

Finding Relations in Climate Data Sets

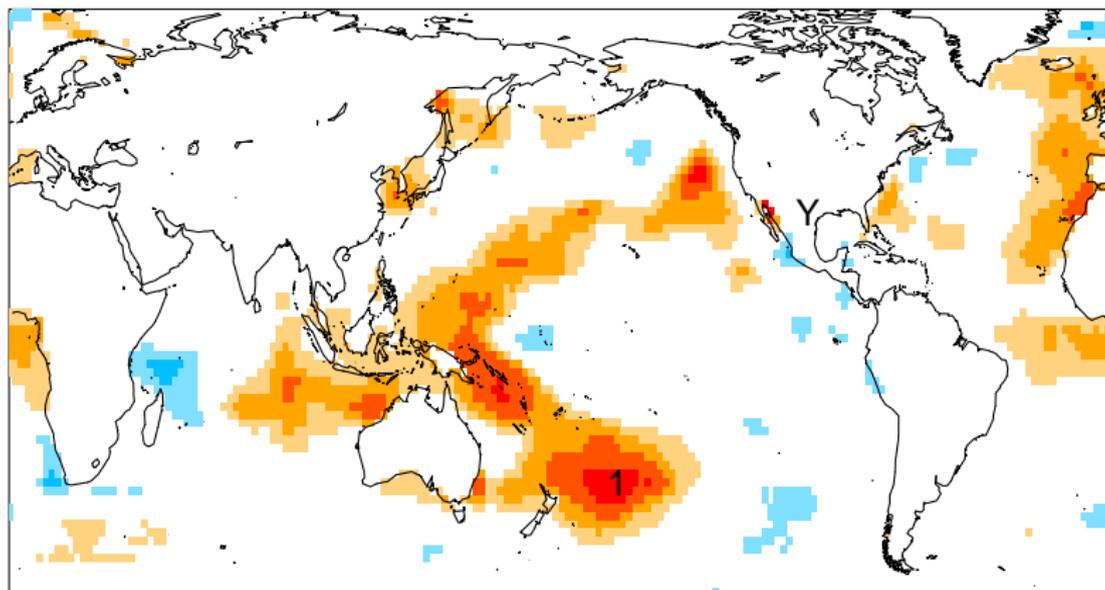
Timothy DelSole

George Mason University, Fairfax, Va and
Center for Ocean-Land-Atmosphere Studies, Calverton, MD

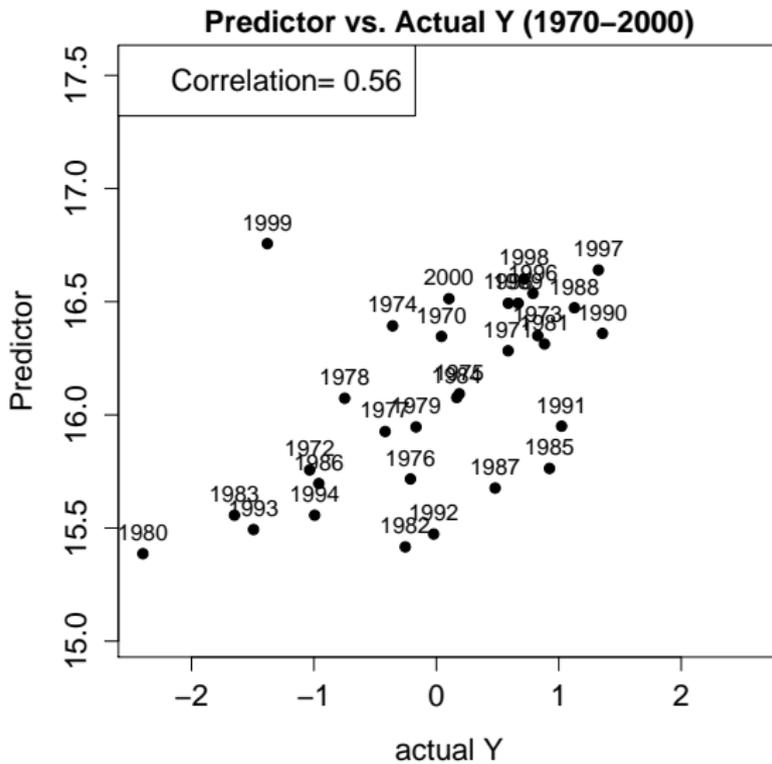
August 5, 2015

collaborators: Michael Tippett, Arindam Bannerjee, Claire Monteleoni

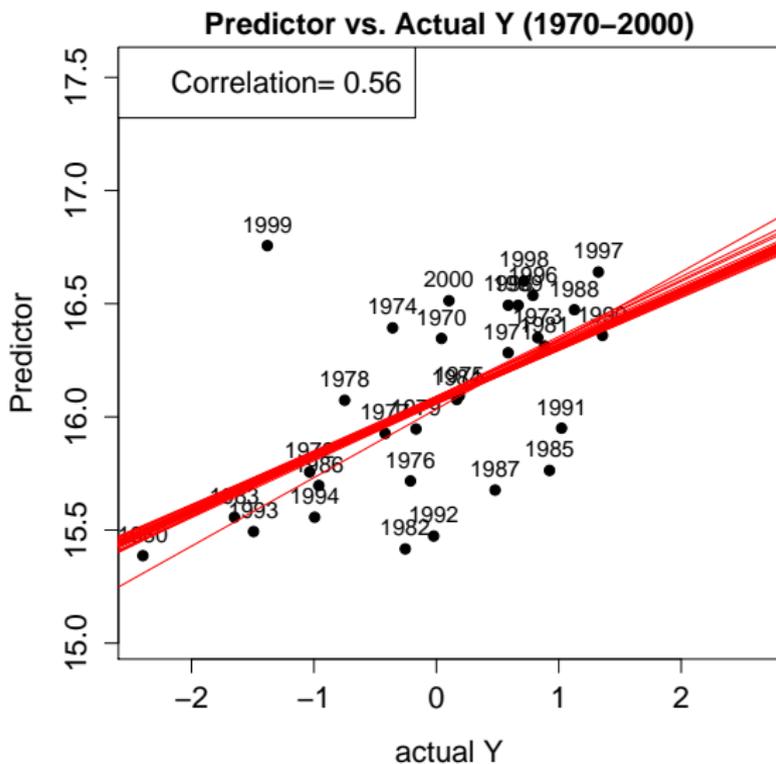
Correlation Map Between X and SST (1970–2000)



Scatter Plot



Leave-One-Out Regression Fits



Surprise! Y is random!

