

Homework #1

by Catalina Jerez

Problem 1

Derivation of the Link Function, Fisher Score, and Fisher Information Matrix for Poisson and Gamma Distributions.

The **Poisson distribution** models the probability of observing y occurrences in a fixed interval, given the rate parameter λ :

$$P(Y = y|\lambda) = \frac{\lambda^y e^{-\lambda}}{y!}, \quad y = 0, 1, 2, \dots$$

where $\lambda > 0$ represents the expected number of occurrences.

1. In Generalized Linear Models (GLMs), we relate the expected value of the response variable to a linear predictor via a **link function** $g(\lambda) = \eta = X\beta$, where $g(\cdot)$ is the link function. For the Poisson distribution, the *canonical link function* is the *log link function*:

$$g(\lambda) = \log(\lambda) \quad \Rightarrow \quad \lambda = e^\eta$$

2. The **Fisher Score** (or Gradient of Log-Likelihood) for a single observation y is the log-likelihood function:

$$\begin{aligned} \ell(\lambda) &= y \log \lambda - \lambda - \log y! \quad (\text{taking the derivative w.r.t } \lambda) \\ \frac{\partial \ell}{\partial \lambda} &= \frac{y}{\lambda} - 1 \end{aligned}$$

Using the transformation $\lambda = e^\eta$, we compute $\frac{\partial \lambda}{\partial \eta} = e^\eta = \lambda$. Then, applying the chain rule:

$$\frac{\partial \ell}{\partial \eta} = \left(\frac{y}{\lambda} - 1 \right) \lambda = y - \lambda$$

3. **Fisher Information Matrix** is the expected negative second derivative of the log-likelihood. Mathematically, $\mathcal{I}(\eta) = -\mathbb{E} \left[\frac{\partial^2 \ell}{\partial \eta^2} \right]$. Computing the second derivative:

$$\begin{aligned}\frac{\partial^2 \ell}{\partial \eta^2} &= -\lambda \quad (\text{using } \mathbb{E}[y] = \lambda) \\ \Rightarrow \mathcal{I}(\eta) &= -\mathbb{E}(-\lambda) \\ \mathcal{I}(\eta) &= \lambda \quad (\text{replacing } \lambda = e^\eta) \\ \mathcal{I}(\eta) &= e^\eta\end{aligned}$$

For a sample of size n , the **total Fisher information** is $\mathcal{I}(\eta) = ne^\eta$.

The **Gamma distribution** is commonly used to model positive continuous data. The probability density function (PDF) is:

$$f(y|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}, \quad y > 0$$

where α is the shape parameter, and β is the rate parameter (inverse of the scale parameter). We use a mean parameterization where $\mathbb{E}[Y] = \mu = \frac{\alpha}{\beta}$.

1. The **canonical link function** for the Gamma distribution is the inverse link function:

$$g(\mu) = \frac{1}{\mu} \quad \Rightarrow \quad \mu = \frac{1}{\eta}$$

2. The log-likelihood function for a single observation is (**Fisher Score**) is:

$$\begin{aligned}\ell(\mu) &= \alpha \log \beta + (\alpha - 1) \log y - \beta y - \log \Gamma(\alpha) \quad (\text{replacing } \beta = \frac{\alpha}{\mu}) \\ \ell(\mu) &= \alpha \log \left(\frac{\alpha}{\mu} \right) + (\alpha - 1) \log y - \frac{\alpha}{\mu} y - \log \Gamma(\alpha) \quad (\text{derivative w.r.t } \mu) \\ \frac{\partial \ell}{\partial \mu} &= -\frac{\alpha}{\mu} + \frac{\alpha}{\mu^2} y \quad (\text{using } \mu = \frac{1}{\eta}) \\ \frac{\partial \mu}{\partial \eta} &= -\frac{1}{\mu^2} \quad (\text{applying the chain rule}) \\ \frac{\partial \ell}{\partial \eta} &= \left(-\frac{\alpha}{\mu} + \frac{\alpha}{\mu^2} y \right) \left(-\frac{1}{\mu^2} \right) \\ \frac{\partial \ell}{\partial \eta} &= \alpha \frac{y - \mu}{\mu^2}\end{aligned}$$

1. For the **Fisher Information Matrix**, we compute the second derivative of:

$$\frac{\partial^2 \ell}{\partial \eta^2} = -\alpha \frac{1}{\mu^2} \quad (\text{using } \mathbb{E}[y] = \mu)$$

$$\mathcal{I}(\eta) = \alpha \frac{1}{\mu^2}$$

Thereby, for a sample of size n , the total Fisher information is $\mathcal{I}(\eta) = n\alpha \frac{1}{\mu^2}$.

Thus, Table 1 presents the canonical link function, the gradient of the log-likelihood (Fisher score), and the Fisher information matrix for both the Poisson and Gamma distributions.

Table 1: Summary of the link function, Fisher score, and Fisher information matrix for Poisson and Gamma distributions.

Distribution	Link Function	Fisher Score	Information Matrix
Poisson	$\log(\lambda)$	$y - \lambda$	e^η
Gamma	$\frac{1}{\mu}$	$\alpha \frac{y - \mu}{\mu^2}$	$\alpha \frac{1}{\mu^2}$

Problem 2a: Spatial map

Figure 1 shows the average annual precipitation in India, overlaid with topographic data to highlight the relationship between elevation and rainfall.

High-altitude regions, such as the Western Ghats and the Himalayan foothills, receive substantial precipitation, often exceeding 3000 mm annually (red-pink zones), due to orographic effects. In contrast, central and northern India experience significantly lower precipitation (<1000 mm annually), characteristic of arid and semi-arid climates.

Rainfall patterns follow distinct spatial trends: intense rainfall occurs along the western coast due to southwest monsoon winds, which gradually weaken inland. Similarly, northeastern states and Himalayan foothills exhibit high precipitation, while peninsular India shows a declining gradient from coastal to interior regions.

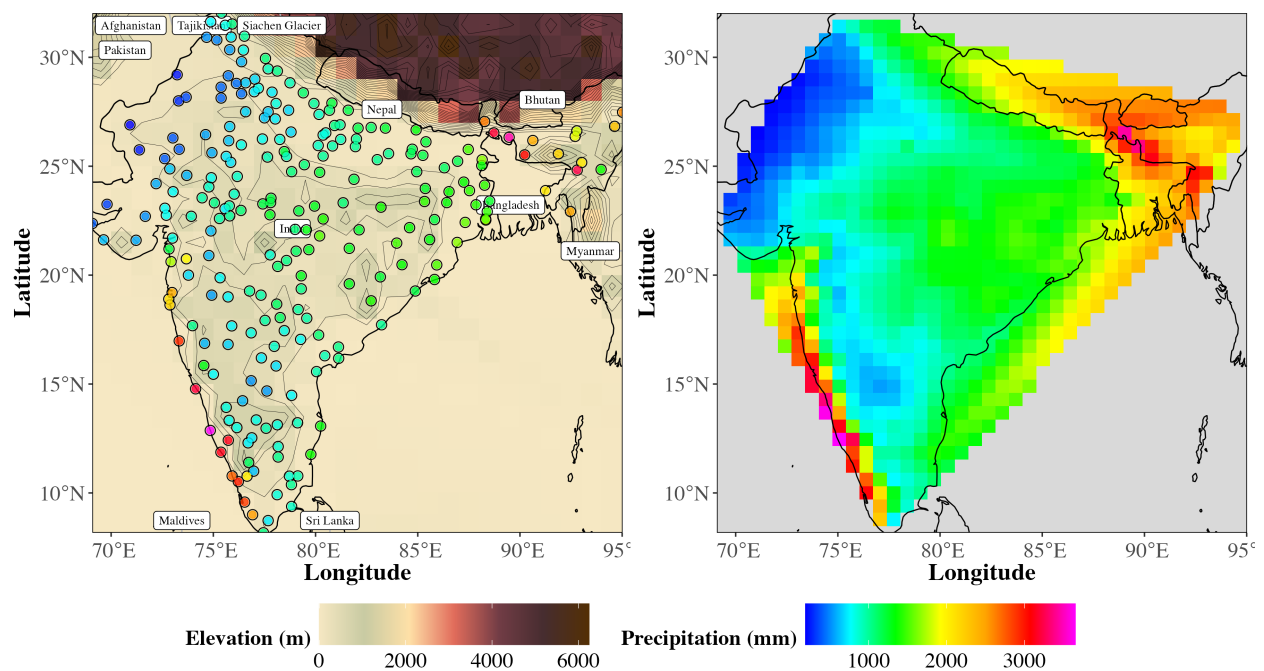


Figure 1: Spatial (left) and Surface (right) map of annual average precipitation along with the topography over India.

Problem 2b: Geodesic Distances

Calculating distances from precipitation stations to the Arabian Sea and the Bay of Bengal was carried out with the **geosphere package** (Hijmans, 2010; Karney, 2013), using the Haversine formula (for spherical distance approximation, Eq. 1) and the Geodesic formula (Vincenty's formula for ellipsoidal distance, Eq. 2) for the calculation of geodetic distances. The minimum distances were calculated from each station to 1) any coast, 2) Arabian Sea and 3) Bay of Bengal.

$$d = 2R \arcsin \left(\sqrt{\sin^2 \left(\frac{\Delta\varphi}{2} \right) + \cos(\varphi_1) \cos(\varphi_2) \sin^2 \left(\frac{\Delta\lambda}{2} \right)} \right) \quad (1)$$

where d is the distance between two points; $R \approx 6378$ Km is Earth's radius; φ_1, φ_2 the latitudes of points; and λ_1, λ_2 the longitudes of points.

$$s = a \tan^{-1} \left(\frac{\sqrt{(\cos U_2 \sin \lambda)^2 + (\cos U_1 \sin U_2 - \sin U_1 \cos U_2 \cos \lambda)^2}}{\sin U_1 \sin U_2 + \cos U_1 \cos U_2 \cos \lambda} \right) \quad (2)$$

where s is the geodesic distance, a is the semi-major axis of Earth (we use $a = 6378137$ for WGS84), U_1, U_2 = reduced latitudes, λ the difference in longitude.

Figure 2 illustrates the spatial distribution of minimum distances calculated using the Haversine formula. The pattern shows that locations closer to the Arabian Sea and the Bay of Bengal have shorter distances, which aligns with monsoonal precipitation influences. However, the Haversine formula assumes a perfect sphere, leading to minor inaccuracies, especially for larger distances.

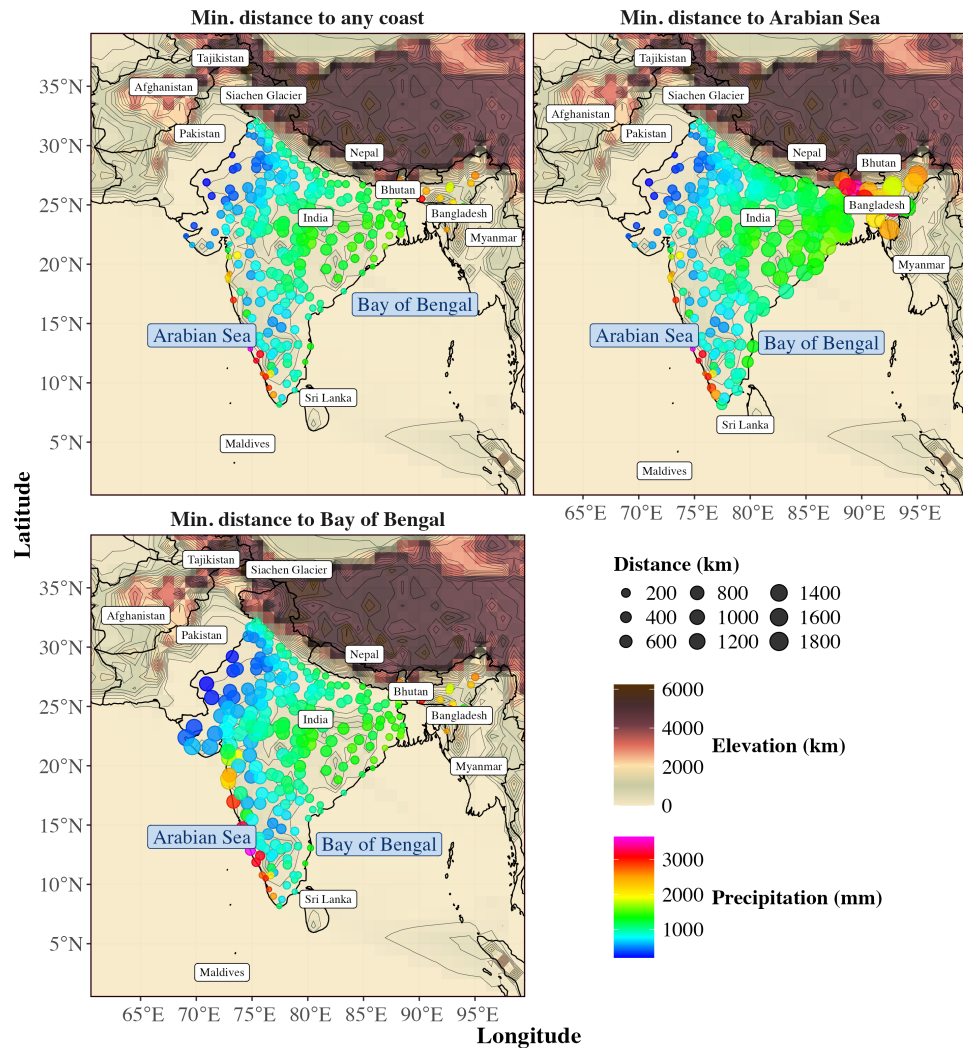


Figure 2: Annual average precipitation and topography over India, and distance calculated with Haversine method from weather stations to Arabian and Bay of Bengal Sea.

Figure 3 shows the geodesic-based distance computation that accounts for the Earth's ellipsoidal shape, making it a more accurate representation of actual distances. The figure demonstrates similar trends as the Haversine method but with slight variations due to the improved accuracy of the geodesic approach.

Figure 4 provides a comparison between the two methods. The differences are more noticeable over longer distances, where the Haversine method tends to slightly underestimate or overestimate distances due to its spherical assumption.

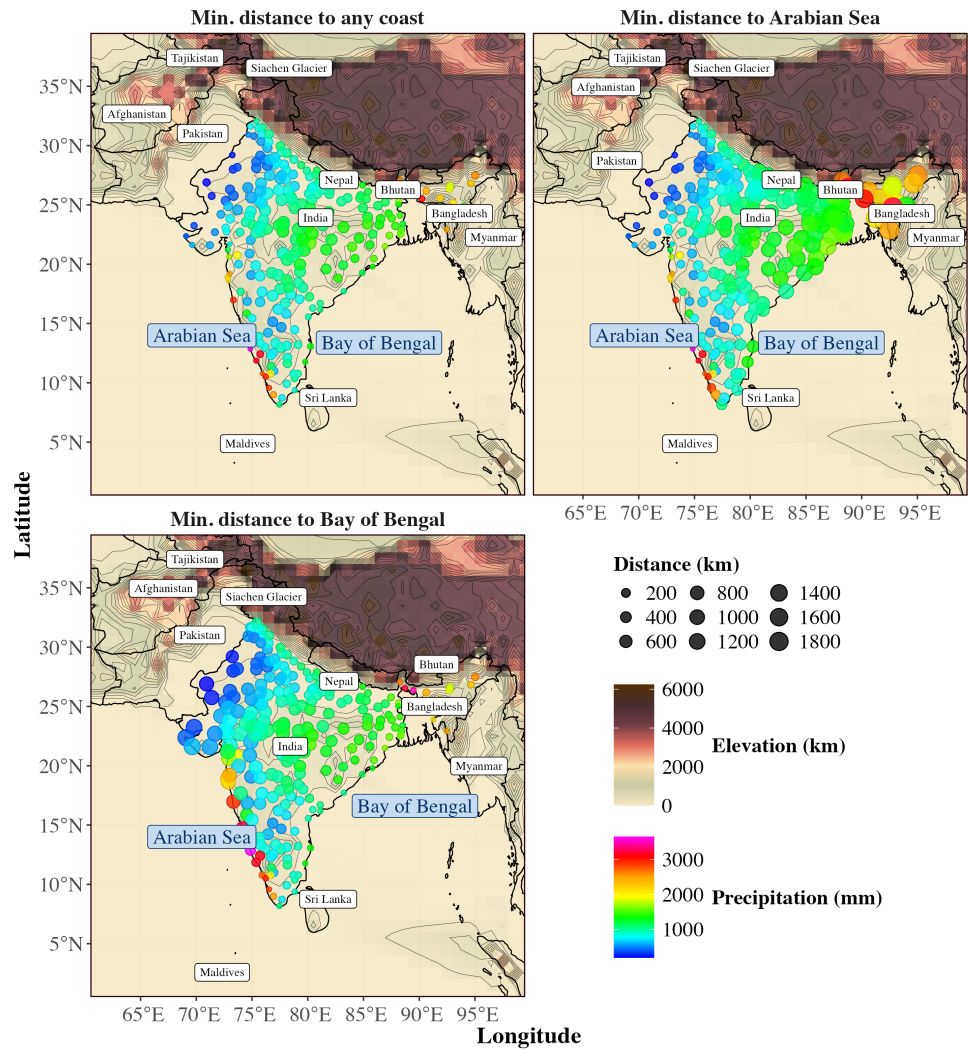


Figure 3: Annual average precipitation and topography over India, and distance calculated with Geodesic method from weather stations to Arabian and Bay of Bengal Sea.

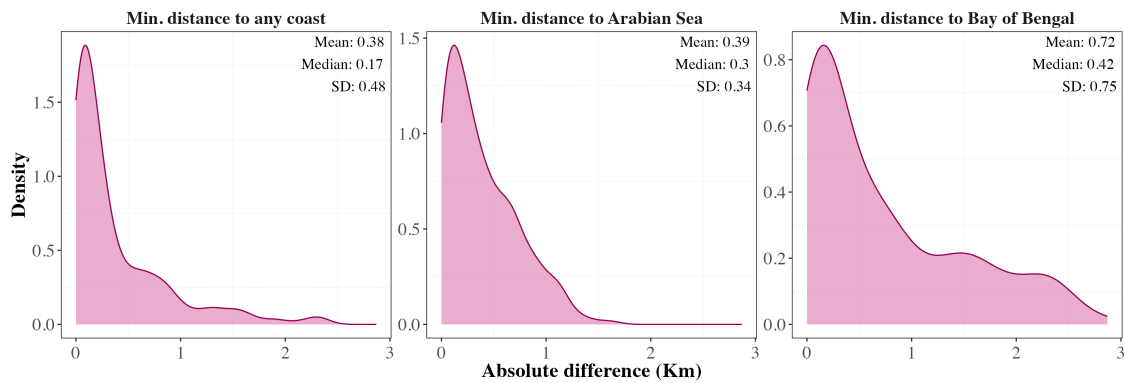


Figure 4: Comparison of Haversine vs. Geodesic Distances

Problem 2c: Underlying precipitation surface under linear model

The goal of this analysis is to estimate the underlying precipitation surface using sparse spatial observations. This is crucial for applications such as floodplain management, natural hazard mitigation, and climate modeling. The study utilizes regression models, cross-validation, and spatial mapping techniques.

Model selection and fitting

The best model was selected based on the Bayesian Information Criterion (BIC). The initial regressors included Latitude (lat), Longitude (lon), Elevation (elev), Distance to the Arabian Sea (mdArabian), and Distance to the Bay of Bengal (mdBengal). A Variance Inflation Factor (VIF) analysis was conducted to check for multicollinearity, and mdArabian and mdBengal were removed due to high VIF values (see Appendix C.2, Figure 5). The final selected predictors are lat, lon, elev, and mdCoast. Additionally, we include the following interaction terms: lld ($\text{lon} \times \text{lat} \times \text{mdCoast}$), ll2d ($\text{lon} \times \text{lat}^2 \times \text{mdCoast}$), l2ld ($\text{lon}^2 \times \text{lat} \times \text{mdCoast}$), and ed ($\text{elev} \times \text{mdCoast}$).

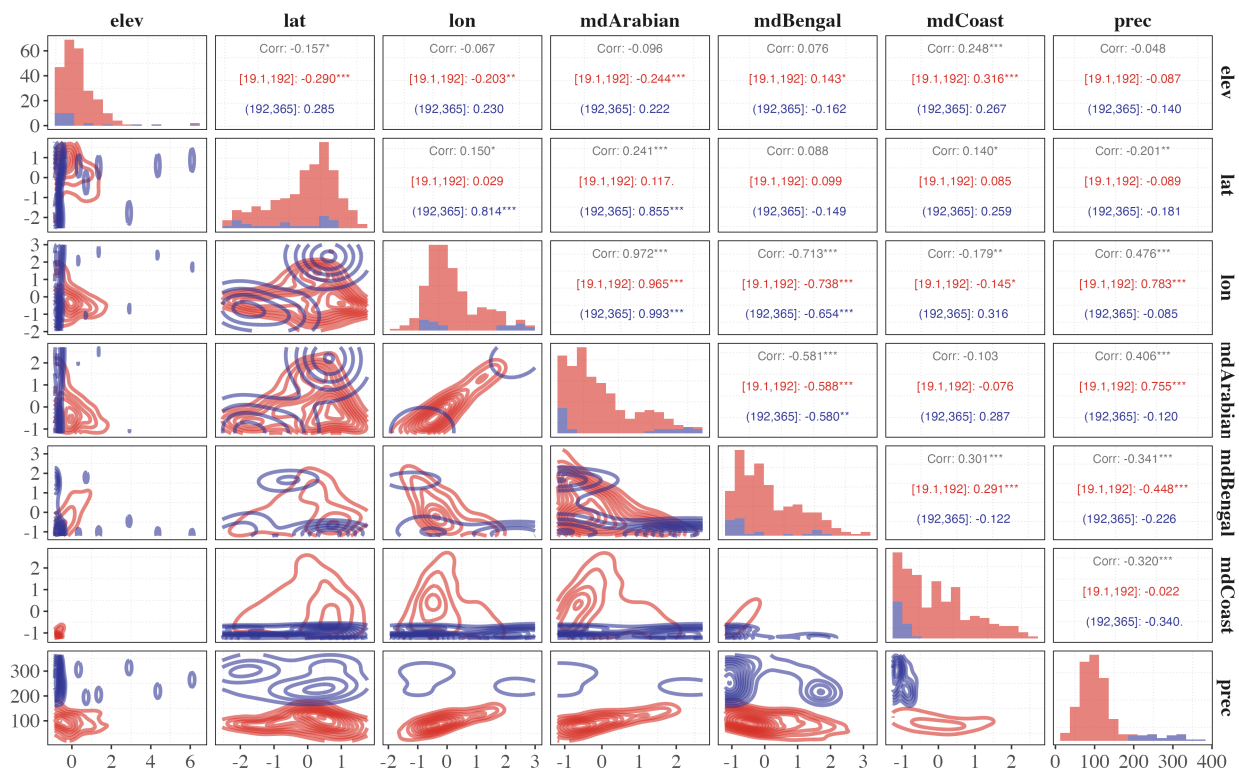


Figure 5: Caption

Scatterplot of observed vs modeled precipitation

The R^2 value of 0.67 (Figure 6) suggests a moderate to strong correlation; however, 33% of the variance remains unexplained, which suggests potential improvements. Besides, the model returns a RMSE of 36.24 cm, which might be large depending on the context (e.g., hazard mitigation).

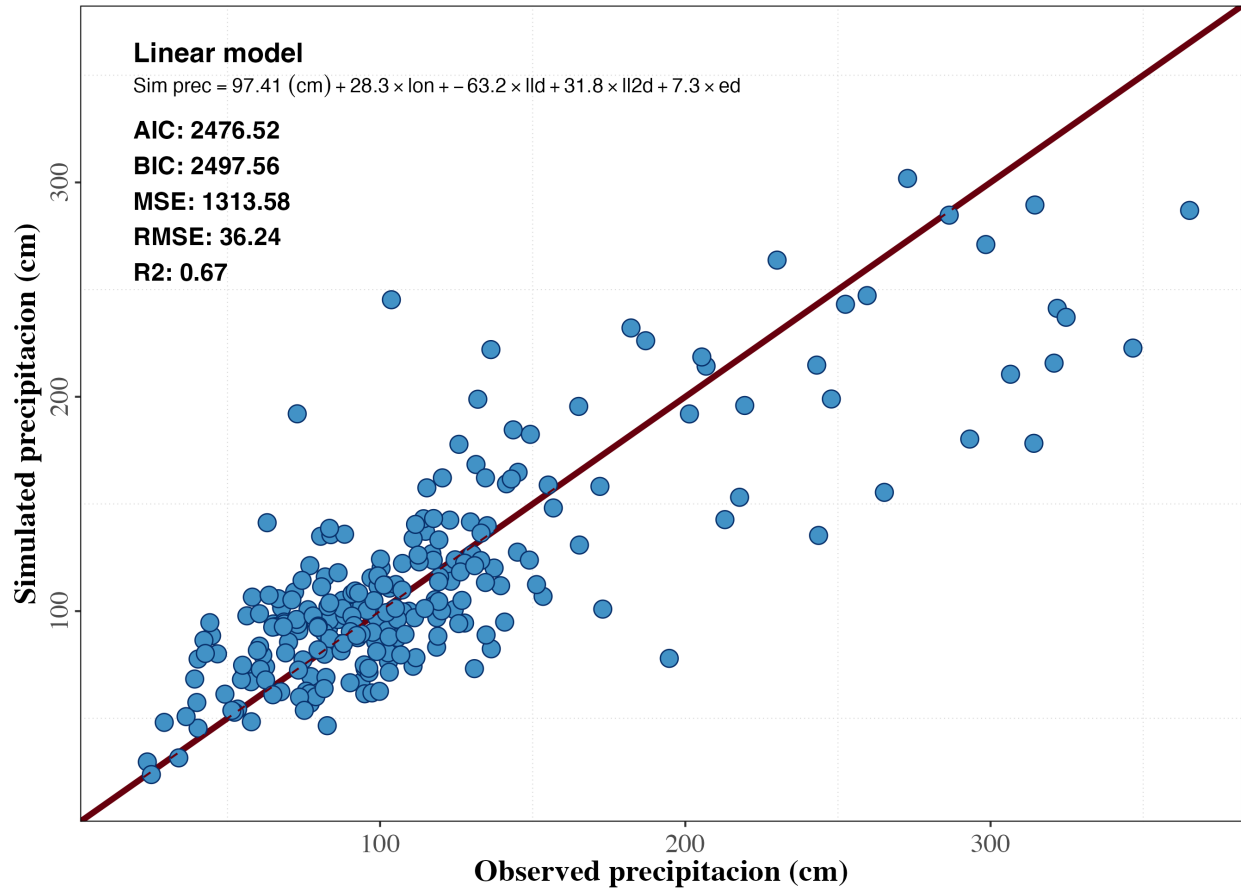


Figure 6: Scatterplot of observed vs modeled precipitation for linear regression model.

Model diagnostics

The ANOVA analysis confirms that the linear model is statistically significant ($p < 0.001$), indicating that at least one predictor contributes significantly to the response variable (precipitation). The results of ANOVA are summarized as: **lon**: significant ($p < 2e-16$, strong influence); **ll2d**: highly significant ($p < 2e-16$, major contribution); **ed**: also significant ($p < 2e-16$, meaningful but smaller effect); **Residuals**: variability remains, suggesting room for model refinement.

Figure 7 provides essential model validation checks:

- a) Normality: the Q-Q plot shows that residuals are approximately normal, though slight deviations exist at the tails.
- b) Histogram: residual distribution appears approximately normal, though some skewness may be present.
- c) Independence: the autocorrelation function (ACF) suggests that residuals are largely uncorrelated, supporting the assumption of independence.
- d) Homoscedasticity: the residuals vs. estimated precipitation plot shows no clear pattern, indicating relatively constant variance.

