 030024, China; li370336892@163.com
* Correspondence: sunxh17@mails.jlu.edu.cn; Tel.: +86-18204318089

**Abstract:** The surface area changes of 151 natural lakes over 37 months in the Yellow River Basin, based on remote sensing data and 21 meteorological indicators, employing spatial distribution feature analysis, principal component analysis (PCA), correlation analysis, and multiple regression analysis, identify key meteorological factors influencing these variations and their interrelationships. During the study period, lake area averages were from 0.009 km$^2$ to 506.497 km$^2$, with standard deviations ranging from 0.003 km$^2$ to 184.372 km$^2$. The coefficient of variation spans from 3.043 to 217.436, indicating considerable variability in lake area stability. Six primary meteorological factors were determined to have a significant impact on lake surface area fluctuations: 24 h precipitation, maximum daily precipitation, hours of sunshine, maximum wind speed, minimum relative humidity, and lakes in the source region of the Yellow River generally showed a significant positive correlation. For maximum wind speed (m/s), 28 lakes showed significant correlations, with five positive and twenty-three negative correlations, correlation coefficients ranging from −0.34 to −0.63, average −0.47, indicating an overall negative correlation between lake surface area and maximum wind speed. For maximum daily precipitation (mm), 36 lakes had 21 showing a positive correlation, indicating a positive correlation between lake surface area and daily precipitation in larger lakes. Furthermore, of the 117 lakes with sufficient data to model, the predictive capabilities of various models for lake surface area changes showcased distinct advantages, with the random forest model outperforming others in a dataset of 65 lakes, Ridge regression is best for 28 lakes, Lasso regression performs best for 20 lakes, Linear model is only best for 4 cases. The random forest model provides the best fit due to its ability to handle a large number of feature variables and consider their interactions, thereby offering the best fitting effect. These insights are crucial for understanding the influence of meteorological factors on lake surface area changes within the Yellow River Basin and are instrumental in developing predictive models based on meteorological data.

**Keywords:** Yellow River Basin; remote sensing data; lake surface area changes; meteorological factors; multiple regression analysis; principal component analysis

## 1. Introduction

Surface water bodies (SWBs), encompassing a wealth of natural lakes and widely distributed artificial reservoirs, bear crucial freshwater resources. SWBs are fundamental to China's geographical environment, ecosystems, and socio-economic development. Statistically, China's surface water bodies store substantial freshwater, playing a decisive role in national water resource security and ecological balance. In the complex hydrological cycle, surface water bodies serve a key role, regulating precipitation runoff, groundwater replenishment, and regional water balance. Additionally, they significantly contribute to carbon cycling [1], sediment and nutrient transport [2], and react markedly to climate change [3] by influencing local climatic conditions [4,5]. SWBs also foster unique and diverse ecosystems [6], providing extensive ecosystem services, such as food and water

supply. The surface area of SWBs is a crucial attribute, as it is the medium through which SWBs interact with many earth system processes and is closely linked to methane emissions, heat flux, and evaporation [7,8].

Lake surface area as the main role of SWBs' features and their driving mechanisms are significant topics in water resource management, ecological environmental protection, and climate change response studies. Recently, with the increase in the number of bands in remote sensing imagery and the growing global demand for water resource research, large-scale extraction of water bodies has become a focal point. Studies by Chen Chen et al., Mo Guifen et al., and Wang Lixuan et al., based on Landsat TM/ETM+ and OLI remote sensing images, revealed the spatiotemporal dynamics of surface water areas in the Altai Mountain ice lakes, the Central Asian countries, and the Sichuan–Tibet transportation corridor glaciers and analyzed their driving forces. They showed that ice lakes in the Altai region are highly sensitive to climate change; surface water area changes in the Central Asian countries are mainly influenced by socio-economic factors, with climate factors having a negligible impact; glacier melting intensity in the Sichuan–Tibet corridor varies significantly at different altitudes, with accelerated glacier area retreat and simultaneous expansion of surface water area in the 4501–5000 and 5001–5500 m elevation ranges [9–11]. Shi Jiancong et al. and Yuan Ruiqiang et al., based on Landsat series data, analyzed the spatiotemporal variations and driving factors of surface water in the Aral Sea basin and Inner Mongolia, showing a decreasing trend in the Aral Sea basin's surface water area, mainly driven by temperature among climate factors, while surface water changes in Inner Mongolia are complex, caused by both climate and human activities, with human activities being the main factor in the reduction of surface water area and lake shrinkage [12,13]. Shunburiji et al., based on 2009–2018 HJ-1A/B remote sensing data, studied the changes in surface water area in various leagues (cities) of the Inner Mongolia Autonomous Region, showing that the surface water area in Inner Mongolia continuously increased from 2009 to 2013 and sharply decreased from 2013 to 2017, with rainfall, runoff, reservoir construction, landfilling, and river diversion all being driving factors of these changes [14].

Existing research shows that the area of surface water bodies (SWBs) varies due to natural and climatic factors, exhibiting significant spatiotemporal differences. The Yellow River Basin is a critical area for China's ecological security and water resource management, providing essential support for agriculture, industry, and the livelihoods of millions. Moreover, the basin showcases a variety of ecological environments and climatic conditions, serving as an ideal natural laboratory for understanding the complex interactions between water bodies and climate change. It includes areas with diverse precipitation patterns, from arid deserts in the upper reaches to more humid climates downstream, offering a unique opportunity to study the impact of different meteorological conditions on lake dynamics. Additionally, the Yellow River Basin is experiencing significant environmental changes due to natural processes and human activities, including climate change, water diversion, and land-use change, further highlighting the need for a comprehensive analysis of its lake ecosystems. Understanding the spatiotemporal variability of lake surface areas in this region and identifying the meteorological factors driving these changes are crucial for developing sustainable water resource management strategies and adapting to climate change impacts. Zhang et al. developed a new combined extraction rule to build an entire annual-scale open-surface water body dataset for 1986–2020 in the Yellow River Basin using all of the available Landsat images [15], and Deng et al. used Landsat series images on the Google Earth Engine (GEE) platform, along with the HydroLAKES and China Reservoir datasets, to establish an extraction process for surface water bodies from 1986 to 2021 in the Yellow River Basin [16].

Based on current SWBs' surface area research, although extensive research has been conducted on the dynamics of SWBs' surface area using remote sensing technology, current studies focus on annual scale changes, with relatively less attention given to the detailed characteristics, storage dynamics, and responses to climate change of SWBs on a monthly scale [17,18].
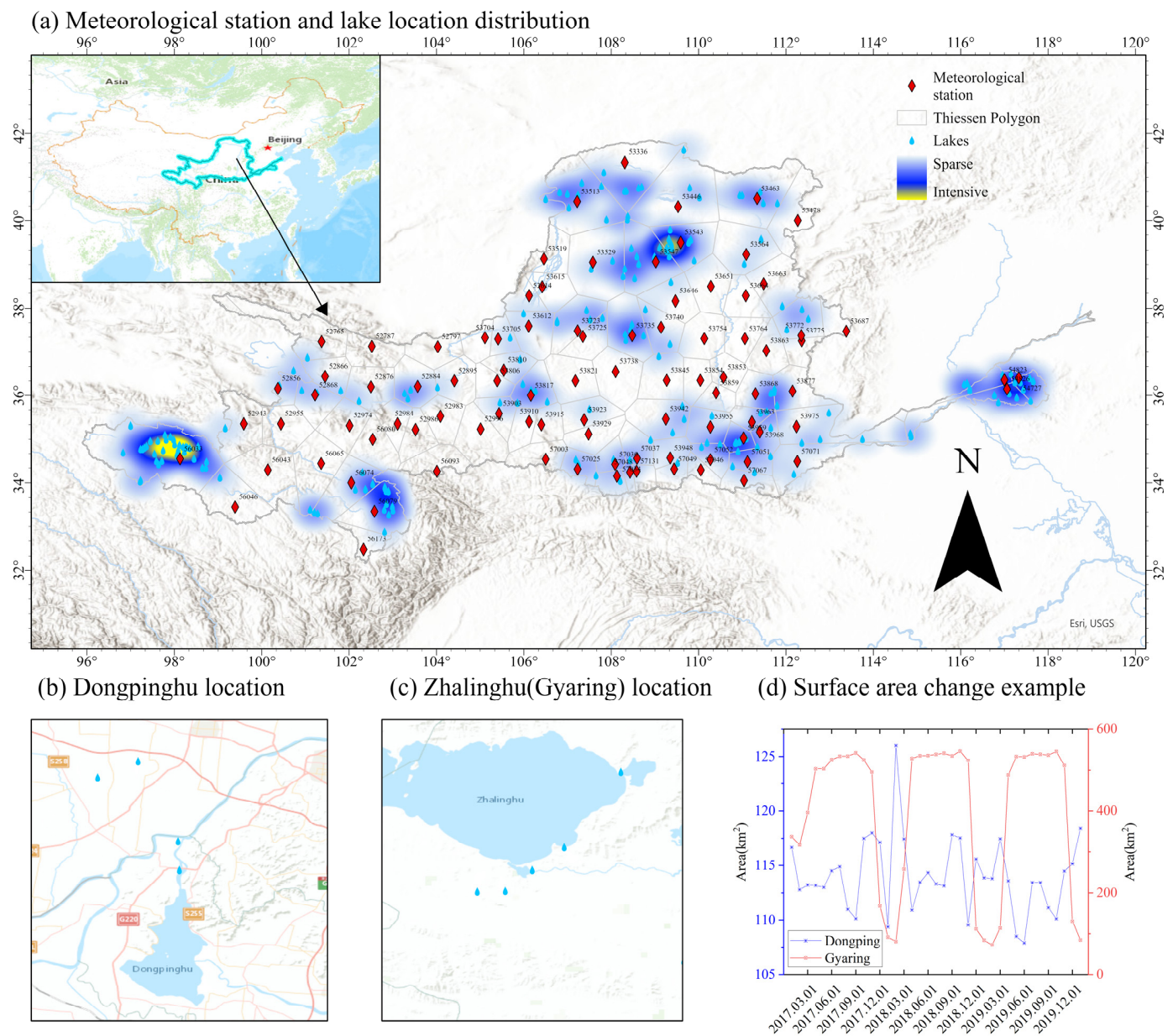
This is mainly attributed to the following limitations of remote sensing imagery: (1) Temporal resolution: Not all remote sensing satellites can provide the necessary temporal resolution to capture changes in lake area each month, and cloud cover and atmospheric conditions can limit the ability to obtain clear images [19]. (2) Cloud cover and atmospheric conditions: one of the main limitations of optical remote sensing is cloud obstruction, which complicates the accurate identification of lake boundaries [20]. (3) Image processing and data gaps: extracting water bodies from satellite imagery requires complex processing steps, and technical issues can lead to data gaps, affecting the creation of continuous monthly time series [21]. (4) Seasonal variability and dynamic water surfaces: changes in lake areas are influenced by seasonal precipitation, evaporation, and human activities, requiring precise remote sensing data and complex hydrological models to capture these changes [22]. (5) Spatial resolution: high-resolution imagery is necessary for monitoring changes in small lakes, but satellites with high-resolution imaging have lower coverage frequencies, which may not support monthly global monitoring [23]. Therefore, monthly scale research on SWBs' surface area based on optical imagery faces significant challenges. The presence of these issues poses significant challenges for current research on monthly scale surface water bodies' (SWBs') area time series based on optical imagery. To address these challenges, synthetic aperture radar (SAR) data, which can penetrate cloud cover, offers a viable alternative for water body detection methods that are not easily obstructed by clouds [24]. This approach can be enhanced by integrating data from multiple satellite systems and utilizing advanced cloud-penetrating radar imagery. To address these challenges, SAR has become a valuable tool because it can penetrate cloud cover [25]. The method of detecting water bodies using SAR data is not affected by cloud cover, allowing researchers to combine data from multiple satellite systems and use advanced radar imagery that can penetrate clouds. For example, researchers have analyzed global monthly scale surface water bodies' (SWBs') area using frequent, high-resolution C-band SAR observations provided by the Copernicus Sentinel-1 mission.

