

A 1 km soil organic carbon density dataset with depth of 20cm and 100cm from 1985 to 2020 in China

Yi Dong[&], Xinting Wang[&], Wei Su^{*}

College of Land Science and Technology, China Agricultural University, Beijing 100083, China

5 *Correspondence to:* Wei Su (suwei@cau.edu.cn)

[&]These authors contributed equally to this work and should be considered co-first authors.

Abstract: Soil organic carbon (SOC) is an important component of the worldwide carbon cycle as a vital indicator of soil quality and ecosystem health, with significant implications for agricultural production and climate change adaptation and mitigation strategies. For tracing the spatiotemporal changes of SOC content in China, this study is aiming at producing the accurate soil organic carbon density (SOCD) products from 1985 to 2020 with the spatial resolution of 1km with depths of 0-20 cm and 0-100 cm. The data sources used in this study include Landsat archives, topographic data, meteorological data, and measured SOCD data. The climate zoning was done for quantifying the climate differences with large area of China and the random forest ensemble learning approach was used for robust SOCD estimation with 8203 samples. Our estimated results show that the zoning model outperformed the global model without climate zoning in estimating SOCD with $R^2=0.63$ and RMSE=1.92 (kg C/m²) for 0-20 cm SOCD estimation and $R^2=0.60$ and RMSE=7.07 (kg C/m²) for 0-100 cm. Comparably, the SOCD estimation using the global model is with $R^2=0.49$ and RMSE=2.24 (kg C/m²) for 0-20 cm SOCD estimation and $R^2=0.48$ and RMSE=7.97 (kg C/m²) for 0-100 cm. Moreover, our SOCD estimated results with 0-20 cm depth are aligned well with independent samples ($R^2=0.71$, RMSE=1.94 kg C/m²) and Xu's dataset ($R^2=0.66$, RMSE=1.75 kg C/m²). Furthermore, the validation of our SOCD estimated results with 0-100 cm depth with independent measurements from Dong et al. (2024) showed strong agreement ($R^2=0.44$, RMSE=5.24 kg C/m²). The comparisons with the published SOC content products including HWSO, SoilGrids250m, and GSOCmap have also shown good consistency, too. Comparably, our predicted SOCD is the best fit with SoilGrids250m products with $R^2=0.72$ and RMSE=1.35 (kg C/m²). Comparisons of model predictions to independent datasets from the 1980s, 2000s, and 2010s in China reveal substantial

connections and demonstrate strong performance over time. The estimated SOCD is available via the Figshare
25 (<https://doi.org/10.6084/m9.figshare.27290310.v1>) (Dong et al., 2024).

1 Introduction

Soils are important because they enable the movement of carbon, energy, and water (Chaney et al., 2019; Crow et al., 2012). The foundation of soil fertility lies in soil carbon, a significant component of terrestrial carbon storage. SOC accounts for more than half of total soil carbon and is an essential component of the soil carbon cycle, which has a major impact on soil
30 fertility and agricultural productivity (Baldock, 2007; Chen et al., 2022). A combination of natural and human forces is placing significant strain on the global SOC reservoir. SOC content estimation has become a hot spot in global climate change due to its close relationship with climate change. The sustainability of agricultural production is threatened worldwide by soil degradation and the loss of intimate relationships. As a largely agricultural country, the distribution and changes of SOC content in China have significant impacts on the global carbon balance and play an important role in
35 balancing global carbon emissions and sequestration (Xu et al., 2018). Moreover, the complicated geography and varied climate in China have resulted in significant regional differences in SOC, leading to the difficulties of SOC content estimation in China.

There has been more and more interest in global, national, and regional SOC content estimation in the last few years (Padarian et al., 2022). In-depth studies to estimate subsurface SOC content estimation, particularly at a regional scale,
40 remain challenging due to the difficulty of data collection, the lack of long-term observations, and the depth dependency of soil carbon sequestration (Padarian et al., 2022). The advancement of digital soil mapping technology opens up new paths for estimating SOC content in large-scale and long-term series (Li et al., 2024). The use of machine learning techniques for digital soil modeling is a common concept in DSM. Compared to traditional mapping methods such as geo-statistics, expert knowledge, and individual representation, machine learning techniques provide a new paradigm for estimating SOC content
45 in large-scale and long-term series. To produce continental-scale SOC-weighted mean maps, Odgers et al. (2012) used an equal-area spline function for soil databases, while Mulder et al. (2016) used a machine learning model with a three-dimensional distribution to estimate SOC content in eastern France. These studies provide evidence for a comprehensive and

accurate understanding of soil properties and their spatial variation. Despite these advances, most digital soil mapping studies have focused on a specific period and the long-term dynamics of SOCD mapping have not yet been developed.

50 Emadi et al. (2020) predicted the SOCD in northern Iran using a sample of 1879 measurements, and Nabiollahi et al. (2019) used a random forest (RF) model to predict the SOCD at 137 sites in Marivan, Kurdistan Province, Iran. However, these studies only focus on local zones. In China, researchers have paid considerable attention to the sequestration potential of SOC storage, but most studies have focused on specific experimental areas or ecosystem types. Fang et al. (2007) estimated the carbon sink of terrestrial vegetation in China. Furthermore, these studies often lack attention to long-term trends and

55 dynamics, resulting in insufficient data sets to fully understand climate change and the impact of human activities on SOCD. At the national level, there is relatively little study on the potential for organic carbon storage across different ecosystem types (O'Rourke et al., 2015). The scarcity and unevenness of SOC data in China, as well as the lack of effective estimation methods, all contribute to the uncertainty of SOC prediction. In addition, the diverse and complex topography in China, as well as the lack of measured SOCD data, have increased the difficulty of SOC content estimation. Previous studies often

60 used the data from inventories of relevant resources to make rough calculations of carbon sinks (Pan et al., 2004). Unfortunately, the spatial continuity and variability of SOC, the spatial differentiation of organic carbon sequestration potential, and the influence of environmental factors have not been considered in previous studies. Especially in western China, there is almost no measured SOC data (Liu et al., 2022), which poses a challenge for understanding terrestrial ecosystems and soil carbon sinks in China. Given these challenges, it is urgent to carry out SOCD mapping and analyze the

65 temporal and spatial changes of SOCD in China.

To produce robust and accurate long-term SOCD products in China, we explore the RF models with climate zoning to predict SOCD in China from 1985 to 2020 and improve the study of SOCD maps for the 0-20 cm and 0-100 cm soil layers in China. The Landsat TM/ETM+/OLI images, topography, meteorology, and soil properties data are used for SOCD mapping in this study. The main contributions of this study can be summarized as follows.

70 (1) A nationwide, long-term soil organic carbon density dataset from 1985 to 2020 with depths of 20cm and 100cm in China is provided in this study.

(2) The machine learning RF models zoned by climate zones in China are developed for SOCD estimation, and the spatial-temporal variability of soil carbon is considered in our SOCD estimation.

(3) The proposed framework provides a comprehensive understanding of SOCD estimation including spectral indices of satellite remote sensing images, digital elevation model (DEM) and its topographic derivatives, meteorological features, and soil properties. The technique offers the potential for SOCD mapping with sufficiently measured SOC content data.

2 Study area and data sources

2.1 Study area

The study area, which extends throughout China, is characterized by complex and diverse terrains including mountains, plateaus, basins, plains, and deserts (Yuan et al., 2023). In addition, China has a large latitude difference from 4°N to 53°N and a large longitude difference from 73°E to 135°E. Therefore, there are obvious differences in precipitation and temperature in the study area, which bring significantly different accumulation processes and spatial patterns of soil carbon (Zheng et al., 2023). In addition, there are various soil types, including red soil, brown soil, black soil, and chestnut calcium soil, which have obvious spatial characteristics in the study area (Shangguan et al., 2014). For these reasons, we developed four different RF models for SOCD estimation for four temperature zones from south to north in China including humid area, semi-humid area, semi-arid area and arid area.

2.2 Data sources

(1) SOC content data

There were 8203 measured SOC content samples in the 1980s, 2000s, and 2010s in China collected for model building and validation of SOCD estimation. The SOC content and soil mass weight data of the 1980s were collected from the profile database of the Second National Soil Survey (1980-1996) (<http://www.geodata.cn>). The SOCD data of the 2000s was collected from the China Terrestrial Ecosystem Carbon Density Dataset (2000-2014) (<http://www.cnern.org.cn/>). The SOC content data of the 2010s was collected from the Soil Attribute Data of the China Soil System Record (2010s) (<https://www.resdc.cn/>), which was measured in the China Soil System Survey Collection and China Soil System Journal Compilation Project. To validate the SOCD estimation results in this study, three independent SOCD datasets were used,

including the measured SOC content data in the Heihe River basin (Song et al., 2016), the measured SOCD data from Xu et al. (2018), and the soil inorganic carbon (SIC) and SOC density dataset from Dong et al. (2024). The SOC content data of the Heihe River basin were collected from the spatio-temporal Tripolar Environmental Big Data Platform (<https://poles.tpdc.ac.cn/zh-hans/>). The measured SOCD data from Xu et al. (2018) focuses on SOC densities and soil carbon storage with a depth of 0–20 cm in various terrestrial ecosystems in China. The data was measured in field campaigns between 2004 and 2014, as well as some unpublished field measurements. The dataset from Dong et al. (2024) provides comprehensive measurements of SOC and inorganic carbon densities across 0-100 cm profiles in Chinese grassland and desert ecosystems, along with key environmental drivers such as climate variables, soil properties (texture, pH, conductivity), nitrogen deposition, and root biomass. This multi-source validation enhances the robustness of our SOCD assessments across different ecosystems and soil depths.

(2) Landsat archives

The time-series archived Landsat 4, 5, 7, and 8 TM/ETM+/OLI images spanning from 1985 to 2020 (Yu et al., 2023) are used for SOCD estimation, which are retrieved from the GEE cloud computing platform (Liu et al., 2024). Preprocessing of Landsat images, including radiometric calibration, atmospheric correction, geometric correction, cloud identification, and spectral index calculating are carried out on the GEE cloud computing platform. Random sampling and statistical regression analysis are performed to determine the calibration coefficients for each band spectral reflectance. Principal major axis regression models are used to normalize the reflectance data for different sensors. Radiometric correction coefficients of different Landsat sensors are calculated (Fig. 1). The spatially overlapping images are combined into one image using the aggregation function, and the combined image dataset is subjected to stitching operations to produce spatially coherent images. A variety of spectral indices were calculated using Landsat images after processing. Spectral indices Normalized Difference Vegetation Index (NDVI), Bare Soil Index (BSI), Enhanced Vegetation Index (EVI), Land Surface Water Index (LSWI), and Soil-Adjusted Vegetation Index (SAVI) were calculated using Landsat images. The formulae for these spectral indices are as follows:

$$NDVI = \frac{\rho_{NIR} - \rho_{Red}}{\rho_{NIR} + \rho_{Red}} \quad (1)$$

$$BSI = \frac{(\rho_{SWIR} + \rho_{Red}) - (\rho_{NIR} + \rho_{blue})}{(\rho_{SWIR} + \rho_{Red}) + (\rho_{NIR} + \rho_{blue})} \quad (2)$$

$$EVI = 2.5 \times \frac{\rho_{NIR} - \rho_{Red}}{\rho_{NIR} + 6 \times \rho_{Red} - 7.5 \times \rho_{blue} + 1} \quad (3)$$

$$SAVI = \frac{\rho_{NIR} - \rho_{Red}}{(\rho_{NIR} + \rho_{Red} + 0.5) \times 1.5} \quad (4)$$

$$LSWI = \frac{\rho_{NIR} - \rho_{SWIR1}}{\rho_{NIR} + \rho_{SWIR1}} \quad (5)$$

Where, ρ_{NIR} is the reference of near-infrared band, ρ_{Red} is the reference of red band, ρ_{blue} is the reference of blue band,

120 ρ_{SWIR} is the reference of short-wave infrared band and ρ_{SWIR1} is the reference of short-wave infrared band 1.

The land cover dataset newly released by Wuhan University (Yang and Huang, 2021) is used in this study. This is the first China Land Cover Annual Data Set (CLCD) derived from Landsat on the GEE platform.

(3) DEM and its topographic derivatives

Terrain is an important factor affecting the formation of soil organic matter. The DEM data is used for SOCD estimation, which is downloaded from the Resource and Environment Science Data Platform of the Chinese Academy of Sciences (<https://www.resdc.cn>) with a spatial resolution of 500 m. Topographic data and its topographic derivatives are extracted from the DEM data. There are four terrain derivatives, including Slope, Aspect, Elevation, and Topographic Wetness Index (TWI), which are calculated using SAGA GIS version 8.0.1 (<https://saga-gis.org/>) (Zhang et al., 2023). The spatial resolution of all raster data was uniformly adjusted to 1000m using resampling techniques to achieve spatial consistency between different datasets.

