# Local versus non-local calibration techniques

Ángel G. Muñoz

with contributions from Simon J. Mason Andrew W. Robertson



## Outline

- 1. Model Output Statistics (MOS): bias correction, calibration.
- 2. Local calibration techniques. Examples.
- 3. Non-local calibration techniques. Examples.
- 4. Exercises: pattern-based calibration with PyCPT

## Outline

- 1. Model Output Statistics (MOS): bias correction, calibration.
- 2. Local calibration techniques. Examples.
- 3. Non-local calibration techniques. Examples.
- 4. Exercises: pattern-based calibration with PyCPT

## Model vs Observations



## Model vs Observations

#### Rainfall Climatology DJF – 1981-1989 – WRF model



## Model Output Statistics

- Because of uncertainties in initial/boundary conditions, unknown or unresolved physical processes and the chaotic nature of the climate system, *models are always subject to error*.
- Part of those errors are systematic, and can be corrected using Model Output Statistics (MOS).
- Other errors are *not* correctible, and it is customary to provide an ensemble forecast to quantify uncertainties. This leads to probabilistic forecasts.

## Model Output Statistics

- Generally speaking, MOS is any postprocessing we conduct on the raw model output (not only bias correction).
- In particular, it refers to statistical postprocessing to correct model errors.
- There are different types of biases, and different methods to correct them! Not a "unique MOS"…
- In this talk, to simplify language, any postprocessing method that corrects/calibrates model outputs are referred to as MOS.

## Calibrated Forecasts

#### Predictor (X)

### Predictand (Y)

### **Calibrated Forecast**



## The NextGen Approach



## Model Output Statistics

- It is common to use Anomaly Correlation Coefficient to assess forecast skill, but it only measures *association*.
- There are a lot of other forecast attributes of interest! (remember Caio's class).



Courtesy of Simon J. Mason

## Model Output Statistics



## Which Biases to take care of?

# Mean & Amplitude Bias



Courtesy of S. Mason

## Which Biases to take care of?

## Local biases may have obvious reasons



Climate can vary dramatically over short distances, especially in the context of precipitation and wind speeds.

## Which Biases to take care of?

# Conditional Bias (errors in patterns of variability)



Important climate features may be displaced in GCMs relative to observations: Systematic spatial biases

Courtesy B. Lyon (or was it Simon?)

## Another example: T2M



Example:

T2M Init: June

> International Research Institute for Climate and Society Earth Institute | Columbia University

60W

0.30

0.60

30W

0.90

90W

-0.30

0.00

Model

Y Spatial Loadings (Mode1)

Y Spatial Loadings (Mode2)

Negative MSSS implies that conditional biases can be LARGE. Need to correct them!

#### MSSS

MSSS: ECMWF Precip Fcst vs CMAP: 1992-2008



CORA

ECMWF Precip Fcst vs CMAP: 1992-2008



The MSSS can be decomposed into the square of the correlation coefficient, together with the squares of the conditional and mean prediction biases (Murphy 1988). If the mean biases are subtracted at the outset, the MSSS consists of only the (CORA)<sup>2</sup> and the squared conditional bias.

 $MSSS = (CORA)^2 - (conditional bias)^2$ 

FIG. 15. Mean square skill score (MSSS) between ECMWF precipitation hindcast and CMAP rainfall data over weeks 1-4.

Li and Robertson, 2015

## Outline

- 1. Model Output Statistics (MOS): bias correction, calibration.
- 2. Local calibration techniques. Examples.
- 3. Non-local calibration techniques. Examples.
- 4. Exercises: pattern-based calibration with PyCPT

- Post-processing gridbox by gridbox so the desired statistical characteristics in the model "match" the observed ones.
- Sometimes considered a "brute force" approach (in the sense that it is usually not informed by physics, just statistics).



- Typically addressing mean and amplitude biases, but there are methods that go well beyond those corrections.
- Matching spatial and temporal resolution of datasets can be important.



#### Let's see some methods (from Manzanas et al., 2019)

All the methods described in this section have been applied gridbox by gridbox considering seasonal interannual series. We use the following notation:  $y_{m,t}$  and  $y'_{m,t}$  denote the original and calibrated values for the ensemble member m at time (season/year) t,  $\hat{y}$  is the average of the ensemble mean  $(\bar{y}_t)$  on all times t,  $\hat{o}$  is the average of the observations on all times t,  $\sigma_f$  is the standard deviation of the complete ensemble (pooling all member interannual time-series) and  $\sigma_o$  is the standard deviation of the observed interannual time-series. Finally,  $\rho$  is the interannual correlation between the ensemble mean and the observational reference.