In summary, this study focuses on the Yellow River Basin in China, integrating previous research that utilized frequent, high-resolution C-band SAR observations from the Copernicus Sentinel-1 mission to analyze the lake surface area data of the Yellow River Basin, along with a meteorological dataset for the region. It aims to reveal the variability of lake water bodies in the Yellow River Basin and their climatic driving factors. The overall goal is to answer the following fundamental Earth science questions: What are the characteristics of monthly scale surface area changes in the Yellow River Basin? What are the meteorological driving factors behind the changes in natural lake surface areas, and how do they each contribute? The research findings are expected to provide valuable insights into the scientific understanding of hydrological and climatic processes in the Yellow River Basin, offering valuable information for policymakers and stakeholders involved in environmental protection and water resource planning in the region.

## 2. Materials and Methods

### 2.1. Study Area

The Yellow River Basin, stretching across the northern part of China, is the second-largest river basin in the country, covering an area of approximately 795,000 square kilometers (Figure 1). It is known for its complex hydrological and climatic systems, playing a crucial role in the ecological balance, agricultural productivity, and water resource management in the region [26,27]. The basin originates from the Bayan Har Mountains in Qinghai Province, flowing through nine provinces before discharging into the Bohai Sea. This extensive journey encompasses diverse climatic zones, ranging from arid and semi-arid climates in the upper and middle reaches to more humid conditions in the lower basin.

(a) Meteorological station and lake location distribution



(b) Dongpinghu location

(c) Zhalinghu(Gyaring) location

(d) Surface area change example



**Figure 1.** Location and distribution of meteorological stations in the Yellow River basin.

Climatically, the Yellow River Basin experiences significant variability, with precipitation patterns markedly changing across its expanse. The upper reaches are characterized by cold and dry conditions with minimal precipitation, while the middle reaches enjoy slightly higher rainfall, critical for agriculture and industry. The lower basin benefits from the East Asian monsoon, receiving the majority of its rainfall during the summer months, which significantly influences the hydrological regime and the availability of water resources.

The Yellow River Basin is home to numerous lakes and reservoirs, which serve as key water sources for irrigation, hydropower, and domestic use. Among these, notable lakes include Dongping, Hongze, Hulun, and the Wuliangsuhai, each playing a pivotal role in the basin's water system [28]. These lakes and reservoirs are not only essential for local water security but also support rich biodiversity and provide vital habitats for various aquatic and terrestrial species.

*2.2. Data*

93 meteorological stations across the Yellow River Basin were selected, covering a total of 37 months from January 2017 to January 2020 (Figure 1a). In total, 21 meteorological indices were analyzed as driving factors. The data were sourced from the China Meteorological Data Sharing Service System: https://data.cma.cn/ (accessed on 31 March 2024). The dataset is developed based on the measured meteorological data of China National Meteorological Information Center, and the data are reliable and in line with the reality.

The monthly surface area data for the lakes in this study were obtained from https://doi.org/10.1029/2022GL098987 (accessed on 31 March 2024) and https://doi.org/10.5281/zenodo.6345234 (accessed on 31 March 2024) [29].Taking Dongping and Gyaring as examples, the available time series of lake surface area is shown in Figure 1b–d. This database was developed using software that processes data from multiple sources, including help from Bruno Collischonn of the Brazilian National Water and Sanitation Agency (ANA) in locating the data: https://www.ana.gov.br/sar0/MedicaoSin (accessed on 31 March 2024), Google Earth Engine (GEE): https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S1_GRD (accessed on 31 March 2024); https://developers.google.com/earth-engine/datasets/catalog/JRC_GSW1_3_MonthlyHistory (accessed on 31 March 2024), and the HydroLAKES dataset: https://www.hydrosheds.org/page/hydrolakes (accessed on 31 March 2024), with the output available at https://doi.org/10.5281/zenodo.6450466 (accessed on 31 March 2024). The data are openly accessible. The database contains a total of 151 lakes in the Yellow River Basin (Figure 1a).

*2.3. Method*

Initially, we utilized GISPro and R language tools for our analysis. Based on meteorological station data and the boundaries of the Yellow River Basin, Thiessen polygons were generated for each weather station. These polygons facilitated the linkage between weather stations and lakes according to the spatial attributes of the Thiessen polygons (Figure 1a), thereby establishing a database for both lake surface area and meteorological time series.

Four characteristic parameters—mean, standard deviation, coefficient of variation, and slope (represents the rate of temporal change in lake surface area over the 37-month study period—were selected to analyze the central tendency, dispersion, and distribution shape of lake changes in the Yellow River Basin. Using the correlation coefficient, the correlation relationship between lake area changes and 21 meteorological factors was established. Based on the correlation coefficient matrix, the Principal Component Analysis (PCA) method was used to extract the main meteorological factors that are significantly correlated [30]. Further, the spatial distribution characteristics of these primary correlated meteorological factors were analyzed, and regression models were established to describe these relationships.

With the surface area of each studied lake as the target variable and selected meteorological factors as independent variables, this study employs four different regression models for predictive analysis, namely linear model (LM), ridge regression, lasso regression, and random forest (RF) [31]. To comprehensively evaluate and compare the performance of each model, this research has chosen the following three commonly used accuracy assessment indicators: mean absolute error (MAE), root mean square error (RMSE), and coefficient of determination ($R^2$).

Mean absolute error (MAE), root mean square error (RMSE) and coefficient of determination ($R^2$) are commonly used performance indicators in statistical analysis and prediction models. Here are their basic equations:

(1)    Mean absolute error (MAE)

MAE is the average of the absolute value of the difference between the observed and predicted values. It can be calculated by the following equation:

$$\text{MAE} = \frac{1}{n}\sum_{i=1}^{n} \mid y_i - \hat{y}_i \mid \qquad (1)$$

where $y_i$ is the $i$ true value, $\hat{y}$ is the $i$ predicted value, and $n$ is the sample size.

(2)    Root mean square error (RMSE)

RMSE is the square root of the mean of the square of the difference between the observed and predicted values. It can be calculated by the following equation:

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2} \tag{2}$$

where the definitions of $y_i$, $\hat{y}$, and $n$ are the same as above.

(3)    Coefficient of determination (R2)

$R^2$ is an index reflecting the goodness of fit of the model, and the closer the value is to 1, the better the model fit. It can be calculated by the following equation:

$$\text{R}^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \overline{y})^2} \tag{3}$$

where $\overline{y}$ is the average of all observations.

The evaluation of the models uses a 10-fold cross-validation method, which divides the dataset into ten equal parts, using nine parts for model training and the remaining one part for testing in a rotating cycle until each part has been used for testing. This process helps to assess the models' generalization ability on unseen data, providing robust insights into their predictive capabilities and potential performance on new observations.