Surprise! Y is random!

- ▶ On a $2.5^\circ \times 2.5^\circ$ grid, ocean surface is about 8000 points.
- ▶ Only about 50 years of observations of global ocean surface.
- ▶ Physics varies with season, so seasons cannot be pooled.

Estimate 8000 parameters using 50-150 data points

This is an example of data fishing.

1. Not all climate scientists understand this issue.
2. For those who do, machine learning “looks” like data fishing.
3. It is useful to challenge climate scientists about how they know what they tell you.
4. Some correlations are trusted.

ENSO Teleconnections?

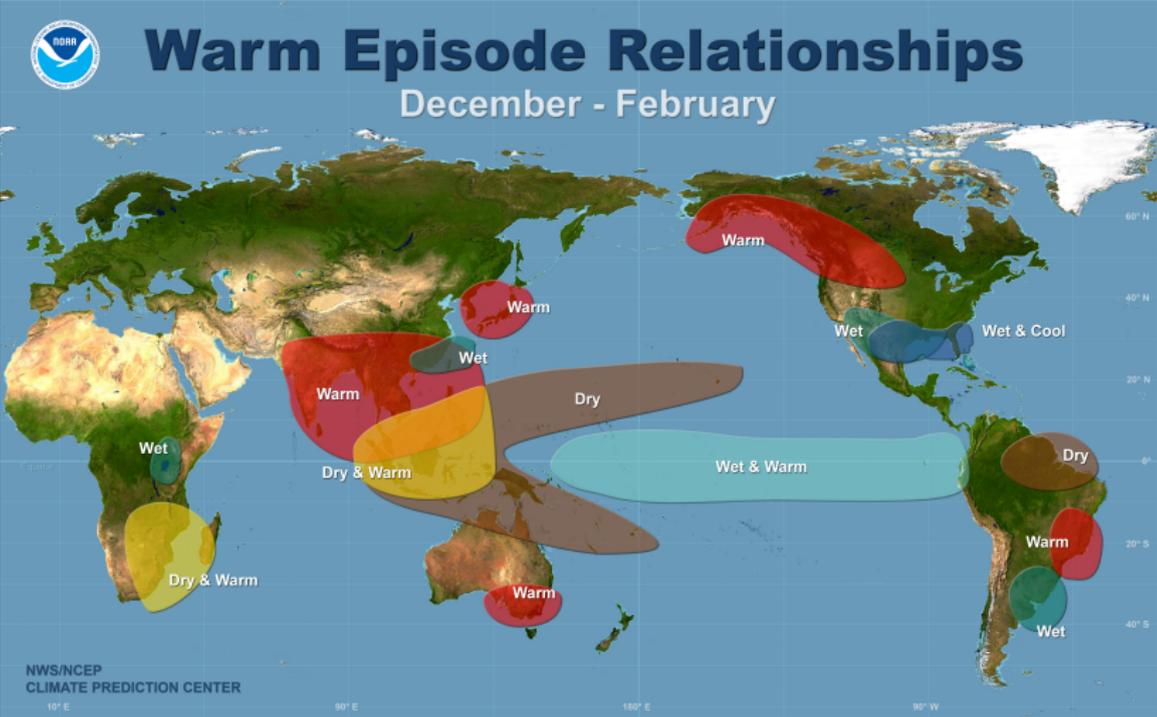


figure http://www.cpc.noaa.gov/products/analysis_monitoring/impacts/warm.gif

Proving an Observed Relation is Real

- ▶ Relation holds in independent data.
- ▶ Relation can be reproduced by climate models.
- ▶ Relation can be understood through simple dynamical models.

ENSO Teleconnection Inferred in 1989 (Precipitation)

SCHEMATIC OF AREAS WITH CONSISTENT HIGH SO INDEX-PRECIPITATION RELATIONSHIP

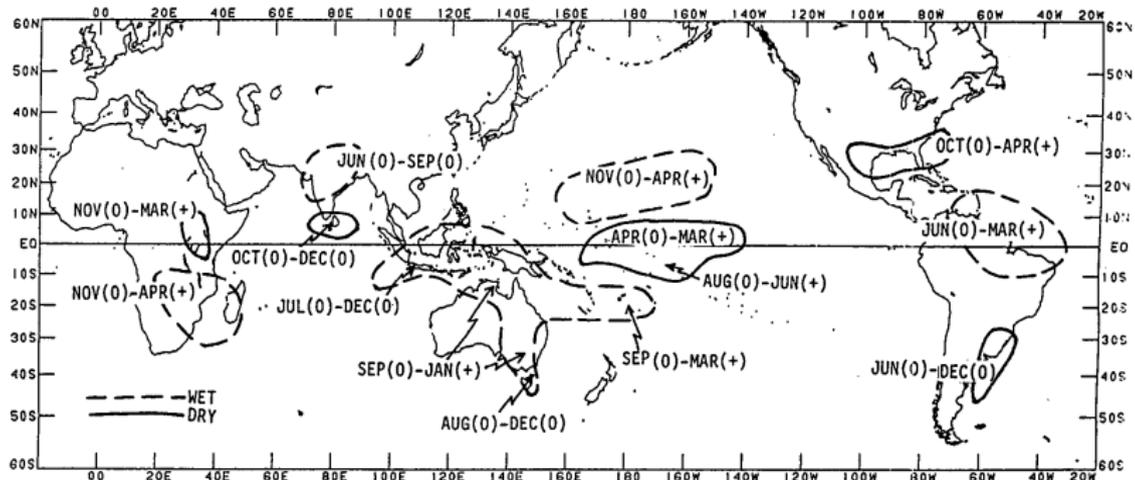


FIG. 18. Schematic representation of the principal areas of high-SO index related precipitation based on the analysis of precipitation composites and time series. Regional maps should be consulted for details.

figure: Ropelewski and Halpert (1989, J. Climate)

Machine Learning Approach: Bet on Sparsity

Frame the estimation such that most parameters are zero.

“Use a procedure that does well in sparse problems, since no procedure does well in dense problems.”

Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*, 2nd edition

Examples of Sparsity

grid points: variable is related to **localized** geographic regions.

matrix rank: variable is related to a **few patterns**.

spectral: variable fluctuates over a **small range of frequencies**

PCs: variable is related to components with **high variance**

variables: variable is related to **certain physical variables**

What is the most appropriate form of sparsity when finding relations in climate data?

Small-scale structure is less trustworthy than large-scale structure in interseasonal climate.

Climate Models

- ▶ Accuracy of numerical solutions of partial differential equations degrades with decreasing length scale.
- ▶ Subgrid scale phenomena are not solved directly, but are parameterized in terms of large scale grid quantities.
- ▶ Surface topography is rarely accurately included, and often induces numerical phenomena that have no counterpart in nature.

Eigenvectors of Laplace Operator provide a natural basis set for representing large spatial scales or long temporal scales.

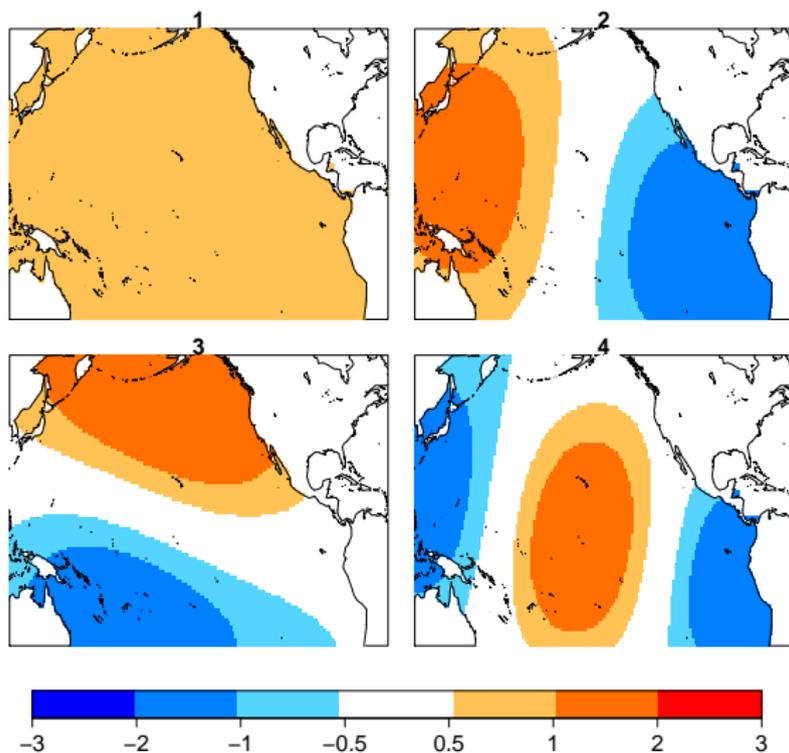
$$\nabla^2\psi = -\lambda^2\psi$$

Circle: Fourier series (λ is the frequency)

Sphere: Spherical Harmonics (λ is the total wavenumber)

Large-scale \rightarrow most amplitudes vanish.

Laplacian Eigenvectors Over the Pacific



DelSole and Tippett (2015, *Journal of Climate*)

Regularized Linear Regression

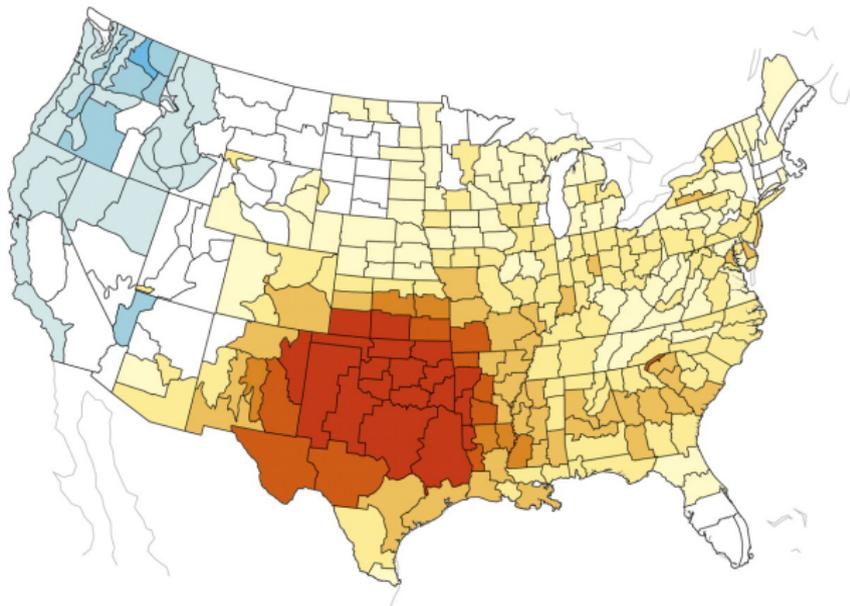
$$\| \underset{\text{predictand}}{\mathbf{y}} - \underset{\text{predictors}}{\mathbf{X}} \underset{\text{weights}}{\mathbf{w}} \|^2 + \lambda R(\mathbf{w})$$

regularizer

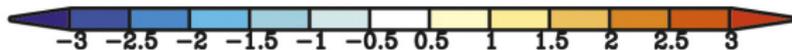
Typical regularizers (choice of R) in machine learning:

- ▶ $w_k = 0$ for $k \geq K$ (Principal Components Regression)
- ▶ L_1 norm ($|\mathbf{w}|$) (LASSO)
- ▶ L_2 norm ($\|\mathbf{w}\|^2$) (Ridge)