Mean and Variance Adjustment (MVA) - e.g., Leung et al. 1999

Mean correction

$$y'_{m,t} = (y_{m,t} - \hat{y}) + \hat{o}$$

Mean and amplitude correction

$$y'_{m,t} = (y_{m,t} - \hat{y})\frac{\sigma_o}{\sigma_f} + \hat{o}$$

Empirical Quantile Mapping (EQM) - e.g., Manzanas et al., 2018



**Figure 3:** The nature of QM: A biased simulated distribution (blue) is corrected towards an observed distribution (black). In the example shown the raw simulated distribution is subject to both a bias of the mean and a bias in variance. The resulting bias-corrected distribution (dashed red) approximates the observed one but is typically not identical to it (e.g. due to the sampling uncertainty during the calibration of the correction function or details of the specific QM implementation). Left panel: Example based on the probability density function (PDF). Right panel: example based on the cumulative distribution function (CDF).

Technical Report MeteoSwiss No. 270

Climate conserving realibration (CCR) - e.g., Doblas-Reyes et al., 2005

$$y_{m,t}' = \rho \frac{\sigma_o}{std(\bar{y}_t)} \bar{y}_t + \sqrt{1 - \rho^2} \frac{\sigma_o}{\sigma_f} (y_{m,t} - \bar{y}_t) + \hat{o}$$

- Corrects forecasts so they have the observed interannual variance
- Preserves inter-annual correlation

Ratio of predictable components (RPC) - e.g., Eade et al., 2014

$$\begin{split} y'_{m,t} &= \rho \frac{\sigma_o}{st d(\bar{y}_t)} (\bar{y}_t - \hat{y}) \\ &+ \sqrt{1 - \rho^2} \frac{\sigma_o}{\sqrt{var(y_{m,t} - \bar{y}_t)}} (y_{m,t} - \bar{y}_t) + \hat{o} \end{split}$$

- Uses ensemble to reduce noise
- Adjust forecast variance to ratio of predictable components in the model matches the observed one

## Correction/downscaling of GCM forecasts

• If we have a set of GCM hindcasts (x) and verifying observations (y), then we can build a regression model to relate them.

$$\hat{\mathbf{y}} = \beta_0 + \beta_1 \mathbf{X}_1$$

- The GCM forecasts becomes the "predictor" x in the regression model, and the observations becomes the "predictand" y
- Regression models trained on GCM hindcasts vs historical data are called "MOS Correction" (for Model Output Statistics)
- Once the regression coefficients have been estimated from hindcasts (e.g., past 20 years), the model can be used to correct a new forecast x for t=2021.

Linear Regression (LR) – e.g., Manzanas et al., 2019

$$o_t = \alpha + \beta \bar{y}_t + \epsilon$$

$$y'_{m,t} = \alpha + \beta \bar{y}_t + \gamma_t (y_{m,t} - \bar{y}_t)$$
  
$$\gamma_t = std(\epsilon_{fit}) \sqrt{1 + 1/n + \frac{(y_t - \bar{y}_t)^2}{(n-1)var(\epsilon_{obs})}}$$

- Simple linear regression between the ensemble mean and observations
- Adjust forecast variance by rescaling the standard deviation of the predictive distribution from the linear fit

Non-homogeneous Gaussian Regression (NGR) – e.g., Gneiting et al., 2005

$$y_{m,t}' = \alpha + \beta(\bar{y}_t - \hat{y}) + \sqrt{\gamma^2 + \delta^2 var(y_t)}(y_{m,t} - \bar{y}_t)$$

- Ensemble mean signal and spread (+constant) are used as predictors for the calibrated forecast mean and variance, respectively.
- Parameters are optimized by minimizing the ensemble CRPS

Extended Logistic Regression (ELR) – e.g., Vigaud, Robertson and Tippet, 2017



Applied at each grid point, using forecast ensemble mean



Manzanas et al 2019

## Outline

- 1. Model Output Statistics (MOS): bias correction, calibration.
- 2. Local calibration techniques. Examples.
- 3. Non-local calibration techniques. Examples.
- 4. Exercises: pattern-based calibration with PyCPT