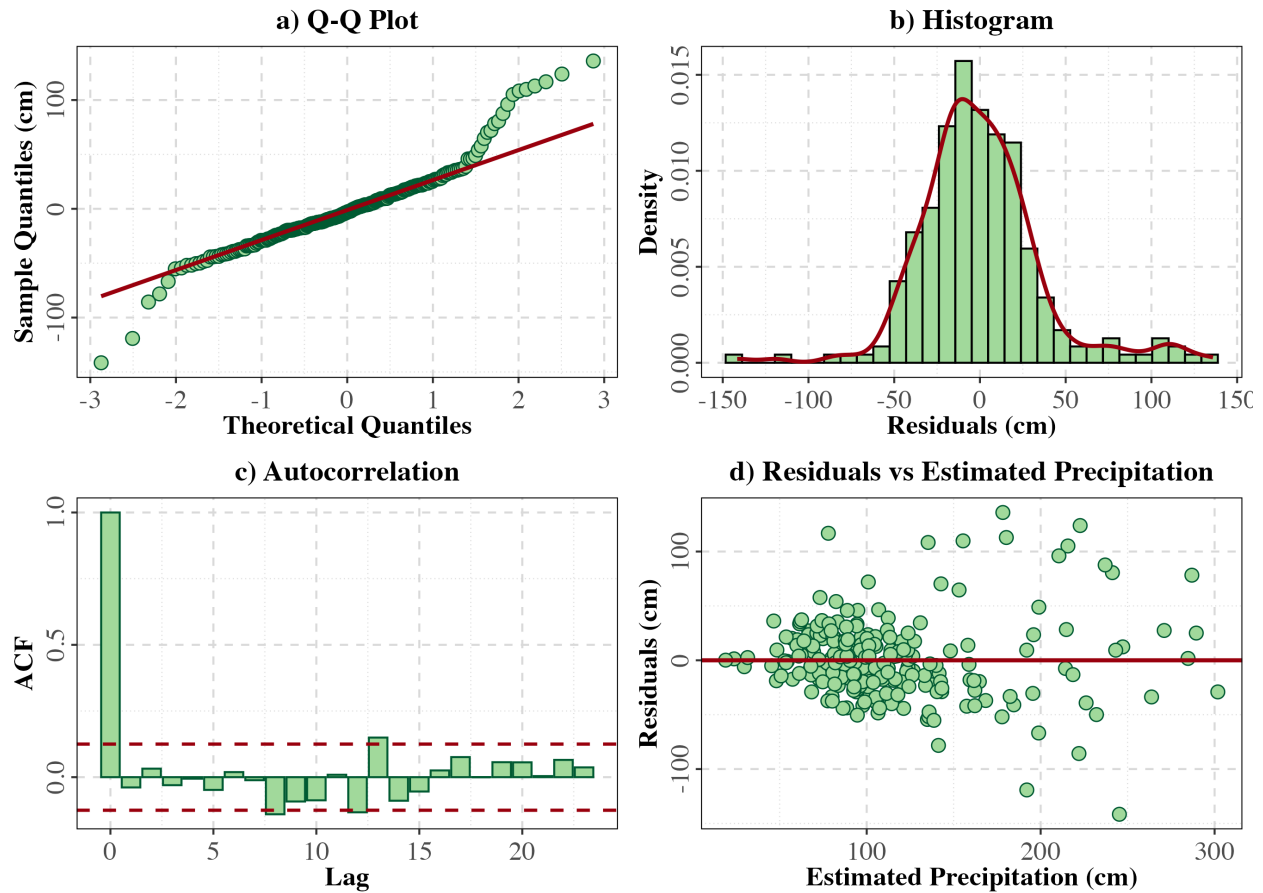


Figure 7: Model diagnostics for linear regression.

Cross-Validation Analysis

Leave-One-Out Cross-Validation (LOOCV) was performed to assess the model's generalization ability. The scatterplot in Figure 8 compares observed and simulated precipitation values from the cross-validated model, showing results consistent with the initial training model. However, the R^2 value (0.63) is slightly lower than the training model's R^2 of 0.67 (Figure 6), suggesting a small drop in predictive performance, which is expected in cross-validation due to reduced bias. In addition, the RMSE increase in 5% (from 36.24 to 38.09), indicating a minor generalization error.

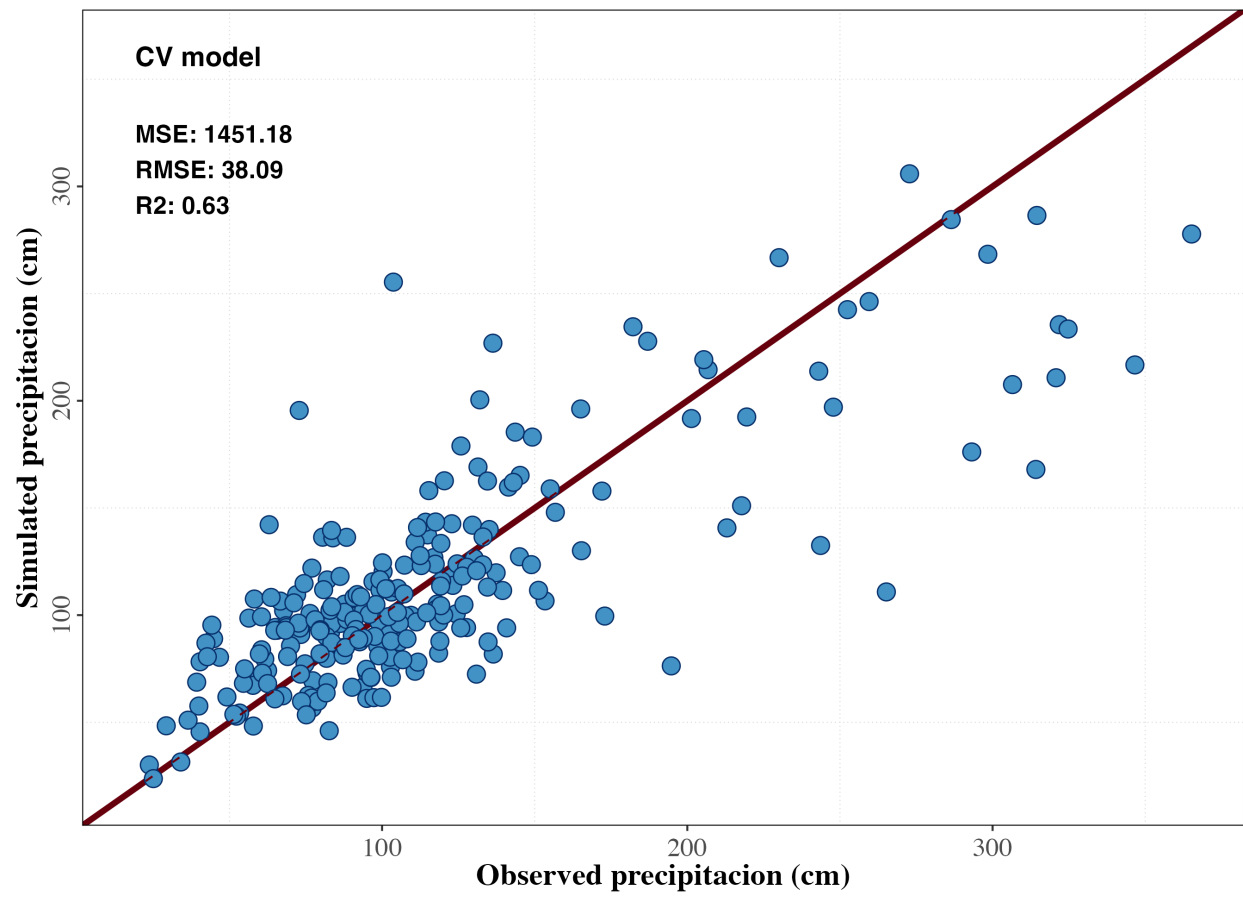


Figure 8: LOOCV Scatterplot.

Model robustness – cropping 10% of observations

To evaluate the robustness of the linear model, 10% of the observations were randomly removed, the model was re-fitted, and predictions were generated for the omitted data. This process assesses how well the model maintains stability when exposed to missing data.

The RMSE boxplot shows an expected increase in error after dropping 10% of the data (left panel in Figure 9), reflecting the model's slight sensitivity to data loss. Some outliers suggest that certain missing observations had a more significant impact on the prediction error.

The correlation between the observed and predicted values shows a wider spread and lower median correlation values (right panel in Figure 9). Many correlation values cluster around zero, suggesting that the model struggles to maintain predictive consistency when observations are removed. Some data points exhibit negative correlation, indicating possible overfitting in certain cases. Despite this, the overall trend remains relatively stable, though the model demonstrates reduced robustness to missing data.

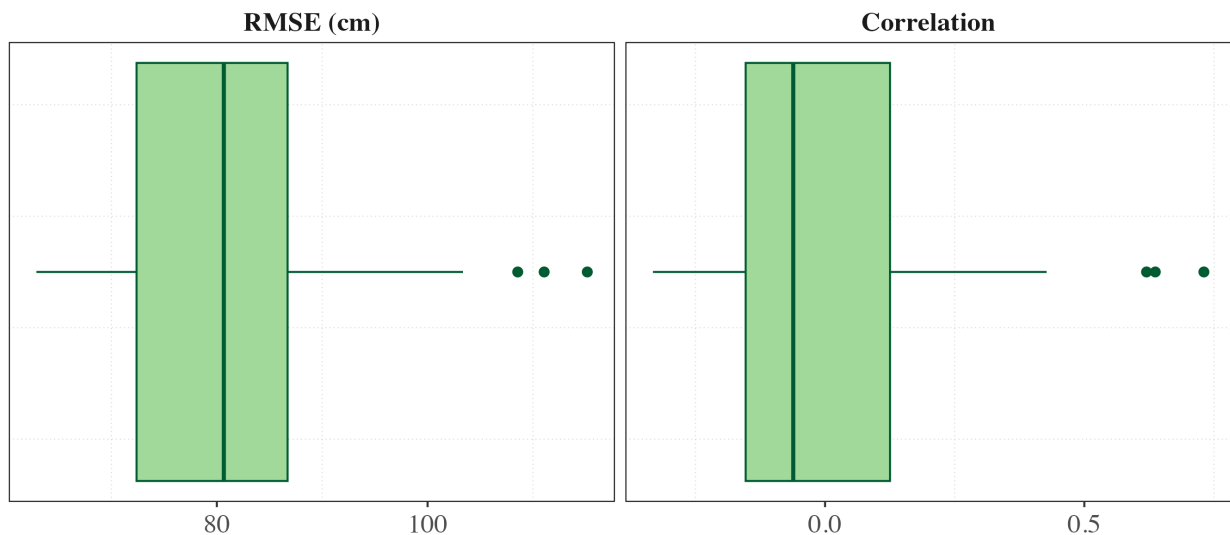


Figure 9: RMSE and correlation boxplots for model robustness analysis.

Spatial mapping

The map on the left of Figure 10 presents the predicted precipitation in India and the right map shows the standard error of predictions.

The predicted precipitation map shows higher precipitation in the northern and coastal regions, aligning with observed patterns (Figure 1). Lower precipitation is seen in central and western India, regions typically associated with drier conditions. The precipitation gradient follows a clear geographical trend, suggesting that the model captures regional climatic variations effectively.

Standard Error map (Figure 11) show lower standard errors are observed in densely sampled regions, indicating higher confidence in predictions. Higher standard errors appear in some inland and less-sampled areas, suggesting potential model uncertainty in regions with fewer observations.

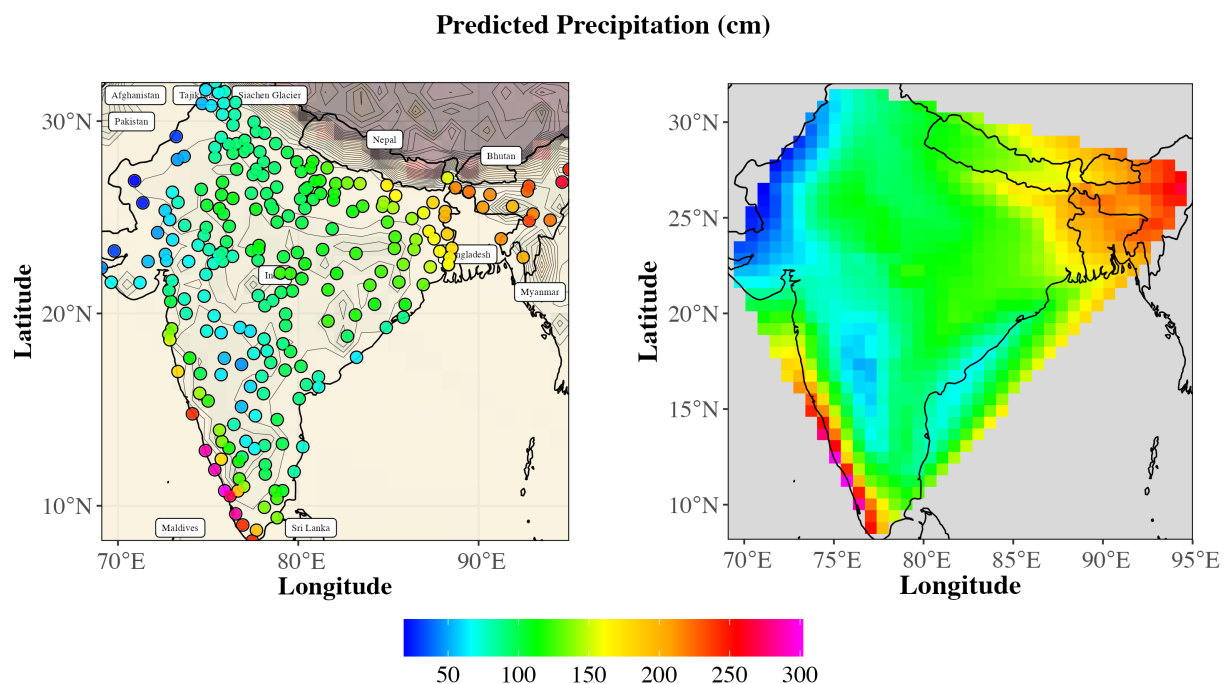


Figure 10: Predicted precipitation surface mapping.

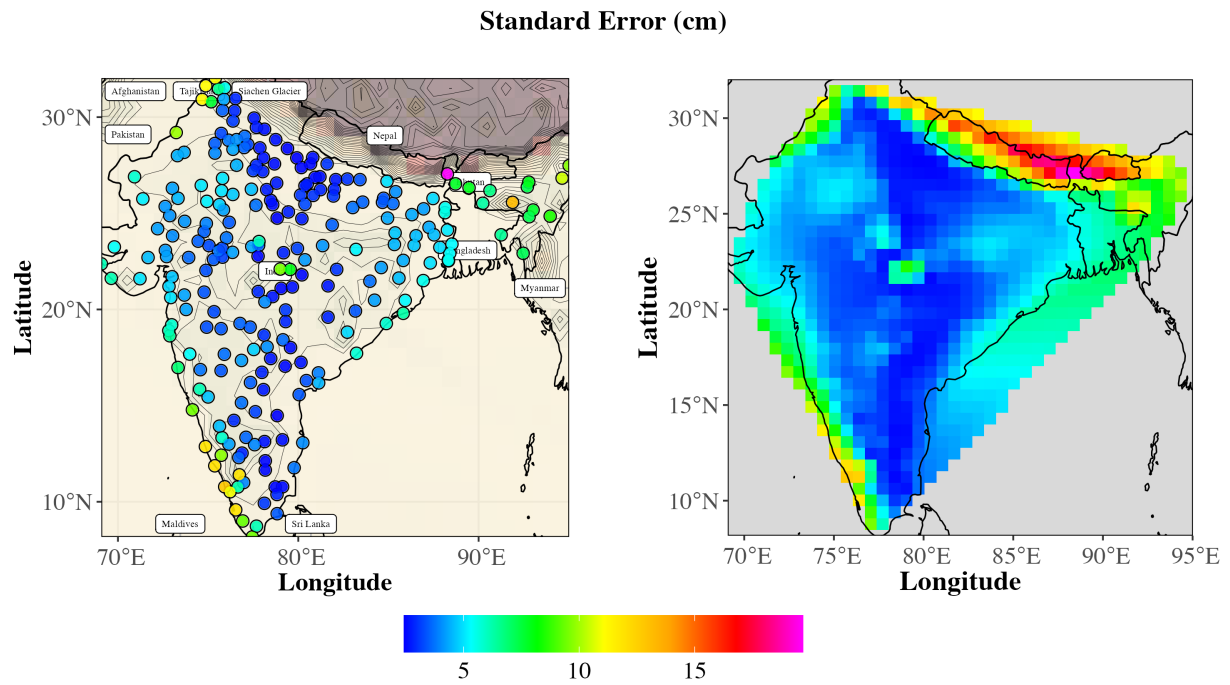


Figure 11: Standard error precipitation surface mapping.

Summary

The linear model was selected based on the lowest BIC, ensuring an optimal balance between predictive power and model complexity. Latitude and longitude emerged as the most significant predictors of precipitation, while elevation had a moderate influence. The distance to the coast (mdCoast) remained an important factor after removing collinear variables.

Model performance & diagnostics:

- The model demonstrated good predictive ability with an R^2 of 0.67 (Training) and 0.63 (LOOCV), indicating a moderate to strong correlation between observed and predicted precipitation.
- Residual analysis suggested that while the model assumptions largely hold, minor deviations in normality and homoscedasticity exist, which could be addressed through further refinements.

Cross-Validation & robustness:

- LOOCV results showed a slight increase in RMSE (from 36.24 to 38.09 cm), reflecting minor generalization error but overall stability.
- Robustness analysis (10% data removal) confirmed model reliability, though some sensitivity to high-leverage observations was observed.
- Correlation remained stable, but RMSE variations suggest that certain influential data points might impact predictions.

Spatial mapping:

- The predicted precipitation map captured key climatic trends, with higher rainfall in northern and coastal regions and lower precipitation in central/western regions.
- The standard error map highlighted higher uncertainty in inland regions, suggesting that improved data coverage in these regions could refine predictions.

Problem 3: Underlying precipitation surface under GLM

Scatterplot of observed vs modeled precipitation

The Generalized Linear Model (GLM) was fitted using the families Gamma, gaussian and binomial, the scatterplot in Figure 12 compares the observed vs. simulated values from the best fitted family.

The AIC and BIC values are slightly higher than the linear model, indicating no major improvement in model selection criteria. The R^2 value (0.67) indicates a strong correlation between observed and predicted precipitation, similar to the linear model. The RMSE of 36.24 cm suggests an error level comparable to the linear regression model. These values suggest that GLM does not provide a significant advantage over the linear approach. The scatterplot follows a strong linear trend, showing that the GLM effectively captures precipitation variations. However, variance increases at higher precipitation levels, which may indicate potential heteroscedasticity.

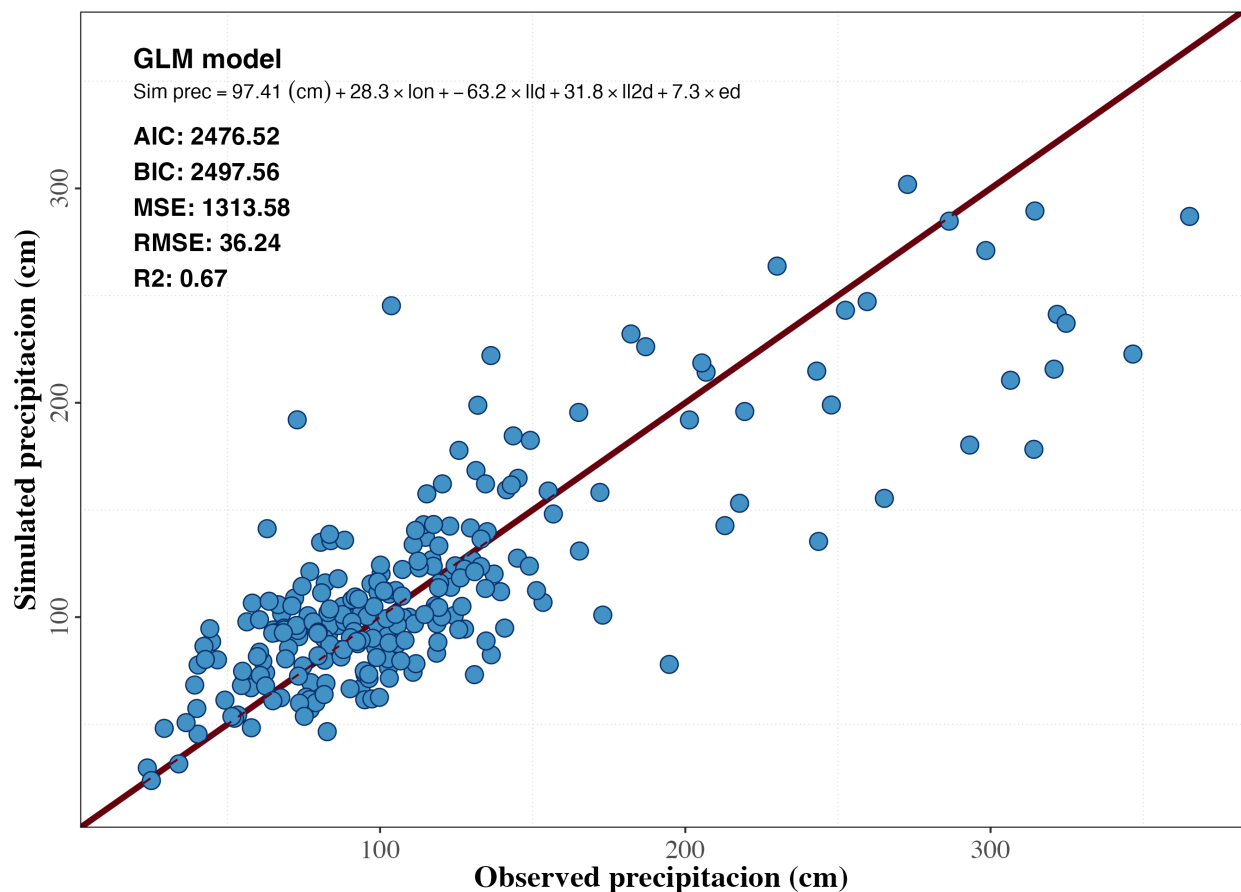


Figure 12: Scatterplot of observed vs modeled precipitation for Generalized Linear Model.

Model diagnostics

To evaluate the performance of the GLM, an ANOVA-based deviance analysis and residual diagnostics were conducted. Figure 13 presents key insights into the model's assumptions and fit quality. The GLM was fitted using a Gaussian distribution with an identity link function. The ANOVA analysis indicates that all predictors (lon, lld, ll2d, ed) are statistically significant ($p < 0.001$ or $p < 0.01$), confirming their contribution to precipitation modeling. Figure 13 provides essential model validation checks:

- a) Normality: residuals in Q-Q plot are approximately normally distributed, but minor deviations at the tails suggest some skewness.
- b) Histogram: the distribution appears roughly normal, though slight asymmetry and heavier tails suggest possible heteroscedasticity.
- c) Independence: low autocorrelation values confirm that residuals are mostly independent, satisfying one of the key model assumptions.
- d) Homoscedasticity: spread of residuals appears random, but a slight funnel shape suggests potential heteroscedasticity (variance increases at higher precipitation values).

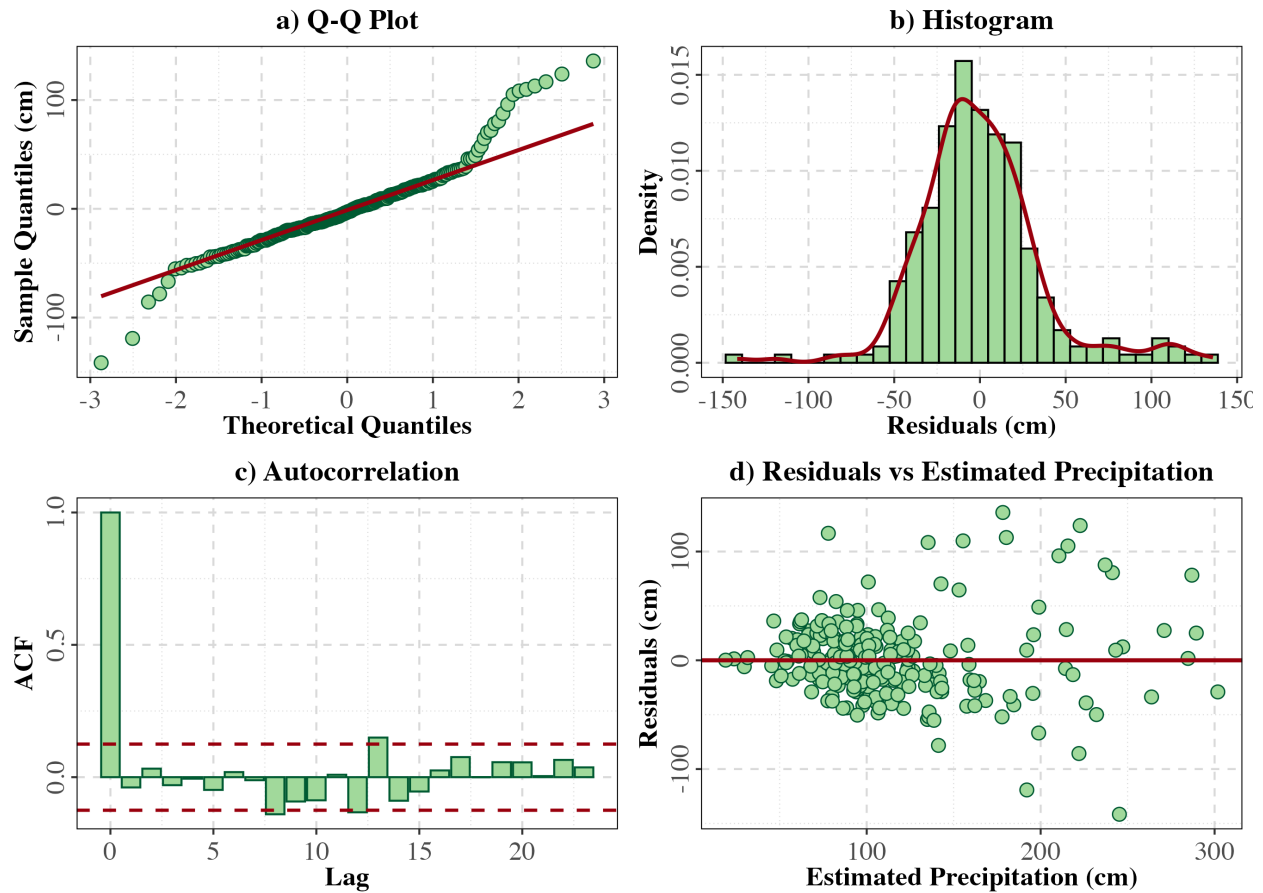


Figure 13: Model diagnostics for GLM regression.

Cross-Validation Analysis

The scatterplot (Figure 14) maintains a strong linear trend, reinforcing that the GLM is capturing overall precipitation patterns effectively. However, variance increases at higher precipitation values, suggesting that GLM predictions are less stable for extreme cases.

LOOCV results for GLM ($R^2 = 0.63$, $\text{RMSE} = 38.38$ cm) are nearly identical to the linear model's LOOCV performance ($R^2 = 0.63$, $\text{RMSE} = 38.09$ cm). This suggests that GLM does not provide a significant advantage over the linear model in terms of predictive accuracy.

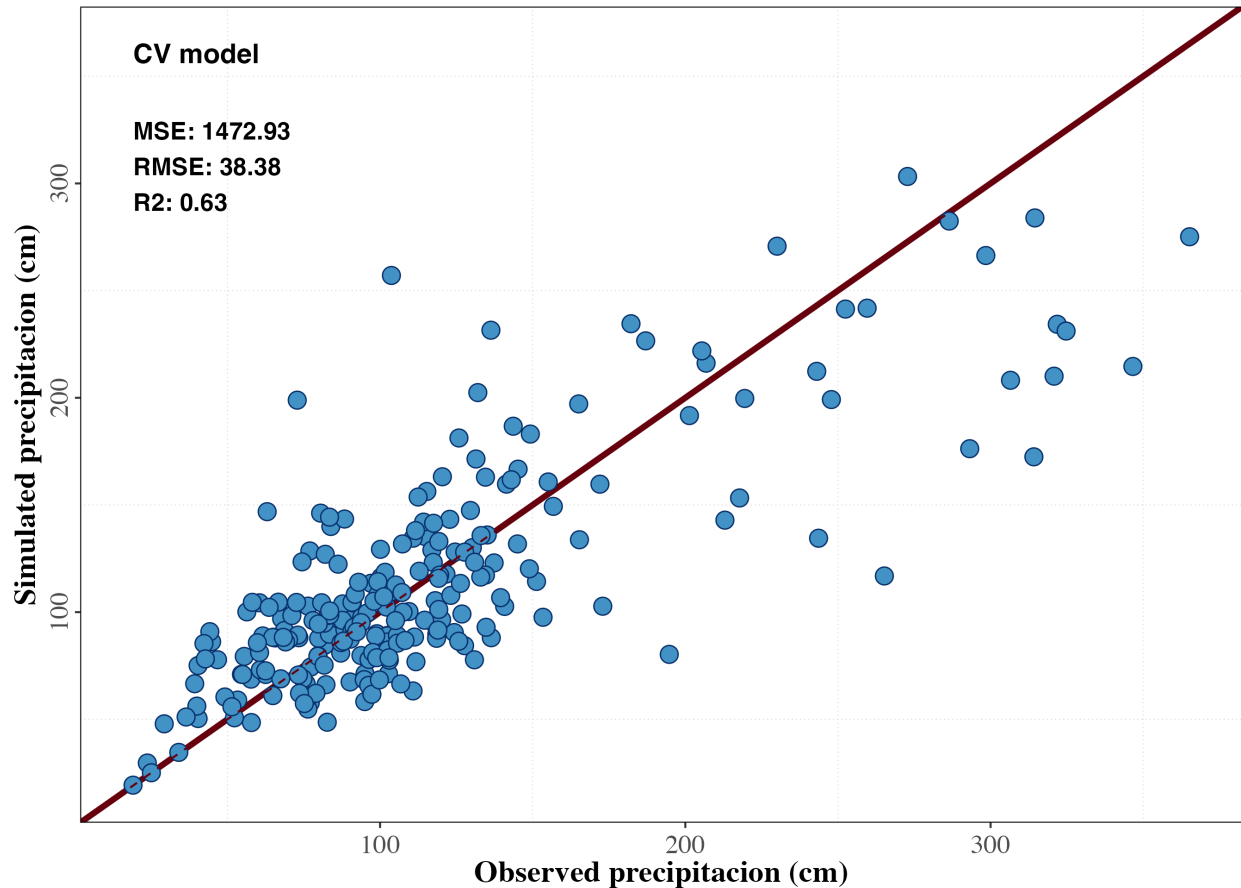


Figure 14: LOOCV Scatterplot.

