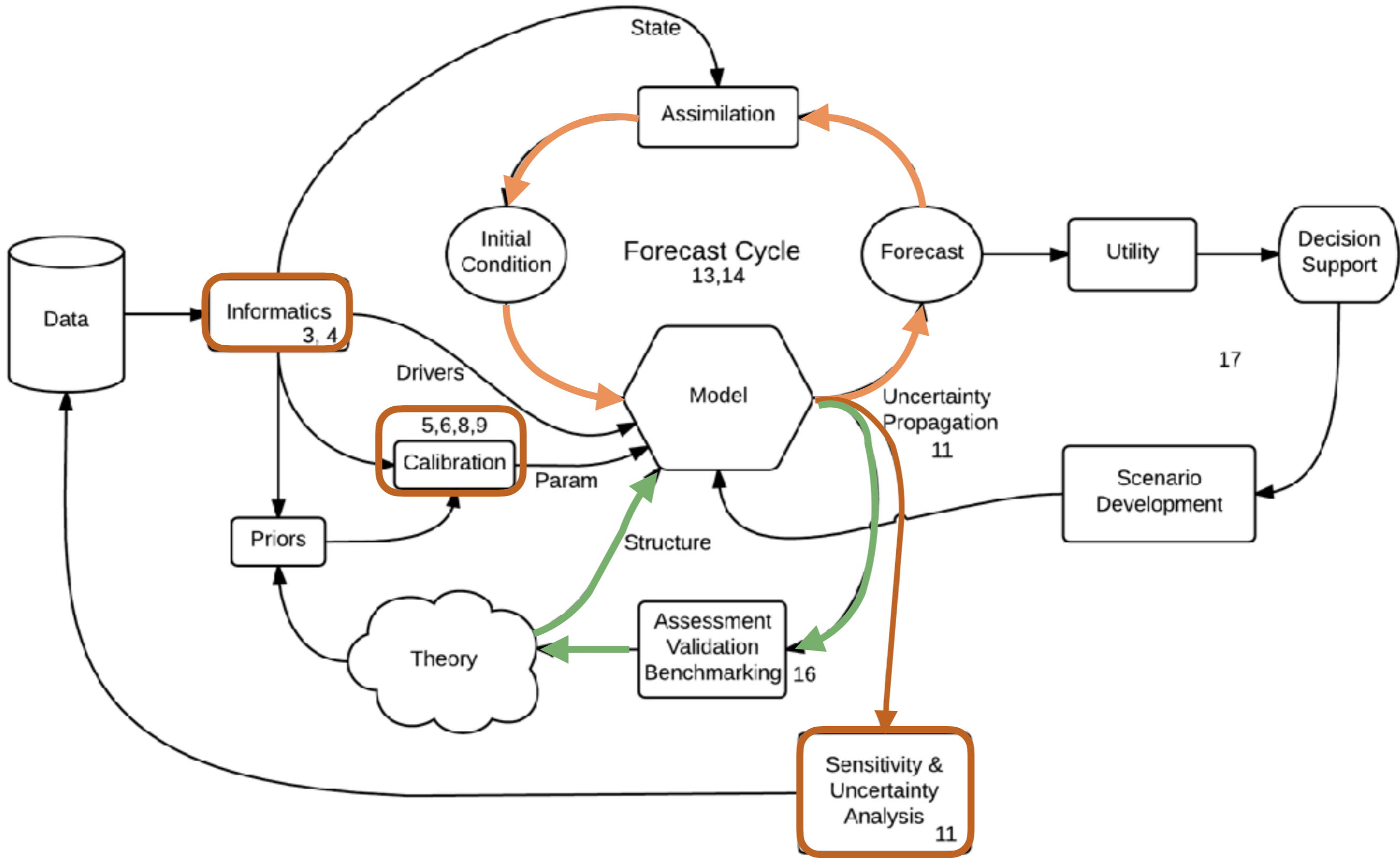


# Assessing Model Performance



- Range of values
- Units
- General pattern in time & space

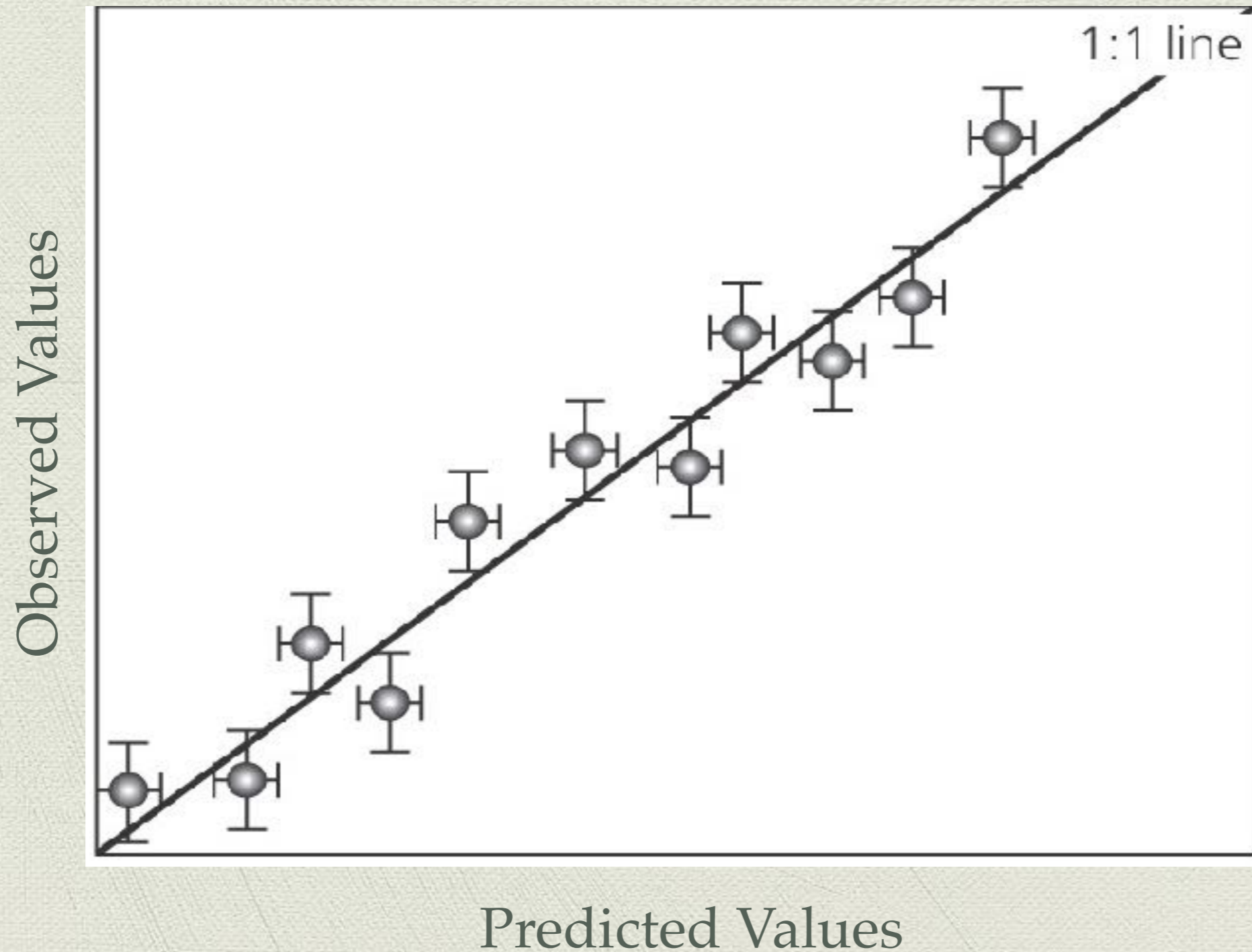


Step 1: Is the model  
output reasonable?