Among these, the average value of RMSE serves as the primary criterion for selecting the optimal model. The model with the smallest average RMSE value is chosen as the final optimal model, aiming to achieve the best predictive performance on the given dataset.

Ultimately, the optimal regression model for each lake is selected, thereby establishing a predictive model library for the surface area of lakes in the Yellow River Basin. This approach ensures the robustness and reliability of the predictive models, enhancing our understanding of the spatiotemporal variations in lake surface areas within the region.
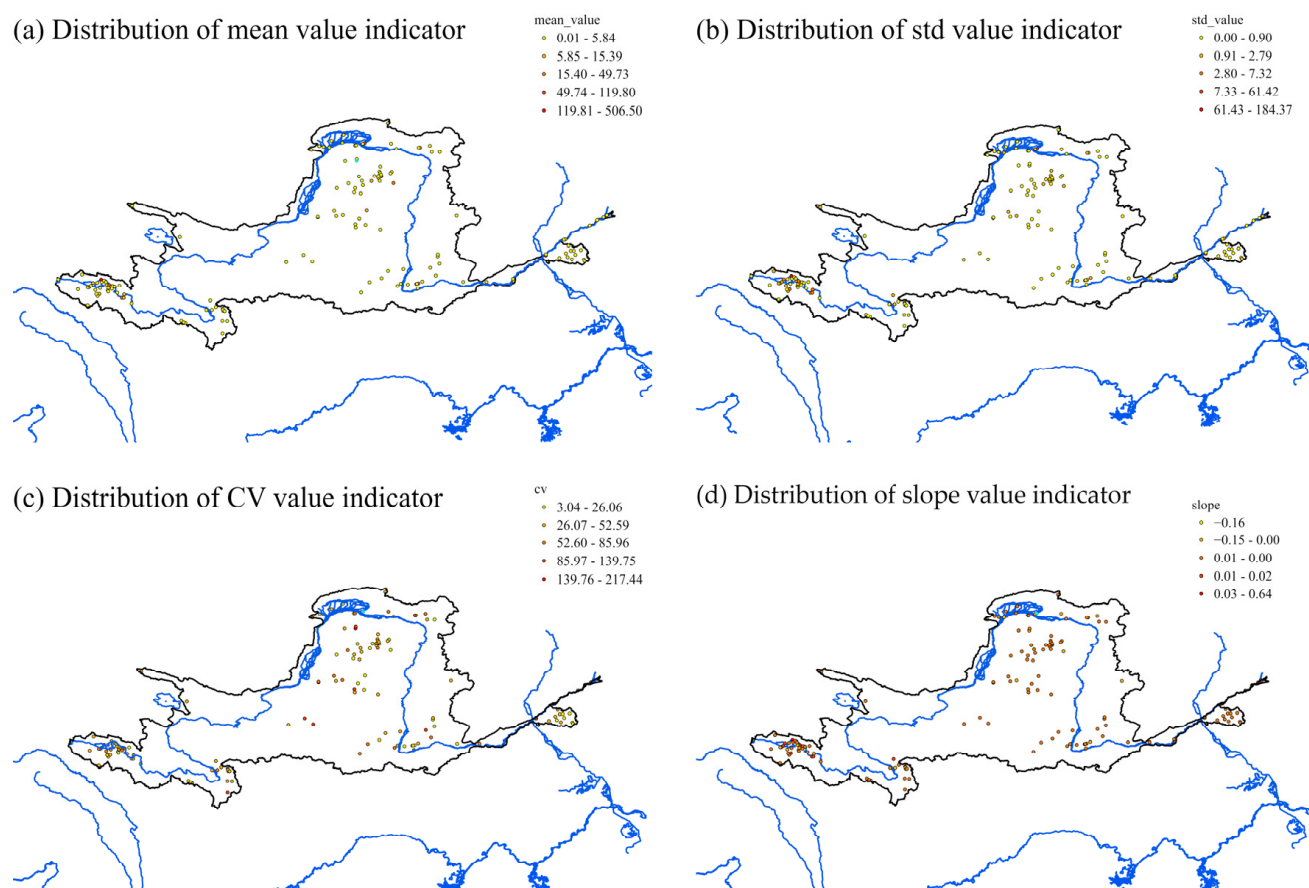
## 3. Results

### 3.1. Surface Area Characteristics of Lakes in the Yellow River Basin

Figure 2 shows the spatial distribution characteristics of the mean surface area, standard deviation (std), coefficient of variation (CV), and slope of 151 natural lakes in the Yellow River Basin over 37 months.

In the Yellow River basin, the average lake area varies from 0.009 km$^2$ (HyLake_ID 174904) to 506.497 km$^2$ (HyLake_ID 1377: Ngoring) as illustrated in Figure 2a. Larger average values are predominantly observed in the western plateau region, exemplified by HyLake_ID 1377: Ngoring and HyLake_ID 1385: Gyaring (403.265 km$^2$), reflecting the characteristic distribution of expansive lakes in the upper reaches of the Yellow River. Conversely, smaller average values are mainly found in the lower Yellow River and eastern areas, which may be associated with higher human activity and lower precipitation levels in these regions.

The standard deviation of lake area ranges from 0.003 km$^2$ (HyLake_ID 174904) to 184.372 km$^2$ (HyLake_ID 1385), with the spatial distribution of standard deviation exhibiting good consistency with the mean values (Figure 2a,b). This consistency suggests that the physical size of lakes (mean surface area) and their variability over time (indicated by standard deviation) may be influenced by similar natural and anthropogenic factors, which affect the stability and variability of lakes on a large scale. For instance, larger lakes within broad geographic regions may be more exposed to the impacts of large-scale climatic pattern changes, such as seasonal fluctuations in precipitation, directly affecting lake surface areas.

(a) Distribution of mean value indicator

mean_value
- 0.01 - 5.84
- 5.85 - 15.39
- 15.40 - 49.73
- 49.74 - 119.80
- 119.81 - 506.50

(b) Distribution of std value indicator

std_value
- 0.00 - 0.90
- 0.91 - 2.79
- 2.80 - 7.32
- 7.33 - 61.42
- 61.43 - 184.37

(c) Distribution of CV value indicator

cv
- 3.04 - 26.06
- 26.07 - 52.59
- 52.60 - 85.96
- 85.97 - 139.75
- 139.76 - 217.44

(d) Distribution of slope value indicator

slope
- −0.16
- −0.15 - 0.00
- 0.01 - 0.00
- 0.01 - 0.02
- 0.03 - 0.64

**Figure 2.** Spatial distribution characteristics of lakes in the Yellow River Basin.

The coefficient of variation, as a key indicator of the stability of lake area changes, spans widely in this dataset from 3.043 (HyLake_ID 1359) to 217.436 (HyLake_ID 172846) (Figure 2c). This range not only reveals the relative stability differences in lake area fluctuations but also reflects the sensitivity of lakes to external environmental changes. Lakes with higher coefficients of variation, such as HyLake_ID 172846, HyLake_ID 173698, and HyLake_ID 174535, exhibit distinct spatial clustering characteristics, primarily concentrated in specific areas of the Yellow River basin: the Mu Us Desert and the Zhengzhou to Bohai segment of the river's lower reaches. This distribution pattern not only reflects the regional characteristics of geographical and climatic factors' impact on lake area changes but also suggests the presence of similar ecological conditions and hydrological dynamics in these regions.

During the study period in the Yellow River basin, the slope of lake area change trends exhibited significant variability, ranging from −0.161 (HyLake_ID 1314: Wu-liang-su) to 0.635 (HyLake_ID 1385: Gyaring) (Figure 2d). This variation unveils the differing trends of lake area expansion or reduction over time within the region, reflecting the unique hydrological and environmental conditions of individual lakes. Specifically, lakes with a positive slope, such as Lake Gyaring (HyLake_ID 1385), demonstrated a noticeable increase in surface area during the observation period. This growth could be closely related to regional increases in precipitation, accelerated snowmelt processes due to rising temperatures, and other changes in the watershed's hydrological cycle. These shifts indicate that some lakes are experiencing accumulations and expansions of water, which could have significant implications for local ecosystems and water resource management. In contrast, lakes with a negative slope, such as Wu-liang-su Lake (HyLake_ID 1314), showed a decreasing trend, possibly indicating water body shrinkage and lake degradation. This reduction could result from the overexploitation of water resources, such as irrigation

and industrial water use, or due to climate change-induced decreases in precipitation and increased evaporation rates. These findings underscore the importance of sustainable water resource management and the urgent need for climate change adaptation strategies.