Model robustness – cropping 10% of observations

The RMSE boxplot (left panel in Figure 15) shows moderate variation, indicating that while the GLM model's performance fluctuates slightly, the error remains within a stable range. Some outliers suggest that certain missing data points had a larger impact, possibly due to influential observations.

The correlation boxplot (right panel in Figure 15) shows a wider spread, with values centered around zero and some negative correlations. This suggests that the relationship between observed and predicted values weakens when 10% of the data is removed, indicating a higher sensitivity to missing data (compared to more flexible models such as GAM). However, the general trend remains stable, reinforcing moderate robustness.

The GLM and linear models exhibit similar robustness, with comparable RMSE and correlation distributions. Both models experience minor sensitivity to high-leverage observations, which may influence predictions. While GLM maintains stability, its sensitivity to certain missing data points could impact its generalization performance.

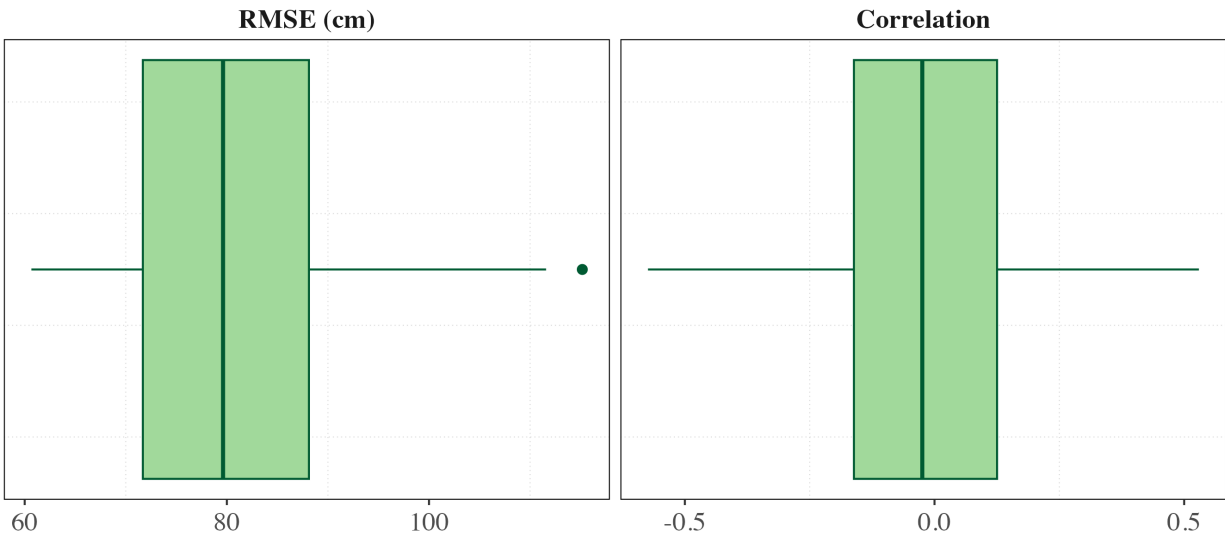


Figure 15: RMSE and correlation boxplots for model robustness analysis.

Spatial mapping

Figure 16 indicates a higher precipitation in northern and coastal regions, aligning well with climatic expectations. Lower precipitation levels are seen in central and western India, consistent with drier climatic conditions. The precipitation gradient appears smooth, suggesting that GLM successfully captures regional variability.

Figure 17 indicates higher standard errors in some inland regions, indicating greater uncertainty in precipitation estimates where fewer observations are available. Lower standard errors in coastal and well-sampled areas suggest more reliable predictions in those regions. Compared to the linear model, the GLM standard error distribution is similar, reinforcing that both models exhibit comparable predictive stability.

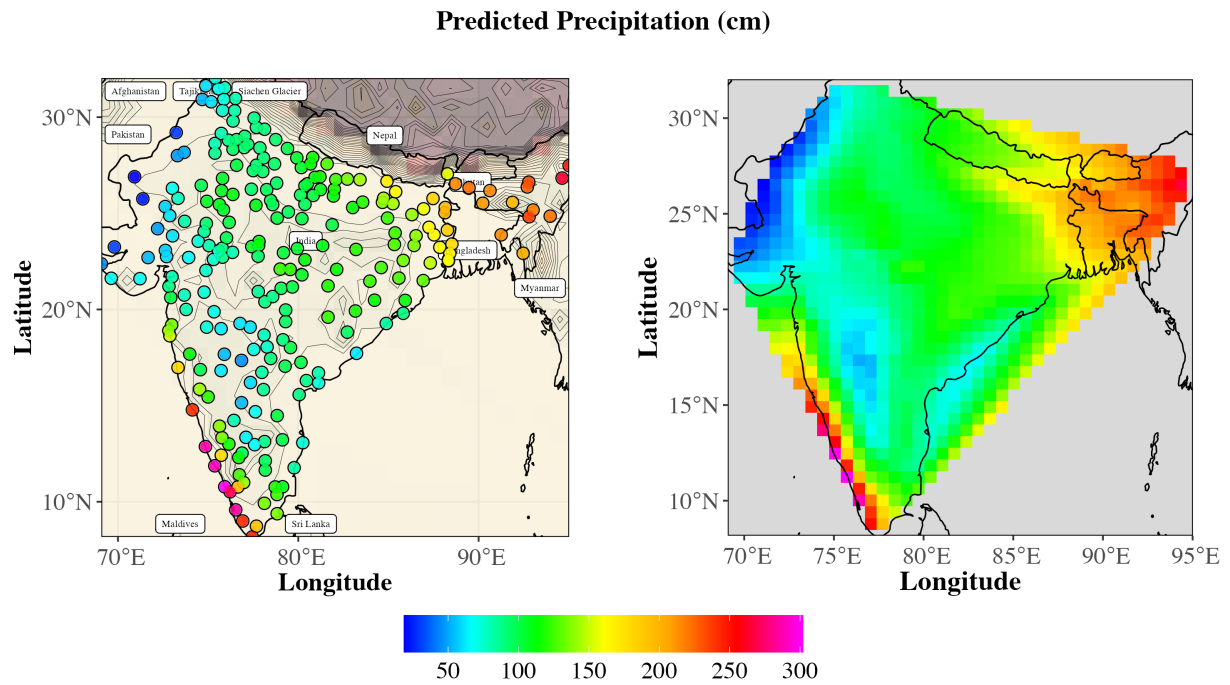


Figure 16: Predicted precipitation surface mapping under GLM.

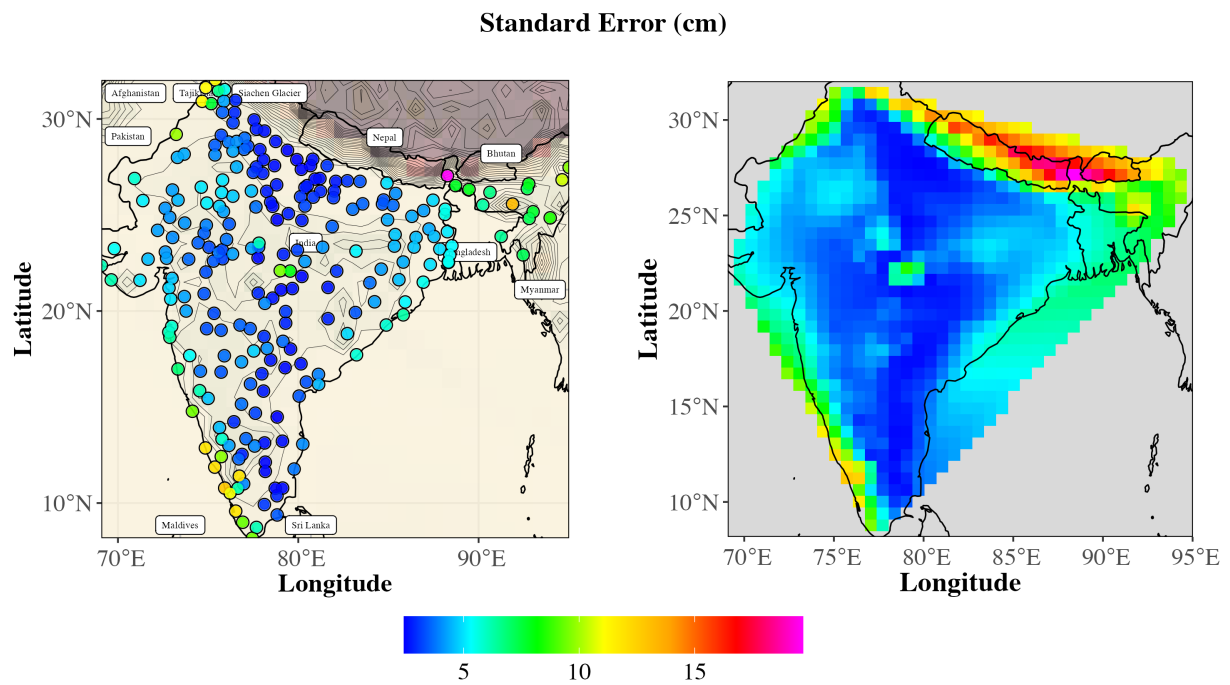


Figure 17: Standard error precipitation surface mapping under GLM.

Summary

The GLM model was evaluated for its predictive accuracy, residual behavior, robustness, and spatial performance. Overall, it demonstrated solid performance, with statistically significant predictors and mostly independent residuals. However, minor deviations in normality and heteroscedasticity indicate areas for potential refinement.

Model performance & diagnostics:

- The GLM fit was strong, with $R^2 = 0.67$ (Training) and $R^2 = 0.63$ (LOOCV).
- Residual analysis showed slight heteroscedasticity, suggesting that transforming the response variable or using a weighted GLM could improve reliability.

Cross-Validation & robustness:

- The RMSE increased slightly (38.38 cm in LOOCV), indicating minor generalization error.
- Predictive performance was nearly identical to the linear model, reinforcing that GLM does not provide significant improvements.
- The model remained stable when 10% of observations were removed, though some sensitivity to influential data points was observed.
- Further leverage diagnostics could identify points that disproportionately affect predictions.

Spatial mapping:

- Predicted precipitation trends align well with known climatic patterns, with higher precipitation in northern and coastal regions.
- Standard error distribution suggests that uncertainty is more influenced by data sparsity than the modeling approach itself.

GLM performs at the same level as the linear model, while the GLM is a viable approach, it does not provide significant gains over the standard linear model, making it a comparable but not superior choice for this dataset.

Problem 4: Underlying precipitation surface under Local GLM

Scatterplot of observed vs modeled precipitation

The scatterplot in Figure 18 exhibits a strong linear trend, showing that the Locfit model closely follows observed precipitation values with fewer deviations. However, variance increases slightly at higher precipitation values, which might indicate some overfitting in extreme cases.

The R^2 value of 0.88 represents a significant improvement over both the linear ($R^2 = 0.67$) and GLM ($R^2 = 0.67$) models, indicating that Locfit captures more variability in precipitation patterns. The RMSE is significantly lower (21.63 cm vs. ~ 36 -38 cm for previous models), suggesting higher prediction accuracy and reduced error variance.

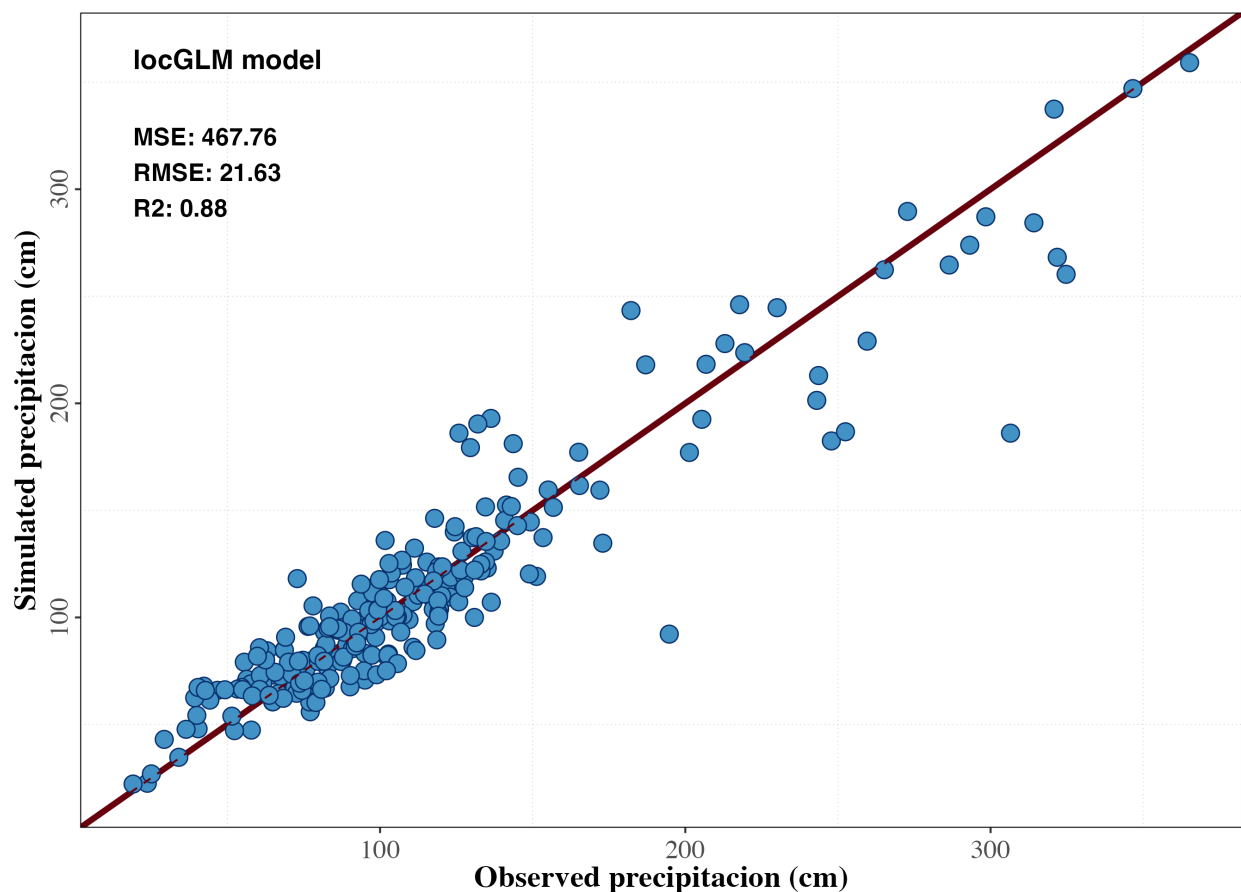


Figure 18: Scatterplot of observed vs modeled precipitation for local GLM.

Model diagnostics

To evaluate the reliability of the Locfit model, residual diagnostics were conducted (Figure 19).

- a) Normality: residuals follow a near-normal distribution, with minor deviations at the tails. The fit is better compared to the Linear and GLM models, indicating improved residual behavior.
- b) Histogram: residuals exhibit a symmetric, bell-shaped distribution, confirming approximate normality. Unlike previous models, fewer extreme residuals are present, further supporting improved fit quality.
- c) Independence: ACF plot shows minimal autocorrelation, confirming that the residuals are mostly independent. This suggests that the Locfit model does not suffer from significant spatial correlation issues.
- d) Homoscedasticity: residuals appear evenly spread, but minor variance increases at higher precipitation levels. Compared to the Linear and GLM models, Locfit exhibits less pronounced heteroscedasticity, meaning error variance is more stable.

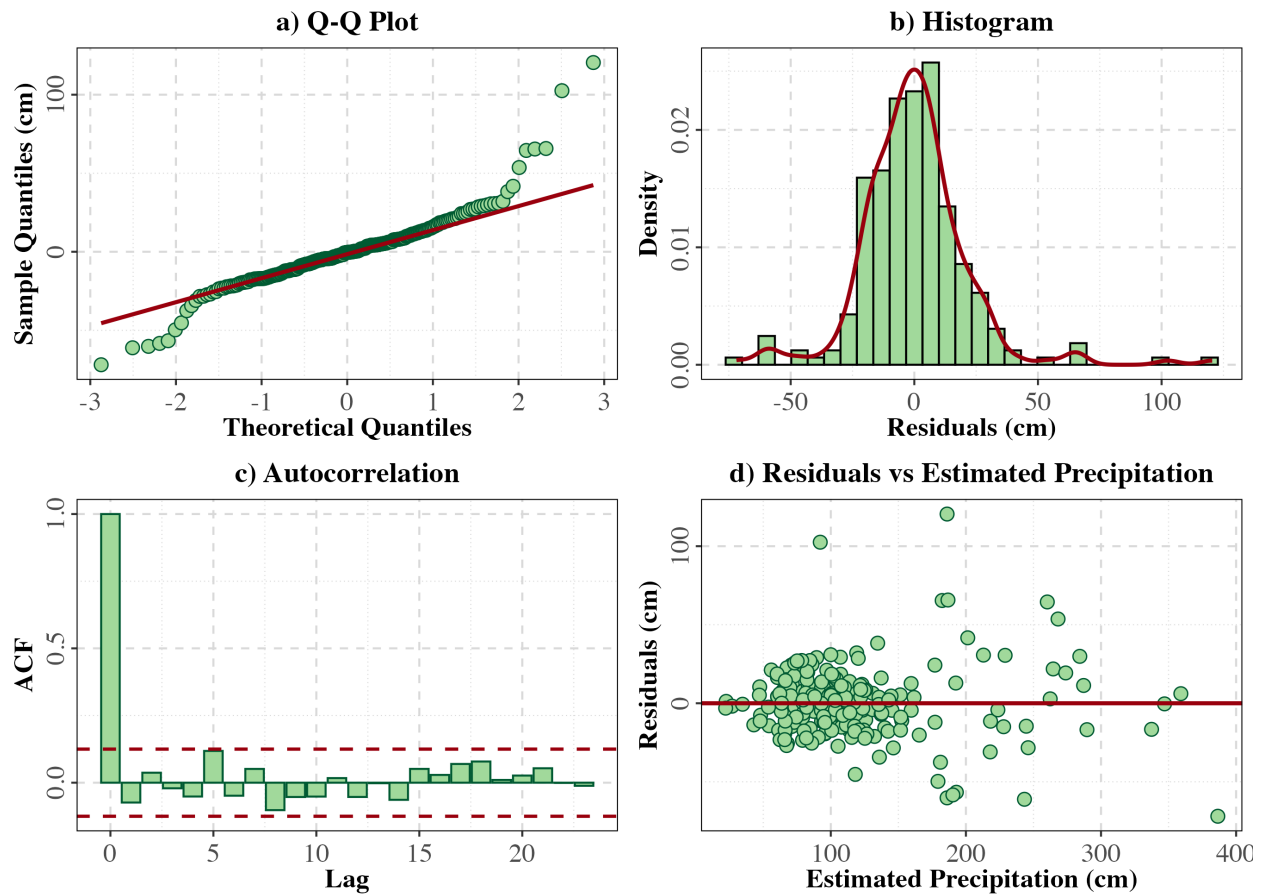


Figure 19: Model diagnostics for local GLM regression.

Cross-Validation Analysis

The scatterplot in Figure 20 shows a weaker correlation compared to training, further confirming that Locfit struggles with generalization. Also, the variance increases at higher precipitation values, highlighting instability in extreme cases.

Drastic drop in R^2 (from 0.88 in training to 0.23 in LOOCV) suggests that Locfit overfits the training data and does not generalize well. The RMSE increased significantly (from 21.63 cm to 55.29 cm), indicating that the model's predictions become much less accurate on unseen data.

Unlike Linear and GLM models, which maintain stable performance in LOOCV, Locfit collapses in predictive power, reinforcing overfitting issues. This suggests that Locfit captures intricate patterns in training data, but fails to generalize to new observations.

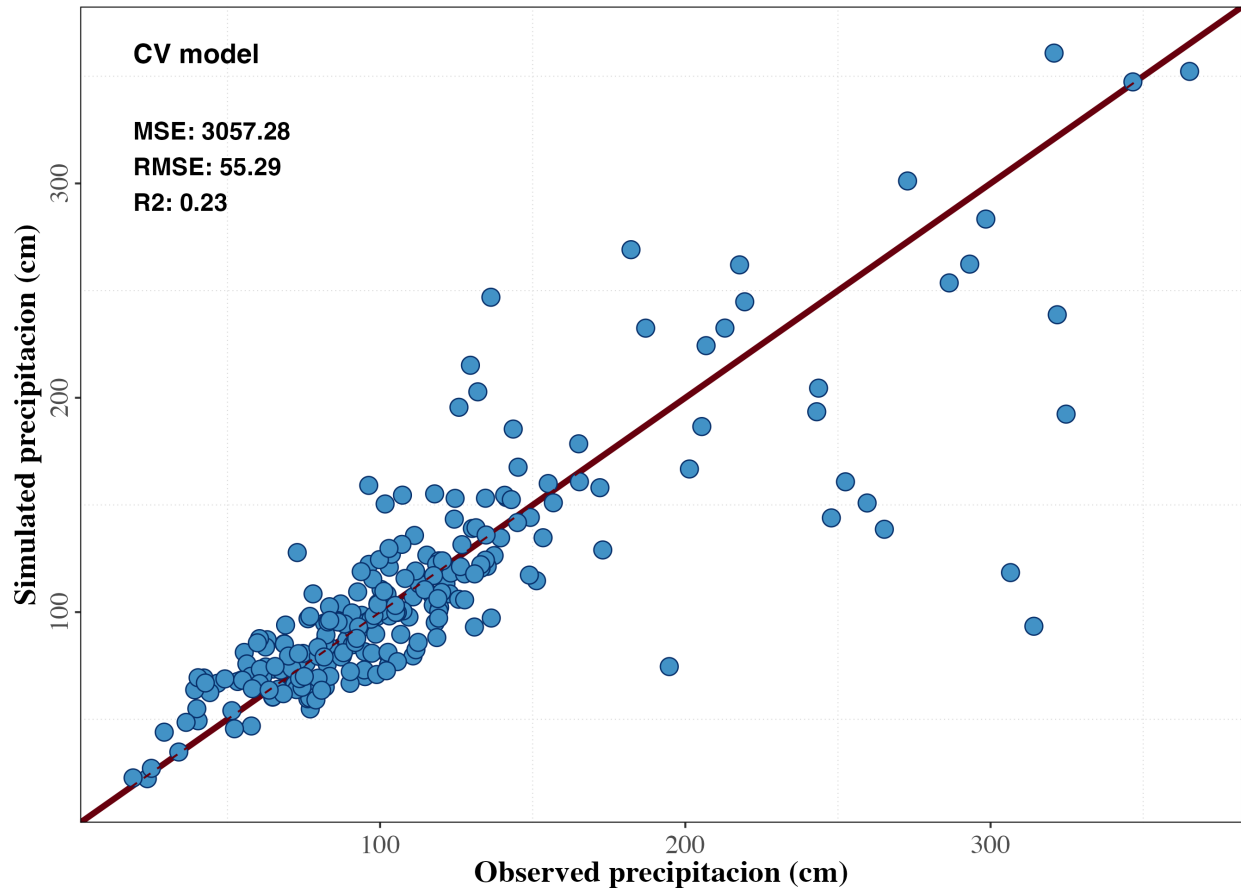


Figure 20: LOOCV Scatterplot.

Model robustness – cropping 10% of observations

The RMSE boxplot shows high variability, with some extreme outliers, indicating that the model's performance is sensitive to missing data (left panel Figure 21). Some extreme RMSE values suggest that certain removed data points had a significant impact, possibly due to overfitting in localized regions. The correlation remains relatively high, but the presence of extreme low values suggests that Locfit may struggle with certain missing data patterns (right panel Figure 21). The wide spread of correlation values indicates greater instability compared to Linear and GLM models.

Unlike Linear and GLM models, which showed consistent robustness, Locfit appears to be more sensitive to missing observations. The higher variation in RMSE and correlation suggests that Locfit may be overfitting certain data points, making it less stable when data is removed. This indicates that while Locfit improves fit quality, its robustness could be enhanced by regularization or alternative smoothing parameters.

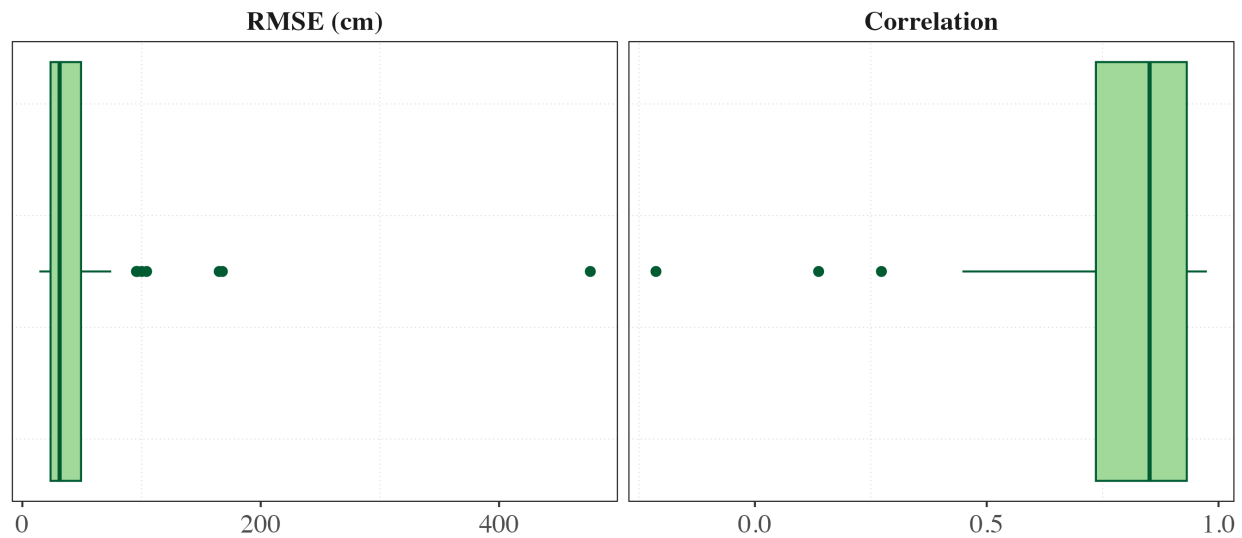


Figure 21: RMSE and correlation boxplots for model robustness analysis.

Spatial mapping

Figure 22 show higher precipitation in northern and coastal regions, aligning with expected climatic patterns. Lower precipitation levels appear in central and western India, consistent with known drier conditions. The precipitation gradient is smooth, indicating that Locfit captures non-linear trends more effectively than linear models.

Figure 23 shows that standard errors are significantly lower overall compared to the Linear and GLM models, suggesting that Locfit provides more stable predictions. Lower standard errors in coastal and well-sampled regions reinforce that Locfit produces more confident estimates in data-dense regions. However, some higher standard errors persist in inland regions, likely due to fewer observations and localized variability.

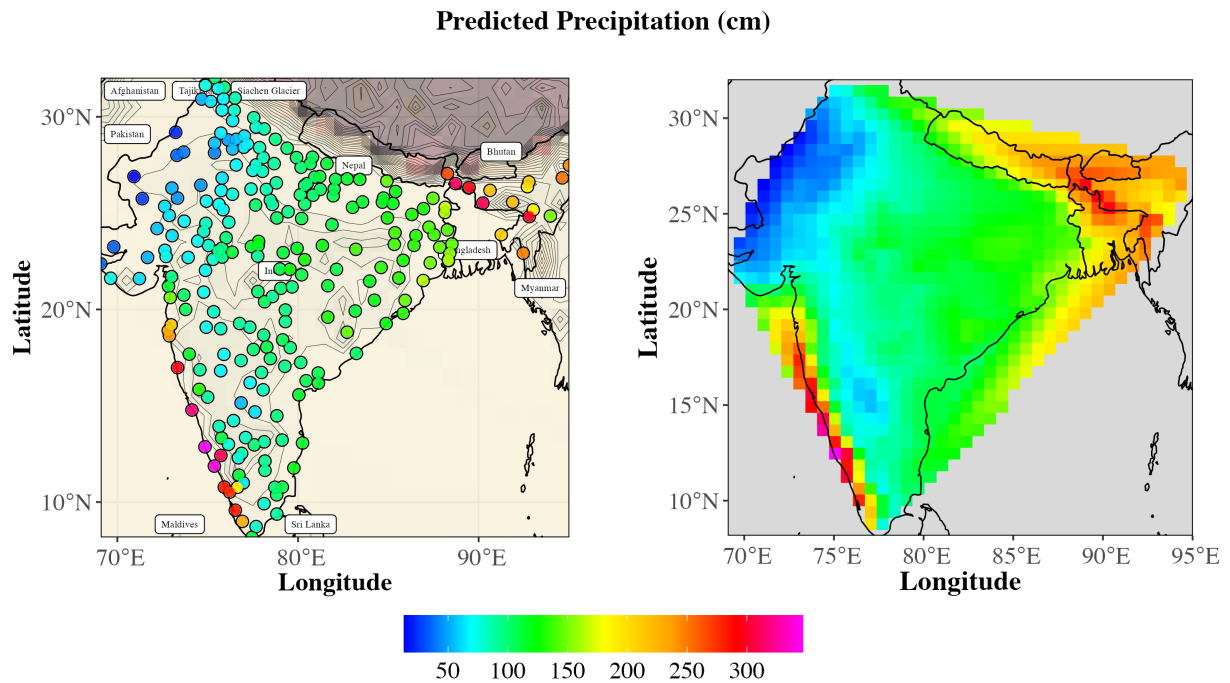


Figure 22: Predicted precipitation surface mapping under Local GLM.

