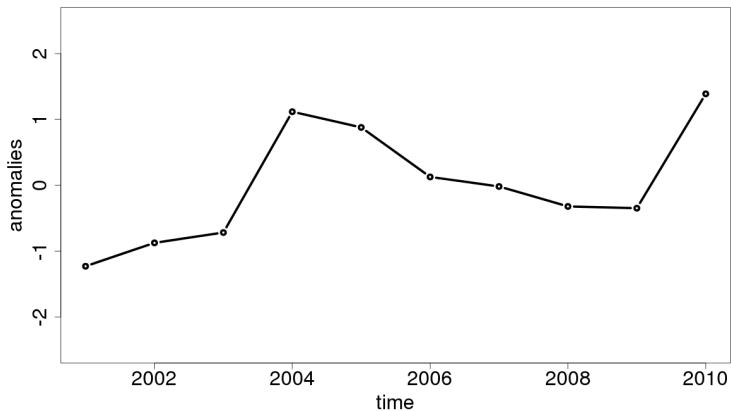


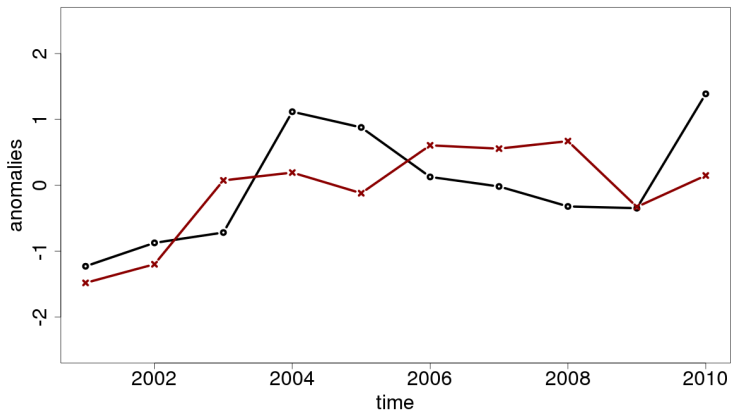
## Calibration of decadal ensemble predictions

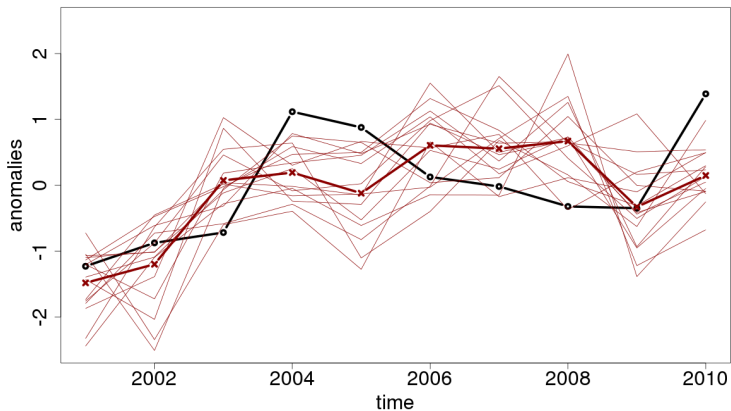
A. Pasternack, H. W. Rust, U. Ulbrich, M. A. Liniger, J. Bhend  
Freie Universität Berlin

Berlin, October, 5<sup>th</sup>, 2016

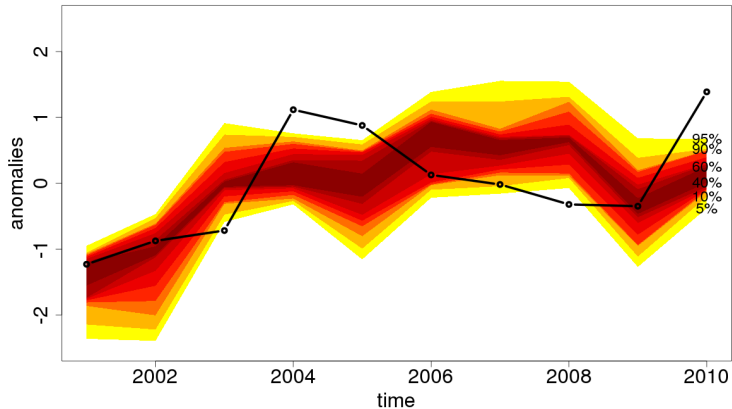
- ▶ Probabilistic forecasts
- ▶ What is a good forecast
- ▶ Re-calibrating an example forecast
- ▶ Tailor re-calibration methods to decadal predictions
- ▶ Apply re-calibration methods to decadal predictions
  - ▶ Validation