Observed 2011 JJA Temperature



Degrees Celsius ($^{\circ}\text{C}$)



from figure 1 of Hoerling et al. 2013, J. Climate

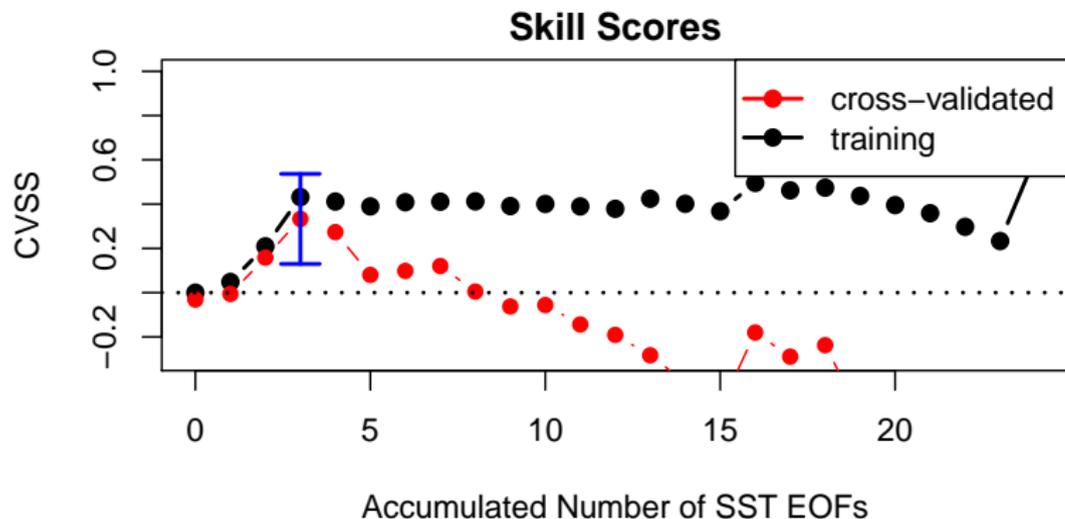
Regression model:

Texas temperature = SST in Pacific * weights + noise

Prediction Measure:

$$\text{Cross-Validated Skill Score} = 1 - \frac{MSE}{\text{variance of Texas Temperature}}$$