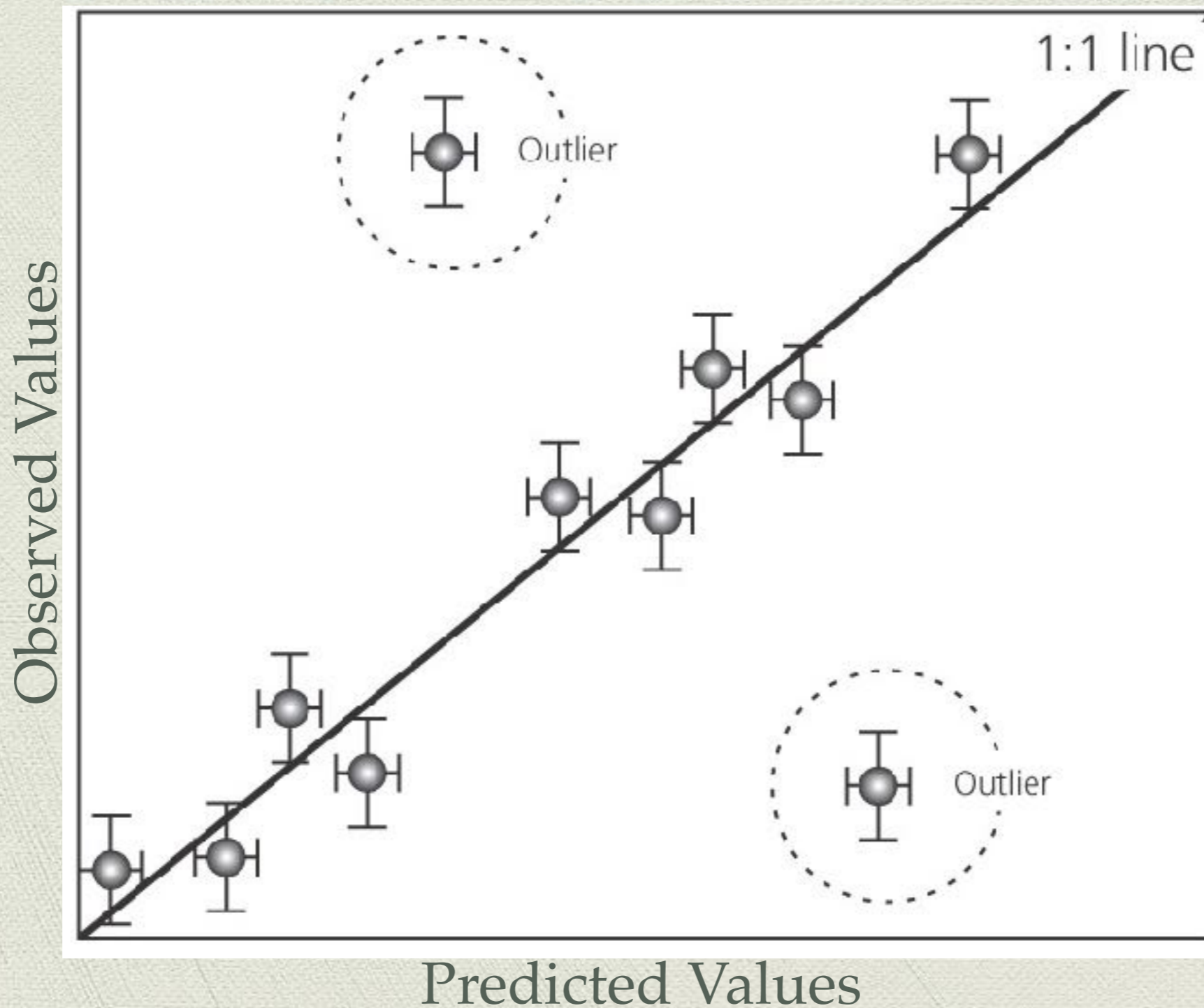


Step 2: Graphical  
comparisons to data

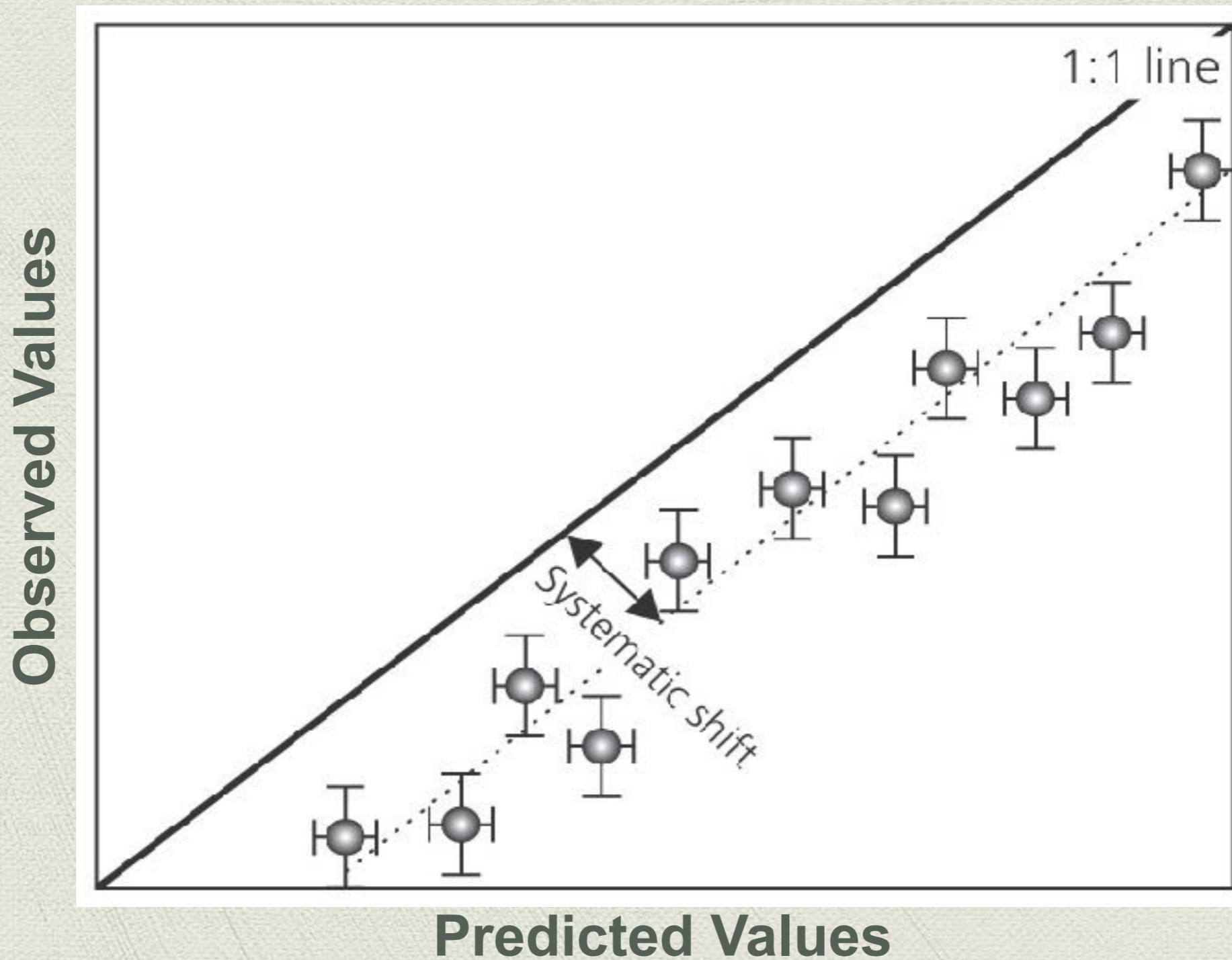
# Accuracy of Prediction



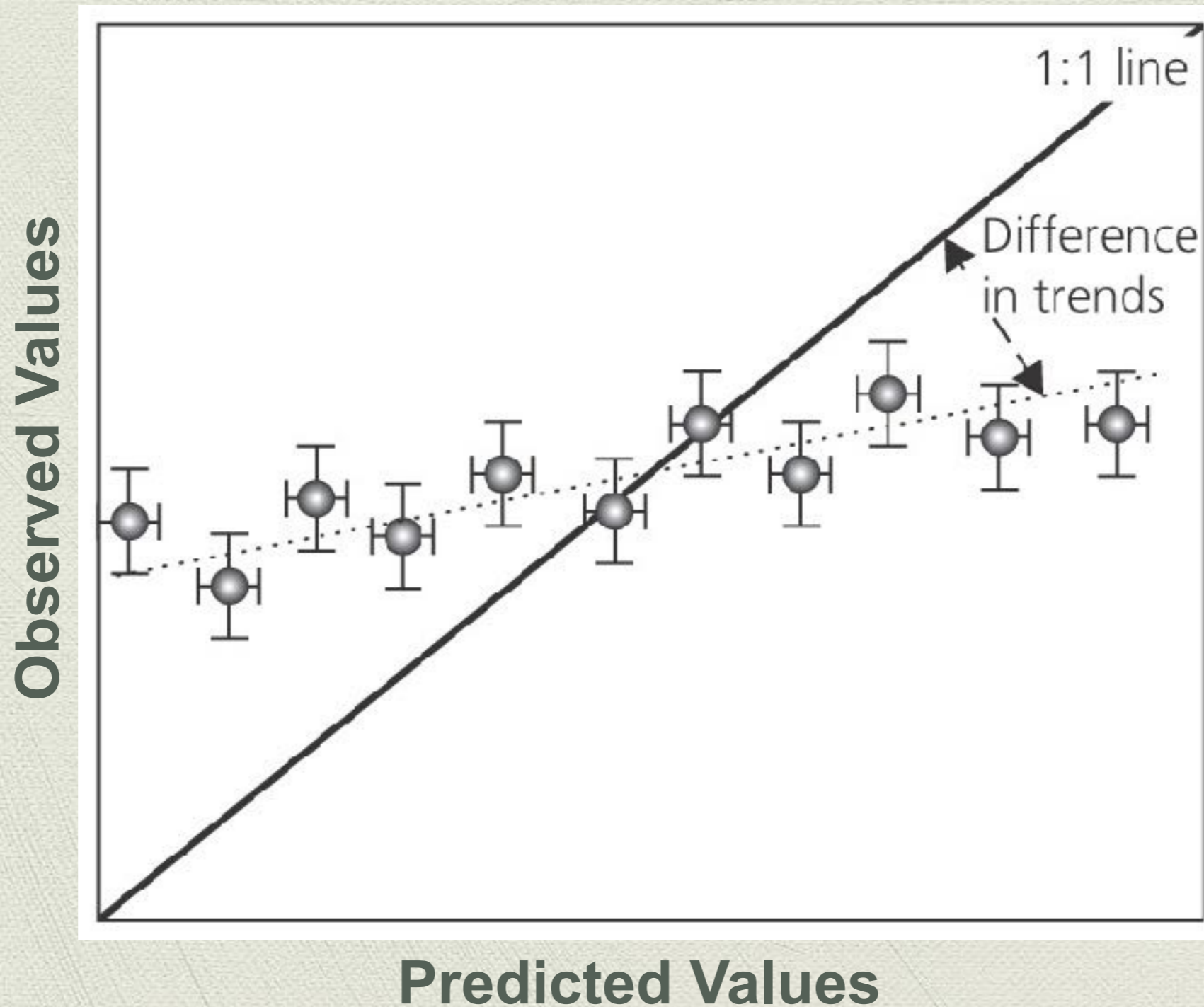
# Identify Outliers



# Assess Biases

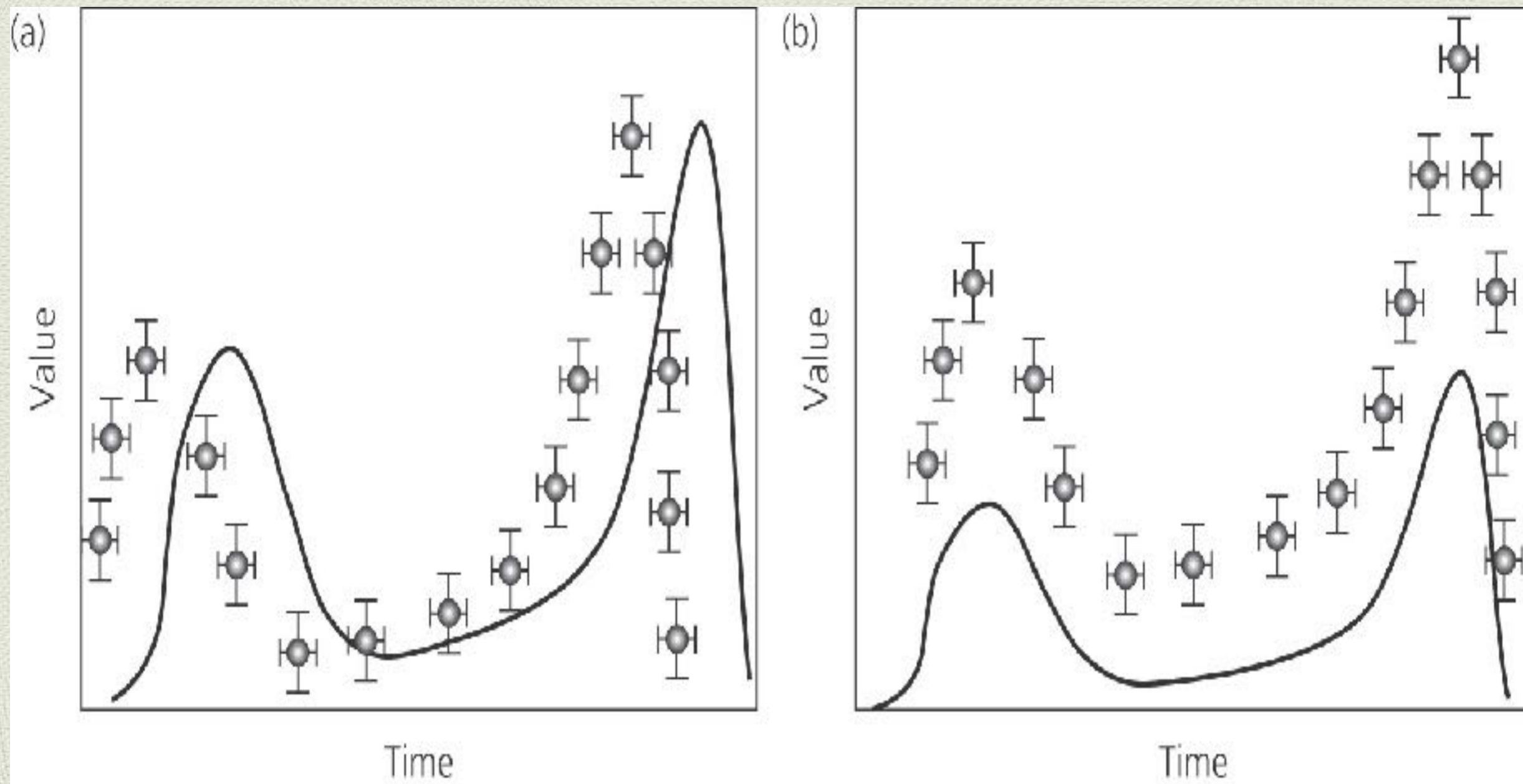