# Probabilistic forecast



# What is a good probabilistic forecast?

*„... an important goal is to maximize sharpness without sacrificing calibration.“ (Wilks, 2011; Gneiting, 2007; Murphy and Winkler, 1987)*

# What is a good probabilistic forecast?

*„... an important goal is to maximize sharpness without sacrificing calibration.“ (Wilks, 2011; Gneiting, 2007; Murphy and Winkler, 1987)*

## Sharpness

Forecasts take a risk, i.e. are frequently different from the climatological value?



# What is a good probabilistic forecast?

*„... an important goal is to maximize sharpness without sacrificing calibration.“ (Wilks, 2011; Gneiting, 2007; Murphy and Winkler, 1987)*

## Sharpness

Forecasts take a risk, i.e. are frequently different from the climatological value?

## Calibration or reliability

Probabilistic forecasts „mean what they say“, e.g. for days with a forecast of 30% chance of rain, we expect a relative frequency of 30% rainy days.

# What is a good probabilistic forecast?

*„... an important goal is to maximize sharpness without sacrificing calibration.“ (Wilks, 2011; Gneiting, 2007; Murphy and Winkler, 1987)*

## Sharpness

Forecasts take a risk, i.e. are frequently different from the climatological value?

## Calibration or reliability

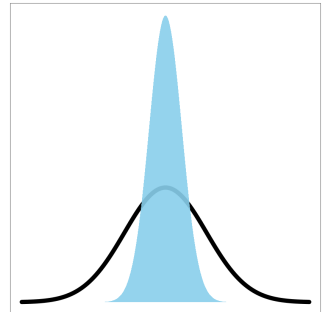
Probabilistic forecasts „mean what they say“, e.g. for days with a forecast of 30% chance of rain, we expect a relative frequency of 30% rainy days.

„Ensemble members are reliable if the MSE between the ensemble mean and observations is identical to the time mean intra-ensemble variance.“ (Palmer et al., 2006)

*„... ensemble distributions typically underestimate the true forecast uncertainty and tend to be overconfident ...“ (e.g. Weigel, 2009)*

*„... ensemble distributions typically underestimate the true forecast uncertainty and tend to be overconfident ...“ (e.g. Weigel, 2009)*

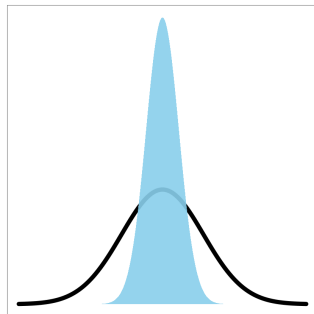
Observed relative frequency distribution  
is broader than forecasted distribution



*„... ensemble distributions typically underestimate the true forecast uncertainty and tend to be overconfident ...“ (e.g. Weigel, 2009)*

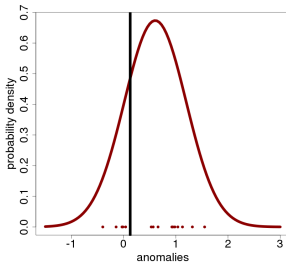
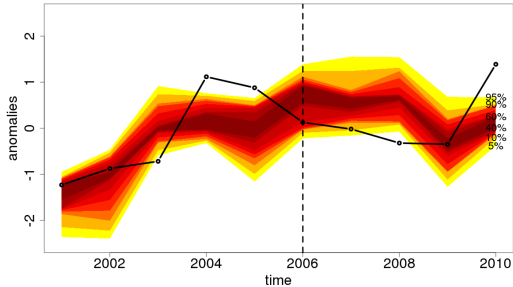
Observed relative frequency distribution  
is broader than forecasted distribution