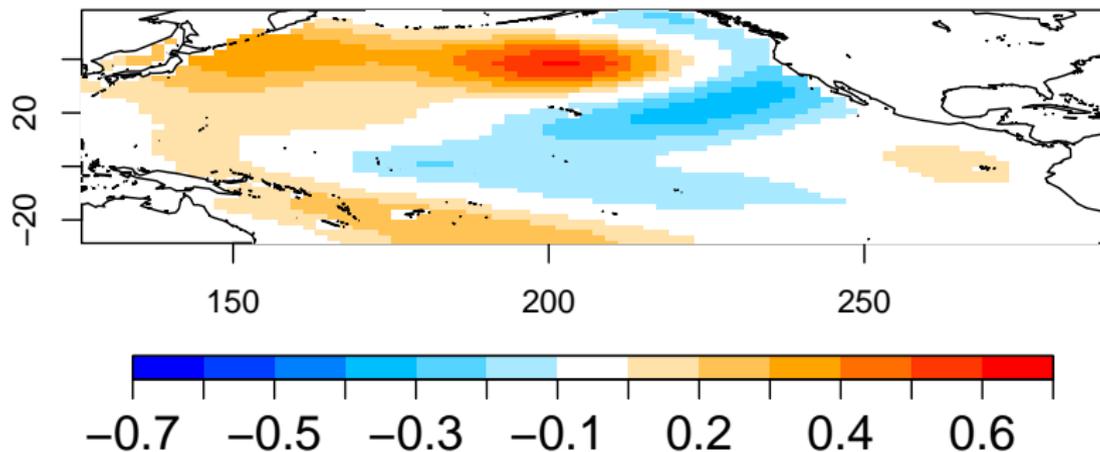
Principal Components Regression



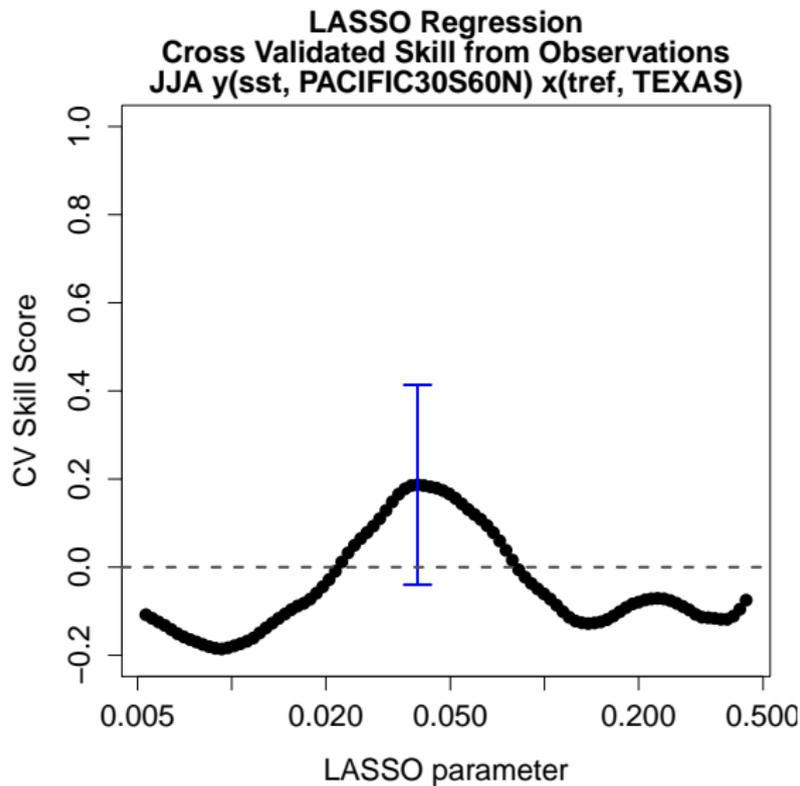
Principal Components Regression

Leading Regression Pattern from Observations

SST, 3 EOFs, JJA

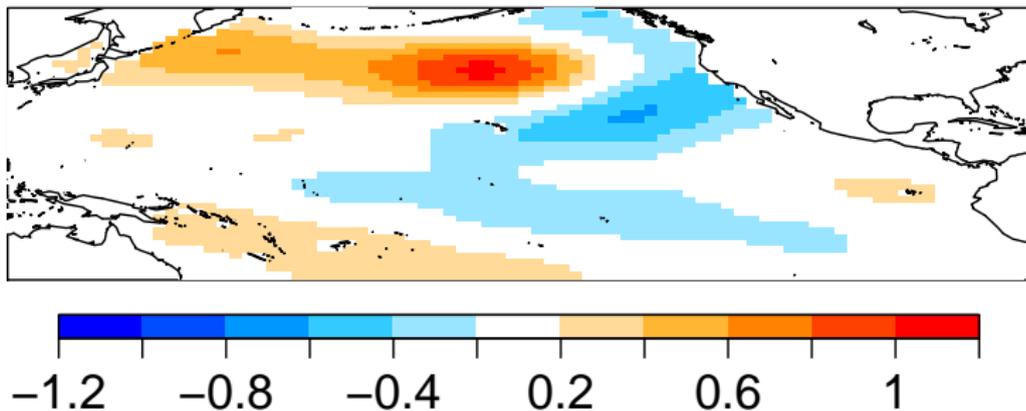


LASSO w/ Laplacians

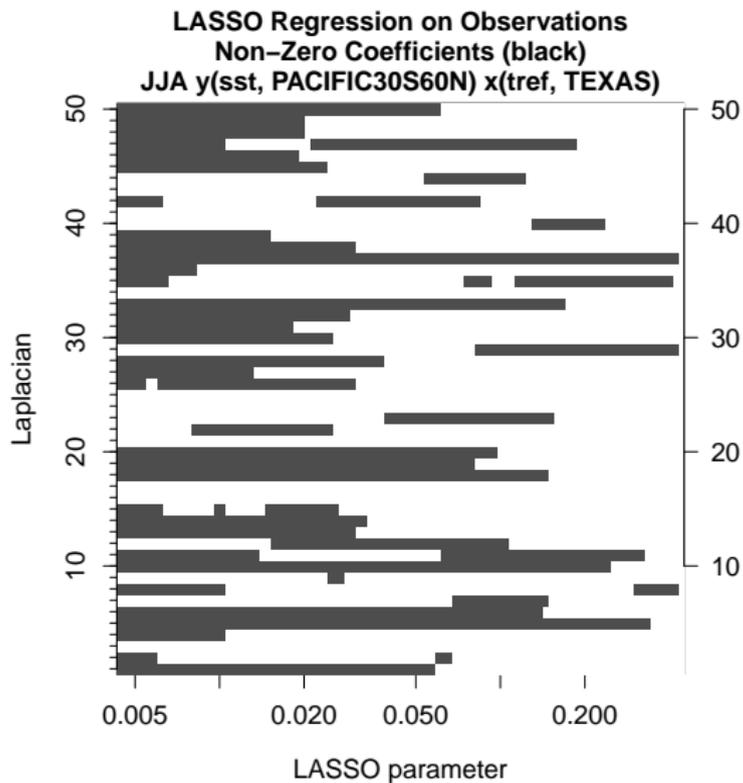


LASSO w/ Laplacians

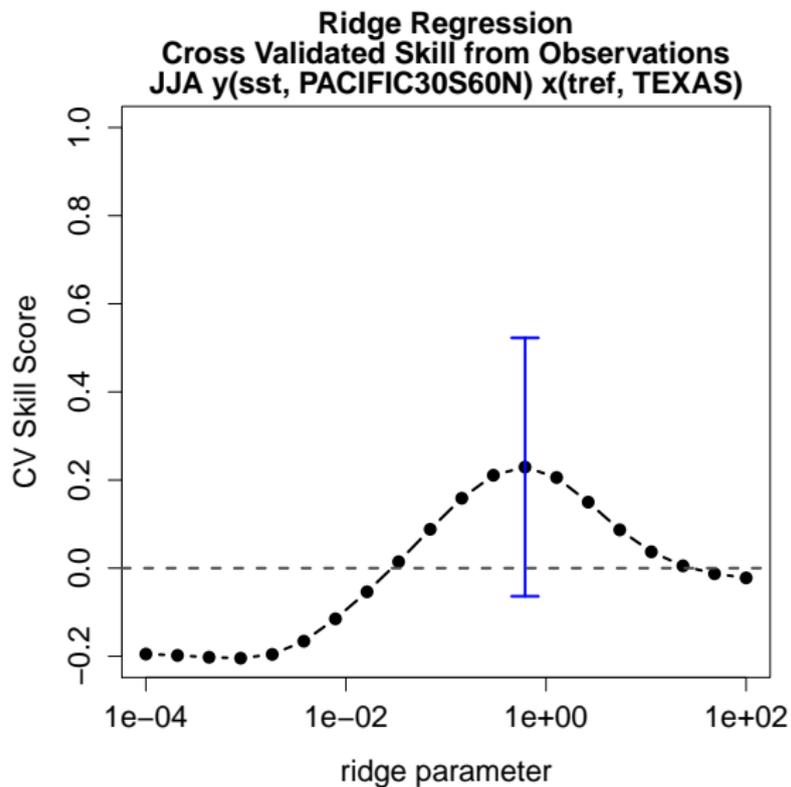
**Leading LASSO Pattern from JJA Observations
SST, LASSO Parameter= 0.039**