130

(4) Meteorological data

The meteorological features including Temperature (Tem), Precipitation (Pre) and Solar Radiation (SR), measured in 2,400 Chinese meteorological stations are used to quantify the effects of meteorological fluctuations. All meteorological data are downloaded from the China Meteorological Data Network (<http://data.cma.cn/>). For spatial consistency, the meteorological data is defined and projected into WGS 84 coordinates. All meteorological point data are interpolated into grid data with 1000 m spatial resolution using the ANUSPLIN program (Padarian et al., 2022).

135

(5) Published soil database

There are four published soil databases used to validate the SOCD estimation results in this study. One is the Harmonized World Soil Database (HWSD v2.0), produced by the International Institute for Applied Systems in Vienna and the Food and Agriculture Organization of the United Nations. There are two soil properties including soil bulk weight and organic carbon content are used for SOCD estimation at depths of 0-20 cm, 20-40 cm, 40-60 cm, 60-80 cm, and 80-100 cm. The SoilGrids250m v2.0 dataset including the soil silt content, sand content, clay content, and organic carbon content data with the spatial resolution of 250 m are downloaded from FAO SoilGrids (<https://soilgrids.org/>) for validation. For spatial consistency, this soil attribute datum is resampled to 1000 m. This soil product with five depth intervals (5, 15, 30, 60, and 100 cm) is used to calculate the soil silt content (Silt), sand content (Sand), clay content (Clay), and organic carbon at 0-20 cm and 0-100 cm (Zhang et al., 2023). Taking the clay content data as an example, the clay content with depths of 0-20 cm and 0-100 cm is calculated as follows:

$$CLY_{020} = \frac{CLY_{05}}{4} + \frac{CLY_{515}}{2} + \frac{CLY_{1530}}{4} \quad (6)$$

$$CLY_{0100} = \frac{CLY_{05}}{20} + \frac{CLY_{515}}{10} + \frac{3}{20} \times CLY_{1530} + \frac{3}{10} \times CLY_{3060} + \frac{2}{5} \times CLY_{60100} \quad (7)$$

Where, CLY_{05} , CLY_{515} , CLY_{1530} , CLY_{3060} , and CLY_{60100} are the clay content (g/kg) at depths of 0-5, 5-15, 15-30, 30-60, and 60-100 cm respectively.

The GSOCmap dataset (<https://www.fao.org/>), which is the first global SOC product led by FAO, is used for validation. GSOCmap is a 1-kilometer soil grid that covers depths ranging from 0 to 30 centimeters. The SOC Dynamics ML dataset in China is now available on the Dryad platform (<https://datadryad.org/>). Using machine learning, the dataset aims to capture the dynamics of SOC and its drivers in different soil horizons in China between the 1980s and 2010s (Li et al., 2022). The dataset contains valuable information such as SOC stocks, carbon fixation rates, and SOC content. While these existing datasets offer broad insights into SOC, our study specifically focuses on refining the estimation of SOCD for precise national-level carbon accounting across multiple historical periods. The organic carbon density with the depth of 20 cm and 100 cm in the 1980s, 2000s, and 2010s in China is used. This study focuses on SOCD, which is different from SOC content. The conversion from SOC content to SOCD is presented in Section 3.1.

3 Methodology

160 3.1 Converting SOC to SOCD with normalized soil depth

The dataset from the 2000s provided pre-calculated SOCD values (derived from SOC, bulk density, and coarse fragments by the original data source), while the 1980s and 2010s data reported SOC content. For consistency, we converted all SOC content data to SOCD using Equation 7. The 1980s SOC content data were from the Second National Soil Survey, and the 2010s SOC content data were from the Soil Attribute Data of China Soil System Record, which had several different soil
165 depths. For the consistency of the measured data, we convert the soil data with different depths into the SOCD with the depth of 0-20 cm and 0-100 cm using the package "mpspline2" v.0.1.3 (Bishop et al., 1999). The default value of 0.1 was used for the spline smoothing parameter lambda. The SOCD (kg C/m^2) is calculated using bulk density (kg/m^3), and coarse fractions percentage (%) provided by SoilGrids 2.0 (Poggio et al., 2021; Zhang et al., 2023).

$$\text{SOCD} = \frac{\text{SOC} \times \text{BD} \times \text{SD}}{100} \times \left(1 - \frac{\text{CF}}{100}\right) \quad (8)$$

Where, SOC is the soil organic carbon content (%), SOCD is the soil organic carbon density, BD is the soil bulk density, SD
170 is the soil depth (cm), and CF is the coarse fractions in a specific soil layer.

3.2 Feature selection for RF modelling

To achieve optimal prediction accuracy for SOCD and to elucidate its underlying mechanisms using RF models, comprehensive feature selection for the numerous potential environmental driving factors is a critical prerequisite (Jiang et al., 2024). This process is instrumental in mitigating model complexity, enhancing computational efficiency, improving
175 model interpretability, and eliminating data redundancy that could adversely affect model performance. In this study, the initial feature set comprised diverse categories of crucial environmental drivers, including remote sensing indices (e.g., NDVI, BSI, etc., derived from Landsat satellite imagery), topographic factors (e.g., elevation, slope, aspect, etc., generated from DEM), climatic factors (e.g., mean annual temperature, mean annual precipitation, etc.), as well as auxiliary soil attributes (e.g., soil type) and other relevant indicators (Fig. 2).

180 Our methodology commenced with a combined approach of correlation analysis, random forest importance ranking, and combinatorial optimization. First, a Pearson correlation matrix was constructed for the initial candidate features, and those

exhibiting high correlation (specifically, where the absolute value of the Pearson correlation coefficient exceeded 0.95) were removed to reduce redundancy. The remaining features, representing a refined set, then underwent an importance assessment and ranking utilizing the RF algorithm. A preliminary RF model was constructed with these features as inputs and SOCD as the target variable, and each feature's importance in predicting SOCD was quantified using Gini importance scores, thereby enabling the preliminary identification of core factors possessing substantial explanatory power for SOCD variation. This iterative procedure ensured the high independence of the selected feature set, preventing information overlap from impairing model performance and interpretability. Finally, to identify the optimal feature combination capable of maximizing model prediction accuracy, an exhaustive combinatorial search was conducted on the 10 most informative features remaining after the initial screening steps. Through a comprehensive evaluation of all possible feature subsets' performance, aiming to maximize the coefficient of determination (R^2), seven key environmental driving factors were ultimately identified as collectively providing the best predictive performance for SOCD: Temperature, Elevation, NDVI, Clay, SR, BSI, and Slope. This rigorous selection process ensures that the chosen feature set effectively characterizes SOCD dynamics while optimizing the model's predictive capability.

The selected features represent fundamental controls on SOCD through their influence on microbial activity (temperature), carbon input (vegetation indices), physical protection (clay content), and soil redistribution processes (slope). This multi-stage selection approach effectively balanced model complexity with predictive power while maintaining the ecological interpretability of the final feature set. The robustness of the selected features was further confirmed through cross-validation, demonstrating consistent performance across different validation datasets.

3.3 Climate zoning in China

Climate zoning is carried out to quantify the differences in temperature and precipitation in China and improve the accuracy of SOCD estimation. China spans a vast geographical area, crossing multiple major climate zones from the eastern coast to the western interior and from the subtropical monsoon climate in the southeast to the temperate continental climate in the northwest. This extensive climatic complexity leads to pronounced regional heterogeneity in soil formation and carbon cycling, which necessitates a zoned approach for accurate SOCD estimation. According to Tang et al. (2018), there are obvious differences in SOCD observed in different climate zones of China for the diverse and complex environmental

factors under warm-temperate climate conditions with a mean precipitation (MAP) threshold of 400 mm and a mean annual temperature (MAT) threshold of 10 °C. To mitigate the interannual variability, the multi-annual average temperature and precipitation are used to classify the climatic differences in China into four subzones including humid areas (MAP \geq 400 mm and MAT \geq 10°C), semi-humid area (MAP \geq 400 mm and MAT \leq 10°C), semi-arid area (MAP \leq 400 mm and MAT \leq 10°C) and arid area (MAP \leq 400 mm and MAT \geq 10°C) (Fig. 3). Soil data and environmental variables are grouped in each subzone, and zonal SOCD estimation models are developed for each subzone with depths of 0-20 cm and 0-100 cm (Fig. 3).

3.4 SOCD estimation using zoned RF models

For the estimation of SOCD, RF models were developed independently across four distinct climate subzones (arid, semi-arid, humid, and semi-humid) and for two soil depths (0-20 cm and 0-100 cm). Within each subzone, the RF model aggregates predictions from numerous decision trees, enhancing forecast stability and accuracy (Wu et al., 2021). This ensemble approach inherently mitigates overfitting, as individual trees are constructed from random subsets of data and features (Sun et al., 2024), thereby significantly improving the generalization of models. Especially, our RF model is conceptualized as a single, unified space-time model, meticulously trained on a comprehensive pooled dataset spanning distinct historical decades (1980s, 2000s, and 2010s). This unified framework, a key novelty of our approach, facilitates consistent SOCD prediction across multiple historical intervals (1985-2020 in five-year increments) for the vast and diverse Chinese region. The methodology effectively leverages the 'space-for-time' principle (Heuvelink et al., 2021) by integrating soil samples collected across these decades into a single training process. This enables the RF model to learn intricate relationships between environmental covariates and SOCD under varying historical conditions, inferring temporal SOCD evolution driven by dynamic factors based on observed spatial patterns.

The RF model inputs, established within the Scikit-Learn framework, comprised both static and dynamic predictors. Dynamic covariates, such as temperature, NDVI, SR, and BSI, were precisely matched to their corresponding five-year mapping periods by utilizing their average values for those intervals (e.g., 1985-1990). Model parameters, including the number of trees, the percentage of randomly selected features, and maximum tree depth, were tuned using a param_dist dictionary to optimize performance during cross-validation. The model's robustness and spatiotemporal capabilities are

underscored by a sophisticated stratified spatiotemporal K-fold cross-validation strategy. This involved spatially stratifying the study area into K independent sub-regions to address autocorrelation and assess generalization to new locations. Critically, temporal stratification ensured proportional representation of samples from all three decades within each spatial fold's training and validation sets, allowing the model to learn complex SOCD change patterns over time. The optimized RF model was subsequently employed to predict SOCD across the entire study area, utilizing measured SOCD values alongside spectral indices from soil properties, Landsat archives, topographic derivatives, and meteorological elements. Model performance and generalization ability were rigorously validated using the coefficient of determination (R^2) and root mean square error (RMSE). The trained model was saved using the joblib library, and the resulting estimations were combined with a geographic coordinate system to generate digital SOCD maps, facilitating the exploration of relationships between SOCD and optimized environmental variables.

4 Results and conclusions

4.1 Statistical analysis of sampling points

The statistics of the measured SOCD are shown in Fig. 4. The SOCD value at 0-20 cm depth shows a range of 0.070 to 22.93 kg C/m² in the 1980s with an average SOCD of 4.12 kg C/m², showing a positive offset and sharp distribution pattern. In the 2000s, average SOCD increased slightly to 4.30 kg C/m² and data variability increased with more measured samples. The number of samples reached its maximum in the 2010s, with the average density decreasing to 4.18 kg C/m², but the maximum reaching 26.58 kg C/m², suggesting that the skewness of the data distribution increased. For the SOCD value with a depth of 0-100 cm, the mean value in the 1980s was 13.23 kg C/m². In the 2000s, the average SOCD decreased significantly to 9.06 kg C/m² and the variability decreased. There is a significant increase in the maximum value of SOCD, with an average of 13.32 kg C/m² and a maximum of 132.92 kg C/m², with a more distorted data distribution and thicker tails (Fig. 4).

Figure 5 shows the geographical arrangement of SOCD data based on Whittaker biomes with depths of 20 cm and 100 cm in the 1980s, 2000s, and 2010s in China. The distribution of samples shows significant regional concentration and geographical variation, with most points concentrated in the northeastern plain, southwestern plateau, hilly zones, and southeastern coastal

zones. There are fewer SOCD samples in northwestern China due to difficult human accessibility, lower vegetation cover, less human activity, and a dry environment. In terms of timing, there are fewer SOCD sample sites in the 1980s. The number of sampling sites increased in the 2000s, particularly in agriculturally developed and densely populated areas.

4.2 Model performance of SOCD estimation

260 To evaluate the model performance of SOCD estimation at depths of 0-20 cm and 0-100 cm, two key indicators were utilized, including the coefficient of determination (R^2) and the Root Mean Square Error (RMSE). R^2 quantifies the proportion of variance in the dependent variable explained by the model, while RMSE assesses the discrepancy between model predictions and estimated results. The precision of RMSE values is further characterized by their 95% confidence intervals (CI), providing insight into the robustness and statistical significance of observed performance differences.