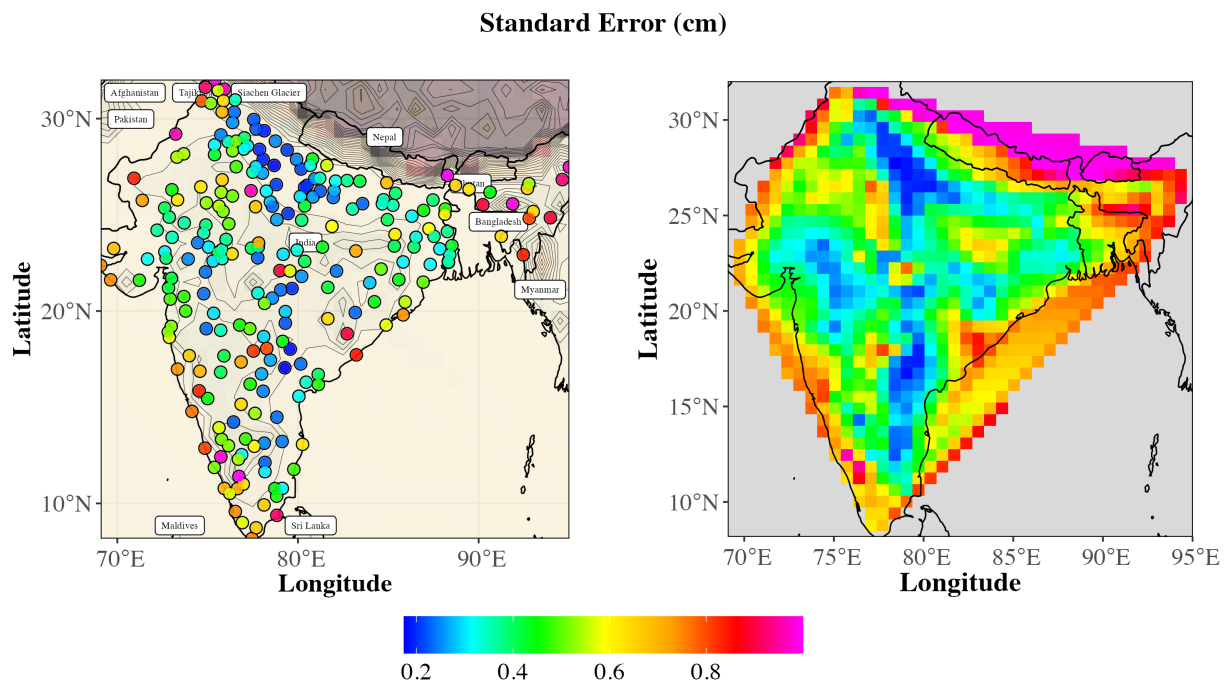


Figure 23: Standard error precipitation surface mapping under Local GLM.

Summary

Locfit significantly outperforms both linear and GLM models, reducing error while improving correlation. The results suggest that Locfit provides a better fit for precipitation modeling, likely due to its ability to capture non-linear relationships.

Model performance & diagnostics:

- Highest predictive accuracy with $R^2 = 0.88$, significantly better than the Linear ($R^2 = 0.67$) and GLM ($R^2 = 0.67$) models.
- RMSE is much lower (21.63 cm vs. ~ 36 cm for previous models), suggesting more precise predictions.
- Best residual behavior, outperforming other models in normality, independence, and homoscedasticity.
- Minor variance increases at high precipitation levels suggest that a transformation or weighted approach could further enhance fit.

Cross-Validation & robustness:

- Weak out-of-sample performance (LOOCV $R^2 = 0.23$, RMSE = 55.29 cm) suggests Locfit may overfit the data.
- Higher RMSE and correlation variability in robustness tests indicate that Locfit is more sensitive to missing data than Linear or GLM models.
- To improve generalization, regularization techniques or a hybrid approach could be explored.

Spatial mapping:

- Locfit significantly improves precipitation predictions and exhibits lower standard errors, suggesting better adaptability to spatial variability.
- Compared to previous models, Locfit better captures precipitation gradients while reducing uncertainty, making it the best-performing model in spatial mapping.

Problem 5: Underlying precipitation surface under GAM

Scatterplot of observed vs modeled precipitation

The Generalized Additive Model (GAM) was applied to predict precipitation, offering a balance between flexibility (Locfit) and generalization (GLM/Linear models). The scatterplot in Figure 24 follows a strong linear trend, with lower variance in residuals compared to Locfit. Unlike this model, variance remains more stable at higher precipitation values, suggesting a better generalization.

The R^2 value of 0.83 is significantly higher than that of the Linear and GLM models ($R^2 = 0.67$) but slightly lower than Locfit ($R^2 = 0.88$). RMSE of 25.27 cm suggests improved predictive accuracy compared to Linear (~ 36 cm) and GLM (~ 36 cm) models, but it is slightly higher than Locfit (21.63 cm).

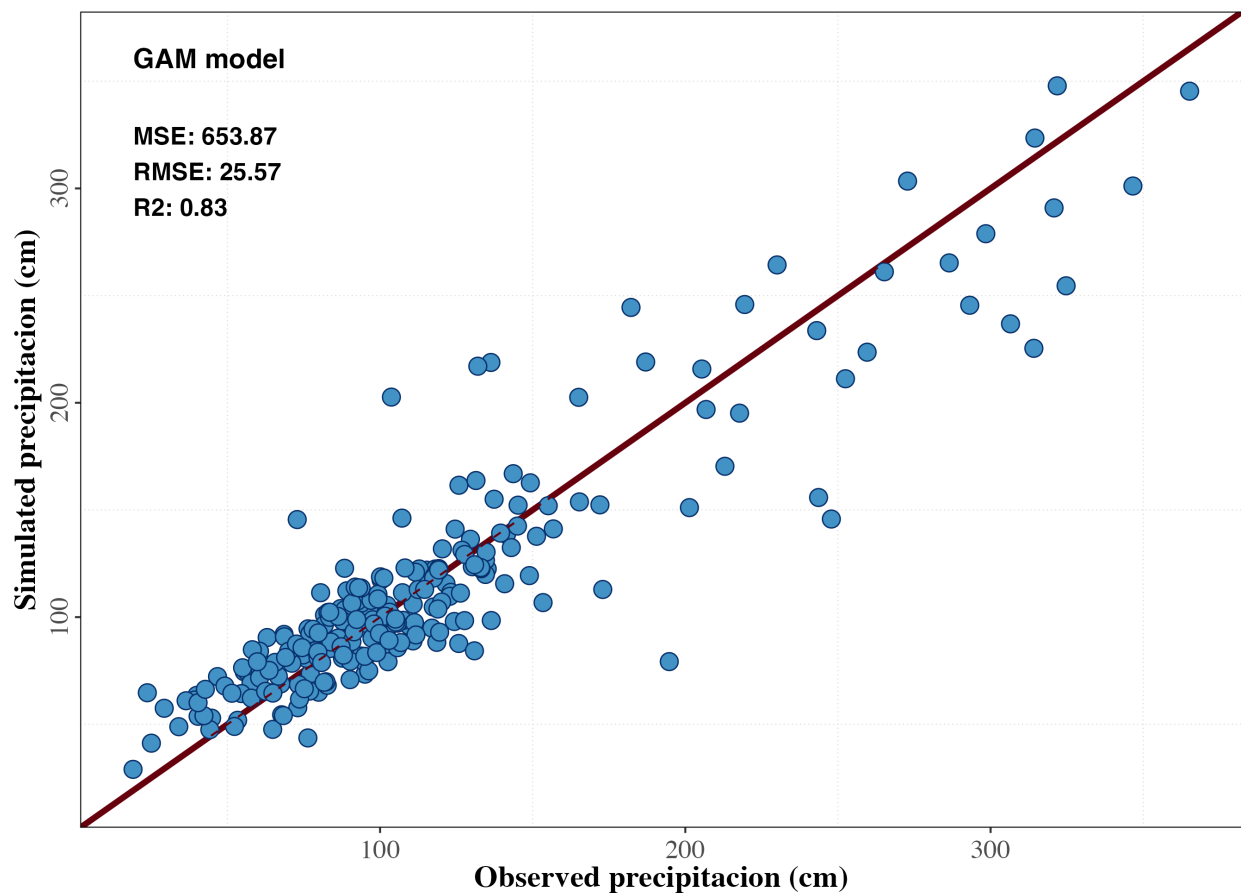


Figure 24: Scatterplot of observed vs modeled precipitation for GAM.

Model diagnostics

ANOVA significance indicates that all smooth terms (longitude, lld, ll2d, ed) are highly significant ($p < 2e-16$), confirming that the non-linear relationships captured by GAM contribute significantly to precipitation predictions. Figure 25 provides essential model validation checks:

- a) Normality: residuals follow a near-normal distribution, though some deviations exist at the tails. Compared to Locfit, GAM has fewer extreme deviations, suggesting a better overall fit.
- b) Histogram: residual distribution is approximately normal, with some minor skewness. The peak of the distribution is sharper than Locfit's, indicating that some residuals are clustered closely around zero.
- c) Independence: ACF plot shows minimal autocorrelation, suggesting that residuals are mostly independent and confirming that GAM does not suffer from major correlation issues.
- d) Homoscedasticity: the spread of residuals appears uniform, but the slight variance increases at higher precipitation values. GAM performs better than Locfit in controlling heteroscedasticity, but some variance issues persist.

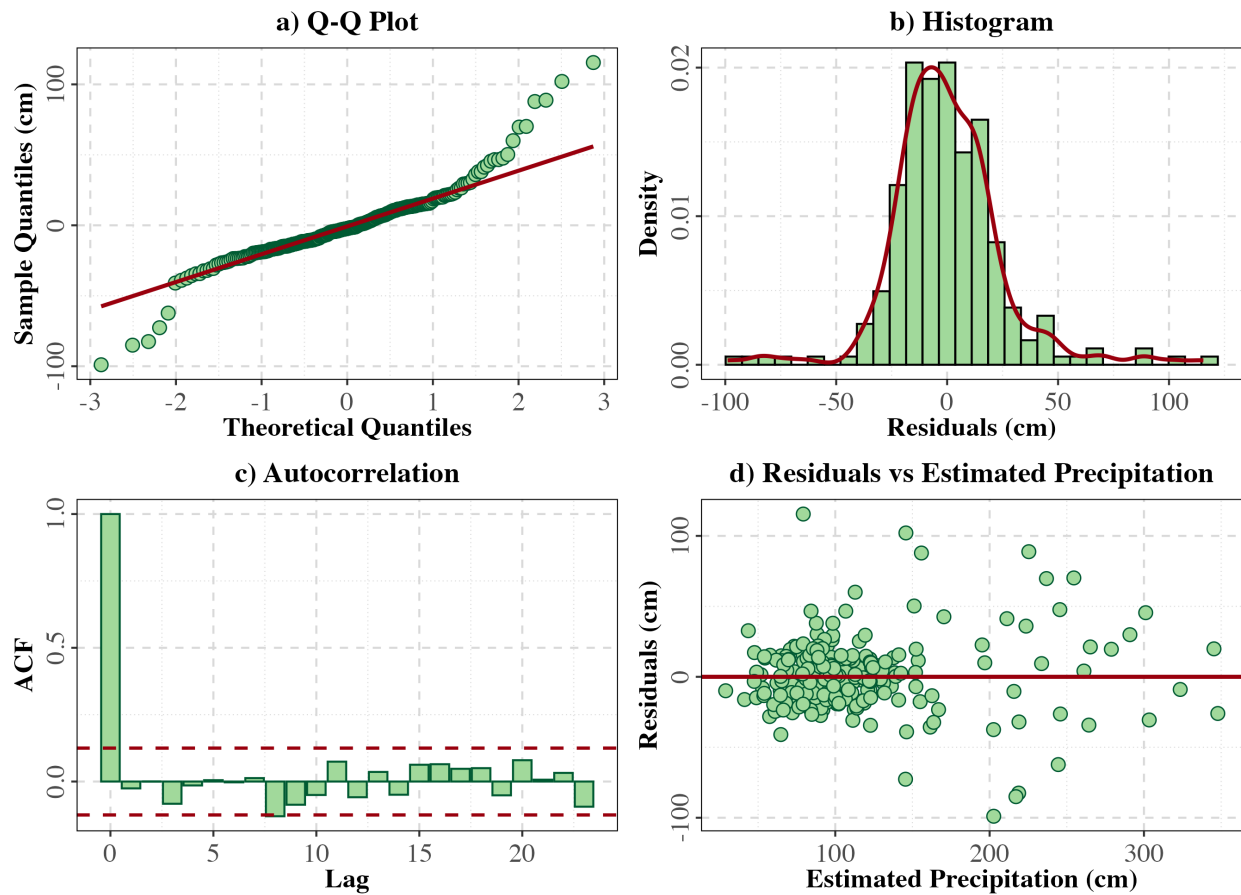


Figure 25: Model diagnostics for GAM regression.

Cross-Validation Analysis

The scatterplot in Figure 26 maintains a strong linear trend, showing that the GAM model generalizes well while still capturing non-linear precipitation patterns. Unlike Locfit, which showed weak LOOCV performance ($R^2 = 0.23$, $RMSE = 55.29$ cm), GAM performs significantly better, proving it is less prone to overfitting.

The R^2 value (0.70) is higher than that of the Linear ($R^2 = 0.63$) and GLM ($R^2 = 0.63$) models, indicating better generalization. $RMSE$ of 34.52 cm is lower than both Linear (36.24 cm) and GLM (36.24 cm) models, meaning GAM maintains stronger predictive accuracy in LOOCV.

GAM outperforms Linear and GLM models in generalization, making it the best-performing model for predictive stability. Locfit, despite its strong in-sample performance, fails to generalize well due to its sensitivity to specific data points. GAM strikes the best balance between flexibility and generalization, making it a strong candidate for precipitation prediction.

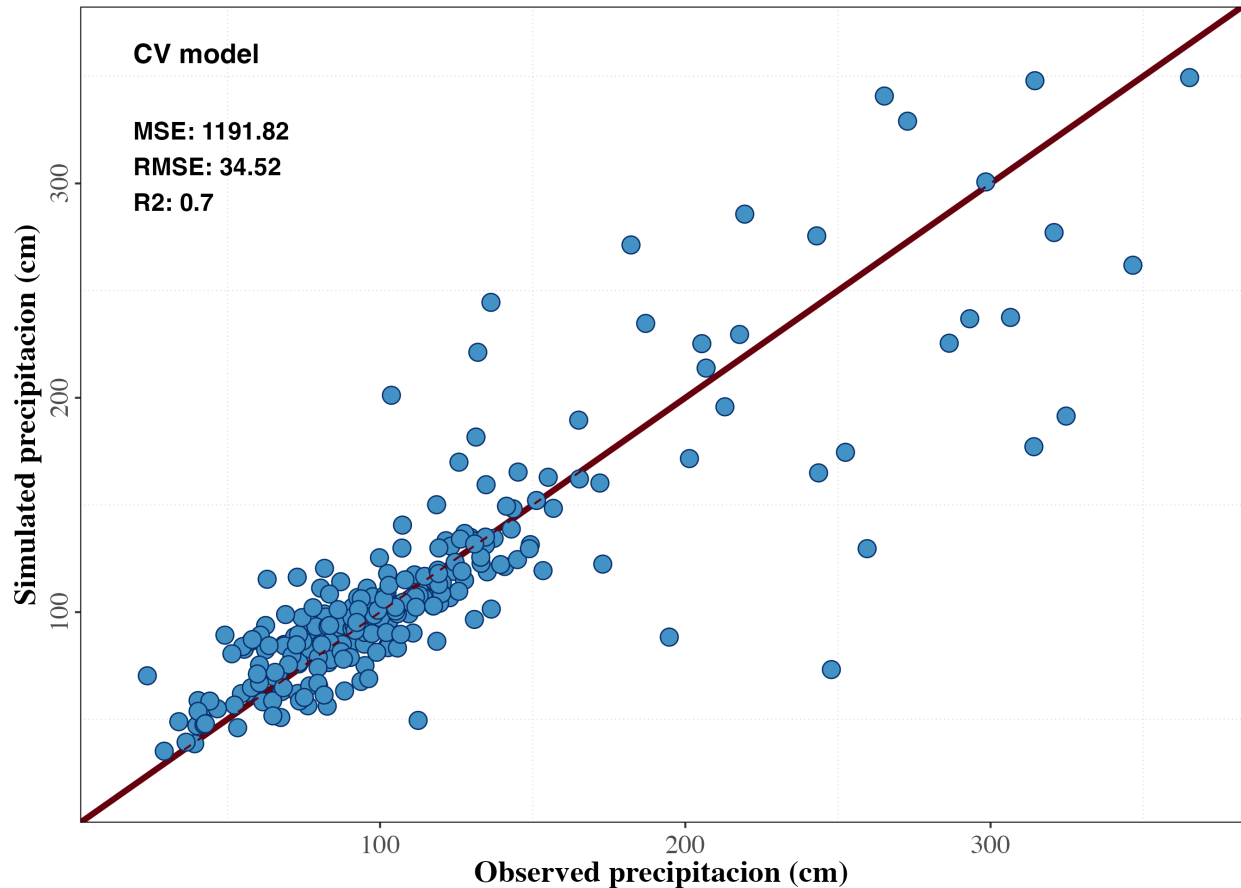


Figure 26: LOOCV Scatterplot.

Model robustness – cropping 10% of observations

The RMSE boxplot shows moderate variation (left panel Figure 27), indicating that GAM remains stable even when 10% of the data are removed. Some outliers suggest that certain missing observations have a larger impact on the model's performance, but overall, GAM maintains relatively low RMSE values. The correlation remains stable (right panel Figure 27), reinforcing that GAM is resilient to missing data, although not necessarily better than Linear or GLM models. Compared to Locfit, which exhibited high sensitivity to data removal, GAM demonstrates significantly better robustness.

GAM proves to be more robust than Locfit, which suffered from extreme RMSE and correlation variability. While GAM performs similarly to Linear and GLM models in terms of robustness, it maintains lower RMSE values, ensuring better predictive performance even when data points are missing. However, its slight sensitivity to specific missing observations suggests that further regularization techniques or alternative smoothing methods may still enhance its stability.

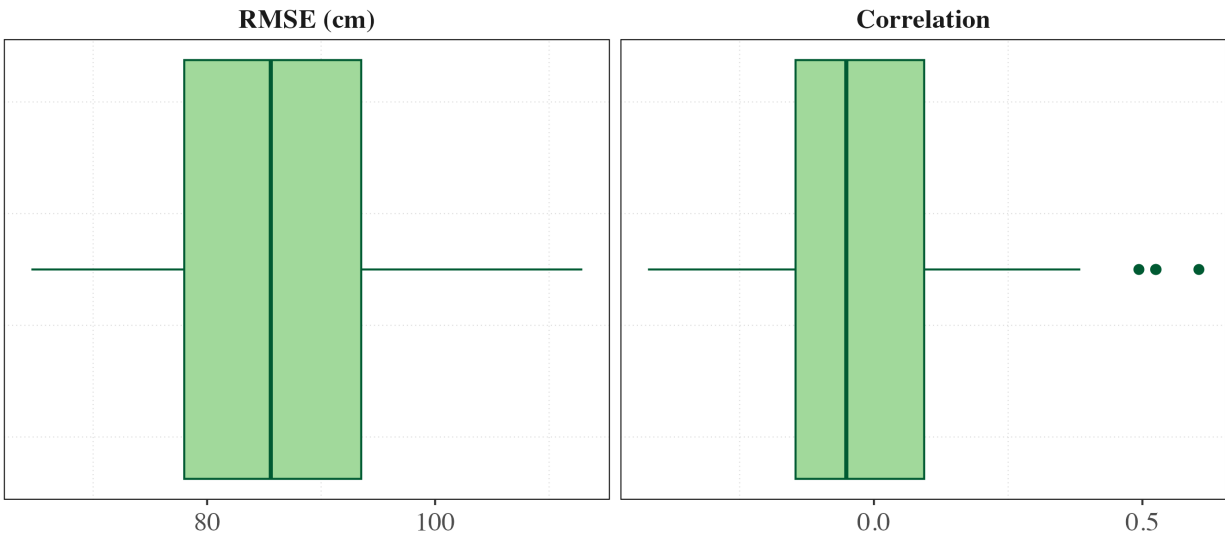


Figure 27: RMSE and correlation boxplots for model robustness analysis.

Spatial mapping

Figure 28 shows higher precipitation in northern and coastal regions, closely aligning with expected climatic patterns. Lower levels of precipitation appear in central and western India, which is consistent with drier conditions. Compared to linear and GLM models, GAM provides a smoother precipitation gradient, effectively capturing non-linear climatic variations.

Figure 29 shows lower standard errors in coastal and well-sampled areas, suggesting higher model confidence in these regions. Higher standard errors in inland regions indicate that precipitation uncertainty remains an issue in data-sparse regions. Compared to Locfit, GAM maintains lower standard errors in most areas, suggesting better generalization and prediction stability.

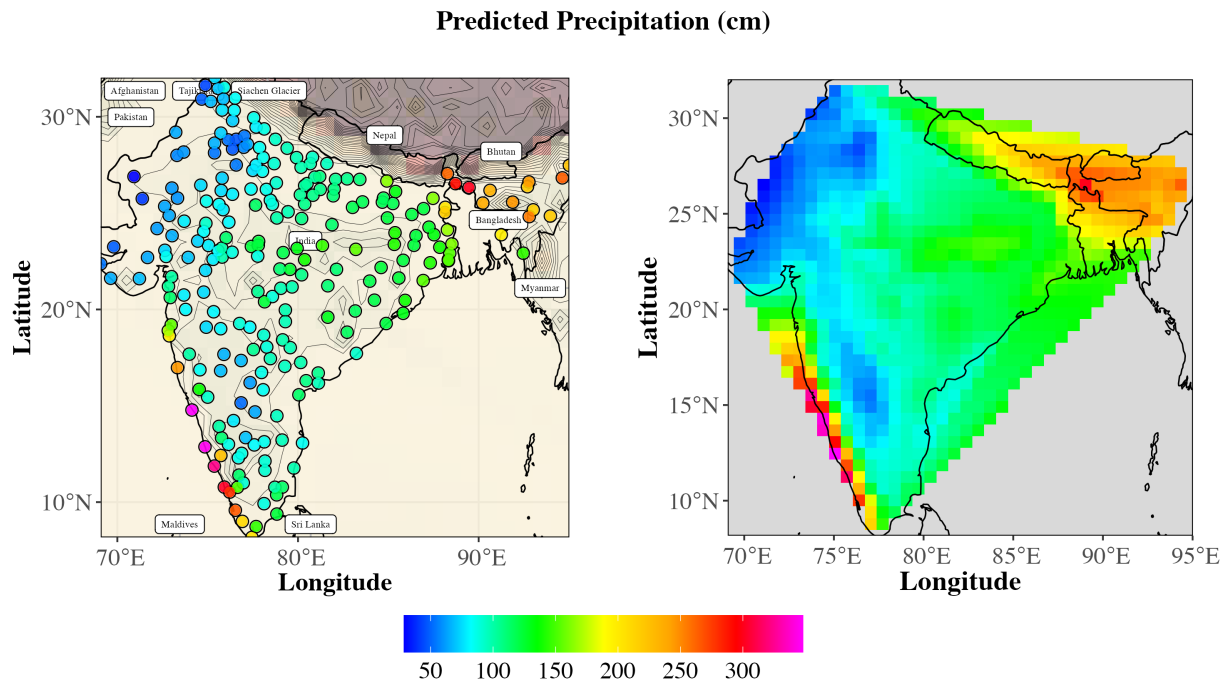


Figure 28: Predicted precipitation surface mapping under GAM.

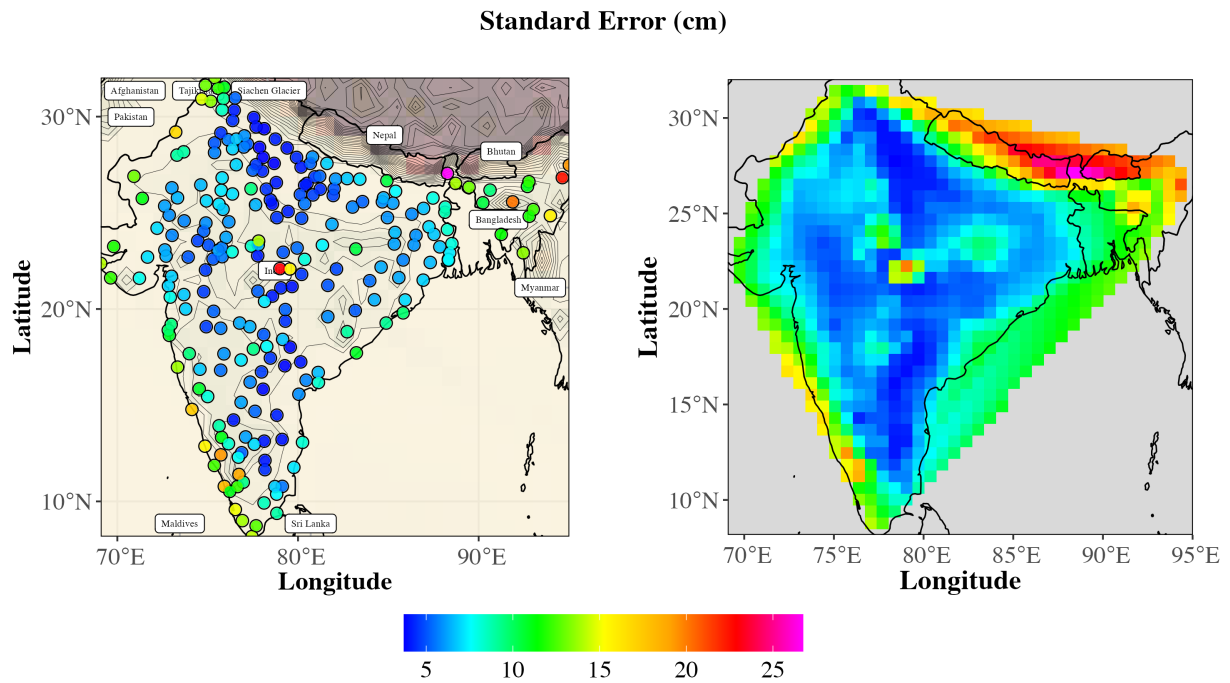


Figure 29: Standard error precipitation surface mapping under GAM.

Summary

The Generalized Additive Model (GAM) outperforms Linear and GLM models, providing better predictive accuracy and generalization. While Locfit remains the most accurate in-sample model ($R^2 = 0.88$), GAM generalizes better (LOOCV $R^2 = 0.70$), making it a strong alternative to Locfit for precipitation modeling.

Model performance & diagnostics:

- GAM achieves high predictive accuracy ($R^2 = 0.83$), significantly improving over Linear and GLM models ($R^2 = 0.67$).
- RMSE is lower (25.27 cm) compared to Linear/GLM (~ 36 cm), but slightly higher than Locfit (21.63 cm).
- Residual analysis confirms good normality, independence, and improved variance control, making GAM more stable than Locfit.

Cross-Validation & robustness:

- GAM maintains strong generalization performance (LOOCV $R^2 = 0.70$, RMSE = 34.52 cm), better than Linear ($R^2 = 0.63$) and Locfit ($R^2 = 0.23$, RMSE = 55.29 cm).
- More robust than Locfit, which showed high RMSE variation when 10% of data was removed.
- GAM exhibits moderate sensitivity to missing data, but performs similarly to Linear/GLM models in robustness while maintaining lower RMSE.

Spatial mapping:

- GAM provides smooth and accurate precipitation predictions, capturing non-linear climatic variations better than Linear and GLM models.
- Lower standard errors than Locfit, making it more stable in spatial predictions.
- Inland standard errors remain slightly high, but GAM still balances accuracy and generalization better than all other models.

Problem 6: Underlying precipitation surface under Hierarchical Spatial Model (HSM) and Kriging

Scatterplot of observed vs modeled precipitation

The binned empirical variogram (Figure 30) shows an increasing trend, where the closer points exhibit lower variance, while the farther apart points show greater variability. The fitted variogram model (dashed red line) indicates a moderate error level (RMSE = 124.58, $R^2 = 0.81$).