→ **adjust ensemble spread**

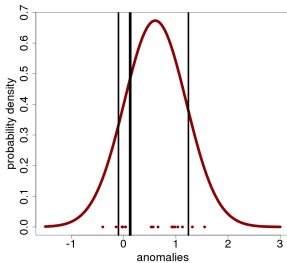
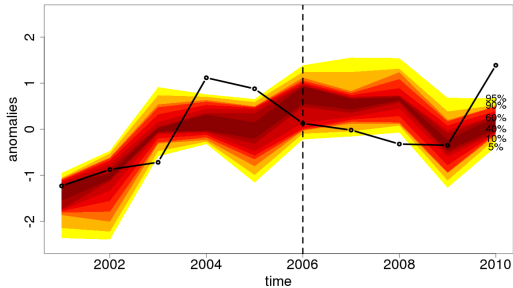


## Re-calibrating an example forecast

# Probabilistic forecast

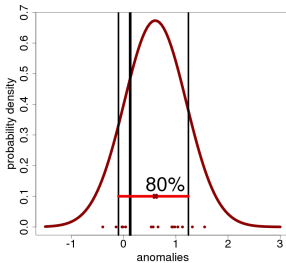
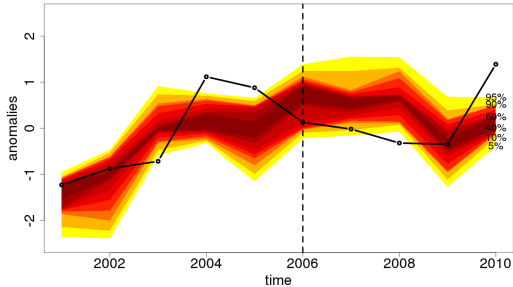


# Probabilistic forecast

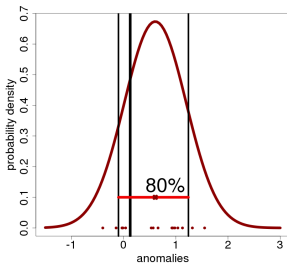
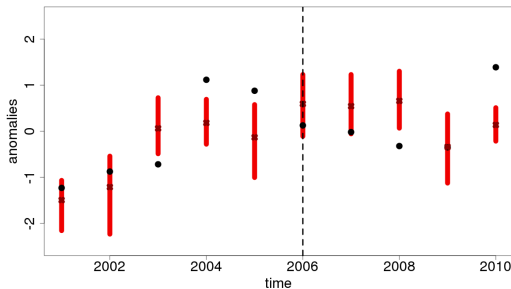




# Probabilistic forecast



# Probabilistic forecast

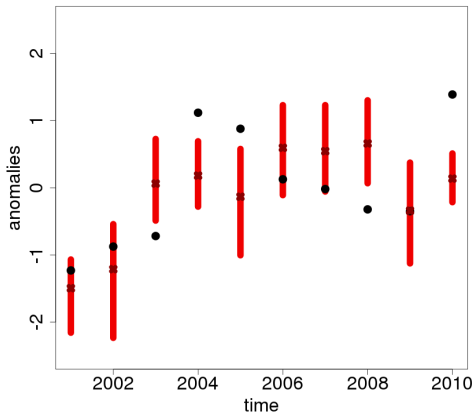


## Reliable forecast:

- ▶ 80% of ens. spread should include 80% of observations
- ▶ only 50% are covered

## What is wrong?

# Re-calibrating an example forecast

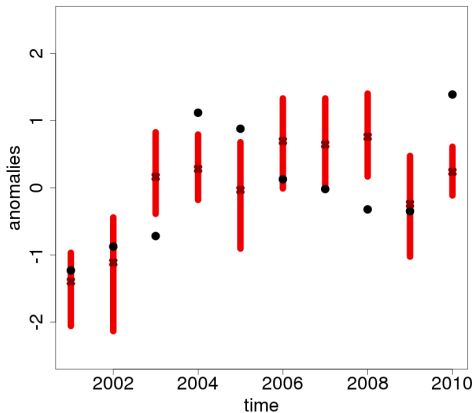


1. forecast is biased (uncond.)