265 As illustrated in Figure 6, the implementation of a climatic zoning strategy significantly improved the model performance for both soil depths compared to global models. For the 0-20 cm SOCD prediction, the global model achieved an accuracy of $R^2=0.49$ and $RMSE = 2.24 \text{ kg C/m}^2$ (95% CI of [2.13, 2.35]). After incorporating climatic zoning, performance significantly improved to $R^2=0.63$ and $RMSE = 1.92 \text{ kg C/m}^2$ (95% CI of [1.85, 1.99]), demonstrating an R^2 increase of 0.14 and an RMSE decrease of 0.32 kg C/m^2 . The non-overlapping confidence intervals for the global and zoned models (95% CI of

270 [2.13, 2.35] vs. [1.85, 1.99]) clearly indicate a statistically significant improvement in RMSE due to climatic zoning. Similarly, for the 0-100 cm SOCD prediction, the global model yielded an $R^2=0.48$ and $RMSE = 7.97 \text{ kg C/m}^2$ (95% CI of [7.34, 8.67]). With climatic zoning, the performance enhanced to $R^2=0.60$ and $RMSE = 7.07 \text{ kg C/m}^2$ (95% CI of [6.49, 7.78]), reflecting an R^2 increase of 0.12 and an RMSE decrease of 0.90 kg C/m^2 . Here again, the distinct confidence intervals ([7.34, 8.67] vs. [6.49, 7.78]) confirm the statistical significance of performance enhancement from zoning. Overall, the

275 SOCD estimation model for the 0-20 cm depth generally exhibited a higher R^2 compared to the 0-100 cm depth model (e.g., peak $R^2=0.63$ for 0-20 cm versus $R^2=0.60$ for 0-100 cm in zoned models), indicating greater complexity in modeling deeper SOC dynamics with available covariates.

Further analysis of model performance within different climate zones revealed distinct patterns for each depth, as detailed in Figure 7. For the 0-20 cm depth, the humid zone showed the highest accuracy with $R^2=0.65$ and $RMSE = 1.77 \text{ kg C/m}^2$ (95% CI of [1.61, 1.93]). The semi-arid zone followed closely with $R^2=0.64$ and $RMSE = 1.77 \text{ kg C/m}^2$ (95% CI of [1.49, 1.63]).

280

The semi-humid zone achieved $R^2=0.63$ and $RMSE = 2.58 \text{ kg C/m}^2$ (95% CI of [2.39, 2.76]). The arid zone showed an $R^2=0.58$ and $RMSE = 1.61 \text{ kg C/m}^2$ (95% CI of [1.17, 2.07]). The higher estimation accuracy observed in humid and semi-arid zones, indicated by generally higher R^2 values and often lower RMSEs with narrower CIs (e.g., arid zone RMSE CI [1.17, 2.07] is distinct from semi-humid RMSE CI [2.39, 2.76]), suggests superior model fit in these regions. This can be attributed to a more even distribution of SOCD measurements in these regions, which reduces the influence of extreme values and facilitates more accurate estimations. Furthermore, environmental factors in these zones are often more consistent with the measured SOCD data used for model training, contributing to improved accuracy.

For the 0-100 cm depth, the humid zone again demonstrated the highest $R^2=0.62$ with $RMSE = 5.44 \text{ kg C/m}^2$ (95% CI of [4.38, 5.63]). The semi-humid zone had an $R^2=0.61$ and $RMSE = 7.70 \text{ kg C/m}^2$ (95% CI of [6.87, 8.52]). The semi-arid zone showed an $R^2=0.60$ and $RMSE = 4.63 \text{ kg C/m}^2$ (95% CI of [4.14, 5.18]). The arid zone had an $R^2=0.54$ and $RMSE = 3.17 \text{ kg C/m}^2$ (95% CI of [2.56, 3.81]). At this deeper layer, while the humid zone still shows the highest R^2 , the arid zone exhibits the lowest RMSE with a relatively narrow confidence interval ([2.56, 3.81]), indicating good precision despite a lower R^2 . This pattern suggests that while overall explained variance might be moderate, the model's predictive error in arid regions for deeper layers is tightly constrained. The explanation might lie in the unique soil moisture conditions of drylands, where more extreme and significant changes lead to sensitive and distinct responses of vegetation and soil microorganisms to water. This high sensitivity can provide a clearer prediction signal, potentially improving accuracy in certain scenarios or for shallower layers. Conversely, the influence of multiple complex factors (e.g., precipitation, temperature, vegetation cover) on soil moisture in wet areas can reduce the predictive power of the model for deeper layers where long-term processes dominate. Our SOCD estimation models successfully capture this feature of SOC accumulation.

Importance analysis of optimized features for SOCD estimation was performed to better understand the contribution of various environmental variables to SOCD estimation. Figure 8 illustrates the hierarchical importance of these features for both depths. For the 0-20 cm depth (Figure 8a), Temperature (Tem) emerged as the most influential variable, contributing 34.41% to the model. It was followed by NDVI (20.3%), SR (12.42%), Elevation (6.92%), Clay (5.2%), BSI (4.19%), and Slope (1.69%). The dominant influence of temperature on SOCD is primarily through its effects on soil microbial activity and respiration. An increase in temperature can accelerate the decomposition of soil organic carbon, but it may also increase

the rate of plant residue decomposition, thereby augmenting carbon return to the soil. The high importance of NDVI underscores the critical role of vegetation health and productivity in contributing to SOCD. Vegetation cover influences the amount of organic matter returned to the soil via litterfall, and root activity affects soil moisture and nutrient cycling. Precipitation also directly affects soil water status, where suitable soil water content is conducive to SOC accumulation, while low moisture content can accelerate SOC decomposition, thus reducing SOCD. For the 0-100 cm depth (Figure 8b), NDVI was the most important feature, contributing 28.3%, followed by Temperature (20.2%), Elevation (17.0%), SR (9.5%), BSI (7.2%), Clay (6.2%), and Slope (0.9%). Elevation is associated with the vertical distribution of surface hydrothermal conditions, which affects the soil formation process and organic carbon distribution properties.

The differences in feature importance between the two depth models highlight the complexity of soil organic carbon dynamics and the necessity of considering depth-specific processes in SOCD estimation models. In the 0-20 cm model, climatic conditions (Temperature) and vegetation cover (NDVI) play a more dominant role, reflecting the direct influence of these factors on topsoil organic carbon accumulation. In contrast, in the 0-100 cm model, soil physical-chemical properties (e. g., Clay) and topographic features (e. g., Elevation, Slope) become more important, indicating that deeper soil layers are influenced by long-term geological processes and soil erosion/deposition dynamics. Understanding these variations allows for a better capture of the spatial and temporal variability of soil organic carbon across different soil layers. All these factors, including climatic conditions, topographic features, vegetation coverage, and soil physicochemical properties, are crucial determinants of SOCD estimation.

4.3 Validation with independent sample points

The SOCD estimation result is validated with independent published SOCD data in the Heihe River basin by Li et al. (2022). The Heihe River basin is a major ecological and agricultural zone in northwest China. There are special geographical and climatic characteristics for the soil carbon accumulation in the Heihe River basin, which are important for exploring soil quality in arid and semi-arid zones. Validation is carried out by comparing measured data in the Heihe River basin with the estimated SOCD in this study. The comparison results show that our estimated SOCD is highly consistent with the measured SOC data from the Heihe River basin (Fig. 9). The estimated SOCD and the measured SOCD have a significant correlation, which is shown by the R^2 value of 0.71, and the RMSE value of 1.94 (kg C/m²) (95% CI of [1.59, 2.27]) for the estimated

result with the depth of 0-20 cm. Additionally, the proposed model demonstrates superior accuracy compared to Li's dataset, which reported an R^2 of 0.60 and an RMSE of 2.27 (kg C/m²) (95% CI of [1.59, 2.98]).

In order to validate with more SOCD samples data with wider ranging area, our estimated SOCD results are compared with the data published data by Xu et al. (2018), which is the data on carbon storage of terrestrial ecosystems in China with a depth of 0-20 cm. These samples are widely distributed across the southern Tibet Autonomous Region, Qinghai Province, and eastern Inner Mongolia Autonomous Region. This is very good evidence for validating the robustness, reliability, and generalizability of the SOCD estimation model in this study. The estimated SOCD results are compared with the measured SOC data in the field campaign of Xu et al. (2018). In addition, the field data were compared with 0-20 cm organic carbon density maps generated by a machine learning analysis dataset of SOC dynamics and their drivers in China during 2000-2014. The results of the comparative analysis are encouraging and show high agreement between the estimated SOCD using our developed model and the measured SOC data. Specifically, the R^2 value is 0.66 and the RMSE value is 1.75 (kg C/m²) [95%CI: 1.49, 2.01], which further confirms the accuracy of our SOCD estimation model (Fig. 9). Furthermore, the model outperforms Li's dataset, which yielded an R^2 of 0.54 and an RMSE of 2.04 (kg C/m²) [95%CI: 1.73, 2.34], underscoring the enhanced predictive accuracy of our approach.

Our estimated SOCD results with the depth of 0-100 cm were validated furtherly with independent measurements from Inorganic carbon pools and their drivers in grassland and desert soils (Dong et al., 2024) and compared with the machine learning-derived SOCD simulations by Li et al. (2022). The SOCD dataset of Dong et al. (2024), covering grassland and desert ecosystems across China, provides robust in-situ measurements of SOCD (0-100 cm) alongside critical environmental drivers (e.g., climate, soil properties), offering an ideal benchmark for evaluating model generalizability in arid and semi-arid regions. Comparative analysis revealed that our SOCD estimates achieved significantly better agreement with the independent validation data ($R^2 = 0.44$, RMSE = 5.24 kg C/m²) [95%CI: 4.42, 6.13] than Li et al.'s simulations ($R^2 = 0.31$, RMSE = 5.80 kg C/m²) [95%CI: 4.16, 7.48] (Fig. 9). This demonstrates the superior accuracy of our approach in capturing deep soil carbon dynamics (0-100 cm), particularly in heterogeneous grassland and desert environments. The higher R^2 and lower RMSE values underscore the improved capability of our model to resolve spatial patterns of SOC storage compared to earlier machine learning-based efforts.

4.4 Comparison with published SOCD products

The 1-km-resolution SOCD dataset of China is created in this study, which is compared with the published SOCD products including HWSD v2.0, SoilGrids 250m, and GSOCmap datasets to validate and confirm its accuracy and reliability. The comparison results shown in Fig. 10 show that our produced 1-km-resolution SOCD dataset is largely consistent with published SOCD products, with the highest fit to SoilGrids250m and an R^2 of 0.72, significantly better than that of 0.61 with GSOCmap dataset and that of 0.54 with HWSD dataset. The HWSD v2.0 dataset is jointly published by the Food and Agriculture Organization of the United Nations (FAO) and the International Institute for Applied Systems (IIAS) in Vienna, which provides soil data on a global scale. Unfortunately, its applicability and accuracy are limited in China. The correlation of our SOCD dataset with HWSD is reported with the R^2 value of 0.54 and the RMSE value of 1.42 (kg C/m²). The GSOCmap dataset is led by the FAO and is intended to cover various ecosystems around the world. This is the first global SOC map. The correlation of our SOCD dataset with GSOCmap is reported with the R^2 of 0.61 and the RMSE of 1.32 (kg C/m²). The SoilGrids250m dataset is created using ISRIC's digital soil mapping technology, which is a global soil dataset. The correlation of our SOCD dataset with SoilGrids250m is reported with the R^2 of 0.72 and the RMSE of 1.07 (kg C/m²). Models are more accurate and applicable than global soil databases in capturing SOCD changes in China. This study highlights the need to create and implement region-specific models that utilize current geographic and environmental data to provide a more precise tool for accurately estimating soil carbon reserves.

For time series estimation accuracy, the estimated SOCDs in China are compared with the SOC Dynamics ML dataset in China in the 1980s, 2000s, and 2010s (Fig. 11). The comparison results show that there are significant correlations between estimated SOCDs and measured data with RMSE of 1.04 (kg C/m²), 1.10 (kg C/m²) and 1.09 (kg C/m²) and R^2 of 0.72, 0.74 and 0.73 in the 1980s, 2000s and 2010, respectively. The performance improvement in the later period is mainly due to the increased sample points. With more sample data available, the model has captured the spatial heterogeneity of SOCD more accurately. These comparisons confirm the robustness of the SOCD estimation model in this study and its potential to provide accurate estimates of SOCD. The improvement over time highlights the importance of integrating current data and advanced analytical methods into soil carbon studies.