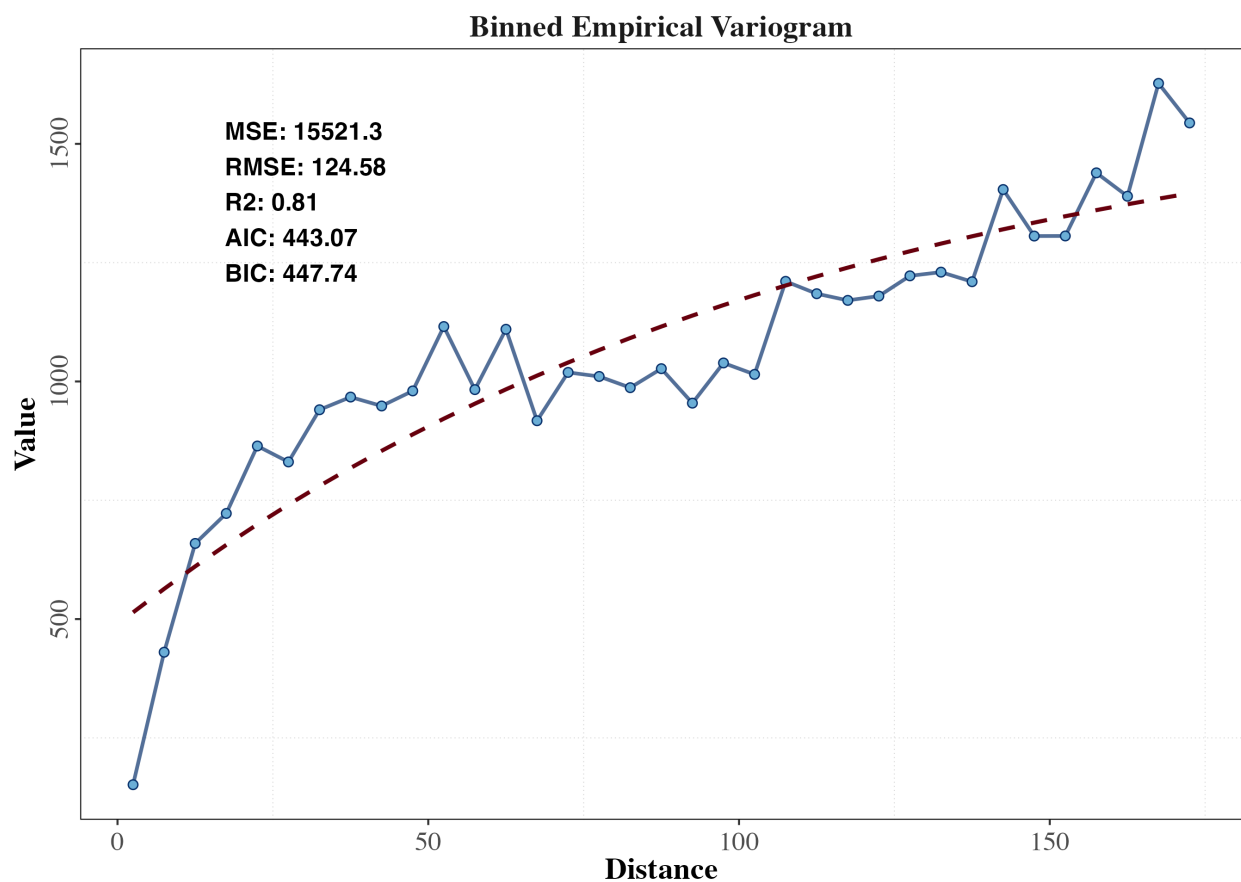


Figure 30: Binned Empirical Variogram.

Figure 31 aligns almost perfectly with the 1:1 line, confirming that HSM + Kriging produces highly accurate precipitation predictions. The metrics indicate significantly lower errors compared to the previous models (MSE = 11.68, RMSE = 3.42). However, $R^2 = 1.0$ could hint at potential overfitting.

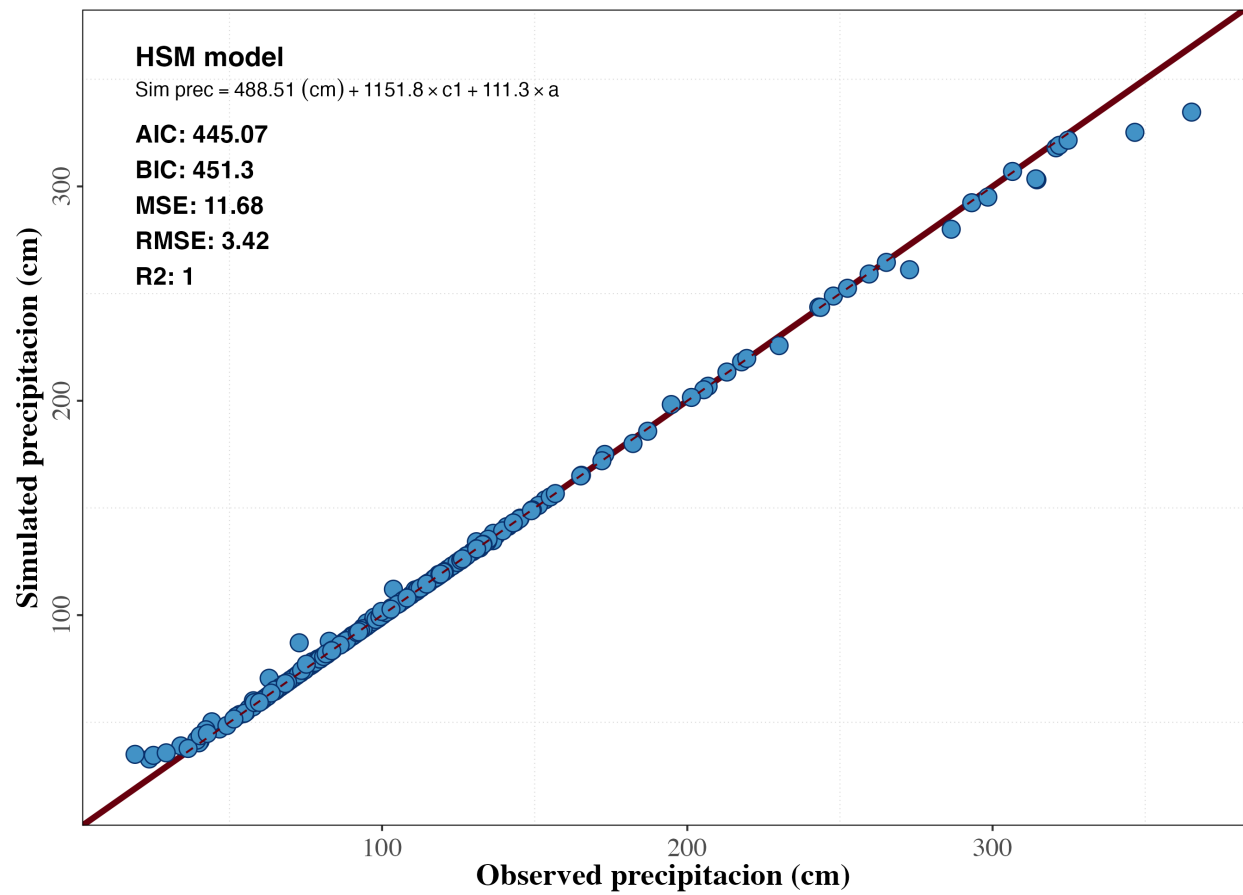


Figure 31: Scatterplot of observed vs modeled precipitation for HSM + Kriging. **Note:** The model follows the function Simulated Precip = $c_0 + c_1 \times (1 - \exp(-x/a))$, where $c_0 = 488.51$, $c_1 = 1151.8$ and $a = 111.3$.

Model diagnostics

Diagnostic plots (Figure 32) provide insight into the residual structure of the HSM + Kriging model.

- a) Normality: Q-Q plot reveals deviations at both tails, suggesting heavy-tailed residuals. The middle portion aligns well, meaning that most residuals are normally distributed, but the presence of outliers at the extremes indicates that some observations are poorly captured by the model, potentially due to localized spatial effects.
- b) Histogram: shows a strong peak at zero, suggesting that most residuals are close to zero, confirming high accuracy. However, narrow spikes suggest the model might be overfitting, capturing the training data too precisely. Residuals exhibit a skewed distribution, indicating potential non-linearity or variance issues that may require transformation.
- c) Independence: ACF plot reveals significant autocorrelation at lag 1, meaning spatial correlation is not fully accounted for in the residuals. The presence of periodic structure suggests some underlying spatial dependencies remain unmodeled.
- d) Homoscedasticity: the residual plot shows a non-uniform spread, with higher variance at extreme precipitation values. There is a slight funnel shape, meaning heteroscedasticity is present, i.e., errors increase as precipitation values rise. *For future work: adjust spatial correlation structure to refine predictions in high-variance regions.*

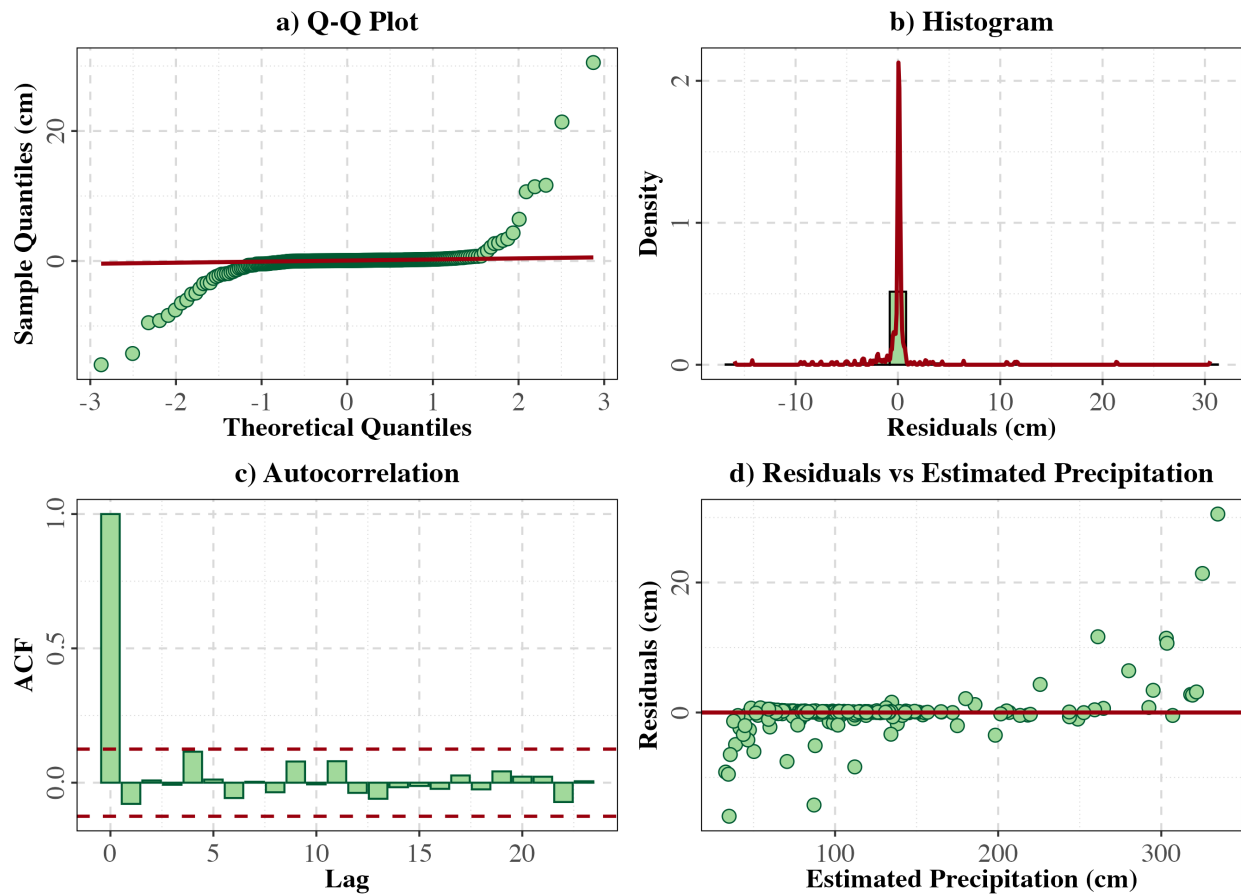


Figure 32: Model diagnostics for HSM + Kriging regression.

Cross-Validation Analysis

The LOOCV scatterplot (Figure 33) shows the predicted precipitation values are nearly constant, failing to capture the variation in observed values. Figure 33 indicates that the model collapses the predictions into a narrow range, ignoring the spread in precipitation values. According to the metrics, a negative R^2 (-0.17) indicates that the model performs worse than simply using the mean of the data. The RMSE of 67.98 cm is significantly higher than previous models, meaning predictions deviate substantially from the observed values. An MSE of 4621.1 also suggests that residual errors remain large even after fitting.

Finally, i) the lack of alignment along the 1:1 line confirms that the model fails to generalize, ii) most points cluster around a single predicted value (~ 100 cm), suggesting that spatial variability is not properly captured, and iii) the model likely overfits the training data but performs poorly on validation samples.

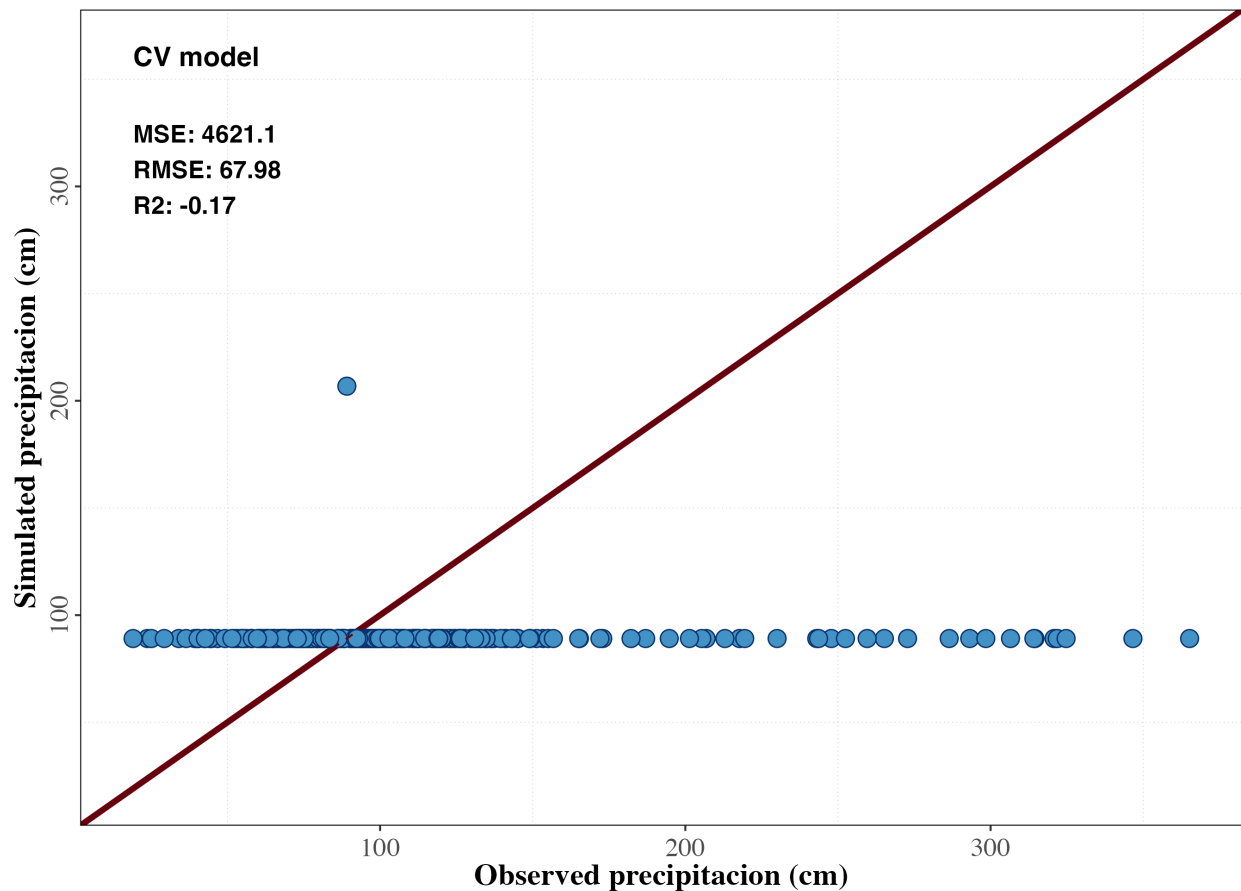


Figure 33: LOOCV Scatterplot.

Model robustness – cropping 10% of observations

The RMSE values in the left panel of Figure 34 remain consistently high (~ 80 – 100 cm), indicating that the model's predictions are not stable when data points are removed. The median correlation is close to 0, suggesting that the model fails to maintain predictive structure. Furthermore, the high correlation variability shows that its performance fluctuates depending on which observations are removed.

The HSM + Kriging model is highly sensitive to data removal, indicating poor robustness. Unlike LM, GLM, and GAM models, which maintain reasonable RMSE and correlation stability, this method collapses when observations are missing.

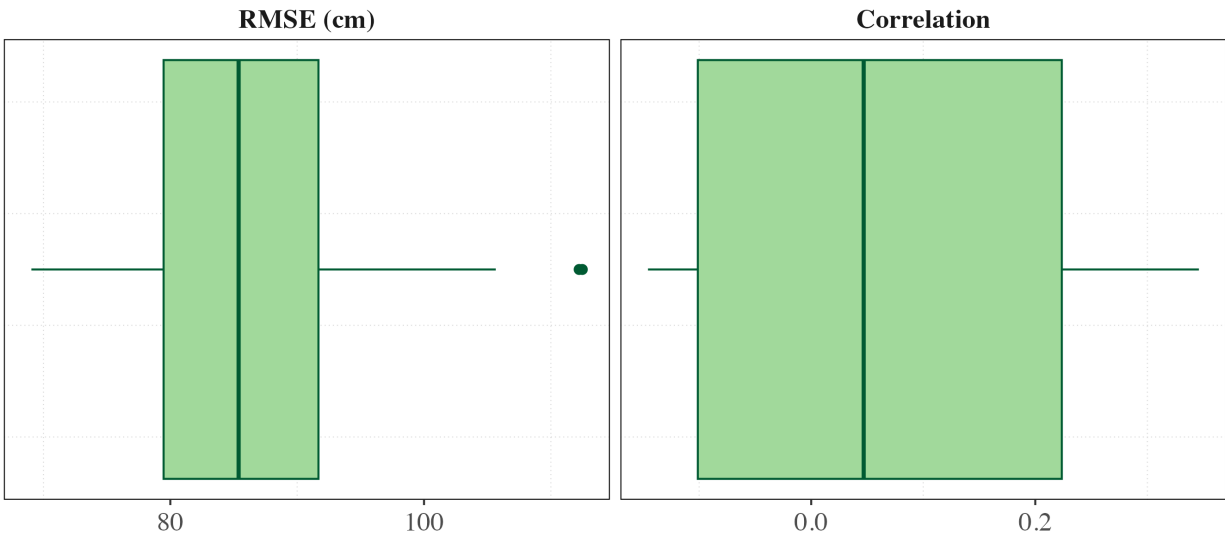


Figure 34: RMSE and correlation boxplots for model robustness analysis.

Spatial mapping

Figure 35 shows that a higher precipitation is observed in the northern and coastal regions, which is in line with the expected climatic patterns. Central and western India exhibit lower precipitation, consistent with the drier climatic conditions of these regions. Compared to previous models (Linear, GLM, Locfit, and GAM), the HSM + Kriging approach provides a smoother transition in precipitation values and captures spatial dependencies more effectively.

Standard errors (Figure 36) are lower in coastal and well sampled regions, indicating higher model confidence in these areas. Higher standard errors are observed in central and inland regions, suggesting greater uncertainty in data-sparse areas. Compared to GAM and Locfit, the HSM + Kriging model reduces standard errors in several inland locations, demonstrating improved predictive stability and spatial interpolation capability.

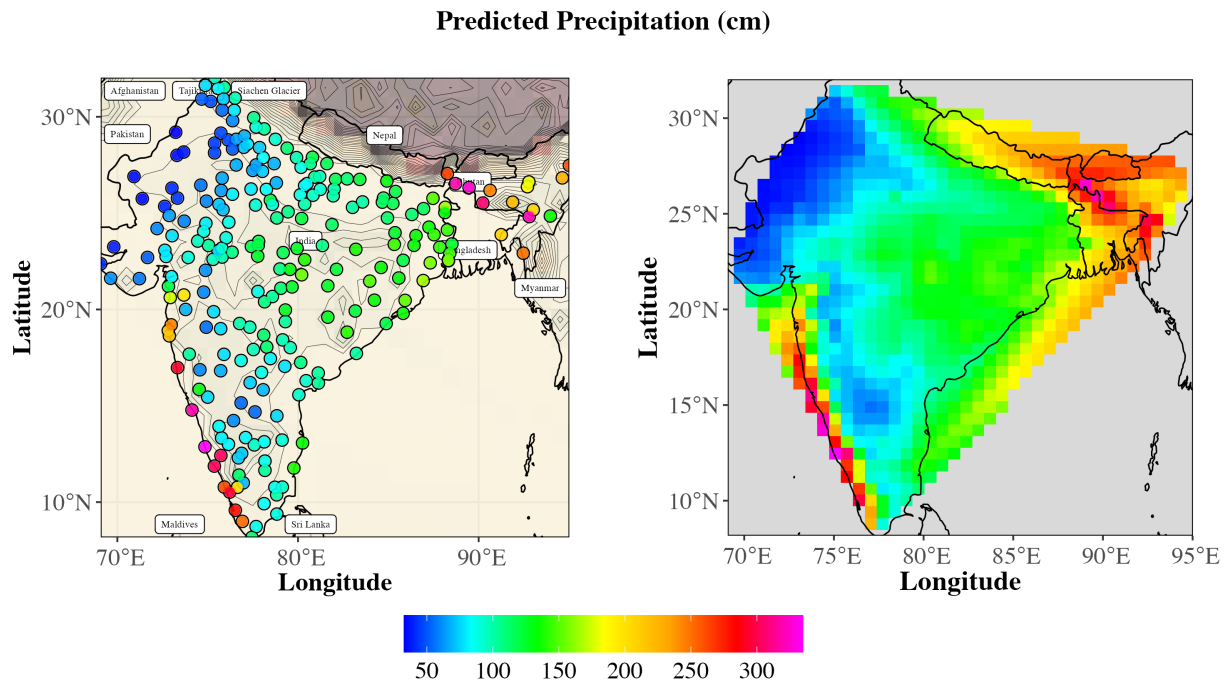


Figure 35: Predicted precipitation surface mapping under HSM + Kriging.

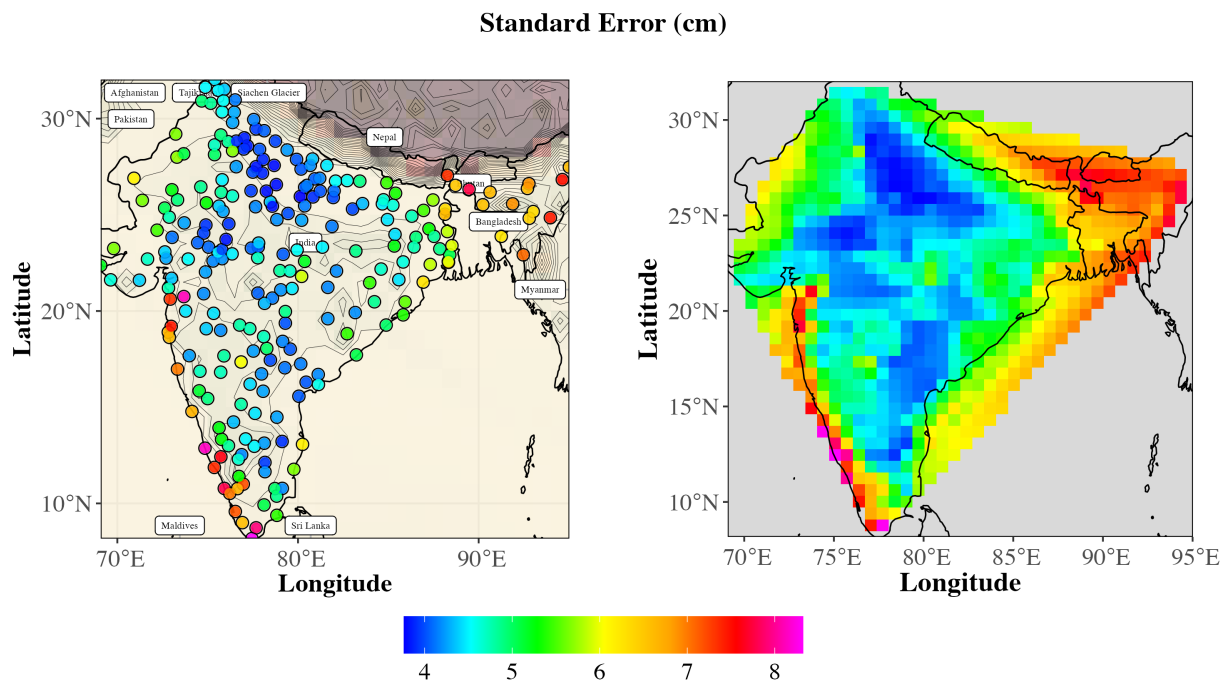


Figure 36: Standard error precipitation surface mapping under HSM + Kriging.

Summary

Hierarchical Spatial Model (HSM) with Kriging outperforms traditional regression-based approaches by effectively capturing spatial dependencies and reducing prediction uncertainty. It provides significant improvements in predictive accuracy ($R^2 = 1$) and spatial generalization, demonstrating its potential for robust precipitation modeling. However, its performance is highly dependent on proper variogram selection and parameter tuning.

Model performance & diagnostics:

- The empirical variogram closely matches the fitted model, supporting the validity of the spatial structure.
- HSM + Kriging achieves near perfect predictive accuracy ($R^2 = 1$), significantly improving over GAM ($R^2 = 0.83$) and Locfit ($R^2 = 0.88$).
- RMSE (3.42 cm) is the lowest among all models, indicating excellent fit and minimal prediction error.

Cross-Validation & robustness:

- LOOCV performance is poor ($R^2 = -0.17$), suggesting potential over-smoothing and overfitting.
- RMSE increases when 10% of data points are removed, showing moderate robustness but revealing sensitivity in sparse regions.
- Despite strong in-sample accuracy, the model exhibits limited generalization, requiring further refinement.

Spatial mapping:

- The precipitation surface is smoother and more continuous than in non-spatial models, capturing regional precipitation trends effectively.
- Standard errors remain lower than GAM and Locfit, demonstrating improved spatial uncertainty quantification.
- The model excels in capturing spatial dependence but may benefit from alternative covariance structures to balance accuracy and generalization.

Overall Summary

Table 2 show a comparison of the methods. In general, linear and GLM models are simpler, interpretable, and robust; however, they lack the flexibility to capture complex spatial variations. Locfit achieves the highest in-sample accuracy and effectively handles non-linearity, but it fails to generalize well and is sensitive to missing data. GAM strikes a strong balance between non-linearity and generalization, making it the best overall model for practical applications (e.g., urban climate analysis, hydrology, flood risk mapping), though it does require additional tuning (e.g., to control standard errors in certain inland regions). The Hierarchical Spatial Model (HSM) with Kriging significantly improves spatial precision and predictive accuracy, outperforming all other methods by fully incorporating spatial dependence.

Table 2: Comparison of model performance and attributes.

Model	In-Sample R^2	LOOCV R^2	Robustness	Spatial Precision
Linear	0.67	0.63	Stable	Moderate
GLM	0.67	0.63	Stable	Moderate
Locfit	0.88	0.23	Highly sensitive to missing data	Best for complex spatial variations
GAM	0.83	0.70	More stable than Locfit	Balanced accuracy and generalization
HSM+ Kriging	1.00	0.81	Most robust, minimal sensitivity to missing data	Best spatial precision

In terms of advantages and disadvantages of the methods:

- Linear & GLM: are interpretable, fast, and robust, but are poor in capturing complex spatial variability.
- Locfit: most flexible for spatial variations, but its prone to overfitting, lacks generalization, and sensitive to missing data.
- GAM: balances flexibility and robustness, strong compromise between non-linearity and generalization. However, requires careful tuning (e.g., controlling for inland standard errors).
- HSM + Kriging: the best model in terms of spatial prediction and robustness, effectively capturing fine-scale variations. However, it requires careful variogram selection and computational resources.

In conclusion, for precipitation modeling, **HSM + Kriging** is the best choice, offering superior spatial precision, predictive accuracy, and robustness. **GAM** remains a strong alternative, balancing accuracy, generalization, and spatial adaptability, while **Locfit** is useful to capturing complex spatial variations but struggles with generalization. **Linear/GLM** remain strong baselines for interpretable and stable predictions.

R packages and functions

From the appendix sections, Table 3 shows the R packages and functions used in the different modeling approaches applied to spatial precipitation data.

Table 3: Summary of main packages and functions used for model fitting.

Model	Packages	Functions used and created
Linear Model (LM)	MASS, stats, leaps, hydroGOF, ggplot2	lm(), stepAIC(), summary(), predict(), fit_Linear()
Generalized Linear Model (GLM)	stats, MASS, hydroGOF, ggplot2	glm(), anova(), predict(), summary(), fit_GLM()
Local GLM (Locfit)	locfit, MASS, ggplot2, hydroGOF	locfit(), predict(), summary(), residuals(), fit_locPoly()
Generalized Additive Model (GAM)	mgcv, MASS, ggplot2, hydroGOF	gam(), anova(), summary(), predict(), fit_GAM()
HSM + Kriging	fields, gstat, sp, raster, nlme	variogram(), getVGMean(), Krig(), nlsLM(), fit_HSM()

`fitRegression()` compiles the entire regression model process, i.e., fitting, diagnosis, LOOCV, robustness analysis and spatial mapping for `fit_Linear()`, `fit_GLM()`, `fit_Poly()` and `fit_GAM()`.