LASSO w/ Laplacians

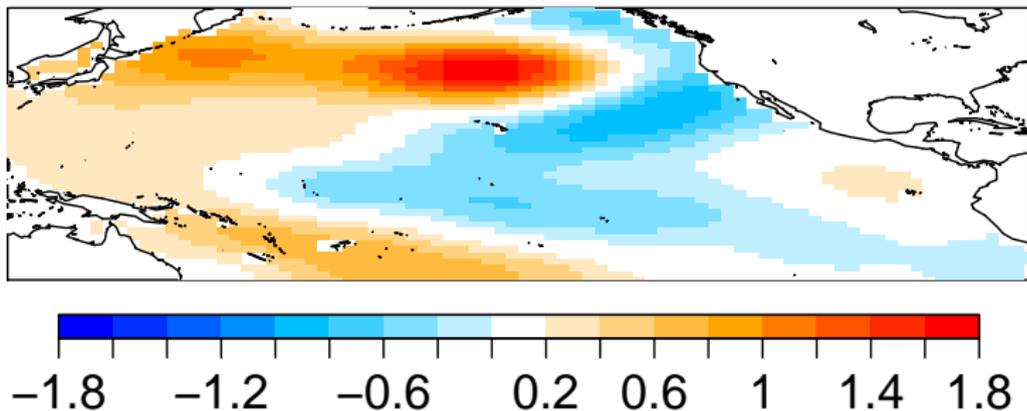


Ridge w/ Laplacians



Ridge w/ Laplacians

Leading Ridge Pattern from JJA Observations
SST, Ridge Parameter= 0.62

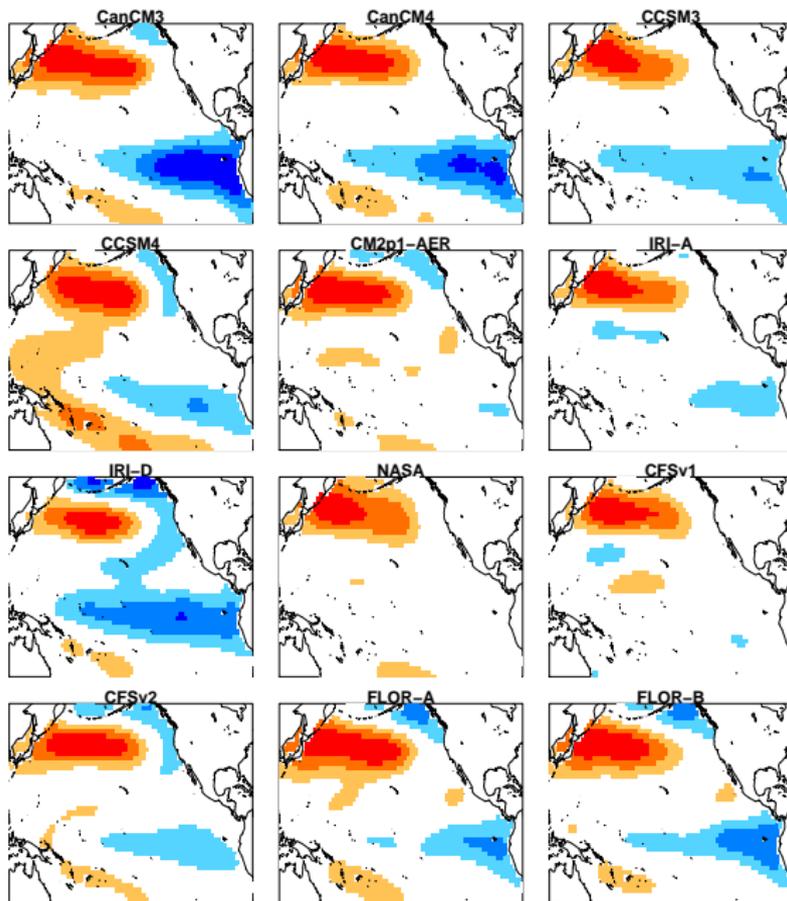


Can Climate Models Reproduce These Relations?

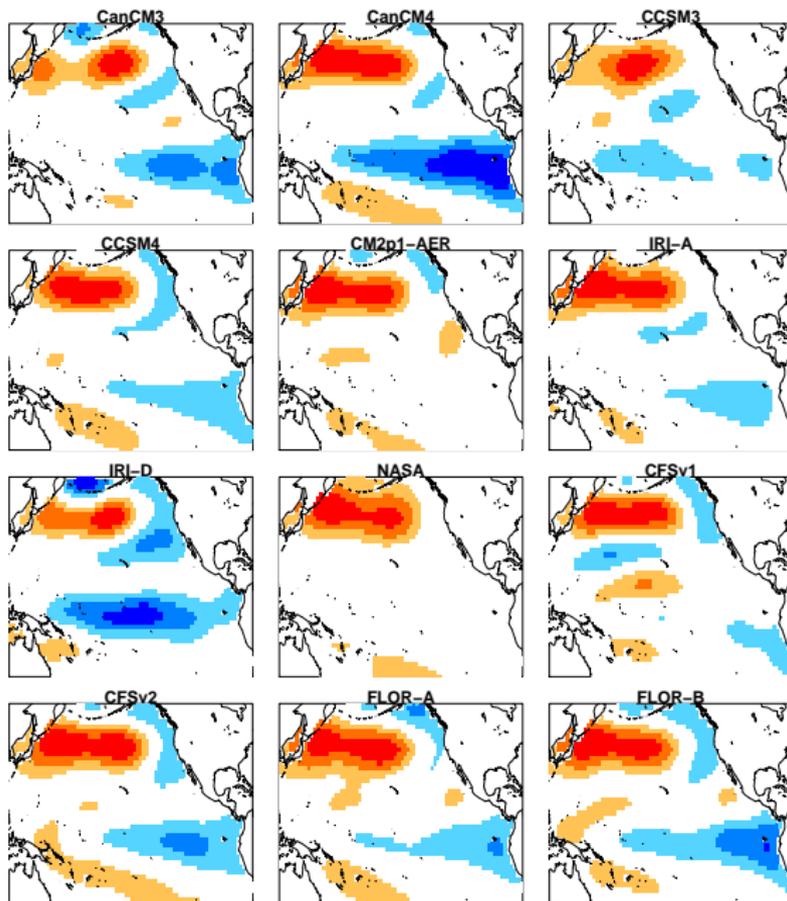
North American Multi-Model Ensemble (NMME)

- ▶ Seasonal-to-intraseasonal predictions by state-of-the-art models
 - CFSv1 NOAA Climate Forecast System version 1
 - CFSv2 NOAA Climate Forecast System version 2
 - CM2p1-AER GFDL Climate Model version 2.2
 - FLOR-A GFDL
 - FLOR-B GFDL
 - NASA NASA Goddard Observing System v5
 - IRI-D IRI-ECHAM4 Direct Coupling
 - IRI-A IRI-ECHAM4 Anomaly Coupling
 - CCSM4 NCAR Community Climate System Model
 - CCSM3 NCAR Community Climate System Model
 - CMC1 Canadian Coupled Climate Model
 - CMC2 Canadian Coupled Climate Model
- ▶ at least 6 ensemble members per model
- ▶ 8-12 month predictions/hindcasts
- ▶ 1982-2014 (not all models available during this period)
- ▶ pool initial start months January - May