380 4.5 Spatiotemporal changing of SOCD in China

The SOCD changes over time from the 1980s to the 2010s are validated in Fig. 12 compared with the published investigations. Fig. 12 reveals that our estimated SOCD results with depths of 0-20 cm (a) and 0-100 cm (b) are falling in the value range of the previous investigations of Ni (2001), Wu et al. (2003), Wang et al. (2004), Xu et al. (2018), Wang et al. (2021), Li et al. (2022), Zhang et al. (2023). We can find that SOCD in China has slightly upward increasing from the 1980s to the 2010s in the 0-20 cm topsoil (Fig. 12a). This resulted from that the topsoil is more susceptible to the direct effects of soil management practices and environmental changes (Oechaiyaphum et al., 2020). However, the estimated SOCD shows an increasing from 1980s to 1990s and keeps stable from 1990s to 2020s in the 0-100 cm deep soil (Fig. 12b). Fig. 13 and Fig. 14 show the spatiotemporal distributions of the estimated SOCD at the 5-year interval from the 1980s to the 2010s. And the regions with high SOCD value in depth of 0-20 cm are in northeast and southwest China with red color in Fig. 13. Comparably speaking, there are the largest area of high SOCD value labeled dark red color bar in period of 2010-2015 (Fig. 13f). From the perspective of longitude, the SOCD distribution shows different pattern, and it is homogeneous in high and low longitudes where the land cover is forest mostly. Conversely, the variance of SOCD is higher in mid-longitude regions where is with distinct land cover types. Similarly, the regions with high SOCD value in depth of 0-100 cm are in northeast and southwest China with green color in Fig. 14. And there are smaller variance of SOCD in high and low longitudes, and there are higher variance of SOCD in mid-longitude regions. There are many driving factors for the changing of SOCD in China. Targeted monitoring and management practices should be implemented for SOCD trends at different soil depths to maximize soil carbon sequestration and continuously improve soil quality.

5 Data availability

The 1 km soil organic carbon density dataset with depths of 20cm and 100cm from 1985 to 2020 in China is currently freely available at <https://doi.org/10.6084/m9.figshare.27290310.v1> (Dong et al., 2024). The data can be imported into remote sensing processing software (e.g., ENVI), standard geographical information system software (e.g., ArcGIS).

6 Conclusions

In this study, a SOCD dataset with a resolution of 1-kilometer resolution and soil depths of 0-20 cm and 0-100 cm is created from 1985 to 2020 in China. The accuracy and validity of this dataset are validated by three independent metrics and data
405 and four types of published global products. The conclusions are as follows.

- (1) The delineation of climatic zones for SOCD estimation modeling has been proven useful for enhancing the precision of the models and effectively addressing the uneven distribution of measured SOC.
- (2) Validation against independent datasets confirmed the robust accuracy of the estimated SOCD. For the 0-20 cm depth, our estimates showed strong agreement with measured data from the Heihe River basin ($R^2=0.71$, $RMSE=1.94$ kg C/m²) and
410 the Xu dataset ($R^2=0.66$, $RMSE=1.75$ kg C/m²). Furthermore, for the 0-100 cm depth, validation against independent measurements from Dong et al. (2024) also indicated strong agreement ($R^2=0.44$, $RMSE=5.24$ kg C/m²). Across all validations with the published datasets of Li et al. (2022), Xu et al. (2018) and Dong et al. (2024), our estimated results showed consistently performance with the published SOC product of Li et al. ($R^2=0.60$ and $RMSE=2.27$ kg C/m² by, $R^2=0.54$), and Xu et al. ($RMSE=2.04$ kg C/m² by, $R^2=0.31$) and Dong et al. ($RMSE=5.80$ kg C/m²).
- 415 (3) Compared to published global products including HWSD, SoilGrids250m, and GSOCmap, the estimated SOCD in this study was consistent and accurate. Comparison with the SoilGrids250m dataset shows the superiority of zoning RF models in capturing variations in SOCD in China with $R^2=0.72$ and $RMSE=1.07$ (kg C/m²).
- (4) The model demonstrated excellent correlation with time series datasets and increased accuracy over time by comparing with the independently measured data from the 1980s, 2000s, and 2010s. This highlights the importance of integrating
420 current data and advanced analytical methods into soil carbon studies. In addition, time series analyses showed the change of SOCD in China over time and at different soil depths, which can be influenced by many reasons such as agricultural management practices, land-use changes, and climate change.

Despite the impressive results of this study, more soil data are required to validate and improve the SOC estimation model. Future studies will focus on the effects of different land management strategies on SOC change as well as the development
425 of more refined models for estimating soil organic carbon. Furthermore, given the uncertainties in existing global SOC

estimates, we urge that future research focus on standardized soil sampling, cross-dataset comparisons, more validation, and global collaboration to improve the accuracy of SOC estimates.

Author contributions

Yi Dong and Xinting Wang collected the SOC data, did the field campaign and SOCD estimation, performed the analysis,
430 and wrote the paper; Wei Su designed the research and revised the manuscript.

Competing interests

The authors declare that they have no conflict of interest.

Acknowledgments

This study was supported by the National Natural Science Foundation of China under the project (No. 42471402), the
435 Intergovernmental International Scientific and Technological Innovation Cooperation Project of National Key R&D Program
(No. 2025YFE0102000), and Beijing Natural Science Foundation (L251053).

References

- Baldock, J. A.: Composition and cycling of organic carbon in soil, in: Nutrient cycling in terrestrial ecosystems, edited, Springer, 1-35, 2007.
- 440 Bishop, T., Mcbratney, A. B., and Laslett, G. M.: Modelling soil attribute depth functions with equal-area quadratic smoothing splines, *Geoderma*, 91, 27-45, [https://doi.org/10.1016/S0016-7061\(99\)00003-8](https://doi.org/10.1016/S0016-7061(99)00003-8), 1999.
- Chaney, N. W., Minasny, B., Herman, J. D., Nauman, T. W., Brungard, C. W., Morgan, C. L., Mcbratney, A. B., Wood, E. F., and Yimam, Y.: Polaris soil properties: 30-m probabilistic maps of soil properties over the contiguous united states, *Water Resour. Res.*, 55, 2916-2938, 2019.
- 445 Chen, S., Arrouays, D., Mulder, V. L., Poggio, L., Minasny, B., Roudier, P., Libohova, Z., Lagacherie, P., Shi, Z., and Hannam, J.: Digital mapping of globalsoilmap soil properties at a broad scale: a review, *Geoderma*, 409, 115567, 2022.
- Crow, W. T., Berg, A. A., Cosh, M. H., Loew, A., Mohanty, B. P., Panciera, R., de Rosnay, P., Ryu, D., and Walker, J. P.: Upscaling sparse ground-based soil moisture observations for the validation of coarse-resolution satellite soil moisture products, *Rev. Geophys.*, 50, 2012.
- 450 Dong, L., Ran, J., Luo, J., Bai, L., Sun, Y., Aqeel, M., Zhang, Y., Wang, X., Du, Q., and Xiong, J.: Inorganic carbon pools and their drivers in grassland and desert soils, *Glob. Change Biol.*, 30, e17536, 2024.
- Dong, Y., Wang, X., and Su, W.: A 1 km soil organic carbon density dataset with depth of 20cm and 100cm from 1985 to 2020 in china, edited, <https://doi.org/10.6084/m9.figshare.27290310.v1>, 2024.
- Emadi, M., Taghizadeh-Mehrjardi, R., Cherati, A., Danesh, M., Mosavi, A., and Scholten, T.: Predicting and mapping of soil
455 organic carbon using machine learning algorithms in northern iran, *Remote Sens.*, 12, <https://doi.org/10.3390/rs12142234>, 2020.
- Fang, J., Guo, Z., Piao, S., and Chen, A.: Terrestrial vegetation carbon sinks in china, 1981–2000, *Science in China Series D: Earth Sciences*, 50, 1341-1350, 2007.
- Heuvelink, G. B., Angelini, M. E., Poggio, L., Bai, Z., Batjes, N. H., van den Bosch, R., Bossio, D., Estella, S., Lehmann, J.,
460 and Olmedo, G. F.: Machine learning in space and time for modelling soil organic carbon change, *Eur. J. Soil Sci.*, 72, 1607-1623, 2021.
- Jiang, L., Fu, H., Zhang, Z., Zhang, H., Zhang, X., Feng, X., Xu, X., Mao, M., and Xie, J.: Synchronously enhancing the strength, toughness, and stress corrosion resistance of high-end aluminum alloys via interpretable machine learning, *Acta Mater.*, 270, 119873, 2024.
- 465 Li, H., Wu, Y., Liu, S., Xiao, J., Zhao, W., Chen, J., Alexandrov, G., and Cao, Y.: Decipher soil organic carbon dynamics and driving forces across china using machine learning, *Glob. Change Biol.*, 28, 3394-3410, 2022.
- Li, T., Cui, L., Wu, Y., McLaren, T. I., Xia, A., Pandey, R., Liu, H., Wang, W., Xu, Z., and Song, X.: Soil organic carbon estimation via remote sensing and machine learning techniques: global topic modeling and research trend exploration, *Remote Sens.*, 16, 3168, 2024.
- 470 Liu, F., Wu, H., Zhao, Y., Li, D., Yang, J., Song, X., Shi, Z., Zhu, A., and Zhang, G.: Mapping high resolution national soil

- information grids of china, *Sci. Bull.*, 67, 328-340, 2022.
- Liu, Z., Chen, G., Wen, Q., Yi, L., and Zhao, J.: Extraction of rocky desertification information based on multi-feature combination optimization and random forest algorithm: a case study of zhaotong city in yunnan province, *Science of Soil and Water Conservation*, 22, 95-105, 2024.
- 475 Mulder, V. L., Lacoste, M., Richer-De-Forges, A. C., Martin, M. P., and Arrouays, D.: National versus global modelling the 3d distribution of soil organic carbon in mainland france, *Geoderma*, 263, 16-34, 2016.
- Nabiollahi, K., Eskandari, S., Taghizadeh-Mehrjardi, R., Kerry, R., and Triantafilis, J.: Assessing soil organic carbon stocks under land-use change scenarios using random forest models, *Carbon Manag.*, 10, 63-77, 2019.
- Ni, J.: Carbon storage in terrestrial ecosystems of china: estimates at different spatial resolutions and their responses to
480 climate change, *Clim. Change*, 49, 339-358, 2001.
- Odgers, N. P., Libohova, Z., and Thompson, J. A.: Equal-area spline functions applied to a legacy soil database to create weighted-means maps of soil organic carbon at a continental scale, *Geoderma*, 189, 153-163,
<https://doi.org/10.1016/j.geoderma.2012.05.026>, 2012.
- Oechaiyaphum, K., Ullah, H., Shrestha, R. P., and Datta, A.: Impact of long-term agricultural management practices on soil
485 organic carbon and soil fertility of paddy fields in northeastern thailand, *Geoderma Reg.*, 22, e00307, 2020.
- O'Rourke, S. M., Angers, D. A., Holden, N. M., and Mcbratney, A. B.: Soil organic carbon across scales, *Glob. Change Biol.*, 21, 3561-3574, 2015.
- Padarian, J., Minasny, B., Mcbratney, A., and Smith, P.: Soil carbon sequestration potential in global croplands, *Peerj*, 10, e13740, 2022.
- 490 Padarian, J., Stockmann, U., Minasny, B., and Mcbratney, A. B.: Monitoring changes in global soil organic carbon stocks from space, *Remote Sens. Environ.*, 281, 113260, 2022.
- Pan, Y., Luo, T., Birdsey, R., Hom, J., and Melillo, J.: New estimates of carbon storage and sequestration in china's forests: effects of age-class and method on inventory-based carbon estimation, *Clim. Change*, 67, 211-236, 2004.
- Poggio, L., de Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., and Rossiter, D.: Soilgrids 2.0:
495 producing soil information for the globe with quantified spatial uncertainty, *Soil*, 7, 217-240, <https://doi.org/10.5194/soil-7-217-2021>, 2021.
- Shangguan, W., Gong, P., Liang, L., Dai, Y., and Zhang, K.: Soil diversity as affected by land use in china: consequences for soil protection, *Scientific World Journal*, <https://doi.org/10.1155/2014/913852>, 2014.
- Song, X., Brus, D. J., Liu, F., Li, D., Zhao, Y., Yang, J., and Zhang, G.: Mapping soil organic carbon content by
500 geographically weighted regression: a case study in the heihe river basin, china, *Geoderma*, 261, 11-22, 2016.
- Sun, Z., Wang, G., Li, P., Wang, H., Zhang, M., and Liang, X.: An improved random forest based on the classification accuracy and correlation measurement of decision trees, *Expert Syst. Appl.*, 237, 121549, 2024.
- Tang, X., Zhao, X., Bai, Y., Tang, Z., Wang, W., Zhao, Y., Wan, H., Xie, Z., Shi, X., and Wu, B.: Carbon pools in china's terrestrial ecosystems: new estimates based on an intensive field survey, *Proceedings of the National Academy of Sciences*,