Ensemble prediction:

$$f_i(t) = \mu(t) + \epsilon_i(t) \quad i = 1 \dots M$$

# Re-calibrating an example forecast



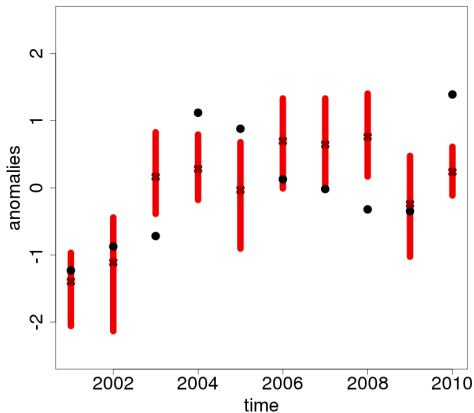
1. forecast is biased (uncond.)  
(shift ensemble mean)

2. forecast is conditionally biased

Ensemble prediction:

$$f_i(t) = \alpha + \mu(t) + \epsilon_i(t) \quad i = 1 \dots M$$

# Re-calibrating an example forecast

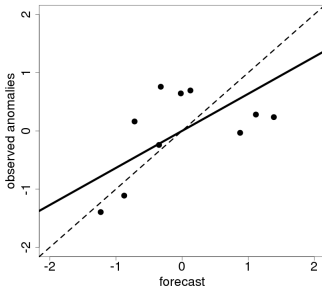


Ensemble prediction:

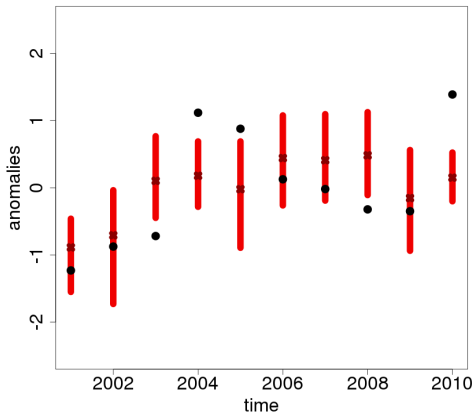
$$f_i(t) = \alpha + \mu(t) + \epsilon_i(t)$$

$$i = 1 \dots M$$

1. forecast is biased (uncond.)  
(shift ensemble mean)
2. forecast is conditionally biased



# Re-calibrating an example forecast

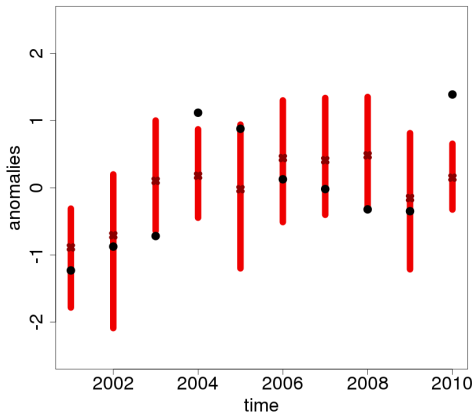


1. forecast is biased (uncond.)  
(shift ensemble mean)
2. forecast is conditionally biased  
(scale ensemble mean)
3. forecast is not reliable

Ensemble prediction:

$$f_i(t) = \alpha + \beta\mu(t) + \epsilon_i(t) \quad i = 1 \dots M$$

# Re-calibrating an example forecast



1. forecast is biased (uncond.)  
(shift ensemble mean)
2. forecast is conditionally biased  
(scale ensemble mean)
3. forecast is not reliable  
(scale ensemble spread)
4. re-calibrated forecast

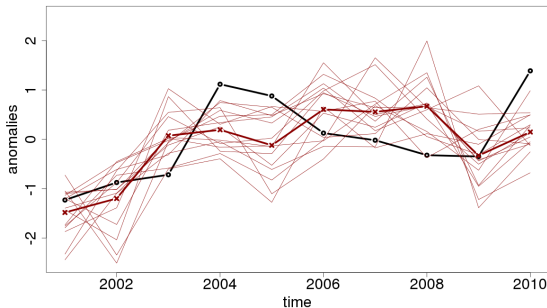
Re-calibrated ensemble:

$$f_i^{Cal}(t) = \alpha + \beta\mu(t) + \gamma\epsilon_i(t) \quad i = 1 \dots M$$