# Non-Local Calibration Techniques Remember our regression equation?

• If we have a set of GCM hindcasts (x) and verifying observations (y), then we can build a regression model to relate them.

$$\hat{\mathbf{y}} = \beta_0 + \beta_1 \mathbf{X}_1$$

- The GCM forecasts becomes the "predictor" x in the regression model, and the observations becomes the "predictand" y
- Regression models trained on GCM hindcasts vs historical data are called "MOS Correction" (for Model Output Statistics)
- Once the regression coefficients have been estimated from hindcasts (e.g., past 20 years), the model can be used to correct a new forecast x for t=2021.

## Varieties of linear regression

- simple regression: a univariate predictor and a univariate predictand:
  y = ax + b
  - multiple regression: two or more predictors, and a single predictand

 $y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n$  (case of n predictors)

-- e.g., Principal Components Regression (PCR)

 multivariate (pattern) regression: two or more predictors, two or more predictands y = Ax + b (matrix A)

Non-local -- e.g., Canonical Correlation Analysis (CCA)

# Canonical Correlation Analysis

- To predict the observed precip anomaly field y from a GCM forecast anomaly field x, we assume the linear relationship y = Ax
- We minimize the regression error squared  $\langle (y Ax)^T(y Ax) \rangle$ by estimating the matrix A from hindcasts:  $A = (yx^T)(xx^T)^{-1}$
- This cannot be done when hindcast series is smaller than the spatial dimension of x and y
- The dimensions of **x** and **y** must be reduced!

# Canonical Correlation Analysis

- Use <u>Principal Component Analysis</u> (PCA) to reduce the spatial dimension of the regression
- Truncate to typically < 10 PCs in x and y
- Data compression! ==> So # of predictors < # of training samples, which makes the multiple regression problem well-posed
- Big bonus: The PC time series are <u>uncorrelated</u> and maximize the variance. This solves the other problem with multiple linear regression when the predictors are correlated.

# Confused? It's always the same idea...

- simple regression: a univariate predictor and a univariate Local predictand: y = ax + b
  - multipl  $\hat{y} = \beta_0 + \beta_1 X_1$ predict ctors)  $\mathbf{v} = \mathbf{a}_0$

a single

-- e.g., Principal Components Regression (PCR)

multivariate (pattern) regression: two or more predictors, two or ٠ more predictands  $\mathbf{v} = \mathbf{A}\mathbf{x} + \mathbf{b} \pmod{\mathbf{A}}$ 

Non-local

-- e.g., Canonical Correlation Analysis (CCA)

## Principal Component Timeseries as predictors



PC 1 of SST anomaly 1971-2006 DJF





Courtesy of Ousmane Nyade

## Non-Local Calibration Examples





## Non-Local Calibration: PCR

 $\hat{\mathbf{y}} = \beta_0 + \beta_1 \mathbf{X}_1$ 



EOF1

## Non-Local Calibration: PCR

 $\hat{\mathbf{y}} = \beta_0 + \beta_1 \mathbf{X}_1$ 



Year

Each location is predicted as a linear combination of the model EOFs

International Research Institute for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY

EOF 2

Week 2

## Non-Local Calibration: CCA

#### Linear combination of observed EOFs ( Linear combination of model EOFs

The CCA modes are linear combinations of the EOFs of x and y y such that their time series are maximally correlated



CCA Mode 2 Week 1 (CFSv2 SubX): Canonical correlation = 0.7865







"Pattern regression" Corrects biases in patterns

How many PCs and CCA modes to retain?







International Research Institute for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY

CCA Mode 1 Week 1 (CFSv2\_SubX): Canonical correlation = 0.8453

## Non-Local Calibration: CCA

#### Linear combination of observed EOFs ( Linear combination of model EOFs

The CCA modes are linear combinations of the EOFs of x and y y such that their time series are maximally correlated



CCA Mode 1 Week 1 (ESRL): Canonical correlation = 0.442

X Spatial Loadings

0

CCA loadings

80°W

24°N

16°N

8°N

\_'3

ESRL







"Pattern regression" Corrects biases in patterns

How many PCs and CCA modes to retain?