- 505 115, 4021-4026, 2018.
- Wang, S., Huang, M., Shao, X., Mickler, R. A., Li, K., and Ji, J.: Vertical distribution of soil organic carbon in china, *Environ. Manage.*, 33, S200-S209, 2004.
- Wang, S., Xu, L., Zhuang, Q., and He, N.: Investigating the spatio-temporal variability of soil organic carbon stocks in different ecosystems of china, *Sci. Total Environ.*, 758, 143644, 2021.
- 510 Wu, H., Guo, Z., and Peng, C.: Land use induced changes of organic carbon storage in soils of china, *Glob. Change Biol.*, 9, 305-315, 2003.
- Wu, J. Y., Lin, Y., Lin, K., Hu, Y. H., and Kong, G. L.: [Predicting prolonged length of intensive care unit stay via machine learning], *Beijing Da Xue Xue Bao. Yi Xue Ban = Journal of Peking University. Health Sciences*, 53, 1163-1170, 2021.
- Xu, L., Yu, G., He, N., Wang, Q., Gao, Y., Wen, D., Li, S., Niu, S., and Ge, J.: Carbon storage in china's terrestrial
- 515 ecosystems: a synthesis, *Sci Rep*, 8, 2806, 2018.
- Xu, L., Yu, G., and He, N.: Changes of soil organic carbon storage in chinese terrestrial ecosystems from the 1980s to the 2010s, *Acta Geographica Sinica*, 73, 2150-2167, 2018.
- Yang, J. and Huang, X.: 30 m annual land cover and its dynamics in china from 1990 to 2019, *Earth System Science Data Discussions*, 2021, 1-29, 2021.
- 520 Yu, Z., Ma, R., Xu, J., Wang, Z., and Hu, M.: A dataset of the tsi of hulun lake in summer, 1986-2020, *China Scientific Data*, 8, 21-25, 2023.
- Yuan, J., Chen, W., and Zeng, J.: Spatio-temporal differentiation of cropland use change and its impact on cropland npp in china, *Journal of Natural Resources*, 38, 3135-3149, 2023.
- Zhang, Z., Ding, J., Zhu, C., Wang, J., Ge, X., Li, X., Han, L., Chen, X., and Wang, J.: Historical and future variation of soil
- 525 organic carbon in china, *Geoderma*, 436, <https://doi.org/10.1016/j.geoderma.2023.116557>, 2023.
- Zhang, Z., Ding, J., Zhu, C., Wang, J., Li, X., Ge, X., Han, L., Chen, X., and Wang, J.: Exploring the inter-decadal variability of soil organic carbon in china, *Catena*, 230, <https://doi.org/10.1016/j.catena.2023.107242>, 2023.
- Zheng, Z., Zhang, K., Shi, J., and Zhang, M.: Analysis of gnss water vapor detection accuracy and temporal sequence characteristics in different climate types in china, *Science of Surveying and Mapping*, 48, 68-77, 2023.

530

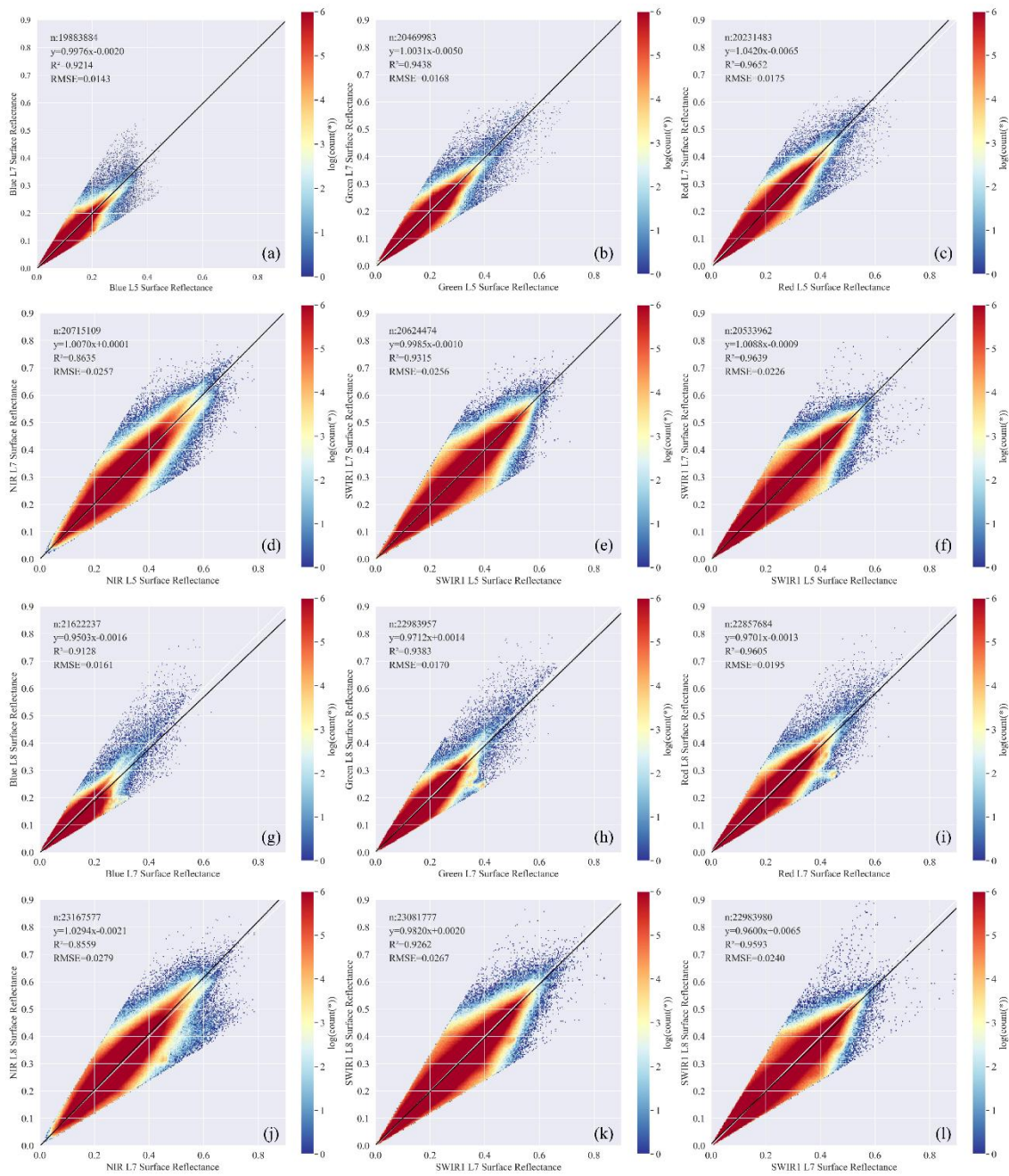


Figure 1. Radiometric normalization coefficients between Landsat 5 TM、Landsat 7 ETM+ (a-f) and Landsat 7 ETM+、Landsat 8 OLI (e-j) sensors for different bands including blue, green, red, NIR, SWIR1, and SWIR2. The radiometric normalization coefficients for each sensor are represented by fitted lines and correlation coefficients, indicating the correlation between the reference of different sensors, and characterizing the spectral response of the sensors in the different wavelength bands.

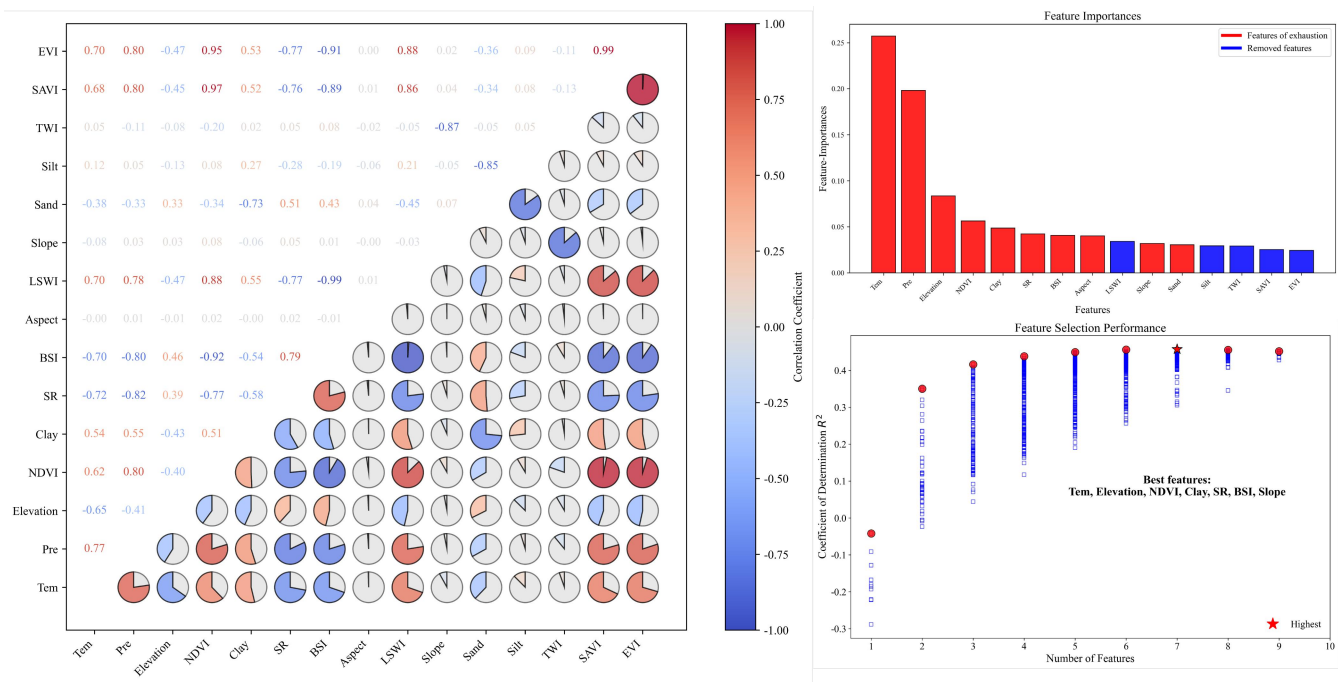
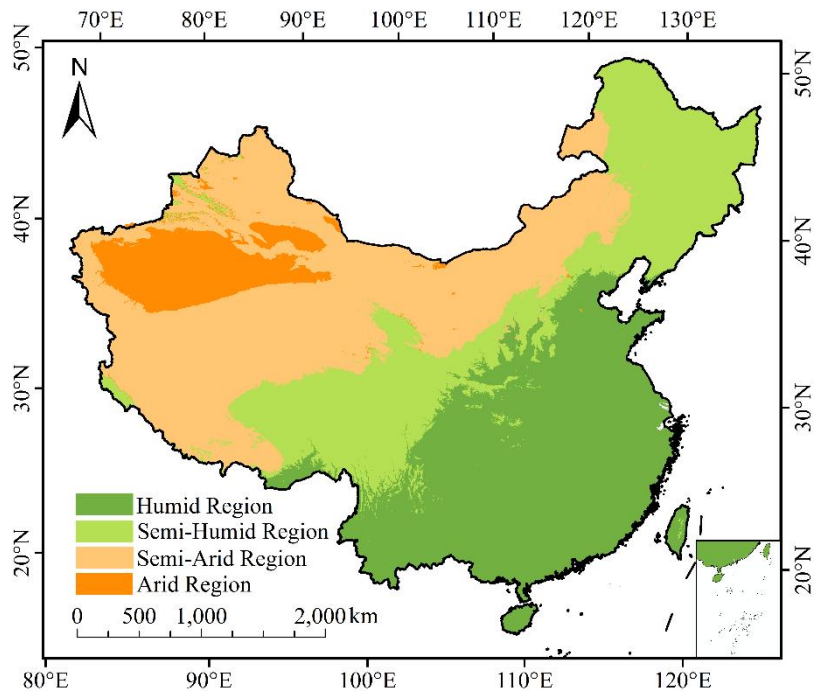


Figure 2. Feature selection process for predicting soil organic carbon density (SOCD). (a) Pearson correlation matrix of top environmental covariates (upper triangle shows correlation coefficients; red=positive, blue=negative), with boxed features indicating the final selected variables. (b) Hierarchical feature importance evaluation combining correlation filtering (removing $|r| > 0.95$), random forest-based ranking (Gini importance), and combinatorial optimization. The optimal feature set (highlighted in bold) comprised seven variables: mean annual temperature (Tem), elevation, NDVI, clay content (Clay), simple ratio index (SR), bare soil index (BSI), and slope, which collectively maximize prediction accuracy (R^2) while maintaining ecological interpretability.



545

Figure 3. Climatic zones for SOCD estimation modeling. Climate zoning comes from the time-series climate data including temperature and precipitation. According to the difference in climate zones, it can be divided into humid, semi-humid, arid, and semi-arid zones.