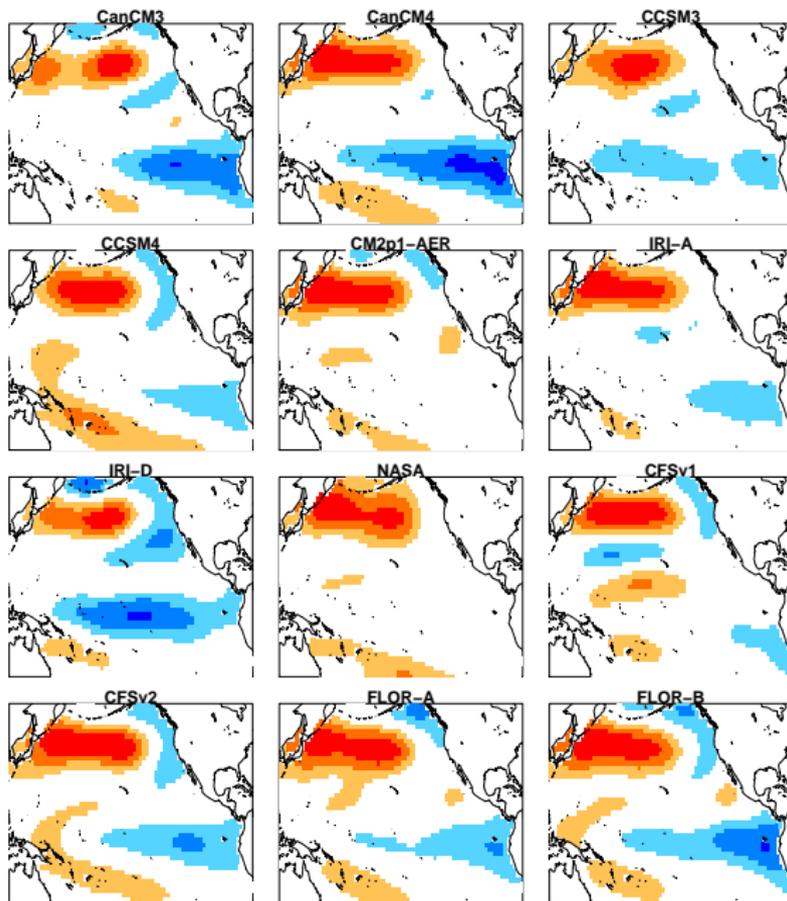
Principal Components Regression



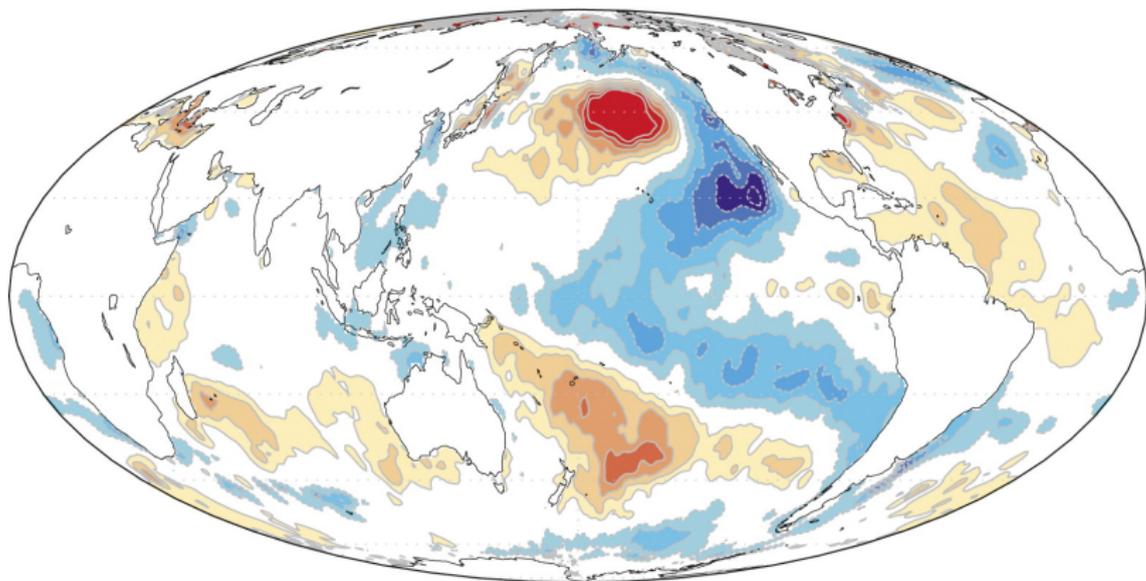
LASSO



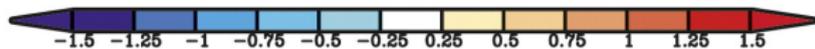
Ridge



Observed SST JJA 2011



Degrees Celsius

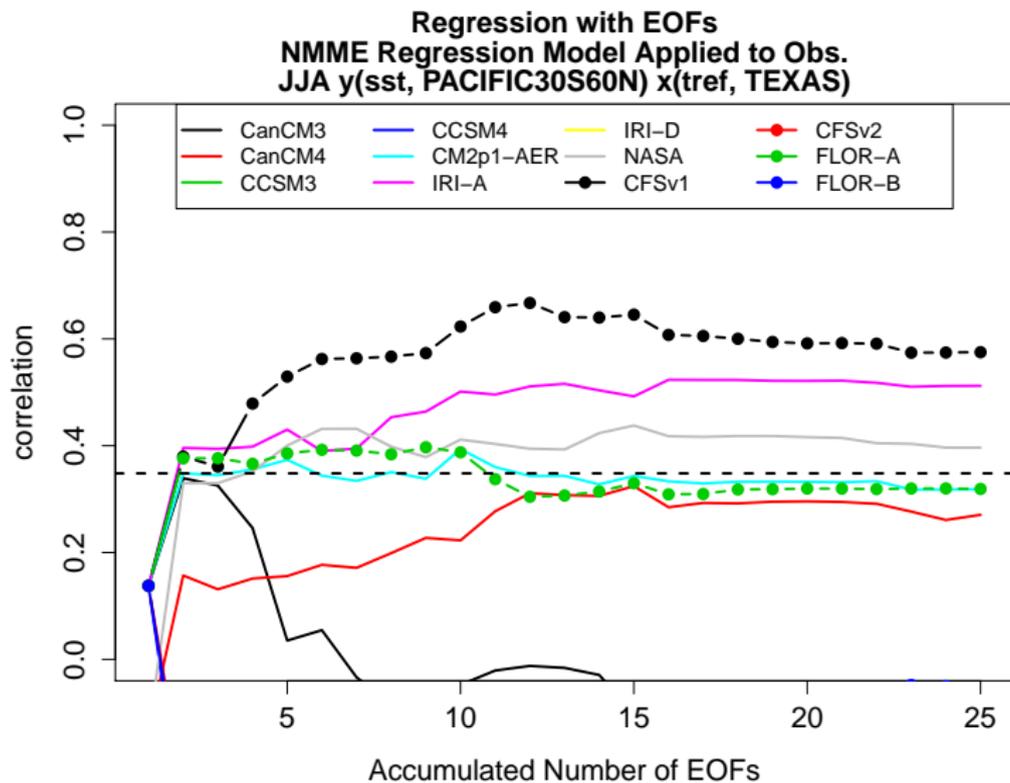


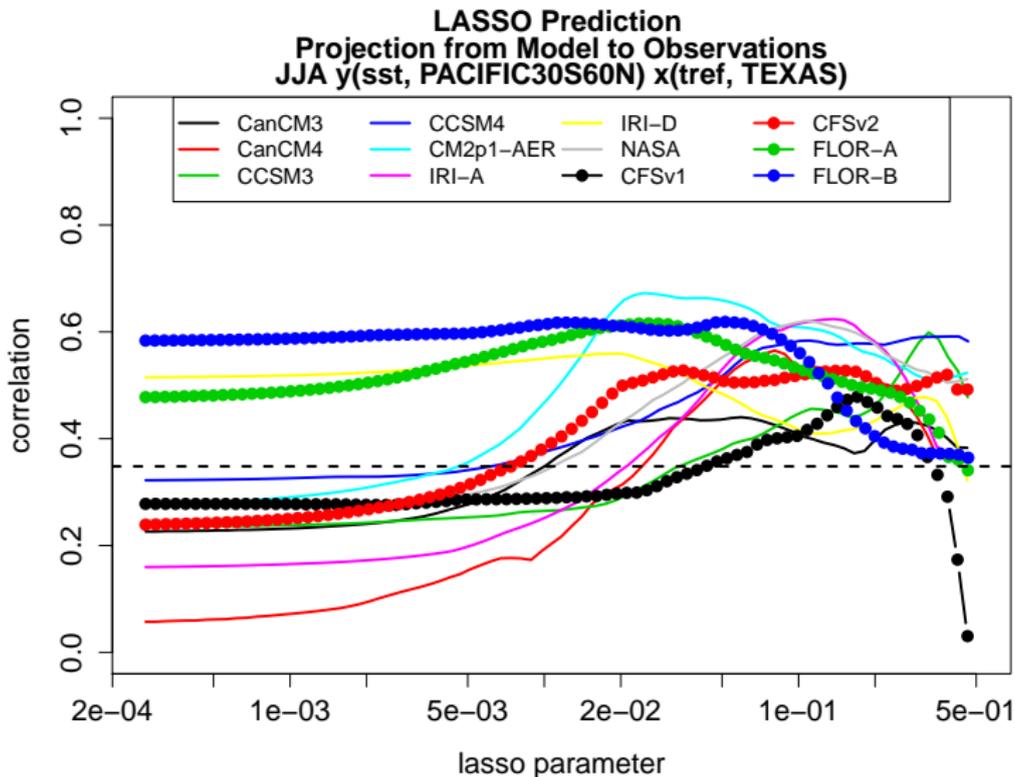
from figure 3 of Hoerling et al. 2013, J. Climate

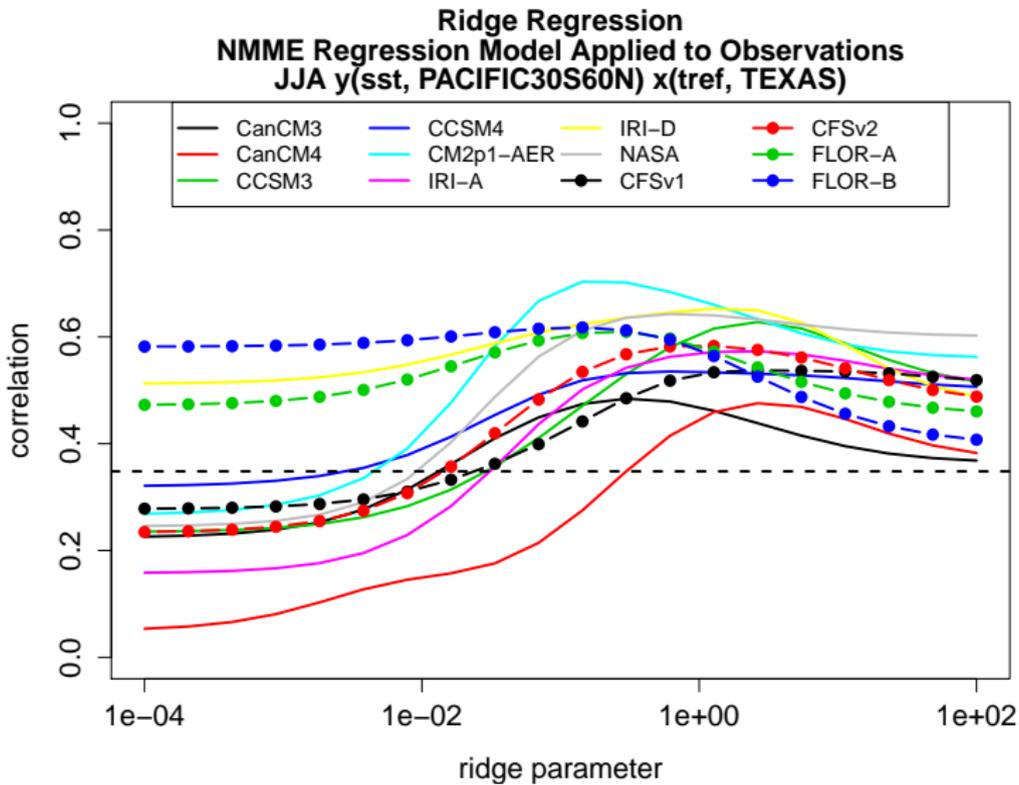
Use relations derived from climate models on observations

$$\begin{array}{ccccccc} \text{Texas temperature} & = & \text{SST in Pacific} & * & \text{weights} & + & \text{noise} \\ \text{observed} & & \text{observed} & & \text{model} & & \end{array}$$

Principal Components Regression







Summary for Finding Relations in Climate Data

1. Data fishing is a serious problem.
2. Convincing evidence of relation involves showing:
 - ▶ relation holds in independent data
 - ▶ relation reproduced by climate models
 - ▶ relation understood from simple dynamical models
3. “Large-scale” principle can be framed as a sparsity problem.
4. N. Pacific is a major predictor of Texas temperature in models and observations.
5. LASSO and Ridge find relations that generalize to independent data better than EOF method.