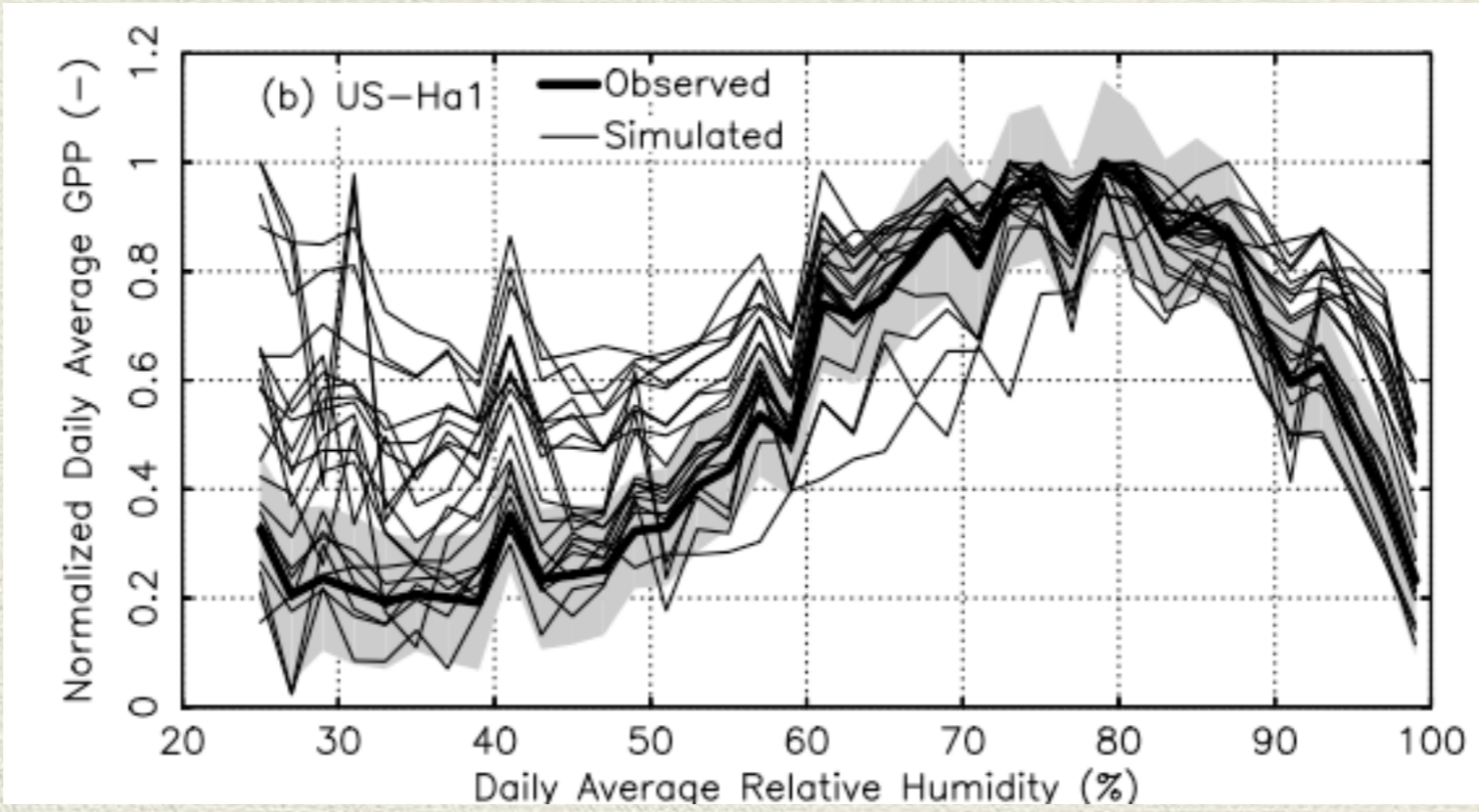
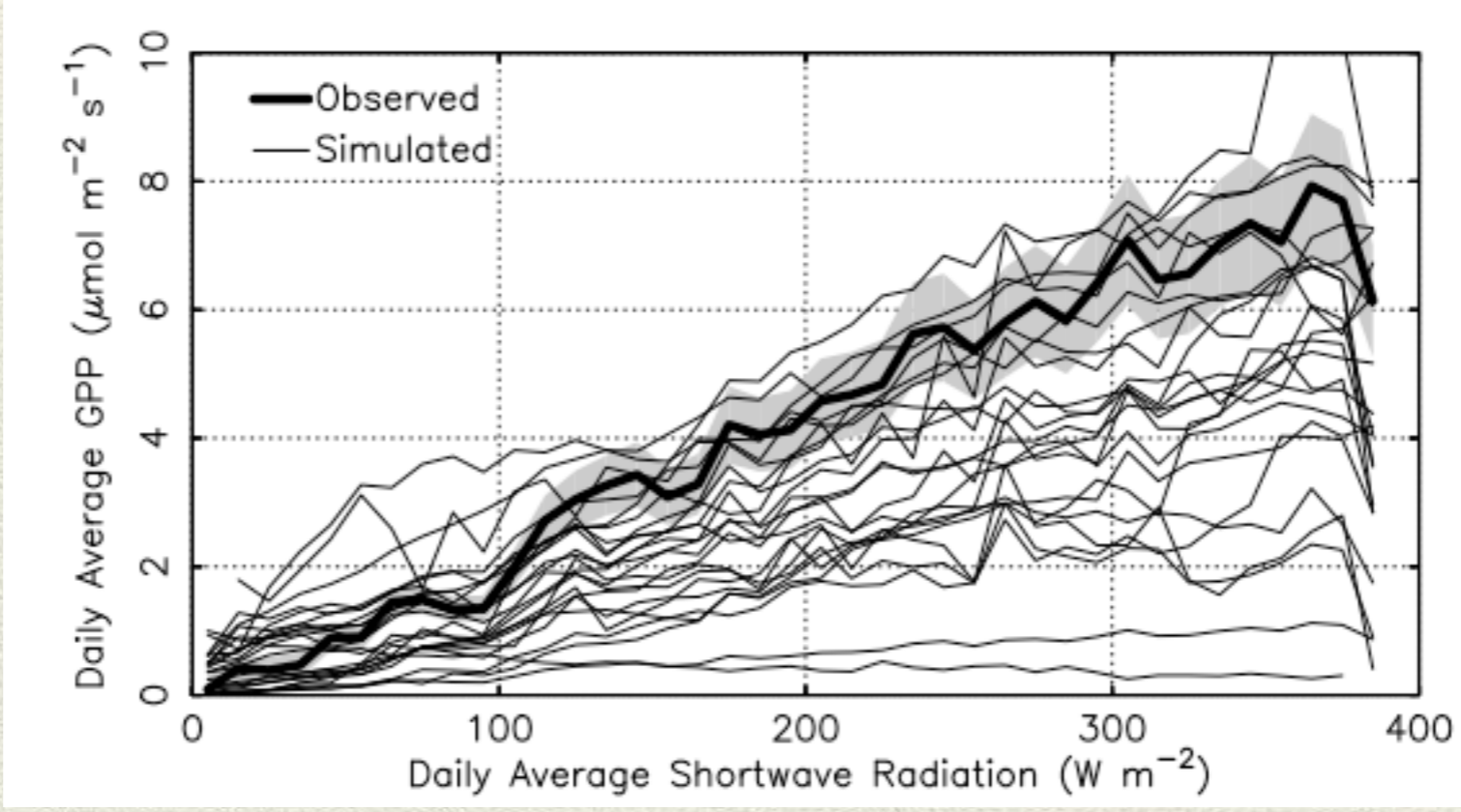
# Miscalibration





# Dynamics & Drivers







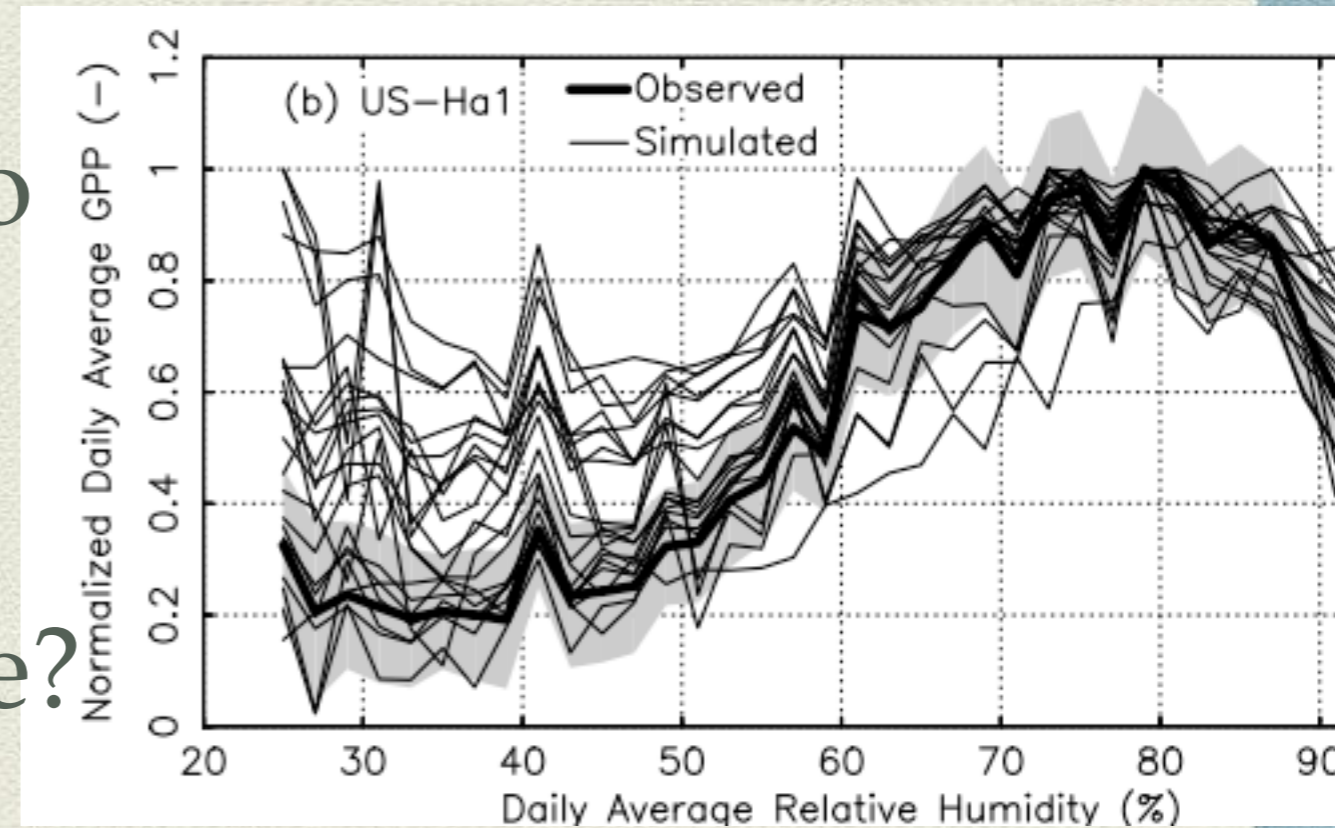
Diagnosing a model is  
Hypothesis Testing

Why would a model fail at low humidity?