Analysis of the relationship models between various indicators and geographic coordinates (Figure 3 and Table 1) indicates that most indicators show either minimal explanation for data variability or lack statistical significance in relation to latitude and longitude. This may suggest that these indicators are less influenced by geographic coordinates or that other unconsidered factors are affecting them. Models of standard deviation against latitude, standard deviation against longitude, and slope against longitude revealed some statistical significance, especially the relationship between standard deviation (Std_value) and longitude, which was most pronounced. This suggests that data variability (standard deviation) and change trends (slope) might vary to some extent across different geographic locations.
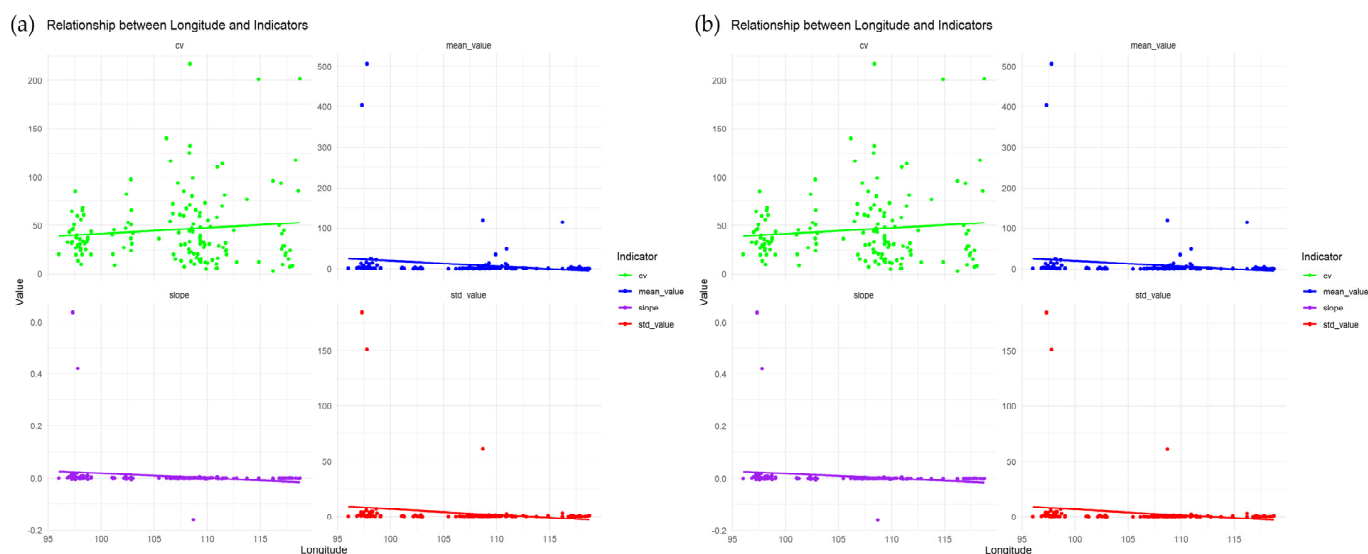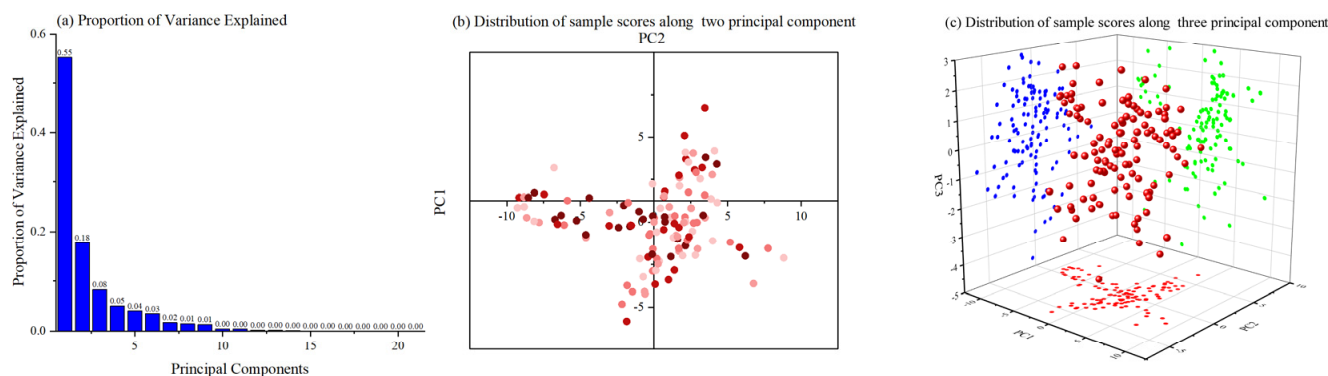


**Figure 3.** Indicators and latitude-longitude relationship curve.

**Table 1.** Relation curve of lake surface area change index and latitude and longitude in Yellow River basin.

| Index | Regression Model | $R^2$ | $p$ |
|---|---|---|---|
| mean_value__lat | Value = 72.727 − 1.691 × Pour_lat | 0.006 | 0.363 |
| mean_value_lon | Value = 149.650 − 1.301× Pour_long | 0.024 | 0.056 |
| Std_value__lat | Value = 22.130 − 0.509 × Pour_lat | 0.004 | 0.046 |
| Std_value_lon | Value = 61.710 − 0.545 × Pour_long | 0.031 | 0.029 |
| CV_value__lat | Value = 4.360 + 1.122× Pour_lat | 0.005 | 0.400 |
| CV_value_lon | Value = −17.618 + 0.591 × Pour_long | 0.011 | 0.203 |
| Slope_value__lat | Value = 0.0127 − 0.003 × Pour_lat | 0.015 | 0.133 |
| Slope_value_lon | Value = 0.197 − 0.002 × Pour_long | 0.033 | 0.025 |

### 3.2. Indicator Selection

Figure 4 presents the results of a principal component analysis (PCA) examining the correlations between various meteorological indicators and lake surface area, including a scree plot of the variance contributions of each principal component to the total dataset variance and a biplot. The analysis reveals that the first principal component (PC1) accounts for 55.16% of the variability in the data, and the second principal component (PC2) captures an additional 17.75%. When considering the three principal components (PC1–PC3) together, they account for 81.217% of the total variation in the dataset (Figure 4a). This significant proportion suggests that these components are sufficient to represent the majority of the information and structure within the dataset, making them pivotal in understanding the underlying patterns.

**Figure 4.** Principal component analysis results of the correlation coefficient matrix.

Figure 4b illustrates the distribution of sample scores along the first principal component (principal component 1, horizontal axis) and the second principal component (principal component 2, vertical axis). The biplot in Figure 4b displays the dispersion of samples primarily along the horizontal axis, captured by PC1, which accounts for the largest proportion of variability. In contrast, PC2 reveals the second-largest proportion of variability along the vertical direction. PC3 offers an additional perspective on the distribution of data points in the depth dimension (Figure 4c).

In terms of component loadings, 20-20 hourly precipitation (mm) and #maximum daily precipitation (mm) exhibit high loadings on PC1, indicating their significant contribution to this component. For PC2, variables such as hours of sunshine and maximum wind speed (m/s) demonstrate high loadings, signifying their pivotal role on this axis. On PC3, minimum relative humidity (%) and average 2-minute wind speed (m/s) show high loadings (Table 2).

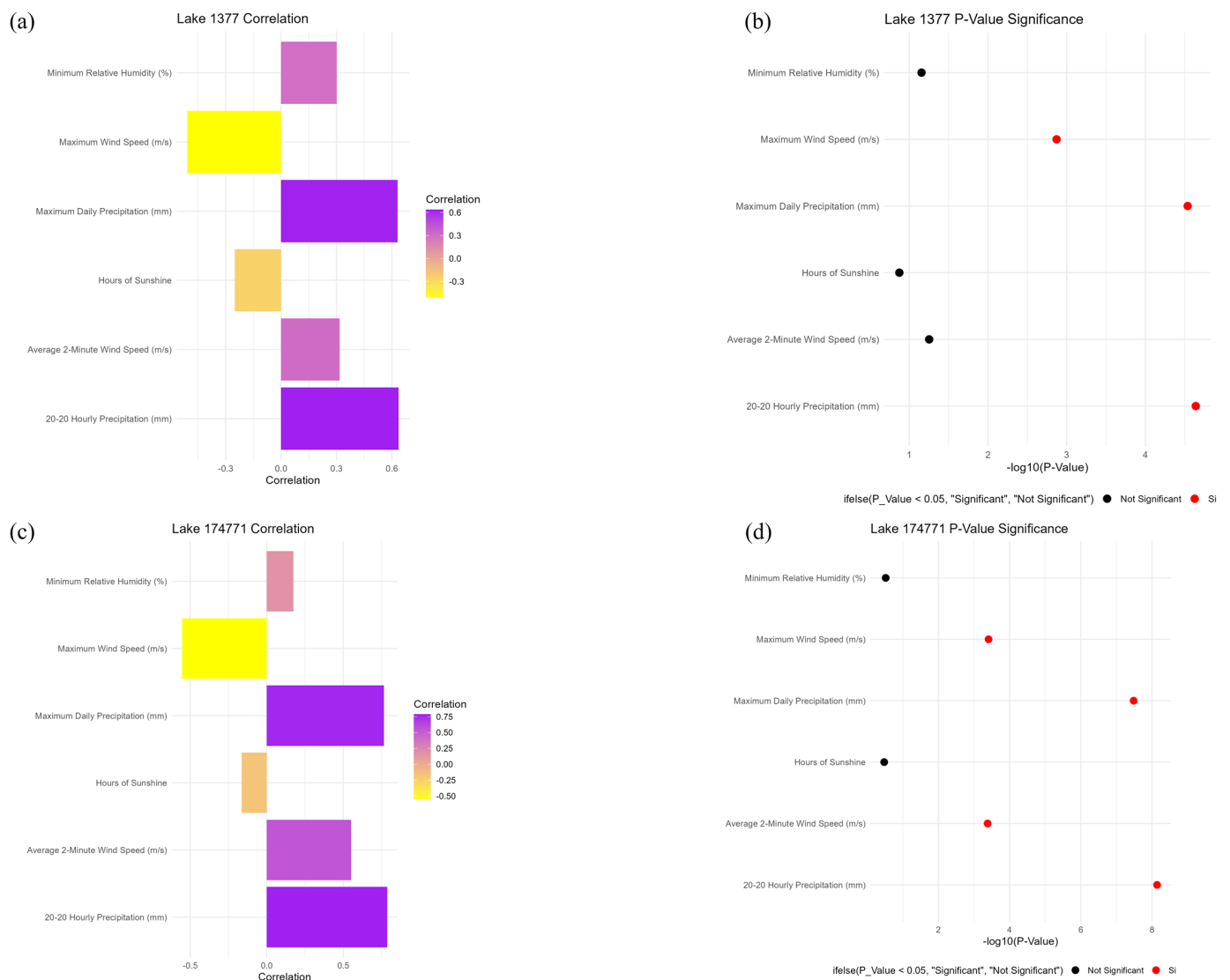**Table 2.** Principal component loading matrix.