MASS is essential for linear and GLM, particularly for variable selection.

locfit is designed for local regression and nonparametric smoothing.

mgcv for fitting GAM.

fields & gstat for spatial interpolation (explored in full in the next assignment).

hydroGOF to verify the model performance.

ggplot2 & ggpubr primary visualization tools for scatterplots, boxplots, spatial maps, and residual diagnostics.

References

- Hijmans, R. J. (2010). geosphere: Spherical Trigonometry. Institution: Comprehensive R Archive Network Pages: 1.5-20.
- Karney, C. F. F. (2013). Algorithms for geodesics. *Journal of Geodesy*, 87(1):43–55.

Appendix A. R code - Problem 2a

```
1 # Code details -----
2 # author: Catalina Jerez
3 # mail: catalina.jerez@colorado.edu
4 # Course: Advanced Data Analysis Techniques
5 # Code: CVEN 6833
6 # Homework #1: Problem 2a) Spatial Map
7 # Department of Civil, Environmental and Architectural Engineering
8 # College of Engineering & Applied Science
9 # University of Colorado Boulder
10
11 # Library -----
12 gc()
13 rm(list = ls())
14
15 # Load necessary libraries
16 spatialAnalysis.Lib = c("sf", "geosphere", "akima")
17 dataManipulation.Lib = c("dplyr", "reshape2", "tidyr")
18 dataVisualization.Lib = c("ggplot2", "ggrepel", "ggpubr",
19                           "scales", "rnatualearth", "rnatualearthdata")
20 list.packages = unique(c(spatialAnalysis.Lib, dataManipulation.Lib, dataVisualization
21                           ↪ .Lib))
22
23 # Load all required packages
24 supply(list.packages, require, character.only = TRUE)
25
26 # Path -----
27 path = "/Users/caje3244/OneDrive - UCB-O365/2025 Spring/CVEN 6833/"
28 path.code = file.path(path, "Codes/")
29 path.plot = file.path(path, "Figures/")
30
31 # Data -----
32
33 # load datasets
34 climatol = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/climatol-
35 ↪ ann.txt")
36 topo = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/india-grid
37 ↪ -topo.txt")
38 rajee.grid = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/Rajeevan
39 ↪ -grid.txt")
40
41 # assign column names
42 colnames(climatol) = c("lon", "lat", "elev", "prec")
43 colnames(topo) = c("lon", "lat", "elev")
44 colnames(rajee.grid) = c("lon", "lat")
45 # colnames(topo) = c("lon", "lat")
46
47 # interpolate elevation (DEM) and precipitation data
```

```

43 interp.elev = interp(x = topo$lon , y = topo$lat , z = topo$elev , extrap = TRUE,
44                     duplicate = "strip")
45 interp.prec = interp(x = climatol$lon, y = climatol$lat, z = climatol$prec, extrap = TRUE,
46                     duplicate = "strip")
47
48 # convert interpolation results to structured data frames
49 df.elev = data.frame(
50   lon = rep(interp.elev$x, times = length(interp.elev$y)),
51   lat = rep(interp.elev$y, each = length(interp.elev$x)),
52   elev = as.vector(interp.elev$z)) %>% drop_na()
53
54 df.prec = data.frame(
55   lon = rep(interp.prec$x, times = length(interp.prec$y)),
56   lat = rep(interp.prec$y, each = length(interp.prec$x)),
57   prec = as.vector(interp.prec$z) ) %>% drop_na()
58
59 # merge elevation and precipitation data
60 df.map = left_join(df.elev, df.prec, by = c("lon", "lat"))
61
62 # a) plot topography with precipitation overlay -----
63 text.size = 15
64 crs.target = st_crs(4326) # WGS84 CRS
65 map.extent = data.frame(xmin = min(df.map$lon),
66                         xmax = max(df.map$lon),
67                         ymin = min(df.map$lat),
68                         ymax = max(df.map$lat))
69
70 # Set elevation and precipitation colors
71 elev.colors = colorRampPalette(c("#f6e8c3", "#c9cba3", "#ffe1a8", "#e26d5c", "#723d46", "
  ↪ #472d30", "#543005"))(1000)
72 prec.colors = colorRampPalette(c("#0000FF", "#00FFFF", "#00FF00", "#FFFF00", "#
  ↪ FFA500", "#FF0000", "#FF00FF"))(1000)
73
74 # Download and clean spatial data
75 countries = ne_countries(scale = "medium", returnclass = "sf") %>% st_transform(crs =
  ↪ crs.target)
76 countries_centroids = countries %>%
77   st_centroid() %>%
78   st_transform(crs = crs.target) %>%
79   st_crop(xmin = min(df.map$lon), xmax = max(df.map$lon),
80          ymin = min(df.map$lat), ymax = max(df.map$lat)) %>%
81   mutate(lon = st_coordinates(.)[,1], lat = st_coordinates(.)[,2])
82
83 spatial.plot =
84   ggplot() +
85   # elevation raster
86   geom_raster(data = df.map, aes(x = lon, y = lat, fill = elev)) +
87

```

```

88 # country, border
89 geom_sf( data = countries, fill = NA, color = "#000", linewidth = 0.4) +
90 geom_rect(data = map.extent, fill = NA, color = "#67001f", size = 0.8,
91           aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax) ) +
92 geom_contour(data = df.map, aes(x = lon, y = lat, z = elev), color = "#000", linewidth =
93             ↪ 0.1, bins = 30) +
94 geom_label_repel(data = countries_centroids, aes(x = lon, y = lat, label = name),
95                 family = "Times", size = 3, fill = "#fff", color = "#000",
96                 box.padding = 0.15, segment.color = "#bdbdbd") +
97 # precipitation points
98 geom_point(data = climatol, aes(x = lon, y = lat), size = 3, shape = 21) +
99 geom_point(data = climatol, aes(x = lon, y = lat, color = prec), size = 2, shape = 19,
100            ↪ alpha = .8) +
101 # coordinate system
102 coord_sf( xlim = c(min(df.prec$lon), max(df.prec$lon)),
103           ylim = c(min(df.prec$lat), max(df.prec$lat)),
104           expand = FALSE) +
105 # scales
106 scale_fill_gradientn( colors = elev.colors, name = "Elevation (m)") +
107 scale_color_gradientn(colors = prec.colors, name = "Precipitation (mm)") +
108 # guides and labels
109 guides(color = guide_colorbar(barwidth = 13, barheight = 2),
110         fill = guide_colorbar(barwidth = 13, barheight = 2)) +
111 labs(x = "Longitude", y = "Latitude") +
112 # theme
113 theme_bw(base_family = "Times") +
114 theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
115       axis.text.y = element_text(size = text.size),
116       axis.title = element_text(size = text.size + 2, face = 'bold'),
117       plot.title = element_text(size = text.size + 2, face = 'bold'),
118       legend.position = "bottom",
119       legend.title = element_text(size = text.size, face = 'bold'),
120       legend.text = element_text(size = text.size),
121       panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
122       panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
123       panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
124       panel.grid.major.x = element_blank(),
125       panel.grid.major.y = element_blank())
126
127 # for surface change the color when prec == 0
128 df.prec = df.prec %>%
129 mutate(prec = ifelse(prec == 0, NA, prec)) # Convert 0 values to NA
130
131 surface.plot =
132 ggplot() +
133 geom_raster( data = df.prec, mapping = aes(x = lon, y = lat, fill = prec) ) +
134 # country, border
135 geom_sf(data = countries, fill = NA, color = "#000", linewidth = 0.4) +

```

```

134 geom_rect(data = map.extent, fill = NA, color = "#67001f", size = 0.8,
135   aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax) ) +
136 coord_sf( xlim = c(min(df.prec$lon), max(df.prec$lon)),
137   ylim = c(min(df.prec$lat), max(df.prec$lat)),
138   expand = FALSE ) +
139 scale_fill_gradientn(colors = prec.colors, na.value = "#d9d9d9",
140   name = "Precipitation (mm)") +
141 guides(fill = guide_colorbar(barwidth = 13, barheight = 2) ) +
142 labs(x = "Longitude", y = "Latitude") +
143 theme_bw(base_family = "Times") +
144 theme(
145   axis.text.x = element_text(size = text.size, vjust = 0.5),
146   axis.text.y = element_text(size = text.size),
147   axis.title = element_text(size = text.size + 2, face = "bold"),
148   plot.title = element_text(size = text.size + 2, face = "bold"),
149   legend.position = "bottom",
150   legend.title = element_text(size = text.size, face = "bold"),
151   legend.text = element_text(size = text.size),
152   panel.border = element_rect(color = "#000000", fill = NA, linewidth = 0.5),
153   panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
154   panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
155   panel.grid.major.x = element_blank(),
156   panel.grid.major.y = element_blank()
157 )
158
159 # join map
160 join.map = ggarrange(spatial.plot, surface.plot,
161   ncol = 2,
162   common.legend = TRUE,
163   legend = "bottom")
164 ggsave(filename = paste0(path.plot, "HW1_FigureP2a_Maps.png"),
165   plot = join.map,
166   width = 12, height = 6.5, bg = "#fff")

```

Appendix B. R code - Problem2b

```
1 # Code details -----
2 # author: Catalina Jerez
3 # mail: catalina.jerez@colorado.edu
4 # Course: Advanced Data Analysis Techniques
5 # Code: CVEN 6833
6 # Homework #1: Problem 2b) distance to sea
7 # Department of Civil, Environmental and Architectural Engineering
8 # College of Engineering & Applied Science
9 # University of Colorado Boulder
10
11 # Library -----
12 gc()
13 rm(list = ls())
14
15 # Load necessary libraries
16 spatialAnalysis.Lib = c("sf", "geosphere", "akima")
17 dataManipulation.Lib = c("dplyr", "reshape2", "tidyr")
18 dataVisualization.Lib = c("ggplot2", "ggrepel", "scales",
19                           "rnatualearth", "rnatualearthdata")
20 list.packages = unique(c(spatialAnalysis.Lib, dataManipulation.Lib, dataVisualization
21                          ↪ .Lib))
22
23 # Load all required packages
24 supply(list.packages, require, character.only = TRUE)
25
26 # Path -----
27 path = "/Users/caje3244/OneDrive - UCB-O365/2025 Spring/CVEN 6833/"
28 path.code = file.path(path, "Codes/")
29 path.plot = file.path(path, "Figures/")
30
31 # Inputs -----
32
33 # load datasets
34 climatol = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/climatol-
35                      ↪ ann.txt")
36 topo = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/india-grid
37                  ↪ -topo.txt")
38 rajee.grid = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/Rajeevan
39                        ↪ -grid.txt")
40
41 # assign column names
42 colnames(climatol) = c("lon", "lat", "elev", "prec")
43 colnames(topo) = c("lon", "lat", "elev")
44 colnames(rajee.grid) = c("lon", "lat")
45
46 # interpolate elevation (DEM) and precipitation data
47 interp.elev = interp(x = topo$lon, y = topo$lat, z = topo$elev, extrap = TRUE,
```



```

43         duplicate = "median")
44 interp.prec = interp(x = climatol$lon, y = climatol$lat, z = climatol$prec, extrap = TRUE,
45                     duplicate = "median")
46
47 # convert interpolation results to structured data frames
48 df.elev = data.frame(
49   lon  = rep(interp.elev$x, times = length(interp.elev$y)),
50   lat  = rep(interp.elev$y, each = length(interp.elev$x)),
51   elev = as.vector(interp.elev$z)) %>% drop_na()
52
53 df.prec = data.frame(
54   lon  = rep(interp.prec$x, times = length(interp.prec$y)),
55   lat  = rep(interp.prec$y, each = length(interp.prec$x)),
56   prec = as.vector(interp.prec$z) ) %>% drop_na()
57
58 # merge elevation and precipitation data
59 df.map  = left_join(df.elev, df.prec, by = c("lon", "lat"))
60 map.extent = data.frame(xmin = min(df.map$lon),
61                         xmax = max(df.map$lon),
62                         ymin = min(df.map$lat),
63                         ymax = max(df.map$lat))
64
65 # ggplot variables
66 elev.colors = colorRampPalette(c("#f6e8c3", "#c9cba3", "#ffe1a8", "#e26d5c",
67                                "#723d46", "#472d30", "#543005"))(1000)
68 prec.colors = colorRampPalette(c("#0000FF", "#00FFFF", "#00FF00", "#FFFF00",
69                                "#FFA500", "#FF0000", "#FF00FF"))(1000)
70 text.size   = 15
71 label.size  = 5
72 crs.target  = st_crs(4326) # WGS84 CRS
73
74 # Functions to calc min dist to coastline -----
75
76 # compute the minimum distance from each station to the coastline
77 MinDistHaversine = function(lon, lat, coast_df) {
78   distances = distHaversine(matrix(c(lon, lat), ncol = 2), coast_df)
79   min(distances) / 1000 # m to km
80 }
81
82 MinDistGeo = function(lon, lat, coast_df) {
83   distances = distGeo(matrix(c(lon, lat), ncol = 2), coast_df)
84   min(distances) / 1000 # m to km
85 }
86
87
88 # Get India's coastline data -----
89 options(warn=-1)
90 # sea's centroid

```

```

91 sea.centroid = data.frame(
92   lon = c(65, 88),
93   lat = c(14, 15),
94   name = c("Arabian Sea", "Bay of Bengal"))
95
96 # download and set crs
97 countries      = ne_countries(scale = "medium", returnclass = "sf") %>%
98   st_transform(crs = crs.target)
99 # calc centroids
100 countries.centroids = countries %>%
101   st_centroid() %>%
102   st_transform(crs = crs.target) %>%
103   st_crop(xmin = min(df.map$lon), xmax = max(df.map$lon),
104           ymin = min(df.map$lat), ymax = max(df.map$lat)) %>%
105   mutate(lon = st_coordinates(.)[,1], lat = st_coordinates(.)[,2])
106 # get coastline
107 countries.coast = countries %>%
108   st_transform(crs = crs.target) %>%
109   st_crop(xmin = min(df.map$lon), xmax = max(df.map$lon),
110           ymin = min(df.map$lat), ymax = max(df.map$lat))
111 coastline.lines = st_cast(countries.coast, "MULTILINESTRING")
112 coast.points    = as.data.frame(st_coordinates(coastline.lines))
113 colnames(coast.points) = c("lon", "lat", "Z", "M")
114
115 # separate coastlines into regions
116 coast.points    = coast.points[, c("lon", "lat")] # drop unnecessary vars
117 coast.arabian   = coast.points %>% filter(lon < 76) # west coast for Arabian Sea
118 coast.bengal    = coast.points %>% filter(lon >= 76) # east coast for Bay of Bengal
119
120 location.labels = c("Min. distance to any coast",
121                    "Min. distance to Arabian Sea",
122                    "Min. distance to Bay of Bengal")
123 # minimum distances by Haversine method -----
124 climaHarv      = climatol
125 climaHarv$mdCoast = mapply(MinDistHaversine, climaHarv$lon, climaHarv$lat,
126                             MoreArgs = list(coast_df = coast.points))
127 climaHarv$mdArabian = mapply(MinDistHaversine, climaHarv$lon, climaHarv$lat,
128                             MoreArgs = list(coast_df = coast.arabian))
129 climaHarv$mdBengal = mapply(MinDistHaversine, climaHarv$lon, climaHarv$lat,
130                             MoreArgs = list(coast_df = coast.bengal))
131 str(climaHarv)
132 # melt for plotting
133 clima.melt      = melt(climaHarv, id = c("lon", "lat", "elev", "prec"))
134 clima.melt$variable = factor(clima.melt$variable)
135 levels(clima.melt$variable) = location.labels
136
137 # spatial map with precipitation, elevation, and distances
138 spatial.plot =

```

```

139 ggplot() +
140 # elevation raster
141 geom_raster(data = df.map, aes(x = lon, y = lat, fill = elev), alpha = .85) +
142 # country, borders
143 geom_sf(data = countries, fill = NA, color = "#000", linewidth = 0.4) +
144 geom_rect(data = map.extent, fill = NA, color = "#67001f", linewidth = 0.8,
145           aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax)) +
146 geom_contour(data = df.map, aes(x = lon, y = lat, z = elev), color = "#000",
147             linewidth = 0.1, bins = 30) +
148 # precipitation points
149 geom_point(data = clima.melt, aes(x = lon, y = lat, size = value),
150           shape = 21, color = "#000") +
151 geom_point(data = clima.melt, aes(x = lon, y = lat, color = prec, size = value),
152           shape = 19, alpha = .8) +
153 # labels
154 geom_label_repel(data = sea.centroid, aes(x = lon, y = lat, label = name),
155                 family = "Times", size = label.size, fill = "#c6dbef", color = "#08306b",
156                 box.padding = 0.3, segment.color = "black") +
157 geom_label_repel(data = countries.centroids, aes(x = lon, y = lat, label = name),
158                 family = "Times", size = label.size-2, fill = "#fff", color = "#000",
159                 box.padding = 0.15, segment.color = "#bdbdbd") +
160 labs(x = "Longitude", y = "Latitude") +
161 # coordinate system
162 coord_sf(xlim = c(min(df.map$lon), max(df.map$lon)),
163          ylim = c(min(df.map$lat), max(df.map$lat)), expand = FALSE) +
164 # scales
165 scale_fill_gradientn(colors = elev.colors, name = "Elevation (km)") +
166 scale_color_gradientn(colors = prec.colors, name = "Precipitation (mm)") +
167 scale_size_continuous(name = "Distance (km)",
168                       breaks = seq(0, max(clima.melt$value, na.rm = TRUE), by = 200)) +
169 # guides
170 guides(color = guide_colorbar(barwidth = 2, barheight = 6, title.position = "right", order
171   ↪ = 3),
172        fill = guide_colorbar(barwidth = 2, barheight = 6, title.position = "right", order
173   ↪ = 2),
174        size = guide_legend(nrow = 3, order = 1)) +
175 # facet
176 facet_wrap(~variable, ncol = 2) +
177 # theme
178 theme_bw(base_family = "Times") +
179 theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
180       axis.text.y = element_text(size = text.size),
181       axis.title = element_text(size = text.size + 2, face = 'bold'),
182       plot.title = element_text(size = text.size + 2, face = 'bold'),
183       legend.position = c(0.75, 0.25),
184       legend.title = element_text(size = text.size, face = 'bold'),
185       legend.text = element_text(size = text.size),

```

```

185
186 strip.text      = element_text(size = text.size, face = "bold"),
187 strip.background = element_blank(),
188
189 panel.border      = element_rect(color = "#000000", fill = NA, linewidth = .5),
190 panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
191 panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
192 panel.grid.major.x = element_blank(),
193 panel.grid.major.y = element_blank()
194
195 # Save
196 ggsave(filename = paste0(path.plot, "HW1_FigureP2b_SpatialMap_DistHarvesine.png"),
197        plot = spatial.plot, width = 10.5, height = 10.5)
198
199 saveRDS(climaHarv, file = paste0(path.code, "RData/hw1_climatol_distHarvesine.rds"))
200
201 # minimum distances by Geodesic method -----
202 climaGeo      = climatol
203 climaGeo$mdCoast = mapply(MinDistGeo, climaGeo$lon, climaGeo$lat,
204                          MoreArgs = list(coast_df = coast.points))
205 climaGeo$mdArabian = mapply(MinDistGeo, climaGeo$lon, climaGeo$lat,
206                          MoreArgs = list(coast_df = coast.arabian))
207 climaGeo$mdBengal = mapply(MinDistGeo, climaGeo$lon, climaGeo$lat,
208                          MoreArgs = list(coast_df = coast.bengal))
209
210 # melt for plotting
211 clima.melt      = melt(climaGeo, id = c("lon", "lat", "elev", "prec"))
212 clima.melt$variable = factor(clima.melt$variable)
213 levels(clima.melt$variable) = location.labels
214
215 # spatial map with precipitation, elevation, and distances
216 spatial.plot =
217   ggplot() +
218     # elevation raster
219     geom_raster(data = df.map, aes(x = lon, y = lat, fill = elev), alpha = .85) +
220     # country, borders
221     geom_sf(data = countries, fill = NA, color = "#000", linewidth = 0.4) +
222     geom_rect(data = map.extent, fill = NA, color = "#67001f", linewidth = 0.8,
223              aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax)) +
224     geom_contour(data = df.map, aes(x = lon, y = lat, z = elev), color = "#000",
225                 linewidth = 0.1, bins = 30) +
226     # precipitation points
227     geom_point(data = clima.melt, aes(x = lon, y = lat, size = value),
228               shape = 21, color = "#000") +
229     geom_point(data = clima.melt, aes(x = lon, y = lat, color = prec, size = value),
230               shape = 19, alpha = .8) +
231     # labels
232     geom_label_repel(data = sea.centroid, aes(x = lon, y = lat, label = name),

```

```

233     family = "Times", size = label.size, fill = "#c6dbef", color = "#08306b",
234     box.padding = 0.3, segment.color = "black") +
235 geom_label_repel(data = countries.centroids, aes(x = lon, y = lat, label = name),
236     family = "Times", size = label.size-2, fill = "#fff", color = "#000",
237     box.padding = 0.15, segment.color = "#bdbdbd") +
238 labs(x = "Longitude", y = "Latitude") +
239 # coordinate system
240 coord_sf(xlim = c(min(df.map$lon), max(df.map$lon)),
241     ylim = c(min(df.map$lat), max(df.map$lat)), expand = FALSE) +
242 # scales
243 scale_fill_gradientn(colors = elev.colors, name = "Elevation (km)") +
244 scale_color_gradientn(colors = prec.colors, name = "Precipitation (mm)") +
245 scale_size_continuous(name = "Distance (km)",
246     breaks = seq(0, max(clima.melt$value, na.rm = TRUE), by = 200)) +
247 # guides
248 guides(color = guide_colorbar(barwidth = 2, barheight = 6, title.position = "right", order
249     ↪ = 3),
250     fill = guide_colorbar(barwidth = 2, barheight = 6, title.position = "right", order =
251     ↪ 2),
252     size = guide_legend(nrow = 3, order = 1)) +
253 # facet
254 facet_wrap(~variable, ncol = 2) +
255 # theme
256 theme_bw(base_family = "Times") +
257 theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
258     axis.text.y = element_text(size = text.size),
259     axis.title = element_text(size = text.size + 2, face = 'bold'),
260     plot.title = element_text(size = text.size + 2, face = 'bold'),
261     legend.position = c(0.75, 0.25),
262     legend.title = element_text(size = text.size, face = 'bold'),
263     legend.text = element_text(size = text.size),
264     strip.text = element_text(size = text.size, face = "bold"),
265     strip.background = element_blank(),
266     panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
267     panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
268     panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
269     panel.grid.major.x = element_blank(),
270     panel.grid.major.y = element_blank())
271
272 # Save
273 ggsave(filename = paste0(path.plot, "HW1_FigureP2b_SpatialMap_DistGeodesic.png"),
274     plot = spatial.plot, width = 10.5, height = 10.5)
275
276 saveRDS(climaGeo, file = paste0(path.code, "RData/hw1_climatol_distGeodesic.rds"))
277
278

```

```

279 # comparison of distance calcs methods -----
280 climaDiff = climatol
281 climaDiff$diffCoast = abs(climaHarv$mdCoast - climaGeo$mdCoast)
282 climaDiff$diffArabian = abs(climaHarv$mdArabian - climaGeo$mdArabian)
283 climaDiff$diffBengal = abs(climaHarv$mdBengal - climaGeo$mdBengal)
284
285 # Summary statistics
286 summaryDiff = data.frame(
287   variable = location.labels,
288   meanDiff = c(mean(climaDiff$diffCoast, na.rm = TRUE),
289                 mean(climaDiff$diffArabian, na.rm = TRUE),
290                 mean(climaDiff$diffBengal, na.rm = TRUE)),
291   medDiff = c(median(climaDiff$diffCoast, na.rm = TRUE),
292               median(climaDiff$diffArabian, na.rm = TRUE),
293               median(climaDiff$diffBengal, na.rm = TRUE)),
294   sdDiff = c(sd(climaDiff$diffCoast, na.rm = TRUE),
295              sd(climaDiff$diffArabian, na.rm = TRUE),
296              sd(climaDiff$diffBengal, na.rm = TRUE)) ) %>%
297   mutate(label = paste0("Mean: ", round(meanDiff, 2),
298                          "\nMedian: ", round(medDiff, 2),
299                          "\nSD: ", round(sdDiff, 2)))
300 print(summaryDiff)
301
302 # melt data
303 clima.melt.diff = melt(climaDiff, id = c("lon", "lat", "elev", "prec"))
304 levels(clima.melt.diff$variable) = location.labels
305
306 density.plot =
307   ggplot(data = clima.melt.diff, aes(x = value)) +
308   geom_density(alpha = 0.6, color = "#8e0152", fill = "#de77ae") +
309   facet_wrap(~variable, nrow = 1, scale = "free_y") +
310   scale_y_continuous(breaks = pretty_breaks(n=4)) +
311   labs(x = "Absolute difference (Km)", y = "Density") +
312   geom_text(data = summaryDiff, aes(x = Inf, y = Inf, label = label),
313             hjust = 1.1, vjust = 1.1, inherit.aes = FALSE, size = 4.5, family = "Times") +
314   # theme
315   theme_bw(base_family = "Times") +
316   theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
317         axis.text.y = element_text(size = text.size),
318         axis.title = element_text(size = text.size + 2, face = 'bold'),
319         legend.position = "none",
320
321         strip.text = element_text(size = text.size, face = "bold"),
322         strip.background = element_blank(),
323
324         panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
325         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
326         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),

```

```
327     panel.grid.major.x = element_blank(),
328     panel.grid.major.y = element_blank())
329
330 # Save plot
331 ggsave(filename = paste0(path.plot, "HW1_FigureP2b_Comparison_HavGeo.png"),
332        plot = density.plot, width = 13, height = 4.5)
333
334 # Save difference data
335 saveRDS(climaDiff, file = paste0(path.code, "RData/hw1_climatol_diff_HavGeo.rds"))
```

Appendix C. Problems 2c-6

C.1. Library

```
1 # Code details -----
2 # author: Catalina Jerez
3 # mail: catalina.jerez@colorado.edu
4 # Course: Advanced Data Analysis Techniques
5 # Code: CVEN 6833
6 # Department of Civil, Environmental and Architectural Engineering
7 # College of Engineering & Applied Science
8 # University of Colorado Boulder
9
10 # Library -----
11 required.lib = c("MASS", "leaps", "dplyr", "reshape2",
12                "mgcv", "locfit", "akima",
13                "psych", "hydroGOF", "fields", "sm", "nlme", "minpack.lm",
14                "ggplot2", "ggpubr", "ggrepel", "purrr",
15                "scales", "maps", "rnatualearth", "rnatualearthdata")
16 supply(required.lib, require, character.only = TRUE)
17 # Plotting -----
18 sptMap.ggplot = function(clima, topography, title.plot = NULL,
19                          crs.target = st_crs(4326), text.size = 14){
20   options(warn=-1)
21   # interpolate elevation (DEM) and precipitation data
22   interp.elev = interp(x = topography$lon, y = topography$lat, z = topography$elev, extrap
23     ↪ = TRUE,
24                       duplicate = "median")
25   interp.prec = interp(x = clima$lon, y = clima$lat, z = clima$prec, extrap = TRUE,
26                        duplicate = "median")
27   # convert interpolation results to structured data frames -----
28   df.elev = data.frame(
29     lon = rep(interp.elev$x, times = length(interp.elev$y)),
30     lat = rep(interp.elev$y, each = length(interp.elev$x)),
31     elev = as.vector(interp.elev$z)) %>% drop_na()
32
33   df.prec = data.frame(
34     lon = rep(interp.prec$x, times = length(interp.prec$y)),
35     lat = rep(interp.prec$y, each = length(interp.prec$x)),
36     prec = as.vector(interp.prec$z)) %>% drop_na()
37
38   # merge elevation and precipitation data
39   df.map = left_join(df.elev, df.prec, by = c("lon", "lat"))
40   map.extent = data.frame(xmin = min(df.map$lon),
41                            xmax = max(df.map$lon),
42                            ymin = min(df.map$lat),
43                            ymax = max(df.map$lat))
```



```

44 # download and clean spatial data
45 countries = ne_countries(scale = "medium", returnclass = "sf") %>% st_transform(crs =
  ↪ crs.target)
46 countries_centroids = countries %>%
47   st_centroid() %>%
48   st_transform(crs = crs.target) %>%
49   st_crop(xmin = min(df.map$lon), xmax = max(df.map$lon),
50     ymin = min(df.map$lat), ymax = max(df.map$lat)) %>%
51   mutate(lon = st_coordinates(.)[,1], lat = st_coordinates(.)[,2])
52
53 # ggplot -----
54 # elevation and precipitation colors
55 elev.colors = colorRampPalette(c("#f6e8c3", "#c9cba3", "#ffe1a8", "#e26d5c", "#723d46",
  ↪ "#472d30", "#543005"))(1000)
56 prec.colors = colorRampPalette(c("#0000FF", "#00FFFF", "#00FF00", "#FFFF00", "#
  ↪ FFA500", "#FF0000", "#FF00FF"))(1000)
57
58 spatial.plot =
59   ggplot() +
60   # elevation raster
61   geom_raster(data = df.map, aes(x = lon, y = lat, fill = elev), alpha = .5, show.legend =
  ↪ F) +
62   # country, border
63   geom_sf( data = countries, fill = NA, color = "#000", linewidth = 0.4) +
64   geom_rect(data = map.extent, fill = NA, color = "#67001f", size = 0.8,
65     aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax) ) +
66   geom_contour(data = df.map, aes(x = lon, y = lat, z = elev), color = "#000", linewidth
  ↪ = 0.1, bins = 30) +
67   geom_label_repel(data = countries_centroids, aes(x = lon, y = lat, label = name),
68     family = "Times", size = 2, fill = "#fff", color = "#000",
69     box.padding = 0.15, segment.color = "#bdbdbd") +
70   # precipitation points
71   geom_point(data = clima, aes(x = lon, y = lat), size = 3, shape = 21) +
72   geom_point(data = clima, aes(x = lon, y = lat, color = prec), size = 2.2, shape = 19,
  ↪ alpha = .9) +
73   # coordinate system
74   coord_sf(xlim = c(min(df.prec$lon), max(df.prec$lon)),
75     ylim = c(min(df.prec$lat), max(df.prec$lat)),
76     expand = FALSE) +
77   scale_x_continuous(breaks = pretty(df.map$lon, n = 3)) + # Move breaks here
78   scale_y_continuous(breaks = pretty(df.map$lat, n = 3)) + # Move breaks here
79
80   # scales
81   scale_fill_gradientn( colors = elev.colors, name = "Elevation (m)") +
82   scale_color_gradientn(colors = prec.colors, name = "Precipitation (cm)",
83     breaks = pretty(range(clima$prec), 5)) +
84   # guides and labels
85   guides(color = guide_colorbar(barwidth = 16, barheight = 1.5),