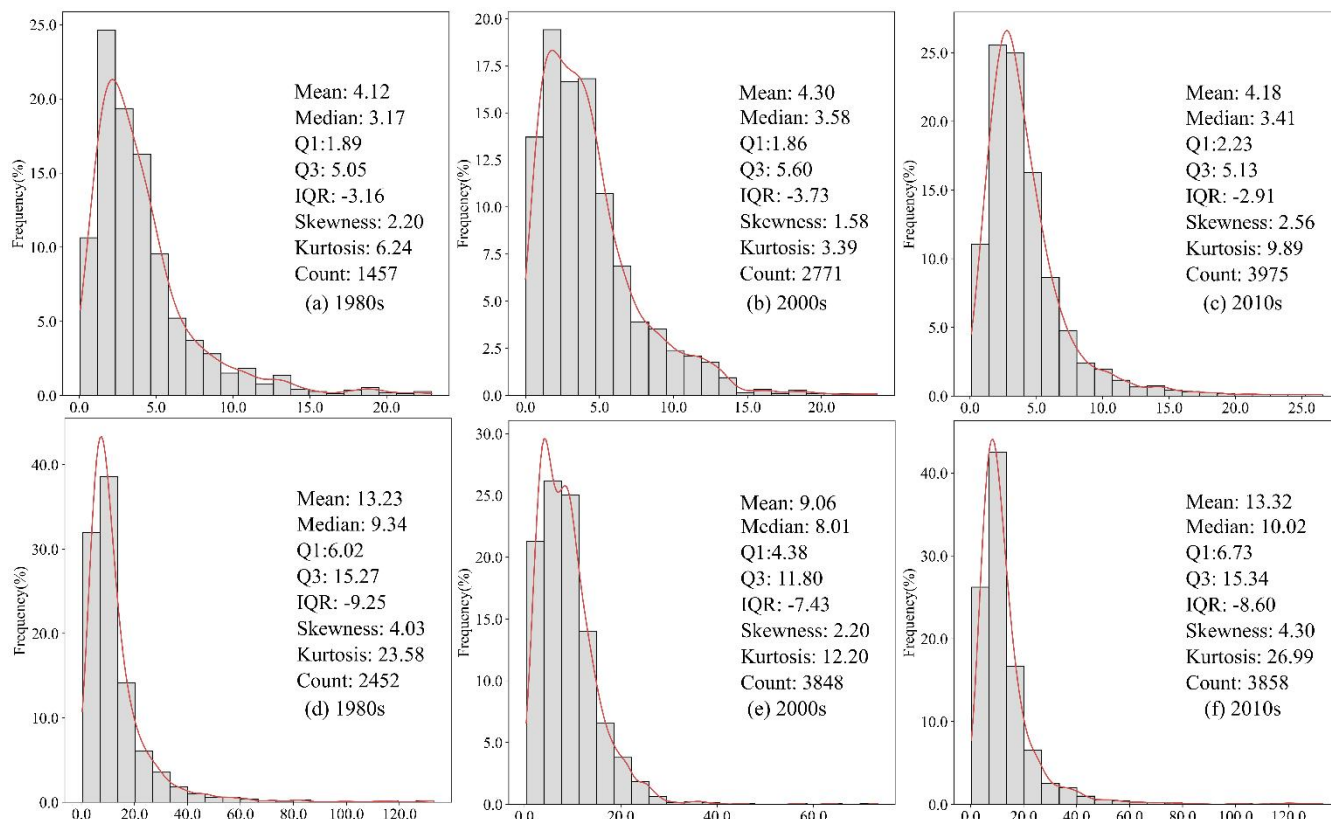
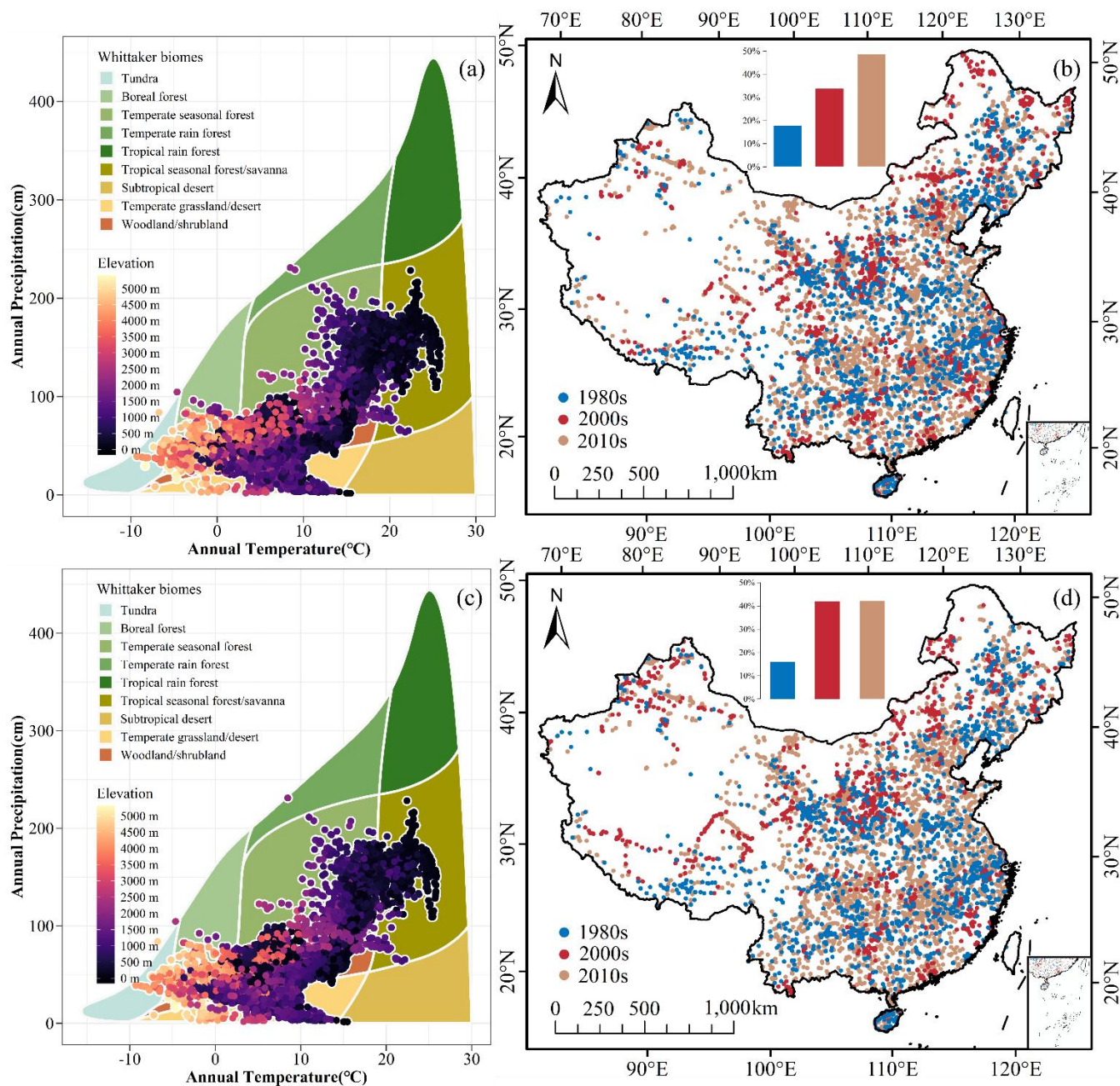


Figure 4. Statistical characteristics of SOCD sample points in different periods. SOCD data with the soil depth of 0-20 cm (a-c) and 0-100 cm (d-f) in China during the 1980s, 2000s, and 2010s are evaluated comprehensively.



555

Figure 5. Spatial distribution of SOC sample points. (a-d) Distribution of the SOCD sampling sites with data used in this study for (b) top 20 cm soil and (d) top 100 cm soil. The distribution of site-level training data is based on Whittaker biomes for (a) top 20 cm soil and (c) top 100 cm soil.

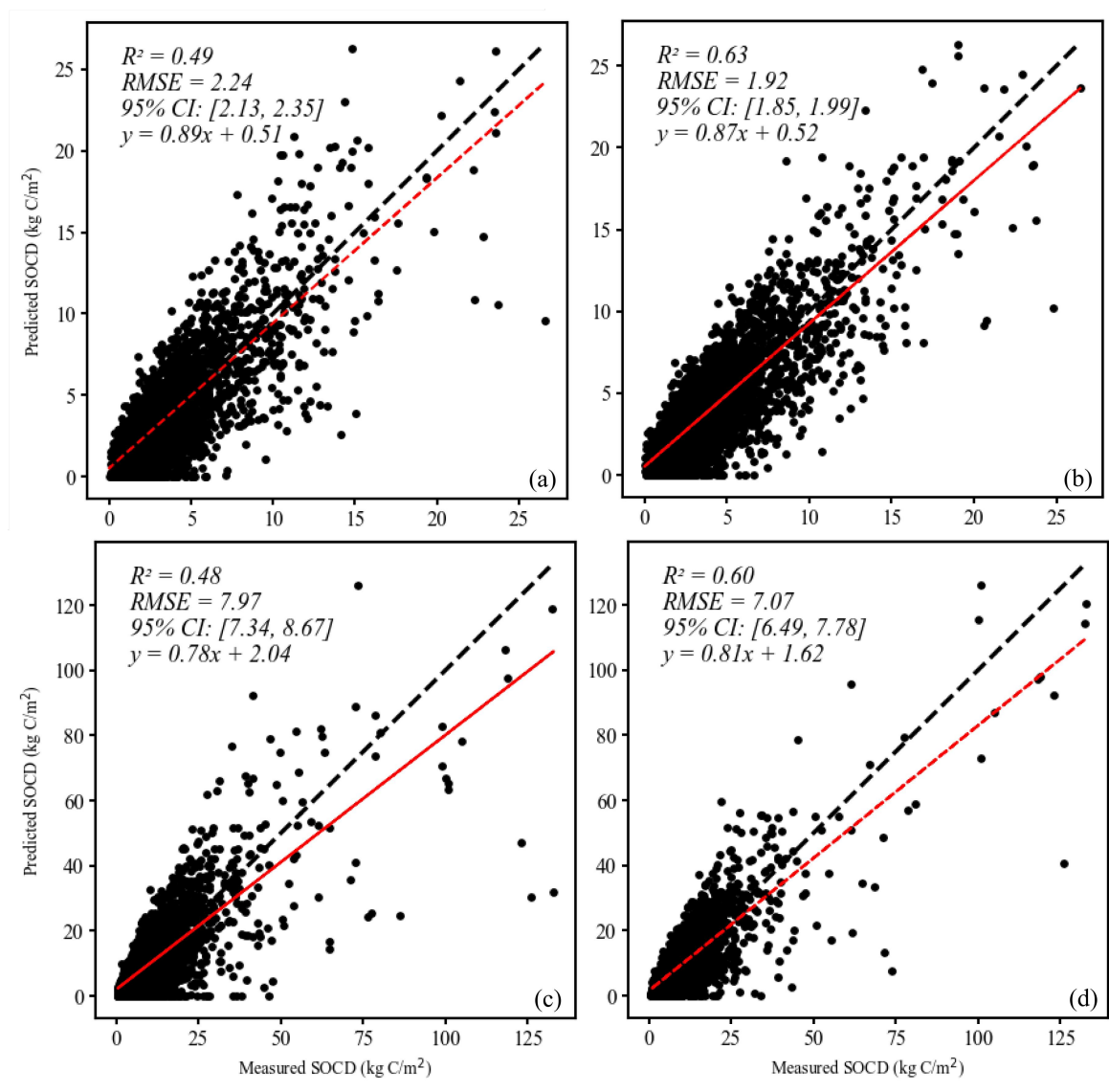


Figure 6. The model performance of global and zoning models with the depth of 0-20 cm and 0-100 cm. The SOCD prediction model of 0-20 cm and 0-100 cm soil depth is evaluated strictly by using a variety of statistical indicators, corresponding to four evaluation results, 0-20 cm global model (a), 0-20 cm regional model (b), 0-100 cm global model (c), and 0-100 cm regional model (d).