| Index | PC1 | PC2 | PC3 |
|---|---|---|---|
| # 20-20 hourly precipitation (mm) | −0.284 * | −0.064 | −0.063 |
| # Maximum wind speed (m/s) | 0.166 | −0.355 ** | 0.027 |
| Maximum wind speed direction(°) | 0.140 | −0.107 | −0.220 |
| # Average 2-minute wind speed (m/s) | −0.172 | −0.123 | 0.454 ** |
| Average temperature (°C) | −0.273 | −0.182 | −0.031 |
| Average atmospheric pressure (hPa) | −0.186 | 0.331 | 0.221 |
| Average vapor pressure (hPa) | −0.282 | −0.121 | −0.082 |
| Average relative humidity (%) | −0.233 | 0.165 | −0.341 |
| Average minimum temperature (°C) | −0.275 | −0.173 | −0.041 |
| Average maximum temperature (°C) | −0.271 | −0.191 | −0.020 |
| Number of days with precipitation ≥ 0.1 mm | −0.276 | −0.086 | −0.120 |
| # Hours of sunshine | 0.020 | −0.409 * | −0.028 |
| Monthly percentage of sunshine (%) | 0.033 | 0.144 | 0.130 |
| Maximum wind speed (m/s) | 0.148 | −0.313 | 0.161 |
| Direction of maximum wind speed (°) | 0.187 | −0.006 | −0.330 |
| # Maximum daily precipitation (mm) | −0.283 ** | −0.064 | −0.052 |
| Minimum temperature (°C) | −0.275 | −0.170 | −0.042 |
| Minimum atmospheric pressure (hPa) | −0.218 | 0.289 | 0.114 |
| Maximum temperature (°C) | −0.268 | −0.198 | 0.011 |
| Maximum atmospheric pressure (hPa) | −0.173 | 0.288 | 0.317 |
| # Minimum relative humidity (%) | −0.109 | 0.245 | −0.529 * |

Note: "#" represents the selected indicators, "*" represents the maximum absolute value, and "**" represents the second largest absolute value.

After evaluating the contributions of individual variables to each principal component, this study identifies the top two variables with the highest loadings from each of the three principal components analyzed, totaling six primary indicators for subsequent spatiotemporal correlation and modeling analyses.

### 3.3. Correlation Analysis

Taking Lake HyLake_ID 1377 (Ngoring, mean area = 506.500 km$^2$) and HyLake_ID 174771 (mean area = 2.960 km$^2$) as examples, the analysis demonstrates significant correlations between meteorological factors and lake surface areas for these distinct environments (Figure 5).
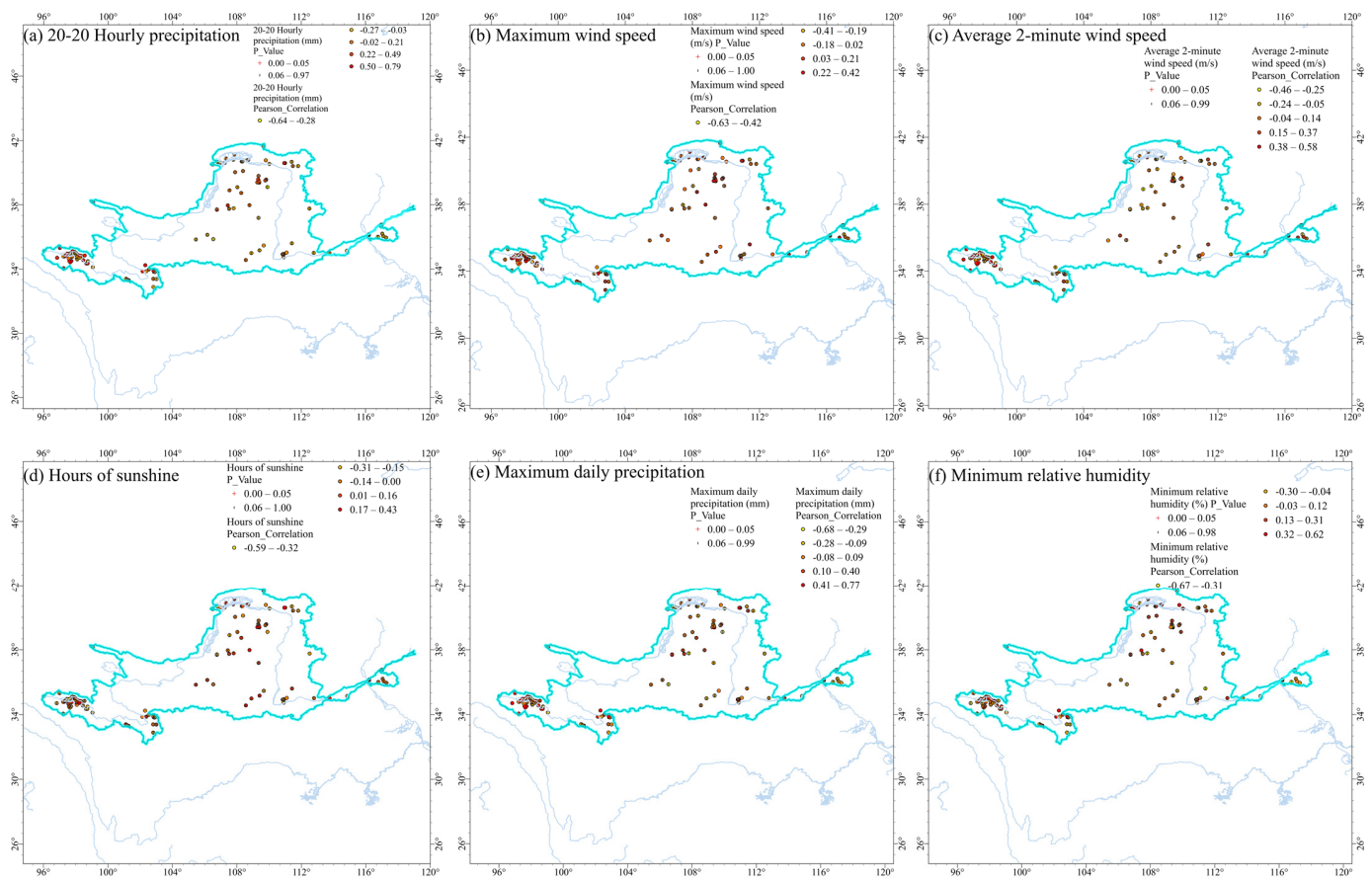


**Figure 5.** Correlation between the changes in surface area and meteorological factors.

For HyLake_ID 1377, there is a strong positive correlation between "20-20 hourly precipitation" and the lake surface area, with a correlation coefficient of 0.6364 and a $p$-value < 0.01, suggesting that increased precipitation is associated with an expansion of the lake surface area. Conversely, "maximum wind speed" shows a moderate negative correlation (correlation coefficient of $-0.508$, $p$-value < 0.01), indicating that higher wind speeds may be associated with a reduction in lake surface area. Other meteorological factors, such as "hours of sunshine" and "maximum daily precipitation", also show strong positive correlations with the lake surface area, with correlation coefficients of 0.6303 ($p$-value < 0.01) and 0.3015 ($p$-value = 0.070), respectively. However, the correlation between "minimum relative humidity" and lake surface area is weaker, with a correlation coefficient of 0.3015 and a $p$-value of 0.070, suggesting a less significant relationship.

Similarly, for HyLake_ID 174771, significant correlations are observed between meteorological factors and lake surface area. "20-20 hourly precipitation" shows a strong positive correlation with the lake surface area (correlation coefficient of 0.788, $p$-value < 0.01), reinforcing the idea that precipitation is a critical factor in lake surface dynamics. Unlike Lake 1377, "maximum wind speed" for Lake 174771 exhibits a slightly weaker negative correlation with the lake surface area (correlation coefficient of $-0.553$, $p$-value < 0.01), which may indicate a different impact of wind on smaller lakes. Furthermore, "hours of sunshine" and "maximum daily precipitation" have substantial positive correlations with the lake surface area, with correlation coefficients of 0.766 ($p$-value < 0.01) and 0.1745 ($p$-value = 0.302), respectively. The correlation between "minimum relative humidity" and lake surface area, similar to Lake 1377, remains weaker and not significant, with a correlation coefficient of 0.175 and a $p$-value of 0.302, indicating minimal impact on the lake's size.

These findings suggest that, despite the considerable size difference between Lakes 1377 and 174771, both lakes exhibit similar trends in the influence of meteorological factors on their surface areas. 20-20 hourly precipitation, maximum daily precipitation and maximum wind speed significantly impact lake surface areas, whereas wind speed shows moderate negative correlations, and relative humidity appears to have minimal effects.