```

```

86     fill = guide_colorbar(barwidth = 16, barheight = 1.5)) +
87     labs(x = "Longitude", y = "Latitude", title = "") +
88     # theme
89     theme_bw(base_family = "Times") +
90     theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
91           axis.text.y = element_text(size = text.size),
92           axis.title = element_text(size = text.size + 1, face = 'bold'),
93           plot.title = element_text(size = text.size + 1, face = 'bold', hjust = .5),
94           legend.position = "bottom",
95           legend.title = element_blank(),
96           legend.text = element_text(size = text.size),
97           panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
98           panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
99           panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
100          panel.grid.major.x = element_blank(),
101          panel.grid.major.y = element_blank())
102
103 # for surface change the color when prec == 0
104 df.prec = df.prec %>%
105   mutate(prec = ifelse(prec == 0, NA, prec)) # Convert 0 values to NA
106
107 surface.plot =
108   ggplot() +
109   geom_raster(data = df.prec, mapping = aes(x = lon, y = lat, fill = prec)) +
110   # country, border
111   geom_sf(data = countries, fill = NA, color = "#000", linewidth = 0.4) +
112   geom_rect(data = map.extent, fill = NA, color = "#67001f", size = 0.8,
113             aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax = ymax)) +
114   coord_sf(xlim = c(min(df.prec$lon), max(df.prec$lon)),
115            ylim = c(min(df.prec$lat), max(df.prec$lat)),
116            expand = FALSE) +
117   scale_fill_gradientn(colors = prec.colors, na.value = "#d9d9d9",
118                        name = "Precipitation (cm)") +
119   guides(fill = guide_colorbar(barwidth = 13, barheight = 2)) +
120   labs(x = "Longitude", y = "Latitude", title = "") +
121   theme_bw(base_family = "Times") +
122   theme(
123     axis.text.x = element_text(size = text.size, vjust = 0.5),
124     axis.text.y = element_text(size = text.size),
125     axis.title = element_text(size = text.size + 2, face = "bold"),
126     plot.title = element_text(size = text.size + 2, face = "bold"),
127     legend.position = "bottom",
128     legend.title = element_text(size = text.size, face = "bold"),
129     legend.text = element_text(size = text.size),
130     panel.border = element_rect(color = "#000000", fill = NA, linewidth = 0.5),
131     panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
132     panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
133     panel.grid.major.x = element_blank(),

```

```

134   panel.grid.major.y = element_blank()
135 )
136
137 # join map
138 join.map = ggarrange(spatial.plot, surface.plot,
139                      ncol = 2,
140                      common.legend = TRUE,
141                      legend = "bottom")
142 join.map = annotate_figure(join.map,
143                            top = text_grob(title.plot, family = "Times",
144                                             face = "bold", size = text.size+2) )
145 return(join.map)
146 }
147 pairs.ggplot = function(data, text.size = 14){
148   ggpairs(
149     data,
150     columns = sort(names(data)),
151     aes(color = cut_interval(data$prec, n=2) ), # categorize prec and plot in cm
152     upper = list(continuous = wrap("cor", size = 3)), # corr in upper panel
153     lower = list(continuous = wrap("density", alpha = 0.6, size = 1.2)), # scatter in lower
154               ↪ panel
155     diag = list(continuous = wrap("barDiag", alpha = 0.6, bins = 15)) # density on
156               ↪ diagonal
157   ) +
158   scale_color_manual("Precipitation (cm per year)", values = scale.color, labels = scale.
159                     ↪ label) +
160   scale_fill_manual("Precipitation (cm per year)", values = scale.color, labels = scale.label)
161                     ↪ +
162   theme_bw(base_family = "Times") +
163   theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
164         axis.text.y = element_text(size = text.size),
165         axis.title = element_text(size = text.size + 2, face = 'bold'),
166         plot.title = element_text(size = text.size + 2, face = 'bold'),
167
168         strip.text = element_text(size = text.size, face = "bold"),
169         strip.background = element_blank(),
170         legend.position = "bottom",
171         legend.title = element_text(size = text.size, face = 'bold'),
172         legend.text = element_text(size = text.size),
173
174         panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
175         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
176         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
177         panel.grid.major.x = element_blank(),
178         panel.grid.major.y = element_blank())
179 }
180
181 diagnostics.ggplot = function(data, text.size = 14){
182   # data is df with "obs", "sim", "res"

```

```

178
179 # are the residuals normal and iid?
180 # 1) normality: qq plot
181 qq.plot = ggplot(data, aes(sample = res)) +
182   stat_qq(color = "#005a32", fill = "#a1d99b", size = 3, shape = 21) +
183   stat_qq_line(color = "#99000d", linewidth = 1) +
184   labs(x = "Theoretical Quantiles", y = "Sample Quantiles (cm)",
185        title = "a) Q-Q Plot")+
186   scale_y_continuous(breaks = scales::pretty_breaks(n = 3)) +
187   theme_bw(base_family = "Times") +
188   theme(axis.text.x      = element_text(size = text.size, vjust = 0.5),
189         axis.text.y      = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
190         axis.title       = element_text(size = text.size, face = 'bold'),
191         plot.title       = element_text(size = text.size+1, face = 'bold', hjust = .5),
192         panel.border     = element_rect(color = "#000000", fill = NA, linewidth = .5),
193         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
194         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
195         panel.grid.major.x = element_line(color = "#d9d9d9", linetype = "dashed"),
196         panel.grid.major.y = element_line(color = "#d9d9d9", linetype = "dashed"))
197
198 # 1) normality: histogram
199 hist.plot = ggplot(data, aes(x = res)) +
200   geom_histogram(aes(y = after_stat(density)), bins = 30, fill = "#a1d99b", color = "
201   ↪ black") +
202   geom_density(color = "#99000d", linewidth = 1) +
203   labs(x = "Residuals (cm)", y = "Density", title = "b) Histogram")+
204   scale_y_continuous(breaks = scales::pretty_breaks(n = 3)) +
205   theme_bw(base_family = "Times") +
206   theme(axis.text.x      = element_text(size = text.size, vjust = 0.5),
207         axis.text.y      = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
208         axis.title       = element_text(size = text.size, face = 'bold'),
209         plot.title       = element_text(size = text.size+1, face = 'bold', hjust = .5),
210         panel.border     = element_rect(color = "#000000", fill = NA, linewidth = .5),
211         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
212         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
213         panel.grid.major.x = element_line(color = "#d9d9d9", linetype = "dashed"),
214         panel.grid.major.y = element_line(color = "#d9d9d9", linetype = "dashed"))
215
216 # 2) independence of residuals: check autocorrelation
217 acf.data = acf(data$res, plot = FALSE)
218 acf.df   = data.frame(lag = acf.data$lag, acf = acf.data$acf)
219 ndata    = length(data$res)
220 conf.int = 1.96 / sqrt(ndata) # confidence interval
221
222 acf.plot = ggplot(acf.df, aes(x = lag, y = acf)) +
223   geom_bar(stat = "identity", color = "#005a32", fill = "#a1d99b") +
224   geom_hline(yintercept = c(-conf.int, conf.int), linetype = "dashed", color = "#99000d",
225             ↪ linewidth = 0.8) +

```

```

224 labs(x = "Lag", y = "ACF", title = "c) Autocorrelation")+
225 scale_y_continuous(breaks = scales::pretty_breaks(n = 3)) +
226 theme_bw(base_family = "Times") +
227 theme(axis.text.x      = element_text(size = text.size, vjust = 0.5),
228       axis.text.y      = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
229       axis.title       = element_text(size = text.size, face = 'bold'),
230       plot.title       = element_text(size = text.size+1, face = 'bold', hjust = .5),
231       panel.border     = element_rect(color = "#000000", fill = NA, linewidth = .5),
232       panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
233       panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
234       panel.grid.major.x = element_line(color = "#d9d9d9", linetype = "dashed"),
235       panel.grid.major.y = element_line(color = "#d9d9d9", linetype = "dashed") )
236
237 # 3) Homoskedasticity (constant variance? residuals vs fitted
238 resid.plot = ggplot(data, aes(x = sim, y = res)) +
239   geom_point(color = "#005a32", fill = "#a1d99b", shape = 21, size = 3) +
240   geom_hline(yintercept = 0, color = "#99000d", linewidth = 1) +
241   labs(x = "Estimated Precipitation (cm)", y = "Residuals (cm)",
242        title = "d) Residuals vs Estimated Precipitation") +
243   scale_y_continuous(breaks = scales::pretty_breaks(n = 3)) +
244   theme_bw(base_family = "Times") +
245   theme(axis.text.x      = element_text(size = text.size, vjust = 0.5),
246         axis.text.y      = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
247         axis.title       = element_text(size = text.size, face = 'bold'),
248         plot.title       = element_text(size = text.size+1, face = 'bold', hjust = .5),
249         panel.border     = element_rect(color = "#000000", fill = NA, linewidth = .5),
250         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
251         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
252         panel.grid.major.x = element_line(color = "#d9d9d9", linetype = "dashed"),
253         panel.grid.major.y = element_line(color = "#d9d9d9", linetype = "dashed") )
254
255 # arrange together
256 ggarrange(qq.plot, hist.plot, acf.plot, resid.plot, ncol = 2, nrow = 2)
257
258 }
259 scatter.ggplot = function(data, model = NULL, title = NULL,
260                          text.size = 14){
261   # data is a dataframe with columns obs and sim.
262
263   # estimate range from df.scatter$sim if range between 0 and 1,
264   # the family its binomial and we need to rescale df.scatter$obs
265   range.obs = range(data$sim, na.rm = TRUE)
266   if(isTRUE(0 <= range.obs[1] & range.obs[2] <= 1)){
267     data$obs = scales::rescale(data$obs, to = c(0, 1))
268   }else{
269     data$obs = ifelse(data$obs <= 0, runif(1, 0.0001, 0.001), data$obs)
270   }
271

```

```

272 # data is df with obs and sim precipitation
273 if(is.null(model)){
274   df.skill = data.frame(
275     metric = c(rownames(melt(evaluation.metrics(data$obs, data$sim)))),
276     value = c(melt(evaluation.metrics(data$obs, data$sim))$value)
277   )
278   skill.text = paste0(df.skill$metric, ": ", round(df.skill$value, 2), collapse = "\n")
279   eq.text = paste0("")
280 }
281 }else{
282   df.skill = data.frame(
283     metric = c("AIC", "BIC", rownames(melt(evaluation.metrics(data$obs, data$sim)))),
284     value = c(AIC(model), BIC(model), melt(evaluation.metrics(data$obs, data$sim))$
285       ↪ value)
286   )
287   skill.text = paste0(df.skill$metric, ": ", round(df.skill$value, 2), collapse = "\n")
288   coefs = coef(model)
289   eq.text = paste0("Sim~prec == ", round(coefs[1], 2), "~(cm)")
290   if (length(coefs) > 1) {
291     terms.text = paste0(round(coefs[-1], 1), "%*%", names(coefs)[-1], collapse = " + ")
292     eq.text = paste(eq.text, "+", terms.text)
293   }
294 }
295
296 ggplot(data = data, aes(x = obs, y = sim)) +
297   geom_abline(slope = 1, intercept=0, color="#67000d", linewidth = 1.5) +
298   geom_point(shape = 21, color = "#08306b", fill = "#4292c6", size = 4) +
299   geom_abline(slope = 1, intercept=0, color="#67000d", linetype = 2) +
300   scale_x_continuous(limits = c(min(data$obs), max(data$obs))) +
301   scale_y_continuous(limits = c(min(data$obs), max(data$obs))) +
302   labs(x = "Observed precipitation (cm)", y = "Simulated precipitation (cm)") +
303   # model equation
304   annotate("text", x = min(data$obs), y = max(data$obs), fontface = "bold",
305     label = paste(title, "model"), hjust = 0, vjust = 1, size = 5) +
306
307   annotate("text", x = min(data$obs), y = .96*max(data$obs), parse = T,
308     label = paste(eq.text), hjust = 0, vjust = 1, size = 3.7) +
309   # model performance
310   annotate("text", x = min(data$obs), y = .9*max(data$obs),
311     label = skill.text, hjust = 0, vjust = 1, size = 4.5, fontface = "bold") +
312
313   theme_bw(base_family = "Times") +
314   theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
315     axis.text.y = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
316     axis.title = element_text(size = text.size + 2, face = 'bold'),
317     panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
318     panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
319     panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),

```

```

319     panel.grid.major.x = element_blank(),
320     panel.grid.major.y = element_blank()
321 }
322 # Skill metrics -----
323 evaluation.metrics = function(actual, predicted) {
324   mse = mean((actual - predicted)^2, na.rm = TRUE)
325   rmse = sqrt(mse)
326   r2 = 1 - sum((actual - predicted)^2) / sum((actual - mean(actual))^2)
327   metrics = c(MSE = mse, RMSE = rmse, R2 = r2)
328   return(metrics)
329 }
330 evaluation.model = function(actual, data, model, local = FALSE) {
331   # predict values
332   predicted = predict(model, data = data)
333
334   # calc metrics
335   mse = mean((actual - predicted)^2, na.rm = TRUE)
336   rmse = sqrt(mse)
337   r2 = 1 - sum((actual - predicted)^2) / sum((actual - mean(actual))^2)
338
339   if(local){
340     metrics = c(MSE = mse, RMSE = rmse, R2 = r2)
341   }else{
342     n.params = length(coef(model)) # parameters in the model
343     n.obs = length(actual) # n obs
344     log.likelihood = -0.5 * n.obs * (log(2 * pi * mse) + 1)
345     aic = -2 * log.likelihood + 2 * n.params
346     bic = -2 * log.likelihood + log(n.obs) * n.params
347     metrics = c(MSE = mse, RMSE = rmse, R2 = r2, AIC = aic, BIC = bic)
348   }
349   return(metrics)
350 }
351
352 # Fitting models (lm, glm, gam, locPoly) -----
353
354 # subset selection using regsubsets
355 subset.select = function(data) {
356   subset.model = regsubsets(prec ~ ., data = data, nbest = 10)
357   subset.summ = summary(subset.model)
358   subset.index = which.min(subset.summ$bic)
359   subset.vars = names(coef(subset.model, subset.index))
360   data[, c("prec", subset.vars[subset.vars != "(Intercept)"])]
361 }
362
363 # fitting functions
364 fit_Linear = function(data, obj.fn = "BIC") {
365   # check if response variable "prec" exists in the dataset
366   if (!"prec" %in% names(data)) {

```



```

367   stop("Error: The response variable 'prec' is missing in the dataset.")
368 }
369
370 # Fitting lm models -----
371 # 0) Basic linear model
372 basic.model = lm(prec ~ ., data = data)
373
374 # 1) AIC-based model selection
375 aic.model = stepAIC(lm(prec ~ ., data = data), direction = "backward",
376                   trace = FALSE)
377 # 2) BIC-based model selection
378 bic.model = stepAIC(lm(prec ~ ., data = data), direction = "backward",
379                   k = log(nrow(data)),
380                   trace = FALSE)
381 # 3) Leaps - all subsets regression
382 leaps.model = regsubsets(prec ~ ., data = data, nbest = length(colnames(data)) - 1 )
383 best.subset = summary(leaps.model)
384 model.index = which.min(best.subset$bic)
385 best.vars   = names(coef(leaps.model, model.index))
386 new.train   = data[, c("prec", best.vars[-1])] # rm intercept
387 # refit the best subset model to get coefficients
388 leaps.model = lm(as.formula(paste("prec ~", paste(best.vars[-1], collapse=" + "))),
389                 data = new.train)
390
391 # Compare and return best model -----
392 cat("\nModel selection comparison:")
393 name.models = c("Linear", "AIC Selection", "BIC Selection", "Leaps")
394 fitted.models = list(basic.model, aic.model, bic.model, leaps.model)
395 names(fitted.models) = name.models
396
397 df.skill     = data.frame(
398   model      = name.models,
399   t(sapply(fitted.models, function(fit) evaluation.model(data$prec, data, fit) ) )
400 )
401 metrics = names(df.skill)[-1]
402
403 # select best model based on BIC values (by default)
404 if (obj.fn == "R2") {
405   best.model = df.skill[which.max(df.skill[[obj.fn]]), ]
406 } else {
407   best.model = df.skill[which.min(df.skill[[obj.fn]]), ]
408 }
409
410 cat("\nModel summary:\n")
411 print(df.skill)
412 cat(sprintf("\nBest link function: %s\n", best.model$model))
413 for (metric in metrics) {
414   cat(sprintf(" %s = %.3f\n", metric, best.model[[metric]]))

```



```

415 }
416
417 # return best model
418 return( fitted.models[[best.model$model]] )
419 }
420 fit_GLM = function(data, family, obj.fn = "BIC") {
421   options(warn = -1)
422
423   if (!"prec" %in% names(data)) {
424     stop("Error: The response variable 'prec' is missing in the dataset.")
425   }
426
427   if (family == "binomial") {
428     data$prec = scales::rescale(data$prec, to = c(0, 1))
429   } else {
430     data$prec = ifelse(data$prec <= 0, runif(1, 0.0001, 0.001), data$prec)
431   }
432
433   links = switch(family,
434     "Gamma" = c("log", "inverse", "identity"),
435     "binomial" = c("logit", "probit", "cauchit"),
436     "gaussian" = c("identity"),
437     stop("Invalid family specified"))
438
439   model.links = list()
440   # j=1
441   for (j in seq_along(links)) {
442     family.func = match.fun(family)
443     cat(paste0("\n----- Fitting GLM with ", family, " + ", links[j], " link\n"))
444     init.model = lm(prec ~ ., data = data)
445     start.values = coef(init.model)
446     basic.model = glm(prec ~ ., data = data, family = family.func(link = links[j]),
447       maxit = 1000, start = start.values)
448
449     model.links[[j]] = basic.model
450   }
451
452   names(model.links) = links
453
454   cat("\nModel link comparison:")
455   df.skill = data.frame(
456     model = links,
457     t(sapply(model.links, function(fit) evaluation.model(data$prec, data, fit)))
458   )
459
460   metrics = names(df.skill)[-1]
461
462   if (obj.fn == "R2") {

```

```

463   best.model = df.skill[which.max(df.skill[[obj.fn]]), ]
464 } else {
465   best.model = df.skill[which.min(df.skill[[obj.fn]]), ]
466 }
467
468 cat("\nModel summary:\n")
469 print(df.skill)
470 cat(sprintf("\nBest link function: %s\n", best.model$model))
471 for (metric in metrics) {
472   cat(sprintf(" %s = %.3f\n", metric, best.model[[metric]]))
473 }
474
475 return(model.links[[best.model$model]])
476 }
477 fit_locPoly = function(data, family, obj.fn = "R2") {
478   options(warn = -1)
479
480   if (!"prec" %in% names(data)) {
481     stop("Error: The response variable 'prec' is missing in the dataset.")
482   }
483
484   if (family == "binomial") {
485     data$prec = scales::rescale(data$prec, to = c(0, 1))
486   } else {
487     data$prec = ifelse(data$prec <= 0, runif(1, 0.0001, 0.001), data$prec)
488   }
489
490   # family = tolower(family) #
491   links = switch(family, # "ident", "log", "logit", "inverse", "sqrt" and "arcsin".
492     "binomial" = c("logit", "probit", "arcsin"),
493     "gamma"    = c("log", "inverse"),
494     "gaussian" = c("ident"),
495     "geom"     = c("log"),
496     "poisson"  = c("log", "sqrt"),
497     stop("Invalid family specified"))
498
499   model.links = list()
500   # j=1
501   for (j in seq_along(links)) {
502     cat(paste0("\n----- Fitting local Polynomial with ", family, " + ", links[j], " link\n"))
503
504     # range for alpha order 1 and 2
505     nvars = ncol(data) - 1
506     alpha1 = round(seq(2*(nvars*1 + 1)/(nrow(data)), 1, by = 0.1), 2)
507     alpha2 = round(seq(2*(nvars*2 + 1)/(nrow(data)), 1, by = 0.1), 2)
508     alpha = c(alpha1, alpha2)
509     n.alpha = length(alpha1)
510

```

```

511 # kern options
512 kern.options = c("bisq", "tcub") # "bisq", "tcub", "trwt", "tria"
513 best.model = NULL
514 best.alpha = NULL
515 best.deg = NULL
516 best.gcv = Inf
517 best.kern = NULL
518
519 capture.output({
520   for (kern in kern.options) {
521     cat(paste0("Testing kernel: ", kern, "\n"))
522     gcv.deg1 = tryCatch({
523       gcvplot(prec ~ ., data = data, maxk = 1000, alpha = alpha1,
524         deg = 1, kern = kern, scale = T,
525         ev = dat(), #lfgrid() for grids, dat for points or none()
526         family = family, link = links[j]
527       )
528     }, error = function(e) NULL)
529
530     gcv.deg2 = tryCatch({
531       gcvplot(prec ~ ., data = data, maxk = 1000, alpha = alpha2,
532         deg = 2, kern = kern, scale = T,
533         ev = dat(),
534         family = family, link = links[j]
535       )
536     }, error = function(e) NULL)
537
538     if (!is.null(gcv.deg1) && !is.null(gcv.deg2)) {
539       gcv.values = order(c(gcv.deg1$values, gcv.deg2$values))
540       deg.selected = ifelse(gcv.values[1] > n.alpha, 2, 1)
541       alpha.selected = alpha[gcv.values[1]]
542       gcv.score = c(gcv.deg1$values, gcv.deg2$values)[gcv.values[1]]
543
544       if (gcv.score < best.gcv) {
545         best.gcv = gcv.score
546         best.deg = deg.selected
547         best.alpha = alpha.selected
548         best.kern = kern
549       }
550     }
551   }
552 })
553 cat(sprintf("Kernel: %s | Alpha: %.3f | Degree: %d\n",
554   best.kern, best.alpha, best.deg))
555 basic.model = locfit(prec ~ ., data = data, alpha = best.alpha,
556   maxk = 100, deg = best.deg, kern = best.kern,
557   ev = dat(),
558   scale = T, family = family, link = links[j]

```

```

559         # lfproc = locfit.robust(prec ~ ., data = data)
560     )
561     model.links[[j]] = basic.model
562 }
563 names(model.links) = links
564
565 cat("\nModel link comparison:")
566 df.skill = data.frame(
567     model = links,
568     t(sapply(model.links, function(fit) evaluation.model(data$prec, data, fit, local = T)))
569 )
570
571 metrics = names(df.skill)[-1]
572
573 if (obj.fn == "R2") {
574     best.model = df.skill[which.max(df.skill[[obj.fn]]), ]
575 } else {
576     best.model = df.skill[which.min(df.skill[[obj.fn]]), ]
577 }
578
579 cat("\nModel summary:\n")
580 print(df.skill)
581 cat(sprintf("\nBest link function: %s\n", best.model$model))
582 for (metric in metrics) {
583     cat(sprintf(" %s = %.3f\n", metric, best.model[[metric]]))
584 }
585
586 return(model.links[[best.model$model]])
587 }
588 fit_GAM = function(data, family, obj.fn = "BIC") {
589     options(warn = -1)
590
591     if (!"prec" %in% names(data)) {
592         stop("Error: The response variable 'prec' is missing in the dataset.")
593     }
594
595     if (family == "binomial") {
596         data$prec = scales::rescale(data$prec, to = c(0, 1))
597     } else {
598         data$prec = ifelse(data$prec <= 0, runif(1, 0.0001, 0.001), data$prec)
599     }
600
601     # GAM equation
602     names.vars = setdiff(names(data), "prec")
603     gam.formula = as.formula(
604         paste("prec ~", paste("s(", names.vars, ")", collapse = " + ")) )
605
606     links = switch(family,

```

```

607     "gaussian" = c("identity", "log", "inverse"),
608     "binomial" = c("logit", "probit", "cauchit"),
609     "Gamma"    = c("inverse", "identity", "log"),
610     "poisson"  = c("log", "identity", "sqrt"),
611     stop("Invalid family specified"))
612
613 model.links = list()
614 # j=1
615 for (j in seq_along(links)) {
616     family.func = match.fun(family)
617     cat(paste0("\n----- Fitting GAM with ", family, " + ", links[j], " link\n"))
618
619     # Fit the GAM model
620     basic.model = gam(gam.formula, data = data, method = "REML",
621                     family = family.func(link = links[j]) )
622
623     model.links[[j]] = basic.model
624 }
625
626 names(model.links) = links
627
628 cat("\nModel link comparison:")
629 df.skill = data.frame(
630     model = links,
631     t(sapply(model.links, function(fit) evaluation.model(data$prec, data, fit)))
632 )
633
634 metrics = names(df.skill)[-1]
635
636 if (obj.fn == "R2") {
637     best.model = df.skill[which.max(df.skill[[obj.fn]]), ]
638 } else {
639     best.model = df.skill[which.min(df.skill[[obj.fn]]), ]
640 }
641
642 cat("\nModel summary:\n")
643 print(df.skill)
644 cat(sprintf("\nBest link function: %s\n", best.model$model))
645 for (metric in metrics) {
646     cat(sprintf(" %s = %.3f\n", metric, best.model[[metric]]))
647 }
648
649 return(model.links[[best.model$model]])
650 }
651
652 # Complete regression function
653 # model.type = "locPoly"
654 # data = climaGeo

```

```

655 # coords = coords
656 # topography = topography
657 # suffix.fig = "HW1/P4/FigP4c_"
658 # family = "gaussian"
659 # path = path
660
661 fitRegression = function(model.type, data, coords, topography, path = getwd(),
662                           suffix.fig = NULL, family = NULL,
663                           text.size = 14) {
664   # figures filenames:
665   figure.out = c("i_scatter", "ii_diagnostics", "iii_loocv", "iv_boxplot",
666                 "v_sptMapPred", "v_sptMapRes")
667   figure.path = paste0(path, "/Figures/", suffix.fig, "FitRegression_",
668                       model.type, "_model_", family, "_", figure.out, ".png")
669
670   # 1) Fit the best regression model -----
671   cat(paste0("\n ----- Fitting ", model.type, " model\n"))
672   model = switch(model.type,
673                 "Linear" = fit_Linear( data ),
674                 "GLM"    = fit_GLM(   subset.select(data), family),
675                 "locPoly" = fit_locPoly(subset.select(data), family),
676                 "GAM"    = fit_GAM(   subset.select(data), family),
677                 stop("Unknown model type") )
678
679   cat(paste0("\n ----- Summary\n"))
680   summary(model)
681
682   # 2) Scatter fitted vs observed values -----
683   preds      = predict(model, data = data, se.fit = T)
684   df.scatter  = data.frame(obs = data$prec, sim = preds$fit)
685
686   if(model.type %in% c("locPoly", "GAM")){
687     scatter    = scatter.ggplot(df.scatter, title = model.type)
688   }else{
689     scatter    = scatter.ggplot(df.scatter, model = model, title = model.type)
690   }
691
692   ggsave(filename = figure.path[1],
693          plot      = scatter,
694          width     = 9, height = 6.5, units   = "in")
695
696   # 3) Diagnostics -----
697   # 3.1) ANOVA
698   if(model.type != "locPoly"){
699     anova.model = anova(model)
700     cat(paste0("\n ----- ANOVA\n"))
701     print(anova.model)
702   }else{