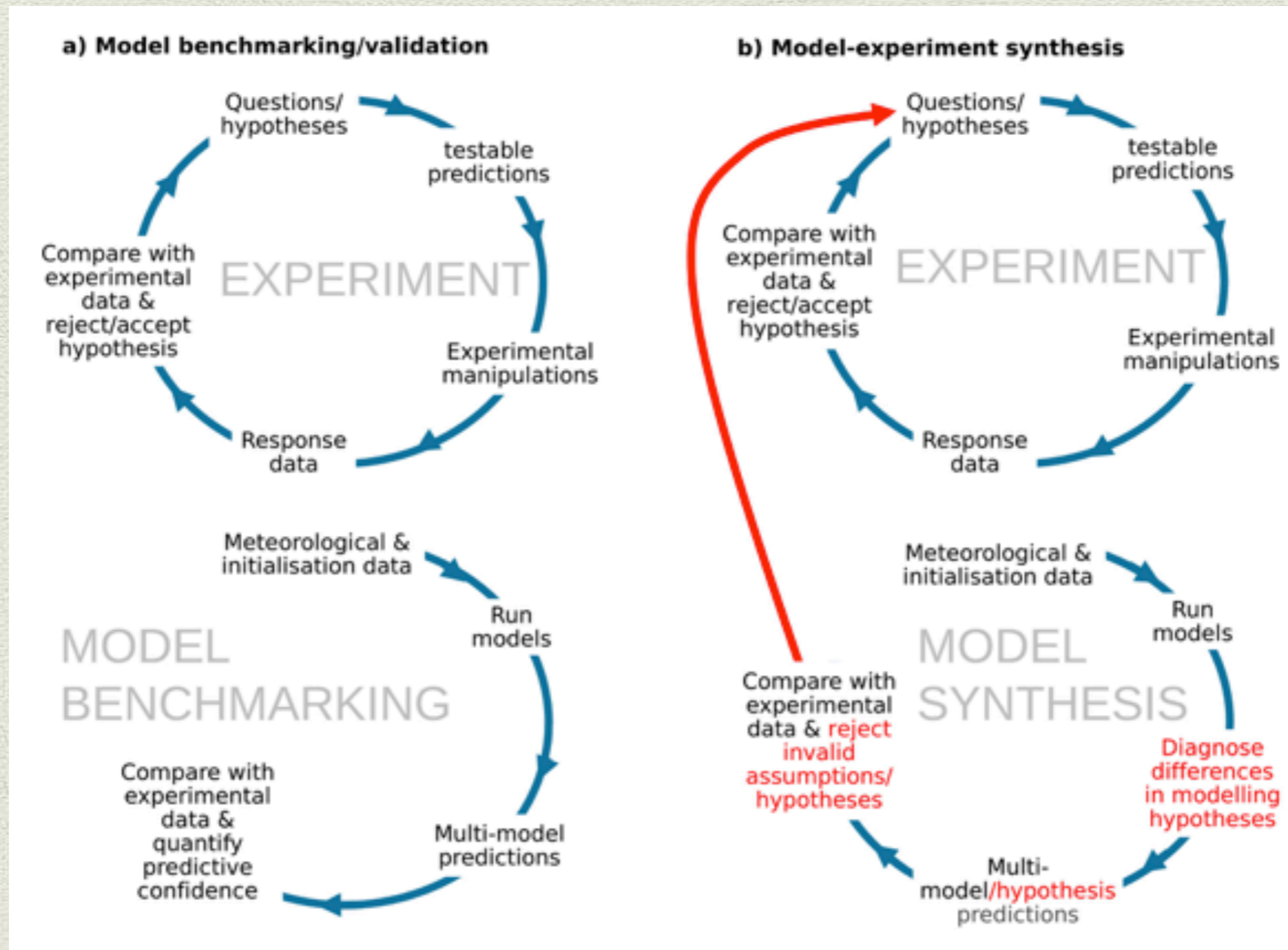
Stomatal sensitivity too low?

Too much soil moisture?

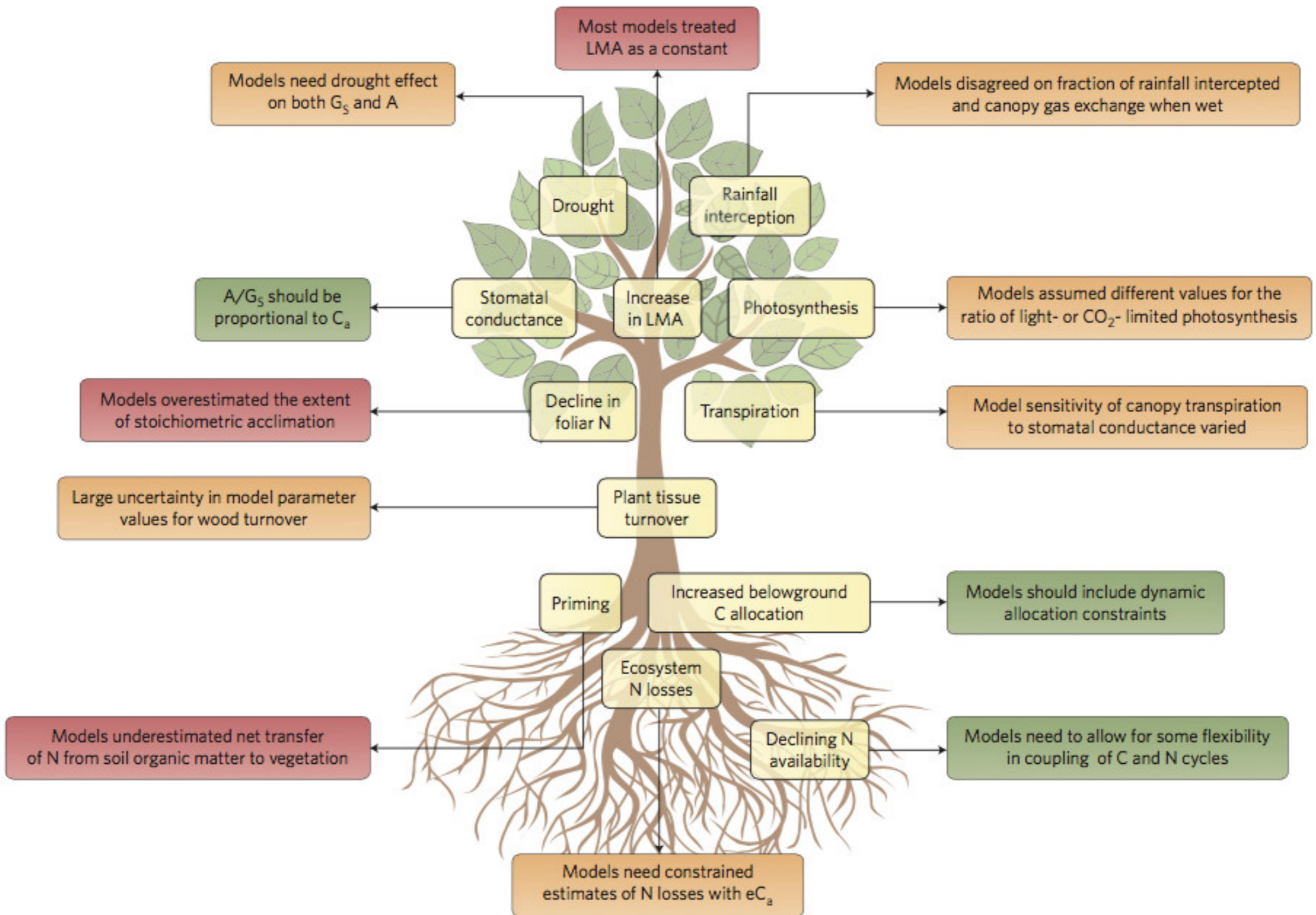
What experiments would I run in the model to test this?



# Focus on key assumptions



Walker et al 2014

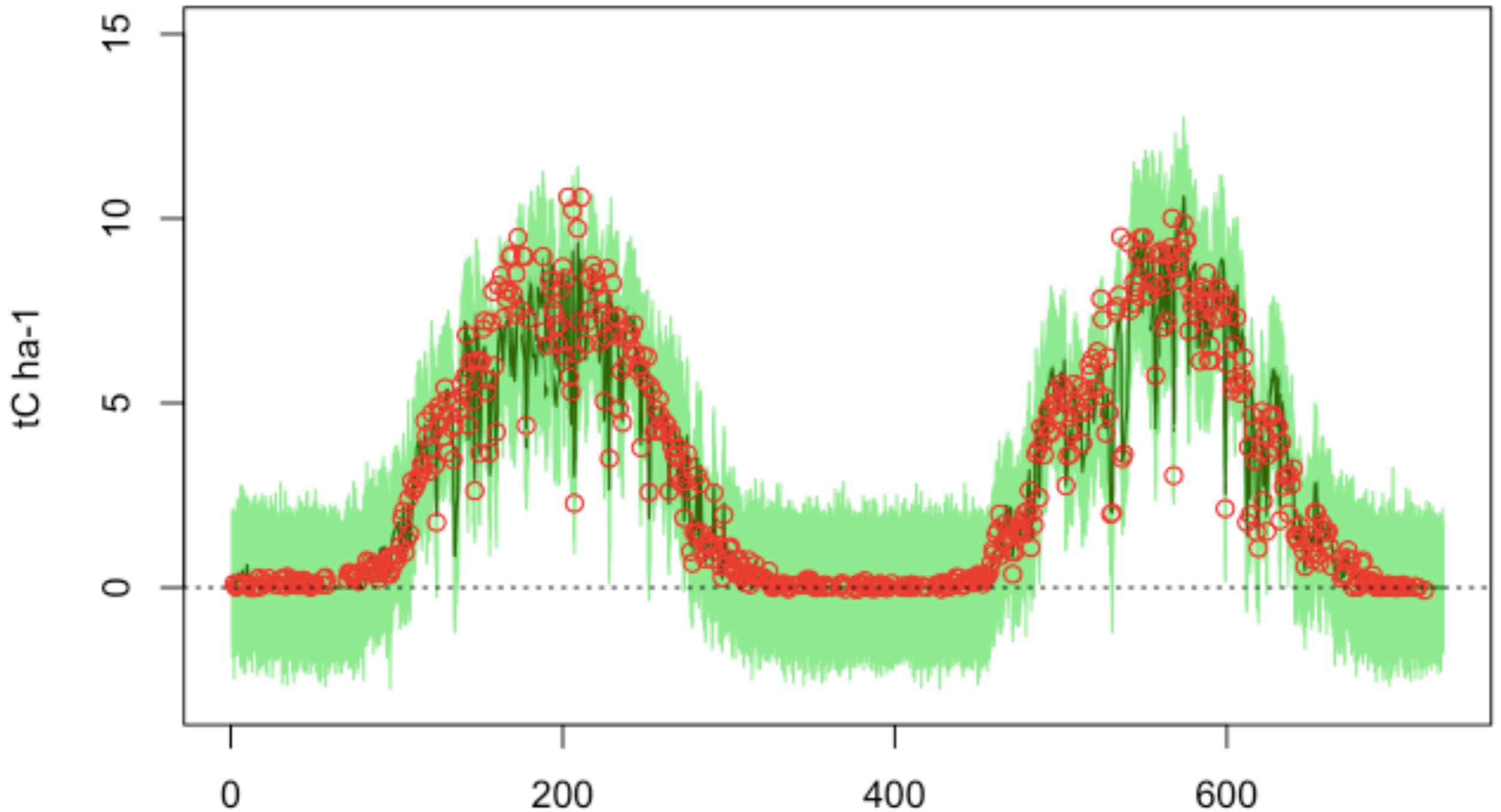


“data simulated under a model should look similar to data gathered in the real world.”