To elucidate the spatial distribution characteristics of the correlation and significance between lake surface areas and meteorological factors in the Yellow River Basin, we created a distribution map showing the correlation and significance between lakes in the Yellow River Basin and changes in meteorological factors, as illustrated in Figure 6. For the 20-20 hourly precipitation (mm), among 118 lakes, 38 exhibited significant correlations, with 22 positively correlated (correlation coefficients ranging from 0.41 to 0.79, average 0.59) and 16 negatively correlated (correlation coefficients ranging from $-0.36$ to $-0.64$, average $-0.43$). Spatially, lakes in the source region of the Yellow River generally showed a significant positive correlation. For maximum wind speed (m/s), 28 lakes showed significant correlations, with five positive and twenty-three negative correlations (correlation coefficients ranging from $-0.34$ to $-0.63$, average $-0.47$), indicating an overall negative correlation between lake surface area and maximum wind speed. For average 2-minute wind speed (m/s) and hours of sunshine, 25 and 22 lakes, respectively, showed significant correlations without a clear pattern. For maximum daily precipitation (mm), 36 lakes had 21 showing a positive correlation, indicating a positive correlation between lake surface area and daily precipitation in larger lakes. For minimum relative humidity (%), 24 lakes exhibited significant correlations but without a discernible pattern.

**Figure 6.** Spatiotemporal distribution of the correlation between lake variations and meteorological factors.
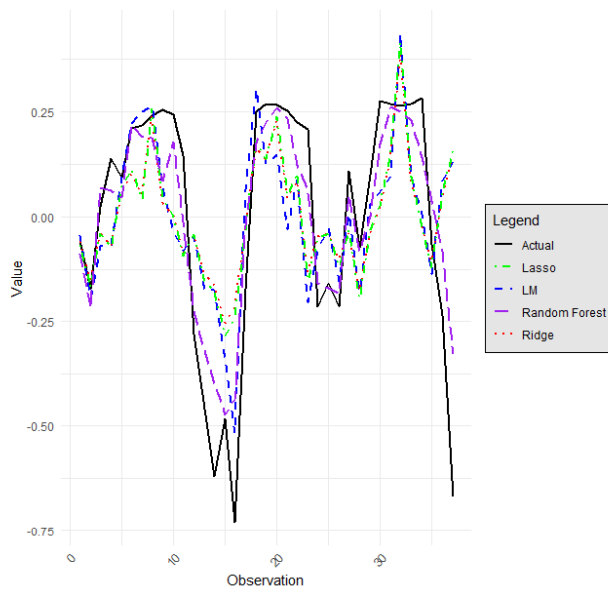
### 3.4. Multivariate Regression Analysis

Taking Lake HyLake_ID 1377 (Ngoring, mean area = 506.500 km$^2$) and HyLake_ID 173250 (mean area = 1.620 km$^2$) as examples, the model analysis results are presented in Figure 7. For HyLake_ID 1377, the lasso model demonstrates the lowest average root mean square error (RMSE) of 0.247, outperforming other models in this analysis. Its consistent and low error rate across different validation sets indicates superior stability and generalization ability. In contrast, for HyLake_ID 173250, the linear regression model achieves the lowest average RMSE of 0.227, marking it as the best performer in this instance. This comparison highlights the variability in model responses across different lakes within the basin, underscoring the importance of model selection tailored to specific lake characteristics.
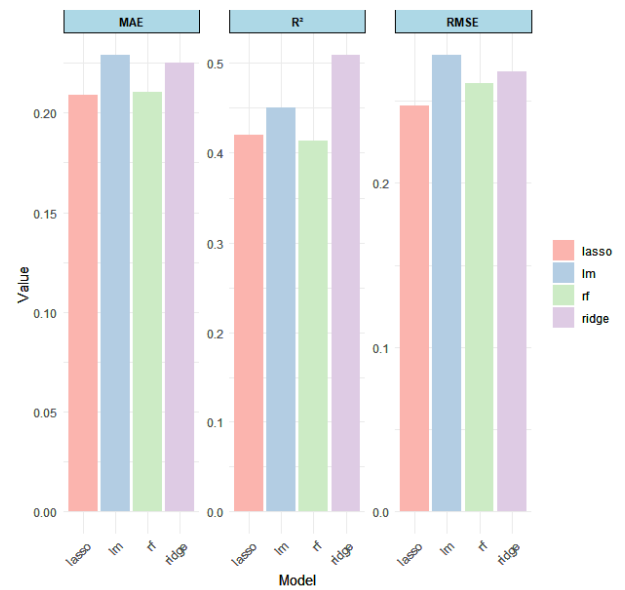
Figure 8 reflects the optimal model fitting results for the surface area of lakes in the Yellow River Basin. The random forest (RF) model performs best in 65 lakes, demonstrating its superiority in dealing with the relationship between the surface area of lakes in the Yellow River Basin and meteorological factors. The strength of the random forest model lies in its ability to handle a large number of input variables and automatically select variables, thus providing deep insights into complex data relationships. Ridge regression is best for 28 lakes, indicating that the introduction of L2 regularization can effectively improve the model's predictive accuracy and stability when data exhibit multicollinearity. Lasso regression performs best for 20 lakes; its use of L1 regularization helps in simplifying variables and enhancing the model's interpretability, which is particularly important in determining the impact of key meteorological factors on the changes in lake surface area. Although the linear model is only best for four cases, it remains the foundation for analyzing linear relationships, providing us with an initial benchmark for model comparison (Figure 7a).

These results suggest that nonlinear models (such as RF) might be more suitable for capturing the complex dynamic relationships between the surface area of lakes in the Yellow River Basin and meteorological factors. The random forest model provides the best fit due to its ability to handle a large number of feature variables and consider their interactions, thereby offering the best fitting effect. Meanwhile, regularized linear models (lasso and ridge) demonstrate robustness in datasets with high multicollinearity, which is crucial for reducing model overfitting and improving prediction accuracy.
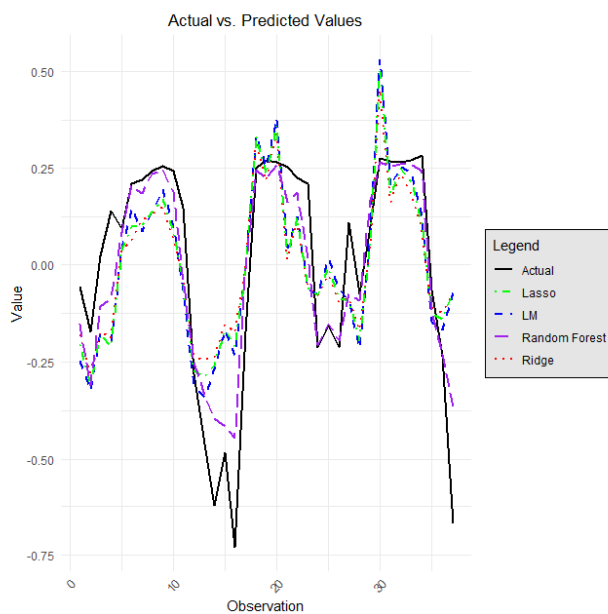
(a) HyLake_ID 1377 actual and simulated value
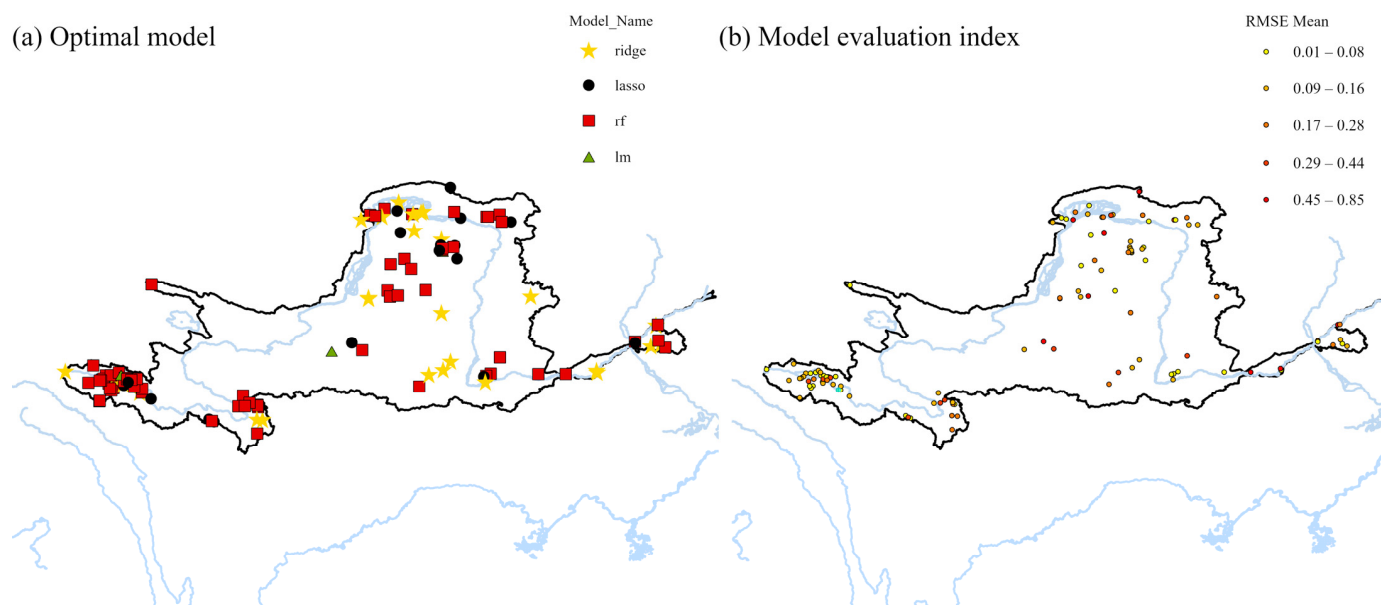
(b) HyLake_ID 1377 model evaluation