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction





## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction

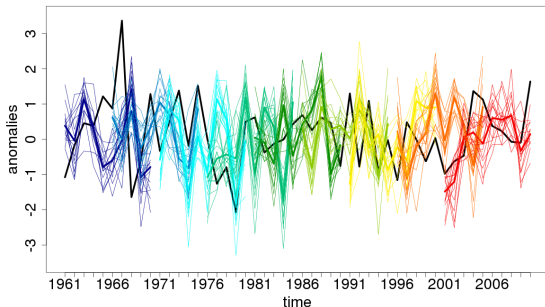
## Tasks for CALIBRATION:

Tailor re-calibration methods to decadal predictions

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction



## Tasks for CALIBRATION:

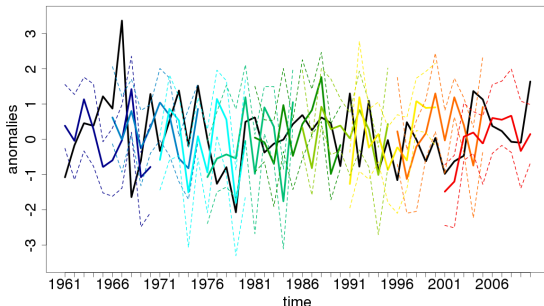
Tailor re-calibration methods to decadal predictions

- ▶ limited number of hindcasts

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction



## Tasks for CALIBRATION:

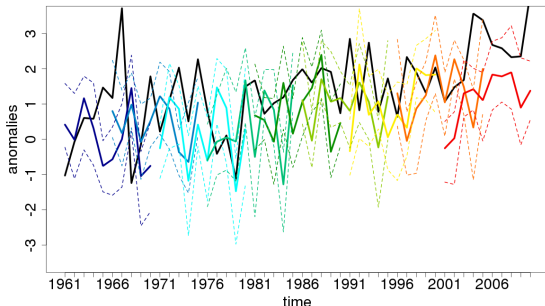
Tailor re-calibration methods to decadal predictions

- ▶ limited number of hindcasts

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction



## Tasks for CALIBRATION:

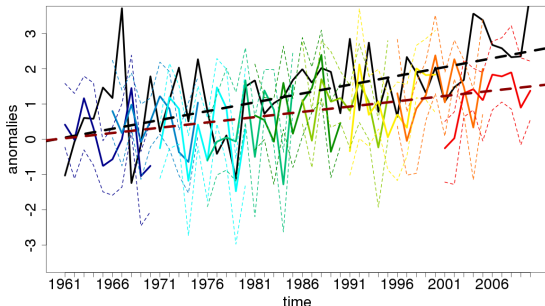
Tailor re-calibration methods to decadal predictions

- ▶ limited number of hindcasts
- ▶ climate trend

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction



## Tasks for CALIBRATION:

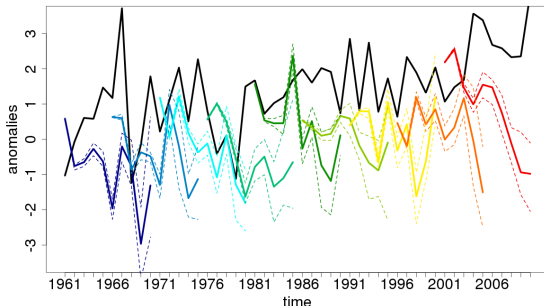
Tailor re-calibration methods to decadal predictions

- ▶ limited number of hindcasts
- ▶ climate trend

## State of the art:

Re-calibration is used in

- ▶ weather prediction
- ▶ seasonal prediction



## Tasks for CALIBRATION:

Tailor re-calibration methods to decadal predictions

- ▶ limited number of hindcasts
- ▶ climate trend
- ▶ dependence on lead years (drift)

## Ensemble prediction

$$f_i(t, \tau) = \mu(t, \tau) + \epsilon_i(t, \tau)$$

$i = 1 \dots M$  ensemble member,  $t$  = start year,  $\tau$  = lead year

with

$$\mu(t, \tau) = E(f_i(t, \tau))$$

## Re-calibrated ensemble

$$f_i^{Cal}(t, \tau) = \alpha(t, \tau) + \beta(t, \tau)\mu(t, \tau) + \gamma(t, \tau)\epsilon_i(t, \tau)$$

find  $\alpha(t, \tau)$ ,  $\beta(t, \tau)$  and  $\gamma(t, \tau)$  such that the ensemble is perfectly calibrated with maximum sharpness