Conn et al 2018

IN THE FITTING, WE ASSUMED IID NORMAL ERRORS

**GPP**

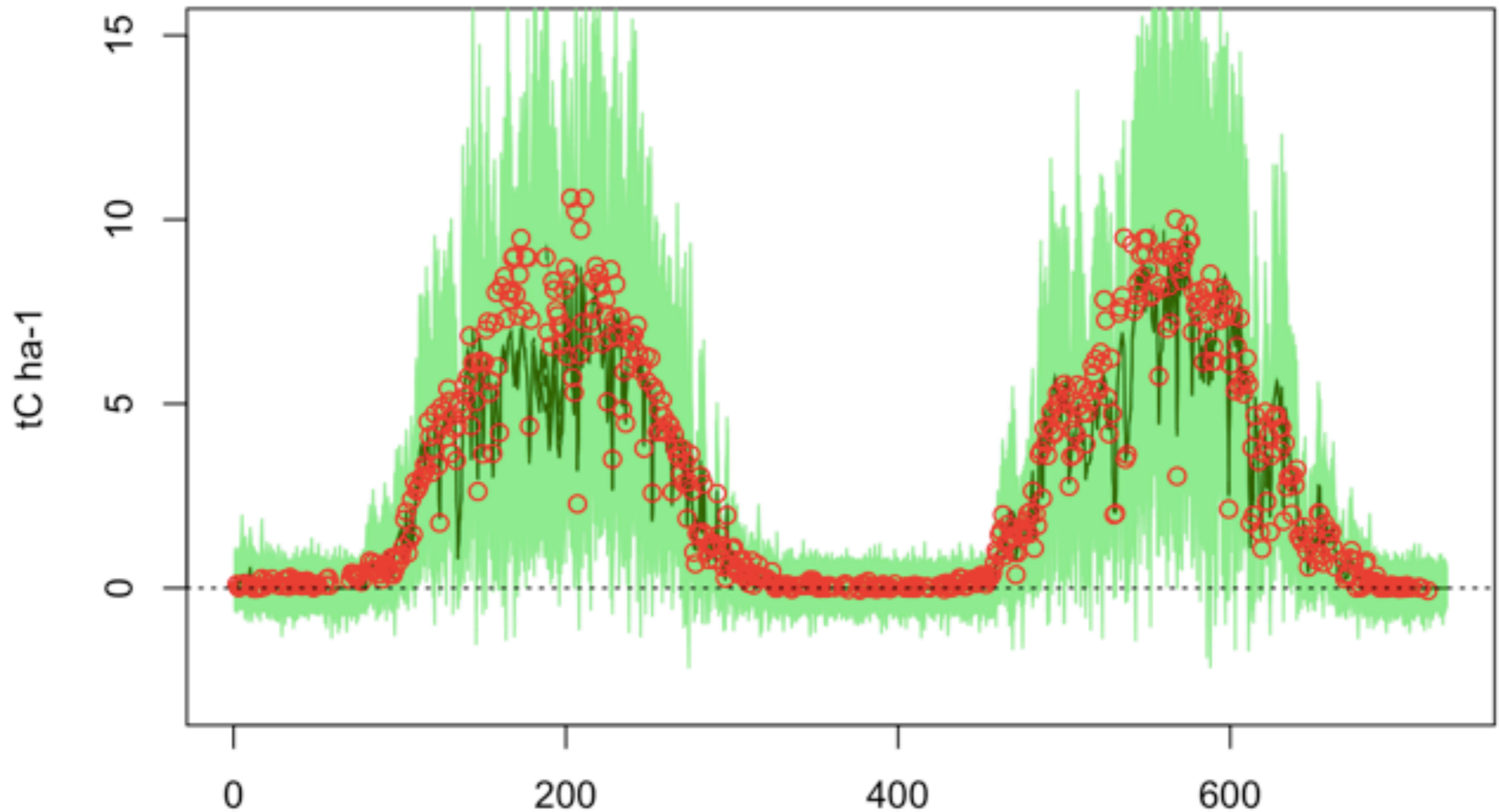


Does that seem like an adequate description of the data?



IN THIS FITTING, WE ASSUMED EXPONENTIAL ERRORS WITH  
NON-CONSTANT VARIANCE

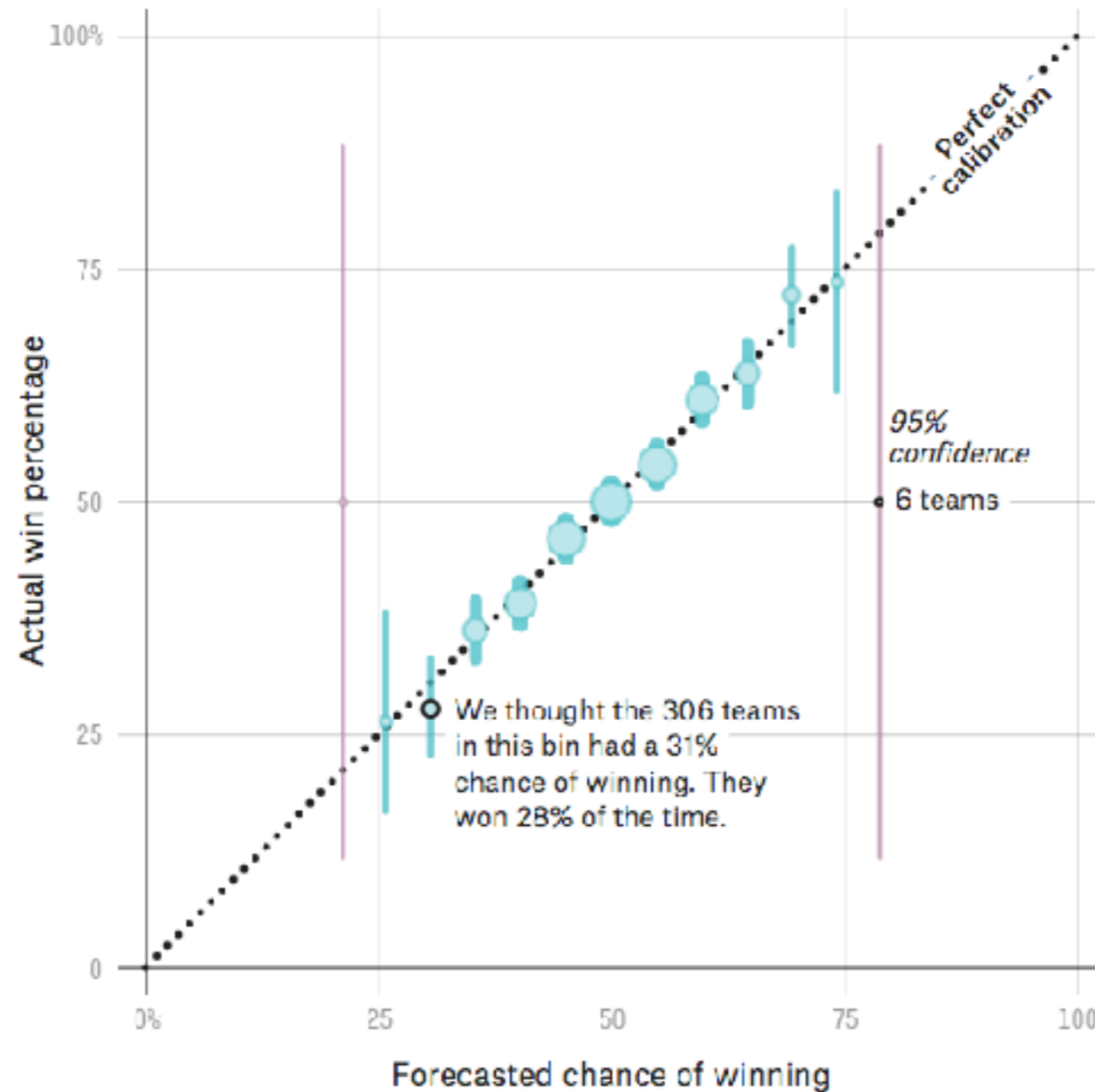
**GPP**



Does that seem like an adequate description of the data?

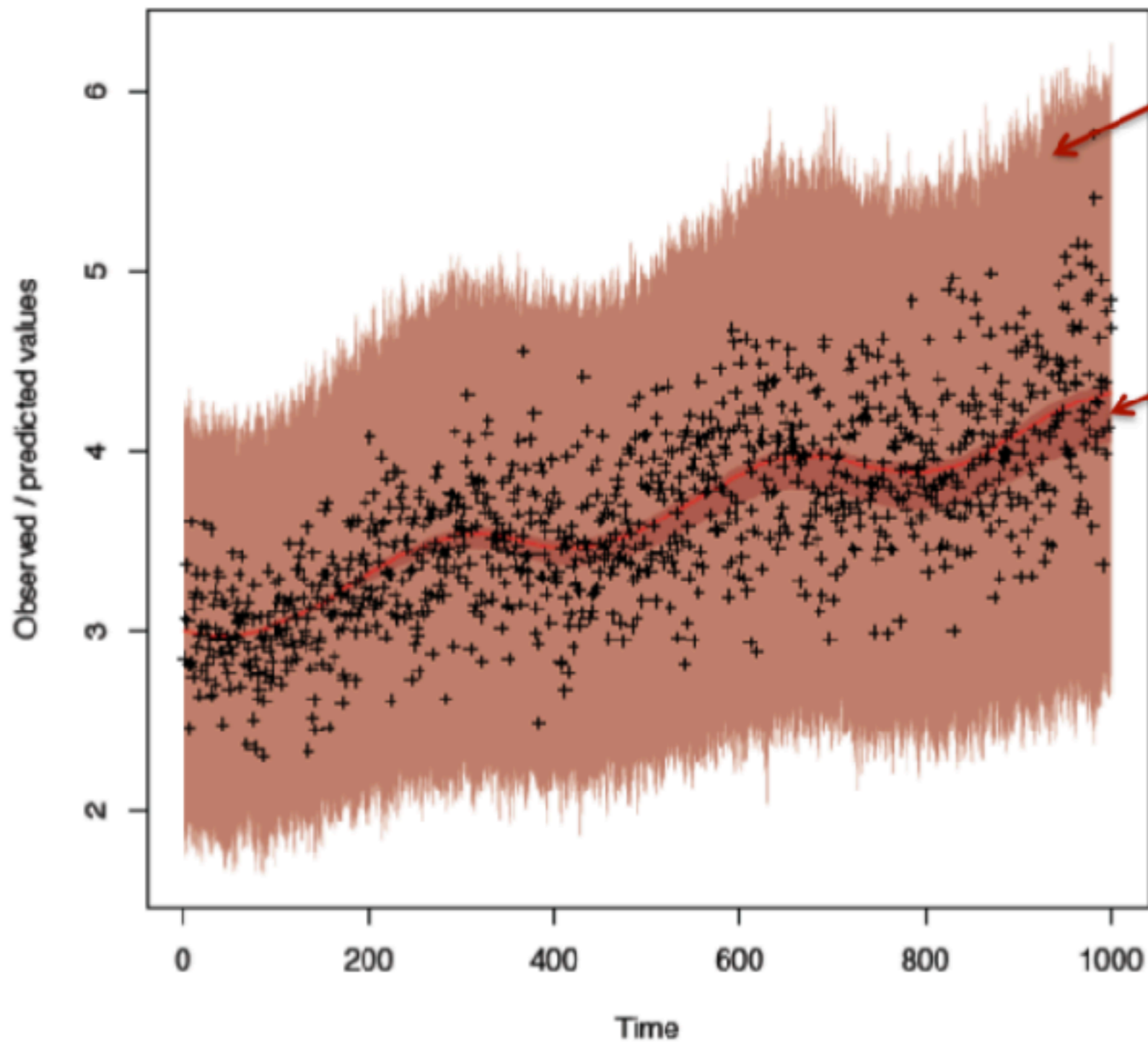
# How Good Are FiveThirtyEight Forecasts?