(c) HyLake_ID 173250 actual and simulated value

(d) HyLake_ID 173250 model evaluation

**Figure 7.** Example of model analysis results.

(a) Optimal model

Model_Name
★ ridge
● lasso
■ rf
▲ lm

(b) Model evaluation index

RMSE Mean
○ 0.01 − 0.08
○ 0.09 − 0.16
○ 0.17 − 0.28
○ 0.29 − 0.44
○ 0.45 − 0.85

**Figure 8.** Model analysis results.

In comparing the fitting results of the average surface area of lakes in the Yellow River Basin across different models, the Lasso model exhibits relatively lower average surface area values, with data points primarily clustered in the lower range and with relatively minor dispersion. The linear model (LM) results show a wider distribution, with data points stretching from values close to 0 up to higher values, though most remain concentrated in the lower range of average surface area. The random forest (RF) model outcomes are scattered across the entire range of values, with a notably higher outlier point visible in the graph, suggesting that the random forest model may predict larger average lake surface areas in certain cases. The ridge regression (ridge) model displays relatively greater variability, with data points concentrated across various average surface area values, including some higher ones (Figure 7b). Overall, there are certain disparities in the fitting effects of the models on the average surface area of lakes in the Yellow River Basin.

## 4. Discussion

This study conducts a comprehensive analysis of the variability of lake surface area in the Yellow River Basin and its influencing factors through the integrated application of spatial distribution analysis, principal component analysis (PCA), correlation analysis, and multiple regression models. Due to the abundance and accessibility of precipitation and temperature data, many studies have identified precipitation and temperature as the primary climatic factors, demonstrating their significant impact on lake surface area [32–34]. Our research further selects the main factors from a large set of meteorological factors, identifying maximum wind speed (m/s), average 2-minute wind speed (m/s), hours of sunshine, and minimum relative humidity (%) as key elements affecting changes in lake surface area. It quantifies the specific impact of these main meteorological elements on lake area and establishes regression models to analyze the combined effects of multiple meteorological factors.

Previous studies on changes in lake surface area often focused on uniform change patterns within a large scale, such as considering the variations over many years, while research on how monthly-scale meteorological conditions affect lake systems based on high-resolution C-band SAR observations from the Copernicus Sentinel-1 is relatively scarce [35,36]. Our study captures the rapid response of lake surface area to seasonal meteorological condition changes. This fine-grained temporal-scale analysis provides a more acute and timely understanding of the hydrological cycle and ecosystem changes in lakes, aiding in a better comprehension of the short-term responses of lakes to climate change.

This discovery of a negative correlation between lake area and maximum wind speed prompts a deeper consideration of how lake ecosystems respond to climatic factors. This negative correlation may reflect the influence of climatic conditions around the lake on the lake's hydrodynamics and hydrological processes. Firstly, the negative correlation could be associated with evaporation from the lake surface. Under conditions of higher wind speeds, the rate of evaporation from the water surface may increase, leading to a decrease in water levels or increased evaporation, which could result in a reduction in lake area, thereby showing a negative correlation with maximum wind speed [37,38]. Secondly, higher wind speeds are often associated with cyclonic systems in the climate, which can be accompanied by precipitation events [39]. During precipitation, lake water levels might rise, causing an expansion of lake area [40]. Therefore, the negative correlation between lake water levels and wind speed might reflect these climate-driven hydrological processes, where lake water level changes are influenced by both wind speed and cyclonic systems. In summary, the negative correlation between lake water bodies and maximum wind speed indicates the sensitivity of lake ecosystems to climatic changes. Further research could explore how lake hydrological processes are affected by climate change and extreme weather events, and how these changes impact the stability and functionality of lake ecosystems.

Analysis of the results from our models shows differences in the capacity to fit relationships between lake surface area and meteorological factors across the Yellow River Basin due to the individual differences in lakes as entities within ecosystems. It is challenging to explain their internal variation rules with a unified regional-scale model. Our study focuses on individual lakes within the research scale, establishing specific surface area change meteorological driving models for different lakes, highlighting the individual differences in each lake as an independent ecosystem. This approach overcomes the problem of unified regional-scale models failing to explain all lakes' internal variation rules. By customizing models for each lake, our method provides a basis for implementing targeted environmental management strategies. This customized approach is more likely to successfully address the specific issues faced by particular lakes, thereby improving resource utilization efficiency and conservation effects.

Despite providing new insights into the variability of lake surface areas in the Yellow River Basin, we also recognize the limitations of our study. Firstly, although we considered multiple meteorological factors in our models, the lack of consideration for anthropogenic factors and other potential influences, such as contributions from groundwater flow and ice/snow melt to lake water volumes, were not covered in this analysis. Secondly, due to limitations in climate models and data acquisition capabilities, predictions of lake area changes under future climate change scenarios remain uncertain. Future research should aim at more refined models and more comprehensive data collection to reveal these complex dynamic processes.

## 5. Conclusions

This study investigates the surface area changes of 151 natural lakes over 37 months in the Yellow River Basin, using spatial distribution feature analysis, PCA, correlation analysis, and multiple regression analysis to reveal the key factors affecting lake surface area changes and their correlation with meteorological factors, and assesses the applicability of different statistical models in specific environments. Key findings include the following:

(1) Key influencing factors: Six main meteorological factors were identified as having a significant impact on lake surface area changes, including 24 h precipitation, maximum daily precipitation, hours of sunshine, maximum wind speed, minimum relative humidity, and average 2 min wind speed. These factors play a decisive role in the dynamics of lake water balance and surface area changes.

(2) The correlation and spatial distribution analysis between lake surface area and meteorological factors showed that lakes in the source area of the Yellow River usually have a significant positive correlation with 24 h precipitation, while most lakes exhibit a

negative correlation with maximum wind speed, reflecting clear spatial differences in the response of lakes in different geographical locations to meteorological changes.

(3) Spatial variability in model performance: In predicting changes in lake surface area in the Yellow River Basin, different models showed their own advantages. The random forest (RF) model performed best across a dataset of 65 lakes, proving its superiority in handling the complex dynamics between lake surface area and meteorological factors in the Yellow River Basin.

In summary, the findings of this study are significant for understanding how meteorological factors influence changes in lake surface area within the Yellow River Basin and provide a scientific basis for lake surface area change prediction models based on meteorological data. Future research needs to further explore other potential environmental and anthropogenic factors and how these factors affect lake surface area changes through complex mechanisms. This study not only provides valuable insights for lake management and conservation in the Yellow River Basin and similar regions globally but also presents new challenges and opportunities for addressing climate change and sustainable water resource management.

**Author Contributions:** Conceptualization, L.T. and X.S.; methodology, L.T. and X.S.; formal analysis, L.T. and X.S.; data curation, L.T. and X.S.; writing—original draft preparation, L.T.; writing—review and editing, L.T. All authors have read and agreed to the published version of the manuscript.