```

```

703   anova.model = NULL
704 }
705
706 # 3.2) Model diagnostics
707 df.scatter$res = df.scatter$obs - df.scatter$sim
708 diagnostics    = diagnostics.ggplot(df.scatter)
709 ggsave(filename = figure.path[2],
710         plot      = diagnostics,
711         width     = 9, height = 6.5, units = "in")
712
713 # 4) LOOCV: Drop-One Cross Validation -----
714 cat(paste0("\n ----- LOOCV\n"))
715 cv.preds = numeric(nrow(data))
716 for (ind in seq_len(nrow(data))) {
717   capture.output({
718     temp.model = switch(model.type,
719                         "Linear" = fit_Linear( data[-ind, ]),
720                         "GLM"    = fit_GLM(   data[-ind, ], family),
721                         "locPoly" = fit_locPoly(data[-ind, ], family),
722                         "GAM"    = fit_GAM(   data[-ind, ], family),
723                         stop("Unknown model type"))
724   })
725   if(model.type != "locPoly"){
726     cv.preds[ind] = predict(temp.model,
727                           newdata = data[ind, , drop = FALSE],
728                           type = "response")
729   }else{
730     cv.preds[ind] = predict(temp.model, data[ind,])
731   }
732 }
733 df.loocv = data.frame(obs = data$prec, sim = cv.preds)
734 loocv.scatter = scatter.ggplot(df.loocv, title = "CV")
735 ggsave(filename = figure.path[3],
736         plot      = loocv.scatter,
737         width     = 9, height = 6.5,
738         units     = "in")
739
740 # 5) 10% holdout test -----
741 cat(paste0("\n ----- Drop 10%\n"))
742 nsamples = .5 * nrow(data) # 50% length of data
743 skill.drop10 = data.frame(iteration = 1:nsamples, rmse = NA, rho = NA)
744
745 for (i in 1:nsamples) {
746   # randomly select 10% holdout test data
747   ind.drop = sample(1:nrow(data), size = round(0.1 * nrow(data)))
748   train.data = data[-ind.drop, ]
749   test.data = data[ind.drop, ]
750 }

```

```

751 capture.output({
752   temp.model = switch(model.type,
753     "Linear" = fit_Linear( train.data),
754     "GLM"    = fit_GLM(   train.data, family),
755     "locPoly" = fit_locPoly(train.data, family),
756     "GAM"    = fit_GAM(   train.data, family),
757     stop("Unknown model type") )
758 })
759 new.pred = predict(temp.model, data = test.data)
760 # metrics
761 skill.drop10$rmse[i] = sqrt(mean((new.pred - test.data$prec)^2))
762 skill.drop10$rho[i] = cor(new.pred[ind.drop], test.data$prec, use = "pairwise.complete.
    ↪ obs")
763 }
764 skill.drop10$iteration = NULL
765 df.drop10 = melt(skill.drop10)
766 levels(df.drop10$variable) = c("RMSE (cm)", "Correlation")
767
768 boxplot.skill =
769   ggplot(data = df.drop10, aes(x = value)) +
770   geom_boxplot(color = "#005a32", fill = "#a1d99b", outlier.size = 2, shape = 21) +
771   scale_x_continuous(breaks = scales::pretty_breaks(n=3)) +
772   facet_wrap(~variable, scales = "free_x") +
773   theme_bw(base_family = "Times") +
774   theme(axis.text.x = element_text(size = text.size, hjust = .5, vjust = .5),
775         axis.title.x = element_blank(),
776         axis.text.y = element_blank(),
777         axis.ticks.y = element_blank(),
778         strip.text = element_text(size = text.size, face = "bold"),
779         strip.background = element_blank(),
780         panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
781         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
782         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
783         panel.grid.major.x = element_blank(),
784         panel.grid.major.y = element_blank())
785 ggsave(filename = figure.path[4],
786       plot = boxplot.skill,
787       width = 9, height = 4,
788       units = "in")
789
790
791 # 6) spatial Mapping -----
792 cat(paste0("\n ----- Spatial Mapping\n"))
793 map.fit = cbind(coords, "prec" = preds$fit)
794 map.res = cbind(coords, "prec" = preds$se.fit)
795
796 # save predicted values

```



```

797 fit.maps      = sptMap.ggplot(map.fit, topography, title.plot = "Predicted Precipitation (
      ↪ cm)")
798 ggsave(filename = figure.path[5],
799       plot      = fit.maps,
800       width     = 10, height = 5.5,
801       units     = "in", bg = "#fff")
802
803 res.maps      = sptMap.ggplot(map.res, topography, title.plot = "Standard Error (cm)")
804 ggsave(filename = figure.path[6],
805       plot      = res.maps,
806       width     = 10, height = 5.5,
807       units     = "in", bg = "#fff")
808
809 # list
810 return( list( model      = model,          # i) best model
811             scatter     = scatter,        # ii) scatterplot of obs and sim prec + 1:1 line
812             anova       = anova.model,    # iii) anova
813             mod.diag    = diagnostics,    # iii) diagnostics: assumption of residuals
814             df.loocv    = df.loocv,      # iv) data loocv
815             cv.scatter  = loocv.scatter,  # iv) loocv scatter + 1:1 line
816             df.drop10   = df.drop10,     # v) drop 10% obs df
817             box.drop10  = boxplot.skill,  # v) boxplot rmse + rho from drop 10%
818             pred.sptMap = fit.maps,      # vi) spatially map for fitted values
819             se.sptMap  = res.maps        # vi) spatially map for SE values
820         )
821     )
822 }
823
824 # NLS functions -----
825
826 # Potential functions
827 # c0 is the nugget effect
828 # c1 is the sill
829 # a is the range parameter
830 exp.model = function(x, c0, c1, a) { # Exponential
831   c0 + c1 * (1 - exp(-x / a))
832 }
833
834 sph.model = function(x, c0, c1, a) { # Spherical (x < a, geostatistics)
835   c0 + c1 * (1.5 * (x / a) - 0.5 * (x / a)^3) * (x < a) + c0 + c1 * (x >= a)
836 }
837
838 gau.model = function(x, c0, c1, a) { # Gaussian (highly continuous data)
839   c0 + c1 * (1 - exp(-(x^2) / (a^2)))
840 }
841
842 pow.model = function(x, c0, c1, p) { # Power (data with no clear sill)
843   if(0 < p & p < 2){

```

```

844   c0 + c1 * (x^p)
845 }else{
846   cat("Error: Enter a valid value for p. Range 0 to 2")
847 }
848 }
849
850 # estimate initial parameters dynamically
851 set.initialParams = function(data) {
852   init.c0 = max(0, round(min(data$value, na.rm = TRUE), 2)) # avoid neg values
853   init.c1 = max(1, round(max(data$value, na.rm = TRUE), 2) - init.c0) # pos range
854   init.a = max(1, round(max(data$distance, na.rm = TRUE), 2) * 0.3)
855   init.p = 1 # power exponent default
856   list(c0 = init.c0, c1 = init.c1, a = init.a, p = init.p)
857 }
858
859 # fit nonlinear least squares models
860 fit_nls = function(data, model) {
861   init.params = set.initialParams(data)
862
863   if (any(is.na(init.params)) || any(!is.finite(unlist(init.params)))) {
864     return(NULL)
865   }
866
867   if(identical(model, pow.model)){ # pow model with p param
868     tryCatch({
869       temp.nls = nlsLM(
870         value ~ model(distance, c0, c1, p),
871         data = data,
872         start = list(c0 = init.params$c0,
873                     c1 = init.params$c1,
874                     p = init.params$p),
875         control = nls.lm.control(maxiter = 2000)
876       )
877       # skill.nls = evaluation.model(data$value, data, temp.nls)
878       # pred.nls = predict(temp.nls)
879       temp.nls
880     }, error = function(e) { return(NULL) })
881   } else{ # e.g., exp, gaussian models with a param
882     tryCatch({
883       temp.nls = nlsLM(
884         value ~ model(distance, c0, c1, a),
885         data = data,
886         start = list(c0 = init.params$c0,
887                     c1 = init.params$c1,
888                     a = init.params$a),
889         control = nls.lm.control(maxiter = 2000)
890       )
891       # skill.nls = evaluation.model(data$value, data, temp.nls)

```

```

892     # pred.nls = predict(temp.nls)
893     temp.nls
894 }, error = function(e) { return(NULL) })
895
896 } # end model condition
897
898 }
899
900 # Hierarchical Spatial Model -----
901 empirical.vgram = function(res, data, type, n.bins = 50) {
902   emp    = vgram(loc = data[, c("lon", "lat")], y = res, type = type, lon.lat = TRUE)
903   emp.bin = getVGMean(x = emp, breaks = pretty(emp$d, n = n.bins, eps.correct = 1))
904   data.frame(
905     distance = emp.bin$centers[-length(emp.bin$centers)],
906     value    = emp.bin$ys[-length(emp.bin$ys)],
907     type     = type
908   )
909 }
910
911 fit_HSM      = function(data, coords, topography,
912                          obj.fn = "BIC", type = "variogram", # covariogram or correlogram
913                          path = getwd(), model.type = "HSM",
914                          suffix.fig = NULL, family = NULL,
915                          text.size = 14) {
916   options(warn = -1)
917   # figures filenames:
918   figure.out = c("0_emp_binned", "i_scatter", "ii_diagnostics", "iii_loocv", "iv_boxplot",
919                 "v_sptMapPred", "v_sptMapRes")
920   figure.path = paste0(path, "/Figures/", suffix.fig, "FitRegression_",
921                        model.type, "_model_", family, "_", figure.out, ".png")
922
923   if (!"prec" %in% names(data)) {
924     stop("Error: The response variable 'prec' is missing in the dataset.")
925   }
926
927   # 1) fit lm model -----
928   lm.model = fit_Linear(data)
929   lm.res   = residuals(lm.model)
930   data$res = lm.res
931
932   # 2) Binned empirical Var, Cov, Cor -----
933   empirical.data = empirical.vgram(data$res, data, type)
934   empirical.data = empirical.data[empirical.data$distance < .5*max(empirical.data$distance),
935     ↪ ]
936   empirical.data = empirical.data %>% filter(!is.na(value) & is.finite(value))
937   bin.type      = unique(empirical.data$type)
938   empirical.data$type = factor(empirical.data$type,
939                               levels = bin.type)

```

```

939
940 # 3) Fitting with NLS -----
941 pot.func      = c("exp", "spherical", "gaussian")
942 fit.mod       = list(exp.model, sph.model, gau.model)
943 names(fit.mod) = pot.func
944 model.nls     = list()
945 # j=1
946 for (j in seq_along(pot.func)) {
947   cat(paste0("\n----- Fitting HSM with ", pot.func[j], " potential function \n"))
948   basic.model  = fit_nls(empirical.data, fit.mod[[j]])
949   model.nls[[j]] = basic.model
950 }
951 names(model.nls) = pot.func
952
953 # 4) NLS comparison -----
954 cat("\nModel link comparison:")
955 df.skill = data.frame(
956   model = pot.func,
957   t(sapply(model.nls, function(fit) evaluation.model(empirical.data$value,
958                                                         empirical.data, fit)))
959 )
960
961 metrics = names(df.skill)[-1]
962
963 if (obj.fn == "R2") {
964   best.model = df.skill[which.max(df.skill[[obj.fn]]), ]
965 } else {
966   best.model = df.skill[which.min(df.skill[[obj.fn]]), ]
967 }
968
969 cat("\nModel summary:\n")
970 print(df.skill)
971 cat(sprintf("\nBest potential function: %s\n", best.model$model))
972 for (metric in metrics) {
973   cat(sprintf(" %s = %.3f\n", metric, best.model[[metric]]))
974 }
975
976 best.skill = as.data.frame(t(best.model[,-1]))
977 colnames(best.skill) = "value"
978 best.skill$metric = row.names(best.skill)
979
980 skill.text = paste0(best.skill$metric, ": ", round(best.skill$value, 2), collapse = "\n")
981 coefs      = coef(model.nls[[best.model$model]])
982 eq.text    = paste0("Sim~ == ", round(coefs[1], 2), "~")
983 if (length(coefs) > 1) {
984   terms.text = paste0(round(coefs[-1], 1), "%*%", names(coefs)[-1], collapse = " + ")
985   eq.text    = paste(eq.text, "+", terms.text)
986 }

```

```

987 empirical.data$pred = predict(model.nls[[best.model$model]], data = empirical.data)
988
989 # Plot empirical data + fit model
990 emp.plot =
991   ggplot(empirical.data) +
992   geom_abline(slope = 0, intercept=0, color = "#000",linetype = 2, linewidth = .5) +
993   # empirical data
994   geom_line(aes(x = distance, y = value), color = "#08306b", linewidth = .9, alpha = .7)
995   ↪ +
996   geom_point(aes(x = distance, y = value), color = "#08306b", fill= "#6baed6",
997             size = 2, shape = 21) +
998   # predicted data
999   geom_line(aes(x = distance, y = pred), color="#67000d", linewidth = 1, linetype = "
1000   ↪ dashed") +
1001   scale_x_continuous(breaks = pretty_breaks(n = 3)) +
1002   scale_y_continuous(breaks = pretty_breaks(n = 3)) +
1003   labs(x = "Distance", y = "Value") +
1004   facet_wrap(~type, scales = "free_y", labeller = as_labeller(c(
1005     variogram = "Binned Empirical Variogram",
1006     correlogram = "Binned Empirical Correlogram",
1007     covariogram = "Binned Empirical Covariance"
1008   ))) +
1009   theme_bw(base_family = "Times") +
1010   theme(axis.text.x = element_text(size = text.size, vjust = 0.5),
1011         axis.text.y = element_text(size = text.size, angle = 90, hjust = .5, vjust = .5),
1012         axis.title = element_text(size = text.size + 2, face = 'bold'),
1013         legend.position = "none",
1014         strip.background = element_blank(),
1015         strip.text = element_text(size = text.size + 2, face = 'bold'),
1016         panel.border = element_rect(color = "#000000", fill = NA, linewidth = .5),
1017         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
1018         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
1019         panel.grid.major.x = element_blank(),
1020         panel.grid.major.y = element_blank())
1021 # model performance
1022 if(type=="variogram"){
1023   emp.plot = emp.plot +
1024   annotate("text", x = .1*max(empirical.data$distance), y = .95*max(empirical.data$
1025   ↪ value),
1026           label = skill.text, hjust = 0, vjust = 1, size = 4.5, fontface = "bold")
1027 }else{
1028   emp.plot = emp.plot +
1029   annotate("text", x = .85*max(empirical.data$distance), y = .95*max(empirical.data$
1030   ↪ value),
1031           label = skill.text, hjust = 0, vjust = 1, size = 4.5, fontface = "bold")
1032 }
1033
1034 ggsave(filename = figure.path[1],

```

```

1031     plot      = emp.plot,
1032     width     = 9, height = 6.5, units   = "in")
1033
1034     ## 5) Kriging on the residuals -----
1035     model      = model.nls[[best.model$model]]
1036     params.fit = as.numeric(coef(model))
1037     ## Check with Krig and predict.krig
1038     krig.fit   = Krig(x = data,
1039                     Y = data$prec,
1040                     theta = params.fit[3],
1041                     lambda = params.fit[1]/params.fit[2],
1042                     m=1)
1043     data$krig.fit = predict.Krig(krig.fit)
1044     data$krig.se  = predictSE(krig.fit)
1045
1046     ## 6) Scatter fitted vs observed values -----
1047     df.scatter    = data.frame(obs = data$prec, sim = data$krig.fit)
1048     scatter       = scatter.ggplot(df.scatter, model = model, title = model.type)
1049     ggsave(filename = figure.path[2],
1050             plot      = scatter,
1051             width     = 9, height = 6.5, units   = "in")
1052
1053     ## 7) Diagnostics -----
1054     df.scatter$res = df.scatter$obs - df.scatter$sim
1055     diagnostics    = diagnostics.ggplot(df.scatter)
1056     ggsave(filename = figure.path[3],
1057             plot      = diagnostics,
1058             width     = 9, height = 6.5, units = "in")
1059
1060     ## 8) LOOCV: Drop-One Cross Validation -----
1061     cat(paste0("\n ----- LOOCV\n"))
1062     cv.preds = numeric(nrow(data))
1063     for (ind in seq_len(nrow(data))) {
1064         capture.output({
1065             temp.model = Krig(x = data[-ind,],
1066                             Y = data$prec[-ind],
1067                             theta = params.fit[3],
1068                             lambda = params.fit[1]/params.fit[2],
1069                             m=1)
1070         })
1071         cv.preds[ind] = predict.Krig(temp.model,
1072                                     newdata = data[ind, , drop = FALSE],
1073                                     type = "response")
1074     }
1075
1076     df.loocv      = data.frame(obs = data$prec, sim = cv.preds)
1077     loocv.scatter  = scatter.ggplot(df.loocv, title = "CV")
1078     ggsave(filename = figure.path[4],

```

```

1079     plot      = loocv.scatter,
1080     width     = 9, height = 6.5,
1081     units     = "in")
1082
1083 # 9) 10% holdout test -----
1084 cat(paste0("\n ----- Drop 10%\n"))
1085 nsamples     = .5 * nrow(data) # 50% length of data
1086 skill.drop10 = data.frame(iteration = 1:nsamples, rmse = NA, rho = NA)
1087
1088 for (i in 1:nsamples) {
1089   # randomly select 10% holdout test data
1090   ind.drop    = sample(1:nrow(data), size = round(0.1 * nrow(data)))
1091   train.data  = data[-ind.drop, ]
1092   test.data   = data[ind.drop, ]
1093
1094   capture.output({
1095     temp.model = Krig(x = train.data,
1096                      Y = train.data$prec,
1097                      theta = params.fit[3],
1098                      lambda = params.fit[1]/params.fit[2],
1099                      m=1)
1100   })
1101   new.pred    = predict.Krig(temp.model, data = test.data)
1102   # metrics
1103   skill.drop10$rmse[i] = sqrt(mean((new.pred - test.data$prec)^2))
1104   skill.drop10$rho[i]  = cor(new.pred[ind.drop], test.data$prec, use = "pairwise.complete.
    ↪ obs")
1105 }
1106 skill.drop10$iteration = NULL
1107 df.drop10 = melt(skill.drop10)
1108 levels(df.drop10$variable) = c("RMSE (cm)", "Correlation")
1109
1110 boxplot.skill =
1111   ggplot(data = df.drop10, aes(x = value))+
1112   geom_boxplot(color = "#005a32", fill = "#a1d99b", outlier.size = 2, shape = 21)+
1113   scale_x_continuous(breaks = scales::pretty_breaks(n=3)) +
1114   facet_wrap(~variable, scales = "free_x") +
1115   theme_bw(base_family = "Times") +
1116   theme(axis.text.x      = element_text(size = text.size, hjust = .5, vjust = .5),
1117         axis.title.x     = element_blank(),
1118         axis.text.y      = element_blank(),
1119         axis.ticks.y     = element_blank(),
1120         strip.text       = element_text(size = text.size, face = "bold"),
1121         strip.background = element_blank(),
1122         panel.border      = element_rect(color = "#000000", fill = NA, linewidth = .5),
1123         panel.grid.minor.x = element_line(color = "#d9d9d9", linetype = "dotted"),
1124         panel.grid.minor.y = element_line(color = "#d9d9d9", linetype = "dotted"),
1125         panel.grid.major.x = element_blank(),

```

```

1126     panel.grid.major.y = element_blank()
1127 ggsave(filename = figure.path[5],
1128     plot    = boxplot.skill,
1129     width   = 9, height = 4,
1130     units   = "in")
1131
1132 # 10) spatial Mapping -----
1133 cat(paste0("\n----- Spatial Mapping\n"))
1134 map.fit    = cbind(coords, "prec" = data$krig.fit)
1135 map.res    = cbind(coords, "prec" = data$krig.se)
1136
1137 # save predicted values
1138 fit.maps   = sptMap.ggplot(map.fit, topography, title.plot = "Predicted Precipitation (
    ↪ cm)")
1139 ggsave(filename = figure.path[6],
1140     plot    = fit.maps,
1141     width   = 10, height = 5.5,
1142     units   = "in", bg = "#fff")
1143
1144 res.maps   = sptMap.ggplot(map.res, topography, title.plot = "Standard Error (cm)")
1145 ggsave(filename = figure.path[7],
1146     plot    = res.maps,
1147     width   = 10, height = 5.5,
1148     units   = "in", bg = "#fff")
1149
1150 # list
1151 return( list( emp.plot = emp.plot, # binned empirical plot.
1152     model    = model,      # i) best model
1153     scatter  = scatter,    # ii) scatterplot of obs and sim prec + 1:1 line
1154     mod.diag = diagnostics, # iii) diagnostics: assumption of residuals
1155     df.loocv  = df.loocv,   # iv) data loocv
1156     cv.scatter = loocv.scatter, # iv) loocv scatter + 1:1 line
1157     df.drop10 = df.drop10,  # v) drop 10% obs df
1158     box.drop10 = boxplot.skill, # v) boxplot rmse + rho from drop 10%
1159     pred.sptMap = fit.maps,   # vi) spatially map for fitted values
1160     se.sptMap  = res.maps     # vi) spatially map for SE values
1161 )
1162 )
1163 }

```

C.2. R code

```

1 # Code details -----
2 # author: Catalina Jerez
3 # mail: catalina.jerez@colorado.edu
4 # Course: Advanced Data Analysis Techniques
5 # Code: CVEN 6833

```



```

6 # Homework #1: Problem 2c-5 fit models linear, glm, gam
7 # Department of Civil, Environmental and Architectural Engineering
8 # College of Engineering & Applied Science
9 # University of Colorado Boulder
10
11 # Library -----
12 gc()
13 rm(list = ls())
14
15 # Load necessary libraries
16 spatialAnalysis.Lib = c("sp", "sf", "geosphere", "geoR")
17 smoothing.Lib = c("mgcv", "locfit", "akima")
18 statisticalAnalysis.Lib = c("MASS", "psych", "hydroGOF", "fields", "leaps", "car")
19 dataManipulation.Lib = c("dplyr", "reshape2", "reshape", "tidyr")
20 dataVisualization.Lib = c("ggplot2", "GGally", "ggpubr",
21                           "scales", "maps",
22                           "rnatuarearth", "rnatuarearthdata")
23 list.packages = unique(c(spatialAnalysis.Lib, smoothing.Lib,
24                          statisticalAnalysis.Lib, dataManipulation.Lib,
25                          dataVisualization.Lib))
26
27 # Load all required packages
28 sapply(list.packages, require, character.only = TRUE)
29
30 # Path -----
31 path = "/Users/caje3244/OneDrive - UCB-O365/2025 Spring/CVEN 6833/"
32 path.code = file.path(path, "Codes/")
33 path.plot = file.path(path, "Figures/")
34
35 # load functions to perform regression
36 source(file = file.path(path.code, "HW1_library.R"))
37
38 # Data -----
39 topography = read.table("http://civil.colorado.edu/~balajir/CVEN6833/HWs/HW-1/india-
    ↪ grid-topo.txt")
40 colnames(topography) = c("lon", "lat", "elev")
41
42 # new climatol data with calculated distances from Geodesic method
43 climaGeo = readRDS(file = paste0(path.code, "RData/hw1_climatol_distGeodesic.rds"))
44 climaGeo = na.omit(climaGeo)
45 climaGeo$prec = climaGeo$prec/10 # mm to cm
46 coords = dplyr::select(climaGeo, c(lon, lat))
47 # scale all vars dependent of lon and lat
48 climaGeo = climaGeo %>%
49   dplyr::mutate(across(!c(prec), scale, center = TRUE, scale = TRUE)) %>%
50   dplyr::mutate(across(!c(prec), as.numeric))
51
52 scale.color = c("#d73027", "#313695")

```

```

53 scale.label = unique(cut_interval(climaGeo$prec, n=2))
54 scatter.prec = pairs.ggplot(climaGeo)
55
56 # ggsave(filename = file.path(path.plot, "HW1_FigureP2c_Scatterplot_v0.png"),
57 #         plot      = scatter.prec,
58 #         width     = 12, height = 7.5, units = "in")
59
60 # Basic analysis -----
61 set.seed(1113)
62 # v0) basic linear model
63 str(data)
64 basic.model = lm(prec ~ ., data = climaGeo)
65 # If the VIF is larger than 1/(1-R2)
66 vif(basic.model) > (1/(1-summary(basic.model)$r.squared)) # check multicollinearity
67
68 cat("Since for some variables VIF is larger than 1/(1-R2), let's remove:
69 mdArabian and mdBengal and just keep the mdCoast (minimun distance to any coast) \n")
70
71 # v1) linear model
72 climaGeo = dplyr::select(climaGeo, -c(mdArabian, mdBengal))
73 basic.model = lm(prec ~ ., data = climaGeo)
74 vif(basic.model) > (1/(1-summary(basic.model)$r.squared)) # check multicollinearity
75
76 scatter.prec = pairs.ggplot(climaGeo)
77 # ggsave(filename = file.path(path.plot, "HW1_FigureP2c_Scatterplot_v1.png"),
78 #         plot      = scatter.prec,
79 #         width     = 12, height = 7.5, units = "in")
80
81 # v2) basic linear + interacterions
82 cat("Cool, for all the variables VIF < 1/(1-R2)
83 Let's add some interaciones between the variables \n")
84
85 climaGeo = climaGeo %>%
86   dplyr::mutate(lld = lon * lat * mdCoast,
87                ll2d = lon * lat^2 * mdCoast,
88                l2ld = lon^2 * lat * mdCoast,
89                ed   = elev * mdCoast
90                )
91
92 basic.model = lm(prec ~ ., data = climaGeo)
93 vif(basic.model) > (1/(1-summary(basic.model)$r.squared)) # check multicollinearity
94
95 scatter.prec = pairs.ggplot(climaGeo)
96 # ggsave(filename = file.path(path.plot, "HW1_FigureP2c_Scatterplot_v2.png"),
97 #         plot      = scatter.prec,
98 #         width     = 12, height = 7.5, units = "in")
99
100 cat("Since for all the variables VIF < 1/(1-R2)

```

```

101 Let's proceed with other regression models \n")
102
103 # Fitting linear model -----
104 fitting.linear = fitRegression(model.type = "Linear", data = climaGeo,
105                               coords = coords, topography = topography,
106                               suffix.fig = "HW1/P2/FigP2c_",
107                               path = path)
108 names(fitting.linear)
109 fitting.linear$anova
110 fitting.linear$box.drop10
111 # Fitting GLM model -----
112 fitting.glm_Gau = fitRegression(model.type = "GLM", data = climaGeo,
113                                coords = coords, topography = topography,
114                                suffix.fig = "HW1/P3/FigP3c_", family = "gaussian",
115                                path = path)
116 fitting.glm_Gau$anova
117 fitting.glm_Gau$se.sptMap
118
119 fitting.glm_Bin = fitRegression(model.type = "GLM", data = climaGeo,
120                                coords = coords, topography = topography,
121                                suffix.fig = "HW1/P3/FigP3c_", family = "binomial",
122                                path = path)
123 fitting.glm_Bin$se.sptMap
124 fitting.glm_Bin$anova
125 fitting.glm_Gamma = fitRegression(model.type = "GLM", data = climaGeo,
126                                   coords = coords, topography = topography,
127                                   suffix.fig = "HW1/P3/FigP3c_", family = "Gamma",
128                                   path = path)
129 fitting.glm_Gamma$scatter
130 fitting.glm_Gamma$cv.scatter
131
132 # Fitting Local Polynomial -----
133 fitting.loc_Gau = fitRegression(model.type = "locPoly", data = climaGeo,
134                                 coords = coords, topography = topography,
135                                 suffix.fig = "HW1/P4/FigP4c_", family = "gaussian",
136                                 path = path)
137
138 fitting.loc_Gamma = fitRegression(model.type = "locPoly", data = climaGeo,
139                                   coords = coords, topography = topography,
140                                   suffix.fig = "HW1/P4/FigP4c_", family = "gamma",
141                                   path = path)
142 fitting.loc_Geom = fitRegression(model.type = "locPoly", data = climaGeo,
143                                  coords = coords, topography = topography,
144                                  suffix.fig = "HW1/P4/FigP4c_", family = "geom",
145                                  path = path)
146 fitting.loc_Geom = fitRegression(model.type = "locPoly", data = climaGeo,
147                                  coords = coords, topography = topography,
148                                  suffix.fig = "HW1/P4/FigP4c_", family = "poisson",

```

```

149                                     path = path)
150 # Fitting GAM -----
151 fitting.GAM__Gau = fitRegression(model.type = "GAM", data = climaGeo,
152                                 coords = coords, topography = topography,
153                                 suffix.fig = "HW1/P5/FigP5c_", family = "gaussian",
154                                 path = path)
155 fitting.GAM__Gau$anova
156
157 fitting.GAM__Bin = fitRegression(model.type = "GAM", data = climaGeo,
158                                 coords = coords, topography = topography,
159                                 suffix.fig = "HW1/P5/FigP5c_", family = "binomial",
160                                 path = path)
161 fitting.GAM__Gamma= fitRegression(model.type = "GAM", data = climaGeo,
162                                 coords = coords, topography = topography,
163                                 suffix.fig = "HW1/P5/FigP5c_", family = "Gamma",
164                                 path = path)
165
166 # Fitting HSM -----
167 # Hierarchical Spatial Model
168 fitting.HSM = fit_HSM(data = climaGeo, coords, topography,
169                       type = "variogram", # covariogram or correlogram
170                       suffix.fig = "HW1/P6/FigP6_", path = path)
171 fitting.HSM$se.sptMap
172 fitting.HSM$scatter

```