MLB games, 2016-18



# Bayesian p-value / prediction interval

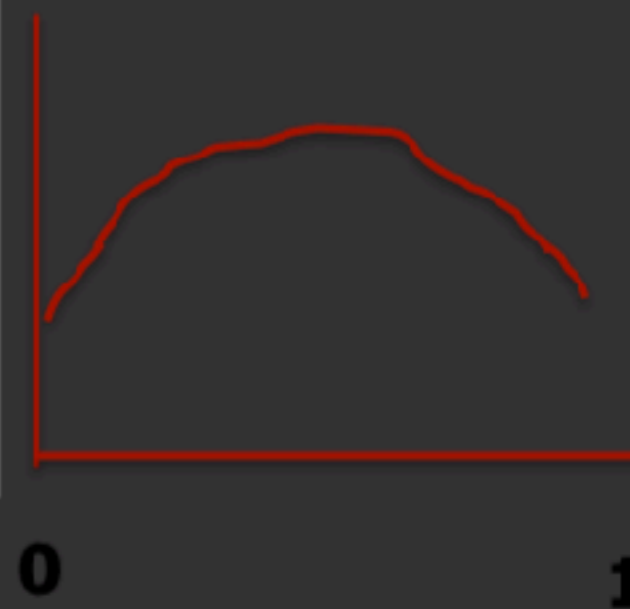
- Posterior predictive distribution is the uncertainty of the „true“ value
- **Prediction interval** is the expected variance of the observed values = PPD + error
  - Shows us what distribution we would expect for the data
- Bayesian p-value is when we use PPD + error to calculate the value of the cdf of the observed data
  - Distribution should be flat (uniform)
  - „Bayesian residuals“



PPD + Error

Posterior Predictive Distribution

Distribution of ecdf values for residuals





# Step 3: Quantitative Skill Assessment

# Error Statistics

◆ Root Mean Square Error (RMSE)

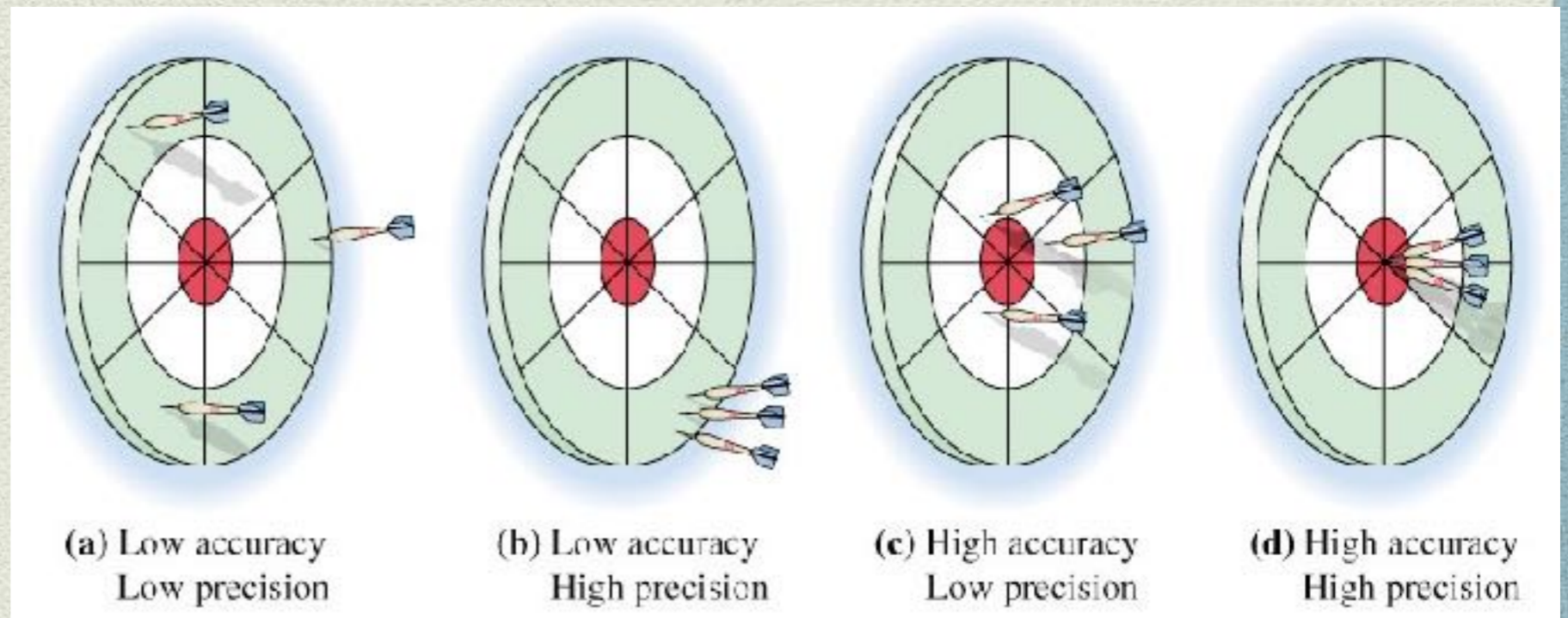
◆ Bias

◆ Correlation (r)

◆ R<sup>2</sup>

◆ Regression  
slope

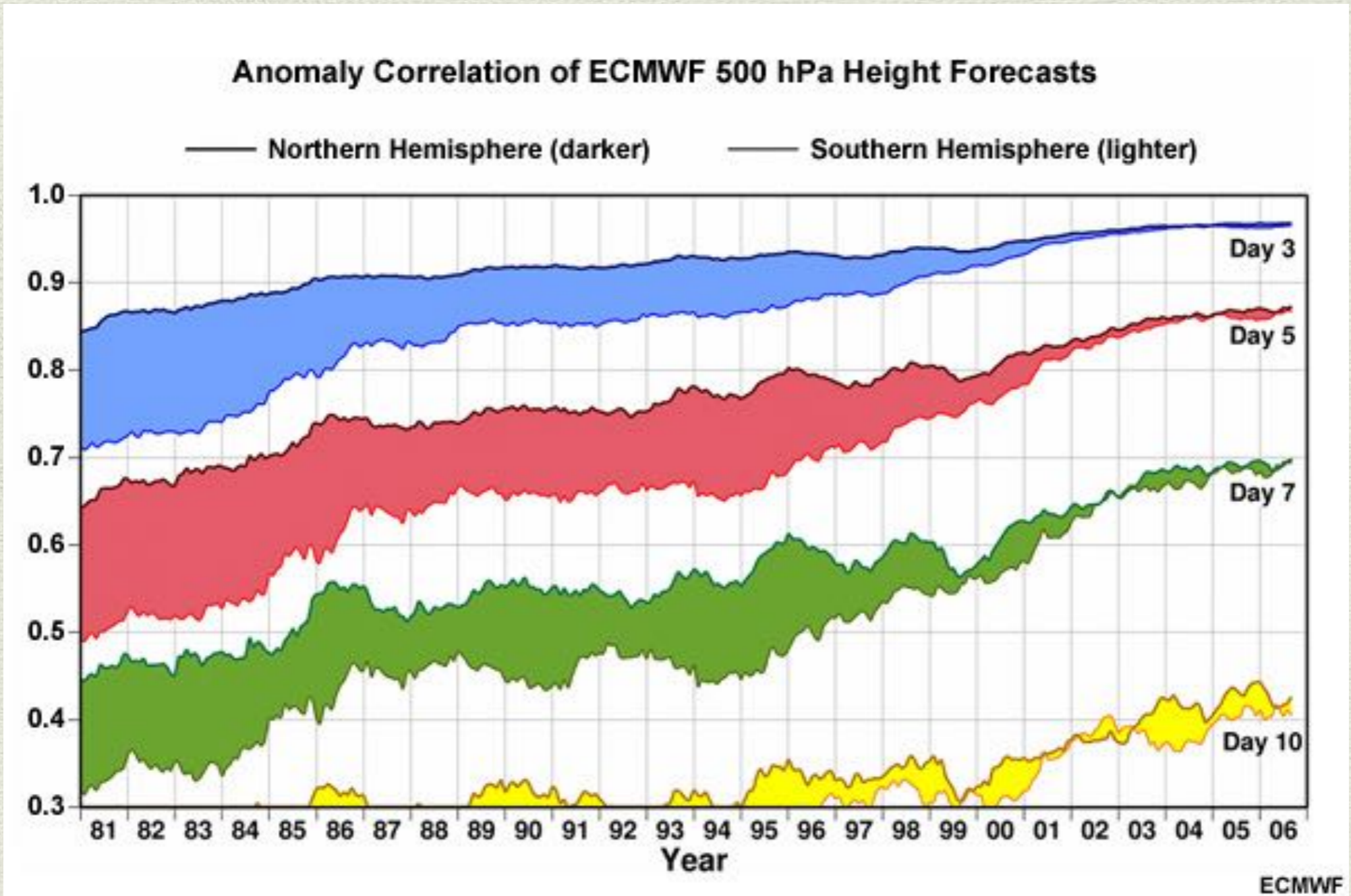
$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - x_i)^2}$$



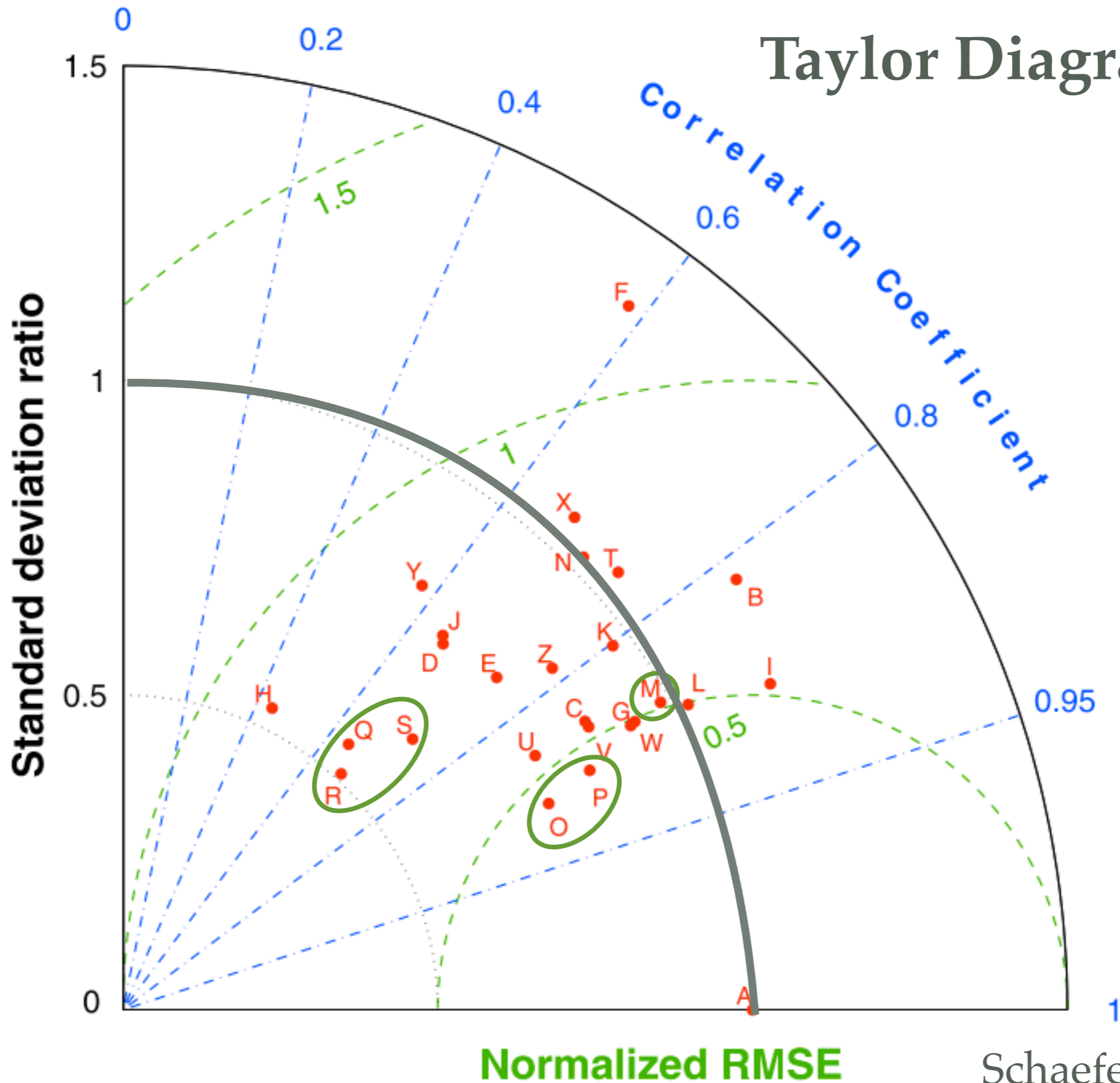
**Proper: based on the metric used for calibration**