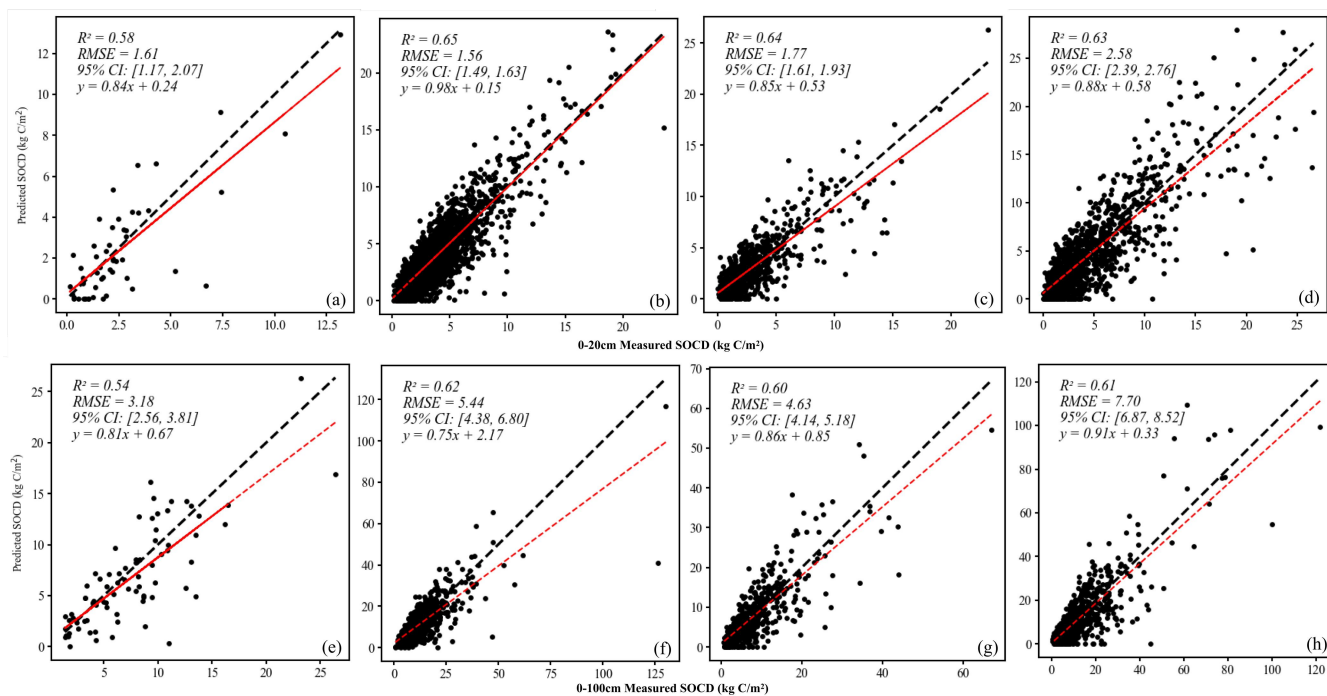


Figure 7. The model performance of different zoning models with the depth of 0-20 cm and 0-100 cm. Panels (a) and (e) depict the model performance for arid regions, where water scarcity is a predominant factor affecting SOC. Panels (b) and (f) illustrate results for humid regions characterized by high moisture availability. Panels (c) and (g) showcase semi-arid regions, where the balance between precipitation and evaporation influences SOC patterns. Finally, panels (d) and (h) display model accuracy in semi-humid regions, which exhibit intermediate conditions between arid and humid environments.

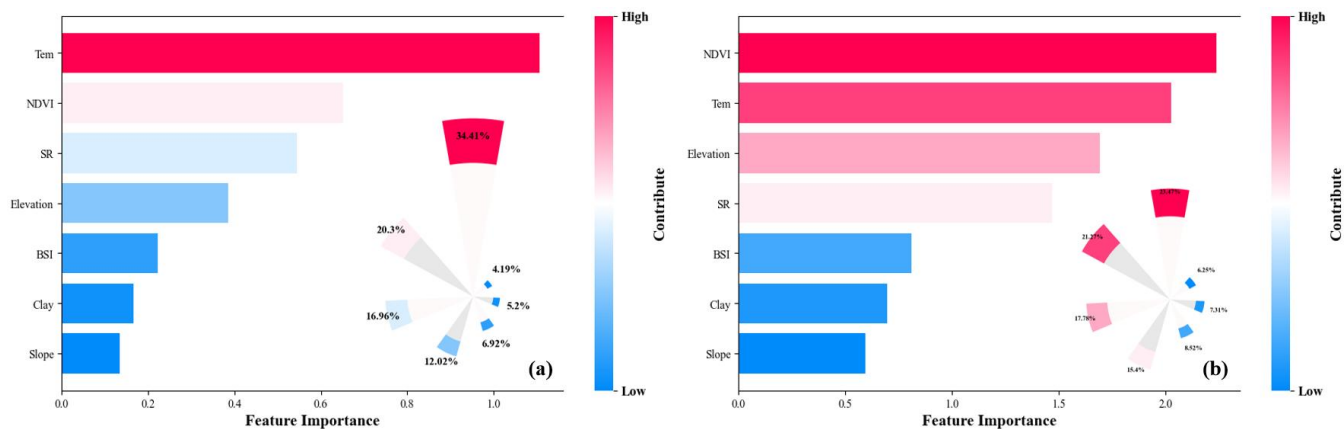


Figure 8. Importance ranking of features for SOCD estimation with the depth of 0-20 cm and 0-100 cm. It reports the contribution of different environmental variables to the SOCD estimation with different soil depths, including feature importance ranking for 0-20 cm depth (a) and feature importance ranking for 0-100 cm depth (b).

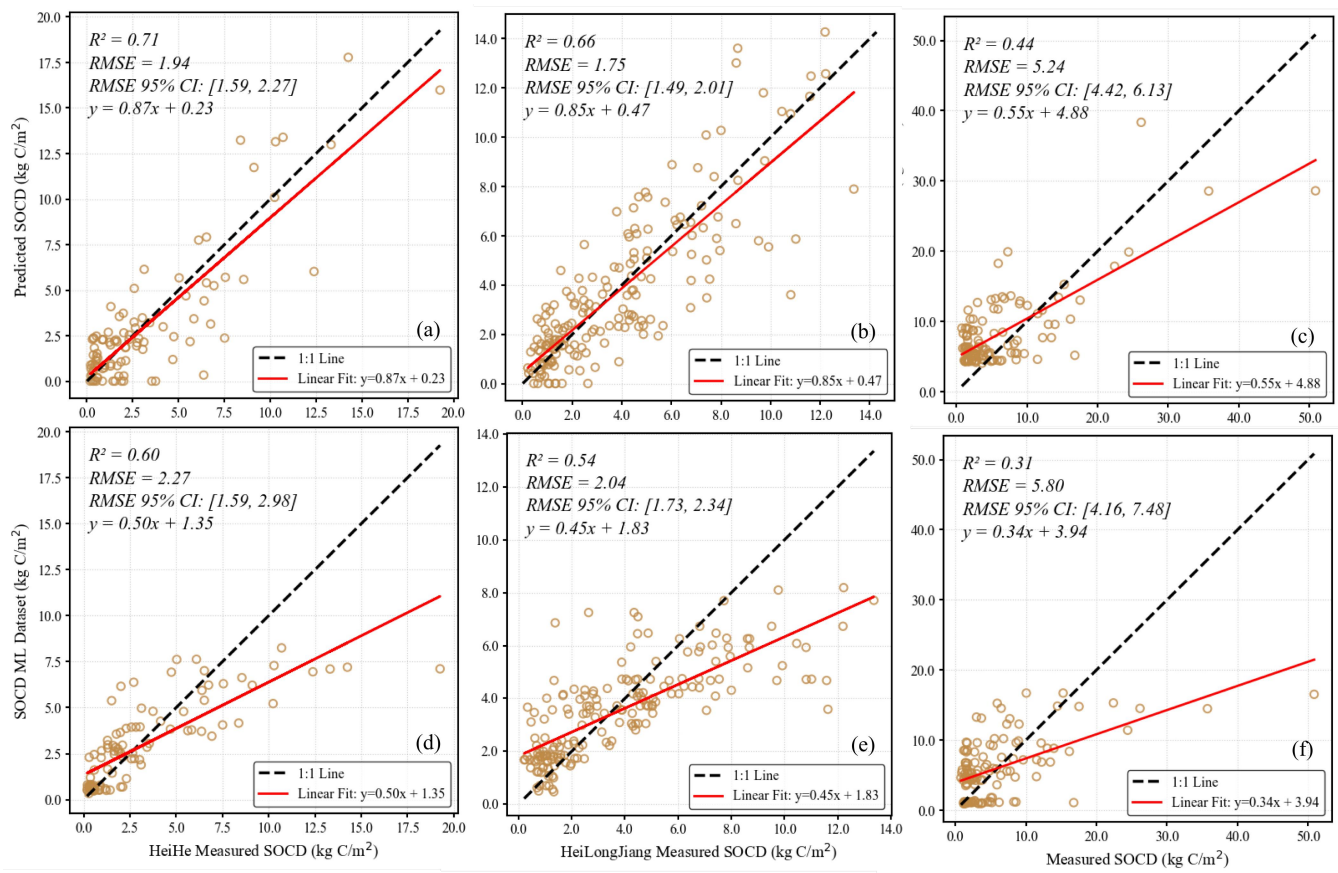


Figure 9. Comparison of predicted and machine learning (ML) derived SOCD with independent measurements at various depths. Panels (a), (b), and (c) display correlations of this study's predicted SOCD, while (d), (e), and (f) show correlations of SOC Dynamics ML dataset. Specifically, (a) and (d) are for 0-20 cm SOCD against Heihe River basin measurements. (b) and (e) compare 0-20 cm SOCD with Xu's published data. (c) and (f) present 0-100 cm SOCD correlations with measurements from Dong et al. (2024) and simulations from Li et al. (2022).

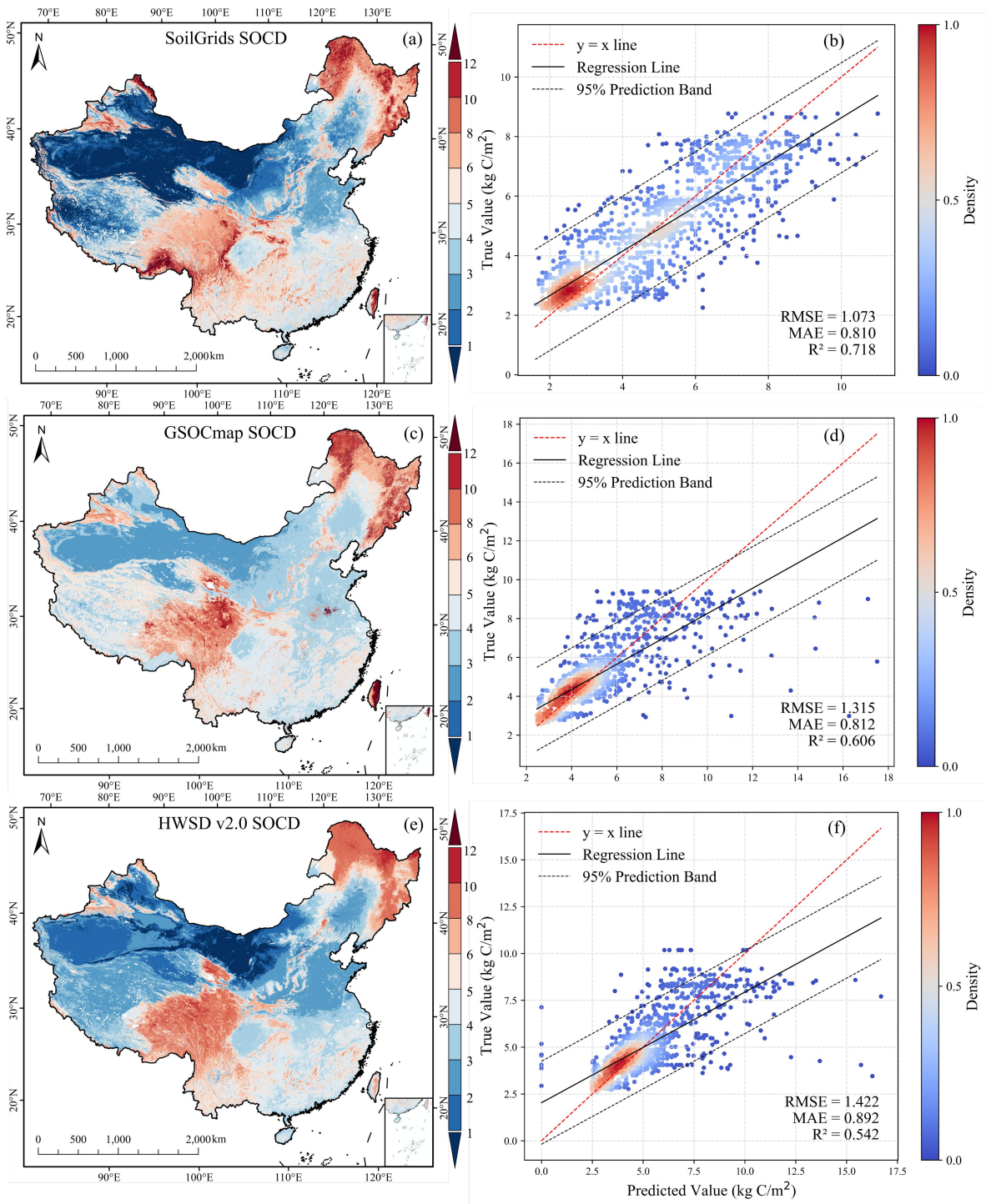


Figure 10. Comparison with three published global products. Our estimated SOCD is compared with the SoilGrids250m (a & b),