**Data Availability Statement:** Data are contained within the article.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Cole, J.J.; Prairie, Y.T.; Caraco, N.F.; McDowell, W.H.; Tranvik, L.J.; Striegl, R.G.; Melack, J. Plumbing the global carbon cycle: Integrating inland waters into the terrestrial carbon budget. *Ecosystems* **2007**, *10*, 172–185. [CrossRef]
2. Harrison, J.A.; Frings, P.J.; Beusen AH, W.; Conley, D.J.; McCrackin, M.L. Global importance, patterns, and controls of dissolved silica retention in lakes and reservoirs. *Glob. Biogeochem. Cycles* **2012**, *26*, GB2037. [CrossRef]
3. Tranvik, L.J.; Downing, J.A.; Cotner, J.B.; Loiselle, S.A.; Striegl, R.G.; Ballatore, T.J.; Weyhenmeyer, G.A. Lakes and reservoirs as regulators of carbon cycling and climate. *Limnol. Oceanogr.* **2009**, *54 Pt 2*, 2298–2314. [CrossRef]
4. Wohlfahrt, G.; Tomelleri, E.; Hammerle, A. The albedo–climate penalty of hydropower reservoirs. *Nat. Energy* **2021**, *6*, 372–377. [CrossRef] [PubMed]
5. Balsamo, G.; Salgado, R.; Dutra, E.; Boussetta, S.; Stockdale, T.; Potes, M. On the contribution of lakes in predicting near-surface temperature in a global weather forecasting model. *Tellus A Dyn. Meteorol. Oceanogr.* **2016**, *68* (Suppl. 1), 15829. [CrossRef]
6. Gleick, P.H. Global freshwater resources: Soft-path solutions for the 21st century. *Science* **2003**, *302*, 1524–1528. [CrossRef]
7. Friedrich, K.; Grossman, R.L.; Huntington, J.; Blanken, P.D.; Lenters, J.; Holman, K.D.; Kowalski, T. Reservoir evaporation in the Western United States: Current science, challenges, and future needs. *Bull. Am. Meteorol. Soc.* **2018**, *99*, 167–187. [CrossRef]
8. Bastviken, D.; Cole, J.; Pace, M.; Tranvik, L. Methane emissions from lakes: Dependence of lake characteristics, two regional assessments, and a global estimate. *Glob. Biogeochem. Cycles* **2004**, *18*, GB4009. [CrossRef]
9. Chen, C.; Zheng, J.H.; Liu, Y.Q.; Xu, Z.L. Spatiotemporal characteristics of glacier lakes in the Altai Mountains of China and their responses to regional climate change over the past 20 years. *Geogr. Res.* **2015**, *2*, 270–284. (In Chinese)
10. Mo, G.F.; Feng, J.Z.; Bai, L.Y.; Wang, Z.M.; Li, H.L.; Yu, T. Spatiotemporal variation characteristics of surface water resources in the arid region of Central Asia from 2001 to 2018. *Geogr. Sci.* **2022**, *1*, 174–184. (In Chinese)
11. Wang, L.X.; Ye, C.M.; Sui, T.B.; Wei, R.L.; Li, H.F. Remote sensing monitoring and coupled analysis of glaciers and surface water in the Sichuan-Tibet transportation corridor. *Bull. Surv. Mapp.* **2023**, *6*, 50–55. (In Chinese)
12. Shi, J.C.; Guo, Q.Z.; Zhao, S.; Su, Y.T.; Shi, Y.Q.; Du, G.; Zhang, L.P.; Zhang, D.; Zhao, Y. Long-term monitoring and analysis of driving factors of surface water changes in the Aral Sea basin. *J. Earth Environ.* **2021**, *5*, 540–548. (In Chinese)
13. Yuan, R.Q.; Qing, S. Study on the spatiotemporal characteristics of surface water in Inner Mongolia from 1990 to 2015. *J. Irrig. Drain.* **2021**, *2*, 136–143. (In Chinese)
14. Shun, Q.S.; Bao, Y.H.; Zhao, W.J. Spatiotemporal variation of surface water area and its influencing factors in Inner Mongolia from 2009 to 2018. *Bull. Soil Water Conserv.* **2021**, *3*, 312–319. (In Chinese)

15. Zhang, Y.; Du, J.; Guo, L.; Fang, S.; Zhang, J.; Sun, B.; Mao, J.; Sheng, Z.; Li, L. Long-term Detection and Spatiotemporal Variation Analysis of Open-Surface Water Bodies in the Yellow River Basin from 1986 to 2020. *Sci. Total Environ.* **2022**, *845*, 157152. [CrossRef]

16. Deng, Y.; Jiang, W.; Tang, Z.; Ling, Z.; Wu, Z. Long-Term Changes of Open-Surface Water Bodies in the Yangtze River Basin Based on the Google Earth Engine Cloud Platform. *Remote Sens.* **2019**, *11*, 2213. [CrossRef]

17. Woolway, R.I.; Kraemer, B.M.; Lenters, J.D.; Merchant, C.J.; O'Reilly, C.M.; Sharma, S. Global lake responses to climate change. *Nat. Rev. Earth Environ.* **2020**, *1*, 388–403. [CrossRef]

18. Lettenmaier, D.P.; Alsdorf, D.; Dozier, J.; Huffman, G.J.; Pan, M.; Wood, E.F. Inroads of remote sensing into hydrologic science. *Water Resour. Res.* **2015**, *51*, 7309–7342.

19. Kuenzer, C.; Knauer, K. Remote sensing of lake dynamics: A review. *Aquat. Sci.* **2013**, *75*, 595–619.

20. Wulder, M.A.; White, J.C.; Nelson, R.F.; Næsset, E.; Ørka, H.O.; Coops, N.C.; Hilker, T.; Bater, C.W.; Gobakken, T. Lidar sampling for large-area forest characterization: A review. *Remote Sens. Environ.* **2012**, *121*, 196–209. [CrossRef]

21. Li, J.; Sheng, Y. An automated scheme for glacial lake dynamics mapping using Landsat imagery and digital elevation models: A case study in the Himalayas. *Int. J. Remote Sens.* **2012**, *33*, 5194–5213. [CrossRef]

22. Song, C.; Huang, B.; Ke, L.; Richards, K.S. Remote sensing of alpine lake water environment changes on the Tibetan Plateau and surroundings: A review. *ISPRS J. Photogramm. Remote Sens.* **2012**, *92*, 26–37. [CrossRef]

23. Woodget, A.S.; Carbonneau, P.E.; Visser, F.; Maddock, I.P. Quantifying submerged fluvial topography using hyperspatial resolution UAS imagery and structure from motion photogrammetry. *Earth Surf. Process. Landf.* **2015**, *40*, 47–64. [CrossRef]

24. Moreira, A.; Prats-Iraola, P.; Younis, M.; Krieger, G.; Hajnsek, I.; Papathanassiou, K.P. A tutorial on synthetic aperture radar. *IEEE Geosci. Remote Sens. Mag.* **2013**, *1*, 6–43. [CrossRef]

25. Strozzi, T.; Wegmüller, U.; Tosi, L.; Bitelli, G.; Spreckels, V. Land subsidence monitoring with differential SAR interferometry. *Photogramm. Eng. Remote Sens.* **2013**, *79*, 901–916.

26. Zhang, L.; Wang, S.; Zhang, J.; Wang, L. Comprehensive assessment of water resources vulnerability in the Yellow River Basin under climate change scenarios. *J. Hydrol.* **2023**, *603*, 127456.

27. Li, M.; Shi, X.; Guo, M.; Liu, Y. Impacts of land-use change on the hydrological processes in the Yellow River Basin, China. *Water Resour. Manag.* **2023**, *37*, 795–812.

28. Cheng, D.; Li, X.; Zhao, T.; Han, M. Satellite-based monitoring of lake area changes and their responses to climate variability in the Yellow River Basin. *Remote Sens. Environ.* **2023**, *270*, 112789.

29. Bonnema, M.; David, C.H.; Frasson, R.P.d.M.; Oaida, C.; Yun, S.-H. The global surface area variations of lakes and reservoirs as seen from satellite remote sensing. *Geophys. Res. Lett.* **2022**, *49*, e2022GL098987. [CrossRef]

30. Turbay, I.; Ortiz, P.; Ortiz, R. Statistical analysis of principal components (PCA) in the study of the vulnerability of Heritage Churches. *Procedia Struct. Integr.* **2024**, *55*, 168–176. [CrossRef]

31. Fernández-Delgado, M.; Sirsat, M.S.; Cernadas, E.; Alawadi, S.; Barro, S.; Febrero-Bande, M. An extensive experimental survey of regression methods. *Neural Netw.* **2019**, *111*, 11–34. [CrossRef] [PubMed]

32. Zhang, G.; Yao, T.; Chen, W.; Zheng, G.; Shum, C.K.; Yang, K.; Piao, S.; Sheng, Y.; Yi, S.; Li, J.; et al. Regional differences of lake evolution across China during 1960s–2015 and its natural and anthropogenic causes. *Remote Sens. Environ.* **2019**, *221*, 386–404. [CrossRef]

33. Zhou, J.; Wang, L.; Zhong, X.; Yao, T.; Qi, J.; Wang, Y.; Xue, Y. Quantifying the major drivers for the expanding lakes in the interior Tibetan Plateau. *Sci. Bull.* **2022**, *67*, 474–478. [CrossRef] [PubMed]

34. Li, X.; Zhang, F.; Shi, J.; Chan, N.W.; Cai, Y.; Cheng, C. Analysis of surface water area dynamics and driving forces in the Bosten Lake basin based on GEE and SEM for the period 2000 to 2021. *Environ. Sci. Pollut. Res.* **2024**, *2024*, 1–14. [CrossRef] [PubMed]

35. Akbas, A. Human or climate? Differentiating the anthropogenic and climatic drivers of lake storage changes on a spatial perspective via remote sensing data. *Sci. Total Environ.* **2024**, *912*, 168982. [CrossRef] [PubMed]

36. Pi, X.; Luo, Q.; Feng, L.; Xu, Y.; Tang, J.; Liang, X. Mapping global lake dynamics reveals the emerging roles of small lakes. *Nat. Commun.* **2022**, *13*, 5777. [CrossRef] [PubMed]

37. Maleki, S.; Mohajeri, S.H.; Mehraein, M.; Sharafati, A. Lake evaporation in arid zones: Leveraging Landsat 8's water temperature retrieval and key meteorological drivers. *J. Environ. Manag.* **2024**, *355*, 120450. [CrossRef]

38. Dibike, Y.; Marshall, R.; de Rham, L. Climatic sensitivity of seasonal ice-cover, water temperature and biogeochemical cycling in Lake 239 of the Experimental Lakes Area (ELA), Ontario, Canada. *Ecol. Model.* **2024**, *489*, 110621. [CrossRef]

39. Reddy, V.M.; Ray, L.K. Concurrent and dynamical interdependency of compound precipitation and wind speed extremes over India. *Atmos. Res.* **2024**, *304*, 107389. [CrossRef]

40. Zhang, Z.; Cong, Z.; Gao, B.; Li, G.; Wang, X. The water level change and its attribution of the Qinghai Lake from 1960 to 2020. *J. Hydrol. Reg. Stud.* **2024**, *52*, 101688. [CrossRef]