**Local: depends on data that could actually be collected**

# Correlation



# Taylor Diagram

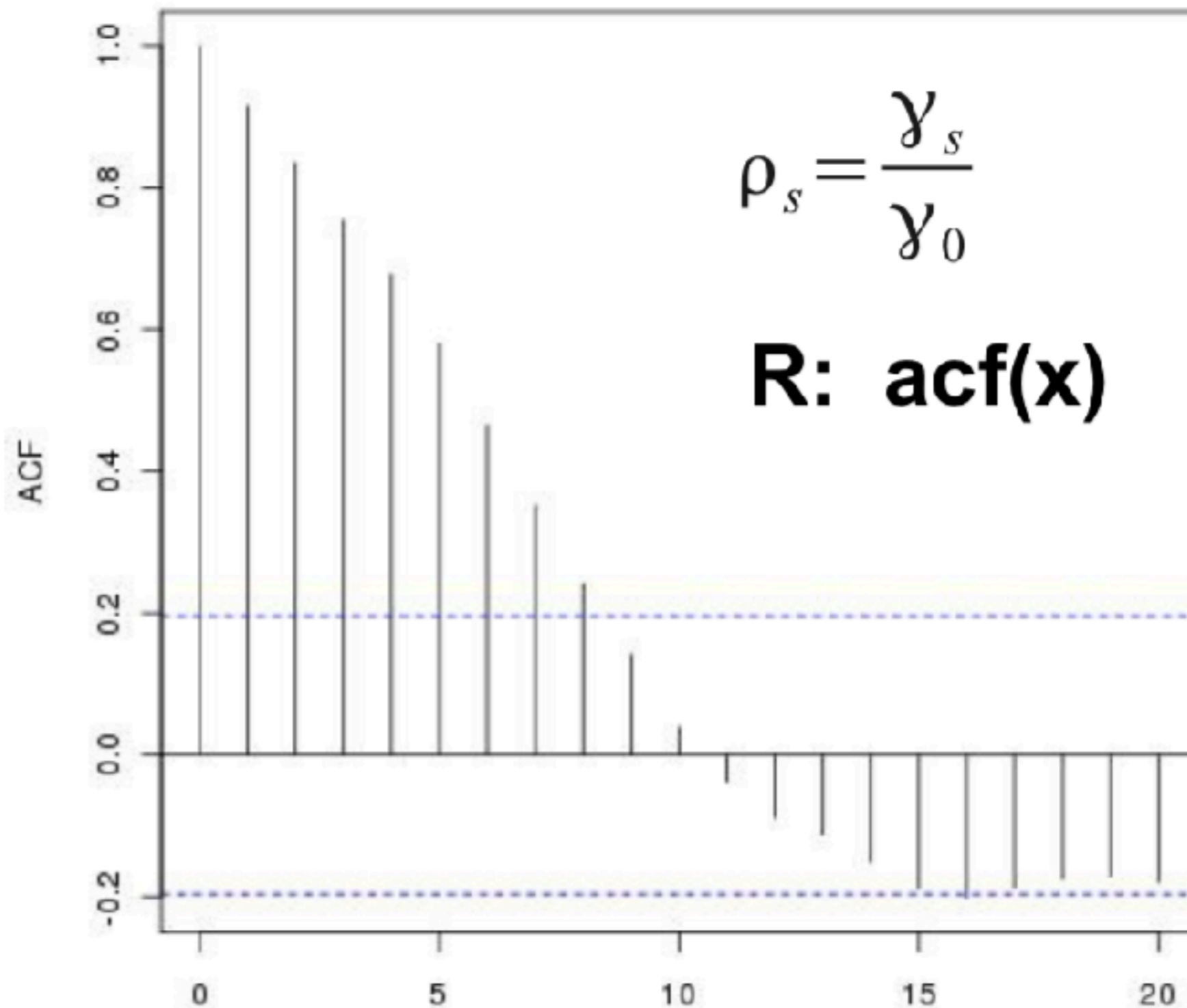


- A OBSERVED
- B AGROIBIS
- C BEPS
- D BIOMEBGC
- E CAN-IBIS
- F CNCLASS
- G DLEM
- H DNDC
- I ECOSYS
- J ED2
- K ISAM
- L ISOLSM
- M LOTEC
- N LPJ
- O MEAN\_ALL
- P MEAN\_DIURN
- Q MODIS\_ALG
- R MODIS\_5
- S MODIS\_5.1
- T ORCHIDEE
- U SIB
- V SIBCASA
- W SIBCROP
- X SSIB2
- Y TECO
- Z TRIPLEX