590 GSOCmap (c & d), and HWSD v2.0 (e & f) datasets.

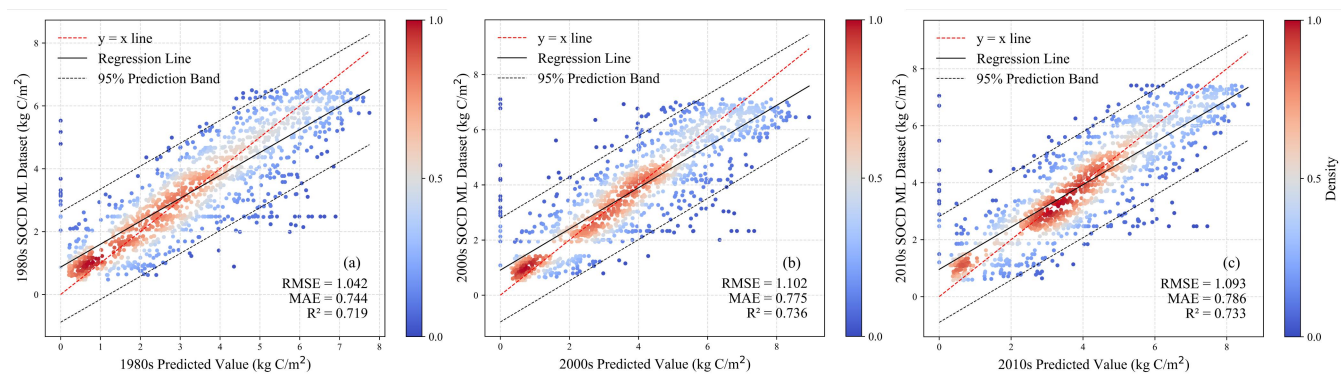


Figure 11. Comparison with the SOC Dynamics ML dataset with a depth of 0-20 cm in China in the 1980s (a), 2000s (b), and 2010s (c).

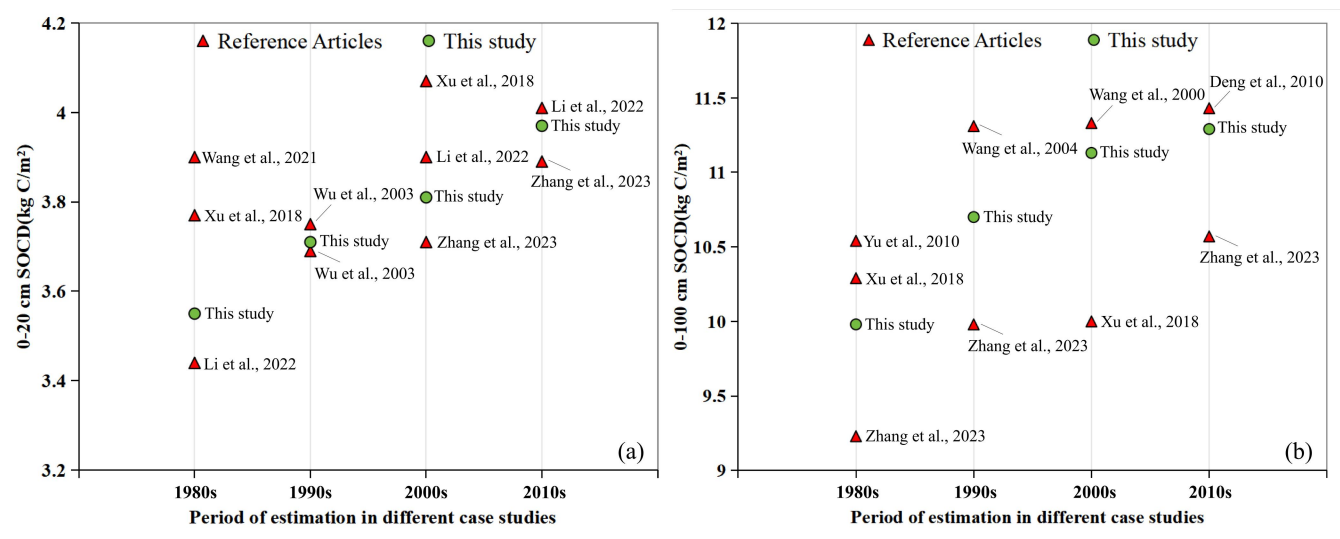


Figure 12. Aggregated results of estimated SOC D with the depth of 0-20 cm (a) and 0-100 cm (b) in China from this study and previous investigations.

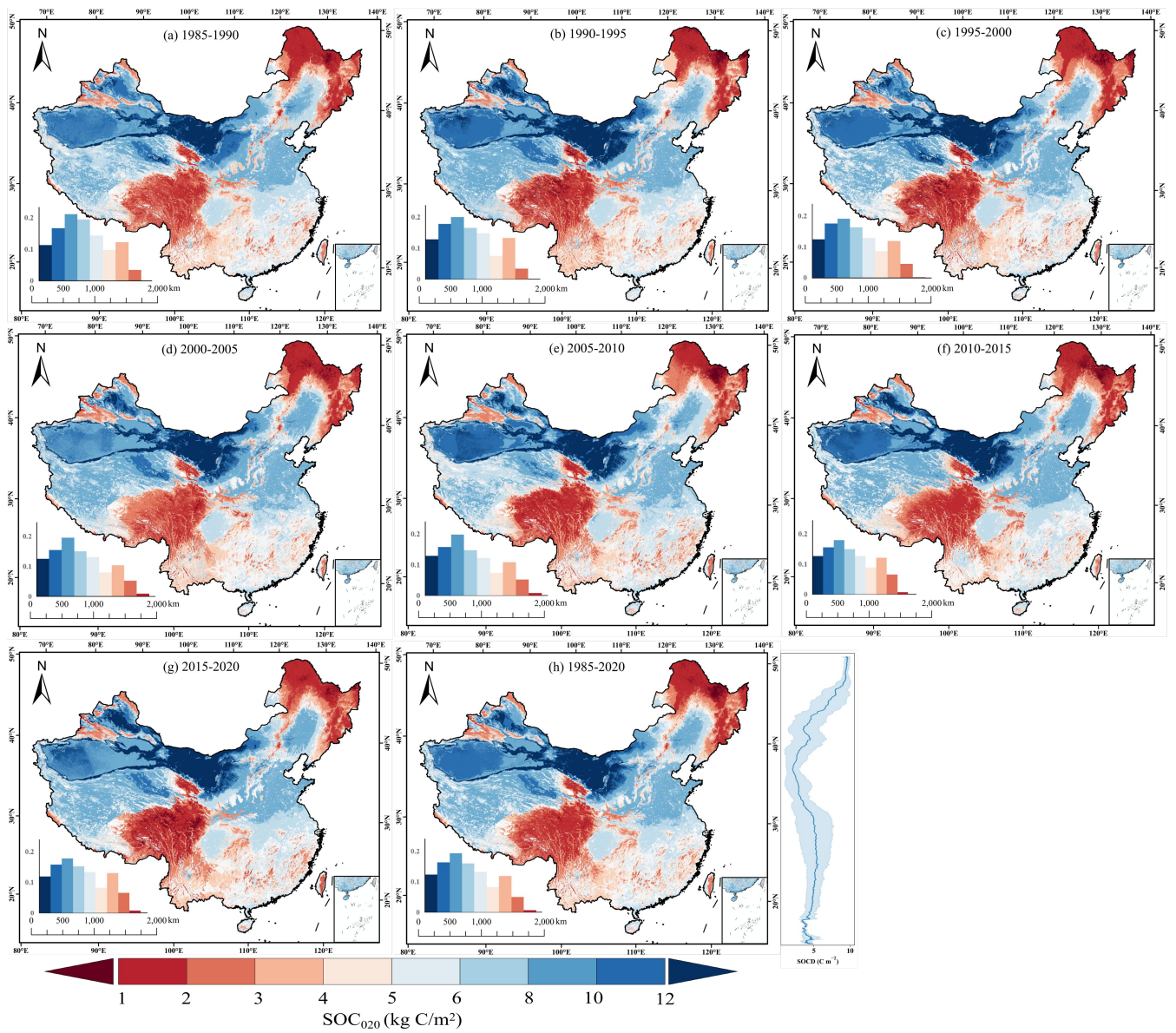


Figure 13. Spatial distribution of estimated SOC_D at a depth of 0-20 cm in 1985-1990 (a), 1990-1995 (b), 1995-2000(c), 2000-2005 (d), 2005-2010 (e), 2010-2015 (f), 2015-2020 (g) and average from 1985 to 2020 (h). The lower left histograms in each panel show the area ratios for different SOC_D levels.

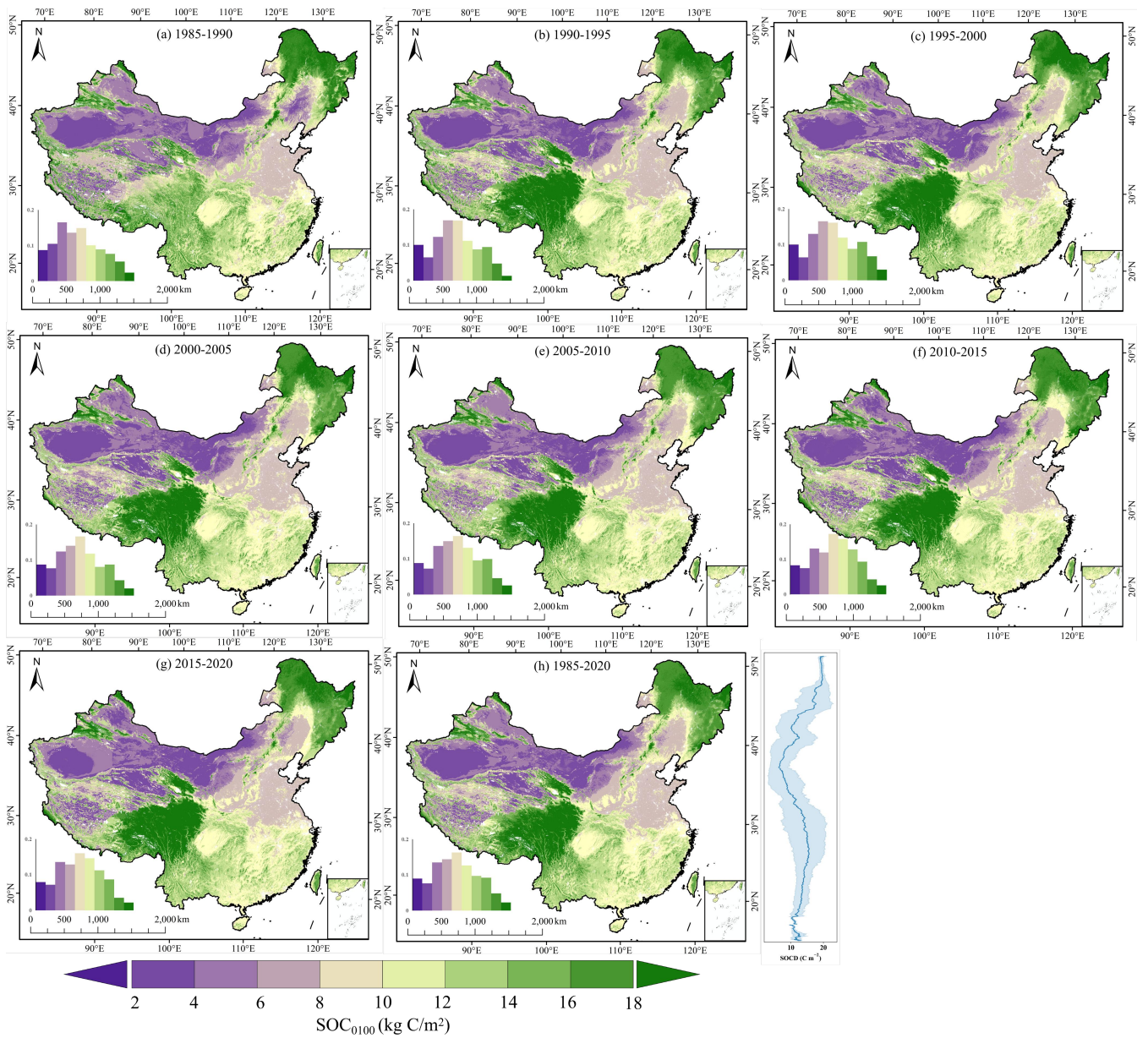


Figure 14. Spatial distribution of estimated SOC_D at a depth of 0-100 cm in 1985-1990 (a), 1990-1995 (b), 1995-2000(c), 2000-2005 (d), 2005-2010 (e), 2010-2015 (f), and 2015-2020 (g) and average from 1985 to 2020 (h). The lower left histograms in each panel show the area ratios for different SOC_D levels.