International Research Institute for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY

CCA Mode 1 Week 1 (CFSv2\_SubX): Canonical correlation = 0.8453



FIG. 10. Raw and MOS-adjusted S2S model forecasts and skill scores for the methods indicated in Table 1. (a)–(e) The heavy rainfall forecast for 1–7 Dec 2015 as odds, defined in Eq. (3) over the target domain. A value greater than 1 indicates that the model forecast greater-than-average odds of rainfall exceeding the 90th percentile. (f)–(j) The IGN defined in Eq. (4), with zero indicating a perfect forecast. (k)–(o) The 2AFC skill score for each grid cell; a value greater than 50 indicates that the model outperforms climatology. Different MOS models except for Raw in (a),(f),(k), which indicates the uncorrected S2S model output. In (top)–(bottom), the grid cells

that observed a 90th percentile exceedance for 1–7 Dec 2015 are outlined in black.

Doss-Gollin et al (2018)

## Outline

- 1. Model Output Statistics (MOS): bias correction, calibration.
- 2. Local calibration techniques. Examples.
- 3. Non-local calibration techniques. Examples.
- 4. Exercises: pattern-based calibration with PyCPT

Before some hands-on examples with PyCPT, let's refresh some ideas…

### Predictive skill depends on space and time scales



Thomson et al. (2018)

## #NextGen: Multi-model calibration, ensemble and verification



## Discrimination – ROC



Courtesy of D. Sydykova

## Discrimination – ROC



A ROC curve represents a classifier with the random performance level. The curve separates the space into two areas for good and poor performance levels.



It shows four AUC scores. The score is 1.0 for the classifier with the perfect performance level (P) and 0.5 for the classifier with the random performance level (R). ROC curves clearly shows classifier A outperforms classifier B, which is also supported by their AUC scores (0.88 and 0.72).

## Discrimination – ROC

ROC curves with equivalent AUC scores



Two classifiers A and B have the same AUC scores, but their ROC curves are different.



Global - Deterministic 2AFC - Model: ECMWF (uncorrected)

Muñoz et al. (2018)



Muñoz et al. (2018)



Muñoz et al. (2018)

## Ignorance Score

The Ignorance Score (IGN), or negative log-likelihood score, of a probabilistic forecast of *n* categories can be written as (Good, 1952; Roulston & Smith, 2002):

$$IGN = -\log_2(p_k) \qquad \qquad k = 1..n$$

and it can be decomposed into reliability, resolution and uncertainty terms:

IGN =	<i>= REL -</i>	(Weijs et al., 2010; Wilks, 2018)		
	calibration	sharpness (if reliable)	obs distribution	

- It measures the information deficit, or ignorance, of a person having a probabilistic forecast but not knowing the actual outcome.
- Units are *bits* of information. IGN=0 means perfect forecast (zero ignorance).
- Each bit of ignorance represents a factor-of-2 increase in uncertainty.
- Related to expected gambling return if used to place proportional bets on the future (cost-loss scenarios).

## Ignorance Score

The Ignorance Score (IGN), or negative log-likelihood score, of a probabilistic forecast of *n* categories can be written as (Good 1952; Roulston & Smith, 2002):

$$IGN = -\log_2(p_k) \qquad \qquad k = 1..n$$

and it can be decomposed into reliability, resolution and uncertainty terms:





- Model: ECMWF
- Rainfall
- Probabilistic Hindcasts
- Obs: CPC Unified
- All initializations available per month (8-9)
- Uncalibrated
- IGN, RPSS, Brier and decompositions, for Week 1-6









for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY



for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY



for Climate and Society Earth Institute | Columbia University



for Climate and Society EARTH INSTITUTE | COLUMBIA UNIVERSITY





# Let's PyCPT!

## Conclusions

- Generally speaking, we *always* need to calibrate (MOS) our raw forecasts.
- Multiple techniques, depending on the desired outcome (different forecast attributes).
- Improving one particular skill metric can decrease skill in other metrics.
- Local and non-local (or pattern-based) calibrations have their own pros and cons.
- Model Output Statistics has the potential to improve forecast skill at subseasonal timescales. In particular, EOF-based MOS methods like Canonical Correlation Analysis (and Principal Component Regression) show clear skill improvement for different regions of the world, both in magnitudes and spatial patterns, but not always.

# Local versus non-local calibration techniques

Ángel G. Muñoz

with contributions from Simon J. Mason Andrew W. Robertson