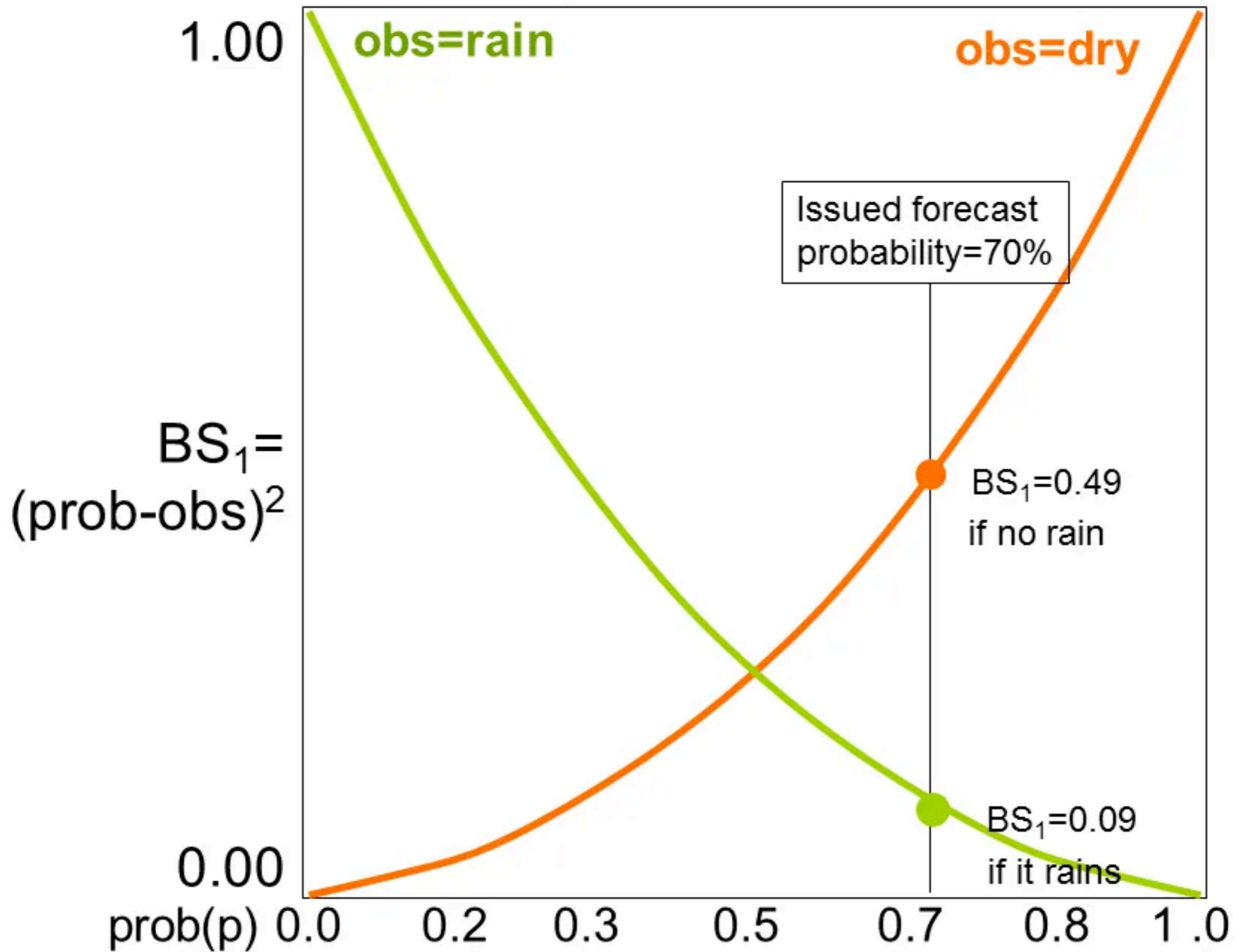
# Autocorrelation

## Correlogram



# Brier score

Contribution  $BS_1$  of one forecast to the total Brier Score



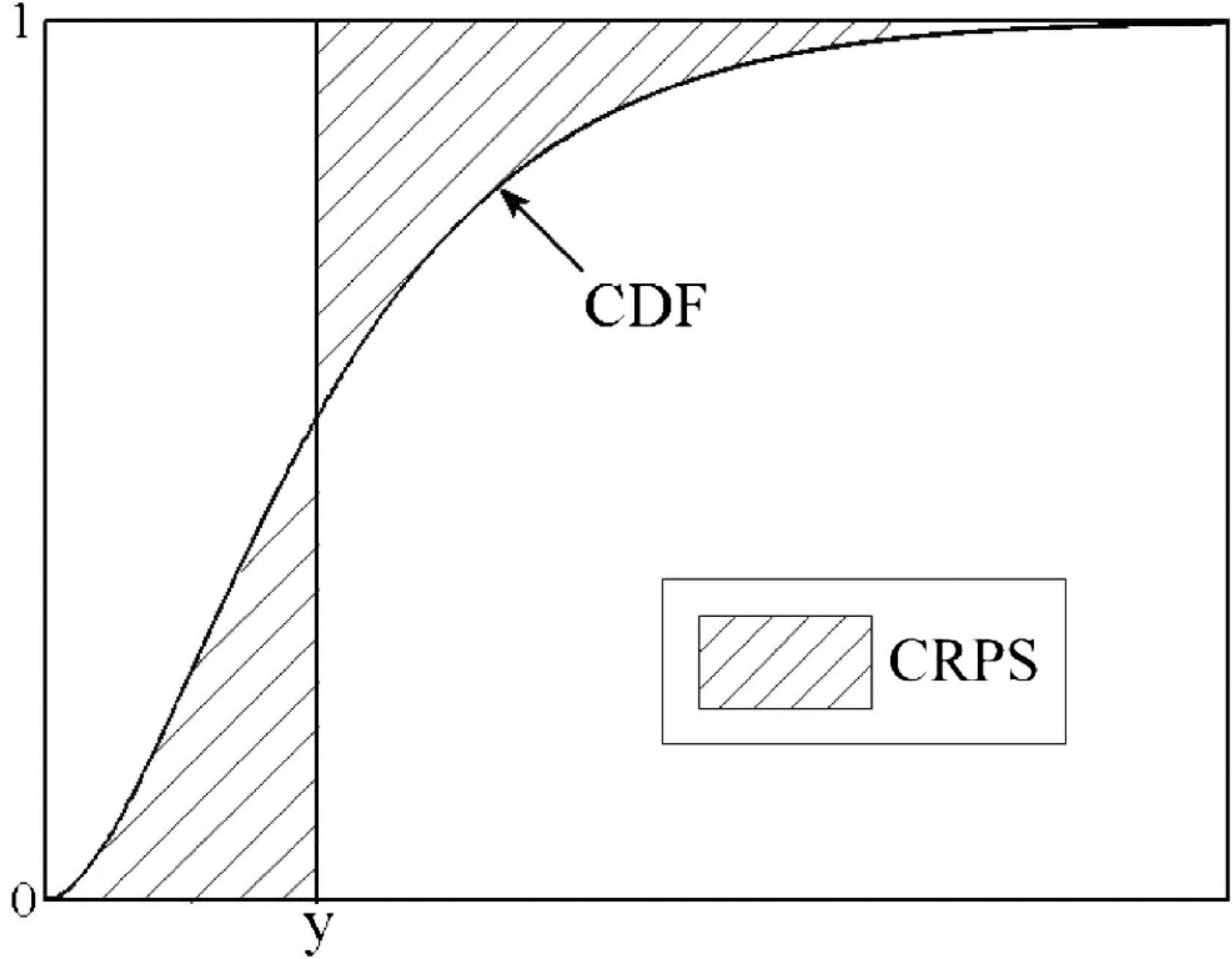
# Continuous Ranked Probability Score

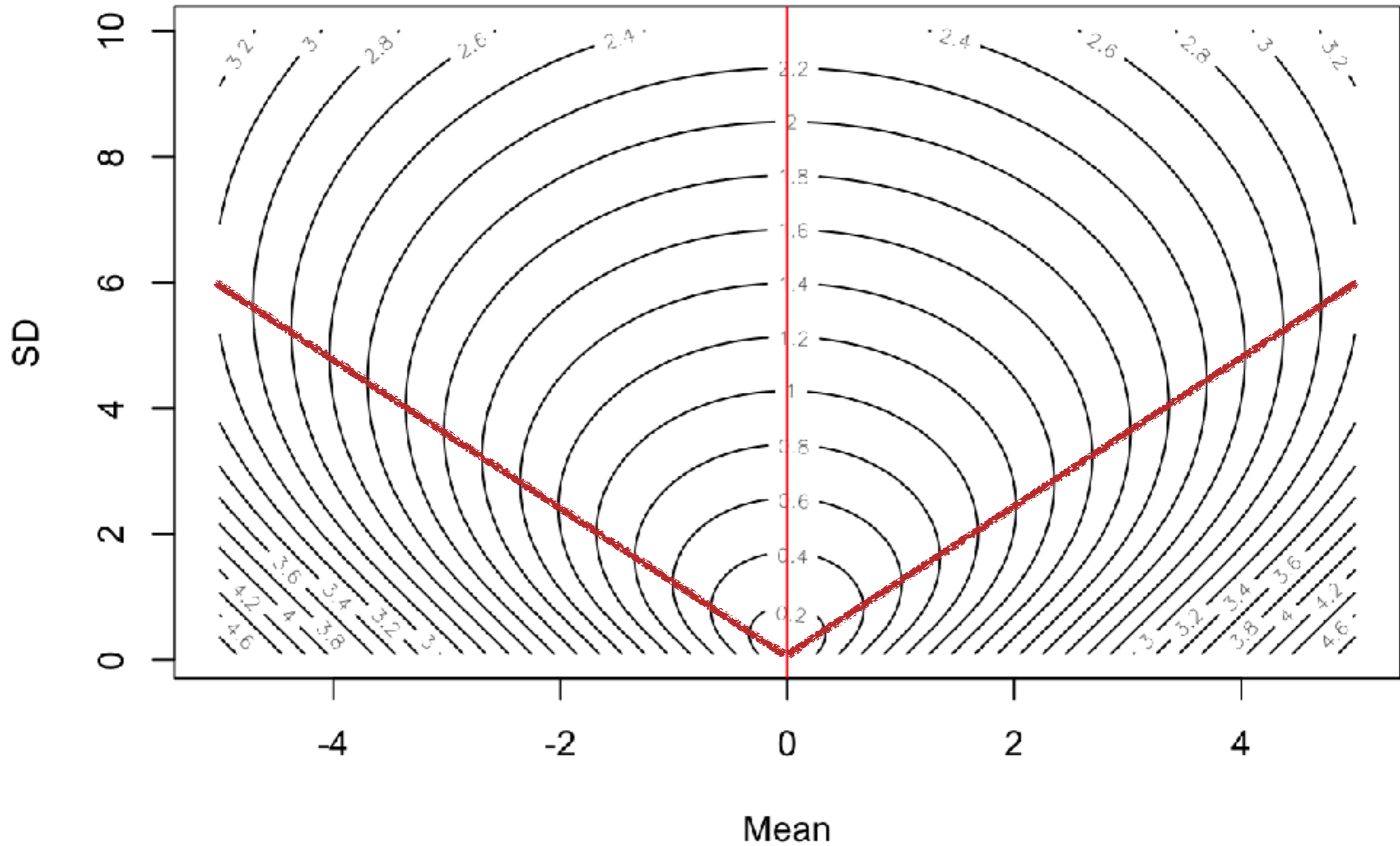
$$\text{CRPS}(F, x) = - \int_{-\infty}^{\infty} (F(y) - \mathbb{1}\{y \geq x\})^2 dy$$

$$\text{CRPS}(\hat{F}_m, y) = \frac{1}{m} \sum_{i=1}^m |X_i - y| - \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m |X_i - X_j|$$

**Ensemble member** ↓ **Data** ↓

**Mean Absolute Error**                      **Penalty for ensemble spread**





<https://github.com/eco4cast/neon4cast-scoring/blob/main/><sub>29</sub>  
CRPS\_example\_JRT.Rmd

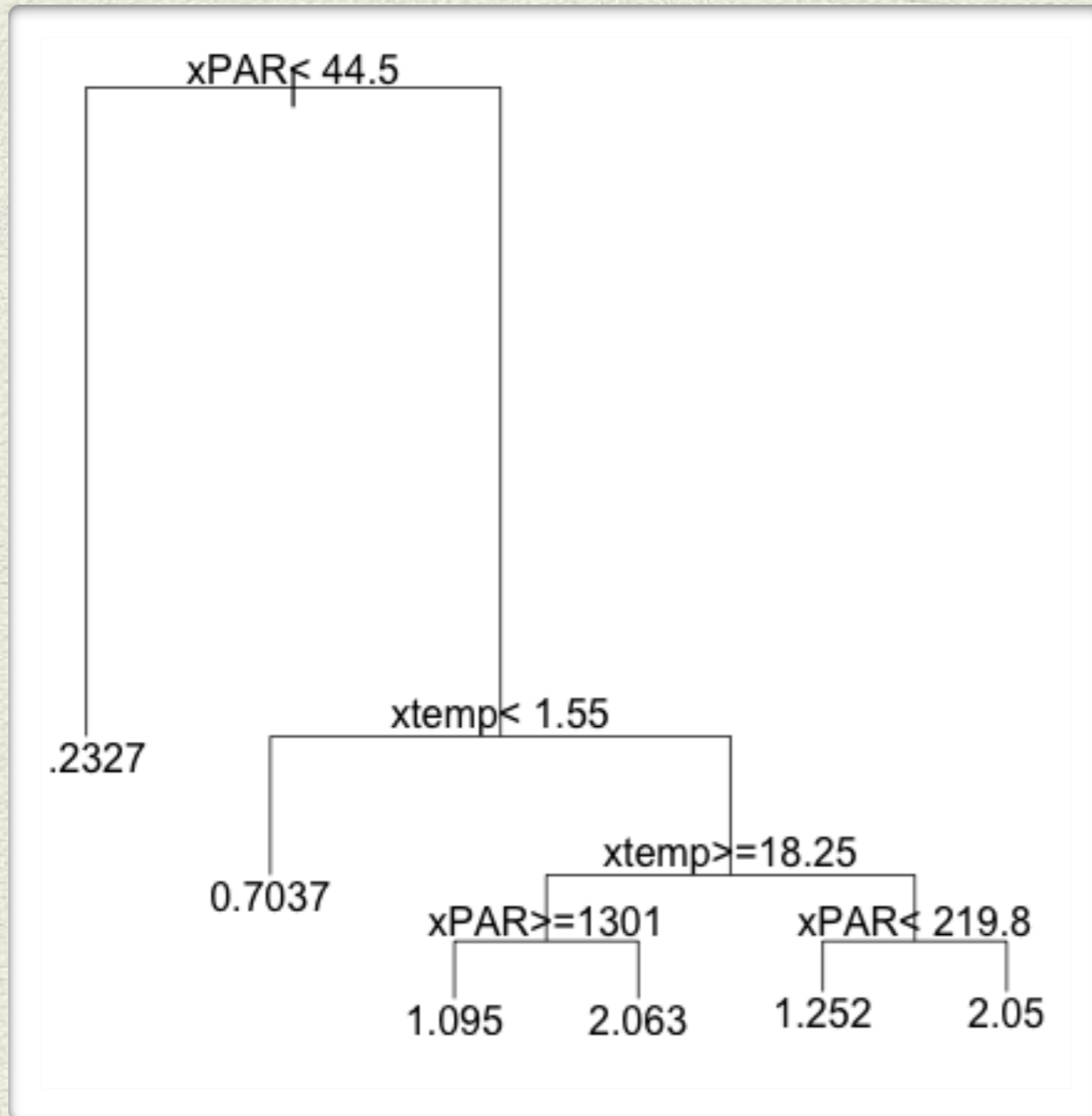
# Data mining the residuals

- ◆ Wide variety of Data Mining algorithms in use
- ◆ Potentially useful for generating hypothesis about when / where model fails
- ◆ “Correct” the forecast
- ◆ Hybrid models increasing

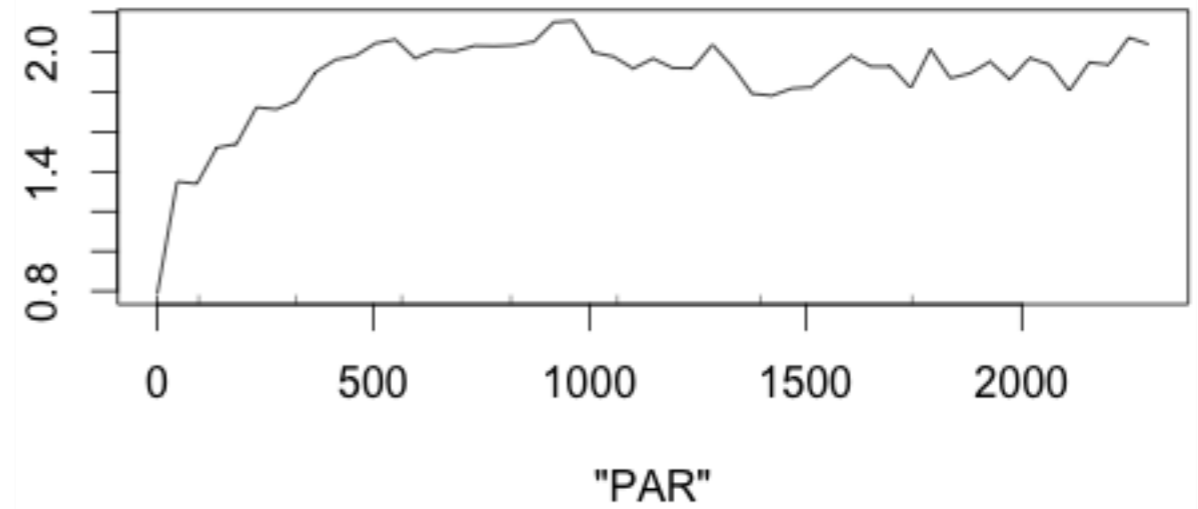
- ◆ CART
- ◆ GAM
- ◆ Random Forests
- ◆ Boosted regression trees
- ◆ Artificial Neural Network
- ◆ Support Vector Machines

# Random Forest

## CART



### Partial Dependence on "PAR"



### Partial Dependence on "temp"