1)  $\alpha$ : bias and drift, 2)  $\beta$ : conditional bias, 3)  $\gamma$ : ensemble spread

Minimize continuous ranked probability score (crps) between model  $f^{Cal}$  and observation  $O$  (*Gneiting et al. 2005*):



Minimize continuous ranked probability score (crps) between model  $f^{Cal}$  and observation  $O$  (Gneiting et al. 2005):

$$crps(F, o) = \int_{-\infty}^{\infty} (F_{f^{Cal}}(y) - F_0(y))^2 dy,$$

Minimize continuous ranked probability score (crps) between model  $f^{Cal}$  and observation  $O$  (Gneiting et al. 2005):

$$crps(F, o) = \int_{-\infty}^{\infty} (F_{f^{Cal}}(y) - F_0(y))^2 dy,$$

$F_{f^{Cal}}(y)$ : CDF derived from  $f^{Cal}$

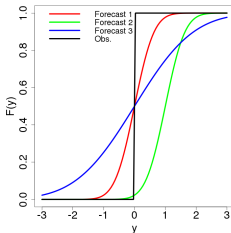
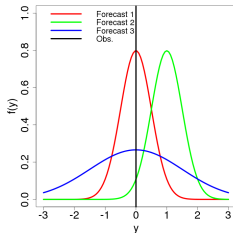
$$F_0(y) = \begin{cases} 0 & y < o \\ 1 & y \geq o \end{cases} : o \text{ is an arbitrary threshold (e.g. observation)}$$

Minimize continuous ranked probability score (crps) between model  $f^{Cal}$  and observation  $O$  (Gneiting et al. 2005):

$$crps(F, o) = \int_{-\infty}^{\infty} (F_{f^{Cal}}(y) - F_0(y))^2 dy,$$

$F_{f^{Cal}}(y)$ : CDF derived from  $f^{Cal}$

$$F_0(y) = \begin{cases} 0 & y < o \\ 1 & y \geq o \end{cases} : o \text{ is an arbitrary threshold (e.g. observation)}$$

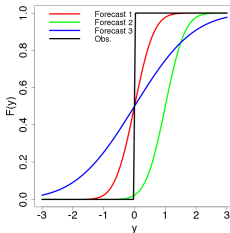
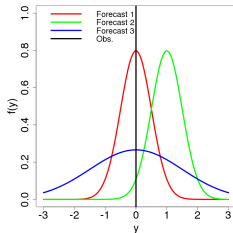


Minimize continuous ranked probability score (crps) between model  $f^{Cal}$  and observation  $O$  (Gneiting et al. 2005):

$$crps(F, o) = \int_{-\infty}^{\infty} (F_{f^{Cal}}(y) - F_0(y))^2 dy,$$

$F_{f^{Cal}}(y)$ : CDF derived from  $f^{Cal}$

$$F_0(y) = \begin{cases} 0 & y < o \\ 1 & y \geq o \end{cases} : o \text{ is an arbitrary threshold (e.g. observation)}$$



$$crps_1 = 0.12, crps_2 = 0.73, crps_3 = 0.35$$

crps measures:

- ▶ reliability
- ▶ sharpness

If the forecast distribution is Gaussian:

$$\longrightarrow f_i^{Cal}(t, \tau) \sim \mathcal{N}(\alpha(t, \tau) + \beta(t, \tau)\mu(t, \tau), \gamma(t, \tau)^2\sigma^2(t, \tau))$$

If the forecast distribution is Gaussian:

$$\longrightarrow f_i^{Cal}(t, \tau) \sim \mathcal{N}(\alpha(t, \tau) + \beta(t, \tau)\mu(t, \tau), \gamma(t, \tau)^2\sigma^2(t, \tau))$$

...the crps simplifies to:

$$crps(\mathcal{N}(\mu, \sigma^2), o) = \sigma \left\{ \frac{o-\mu}{\sigma} [2\Phi(\frac{o-\mu}{\sigma}) - 1] + 2\phi(\frac{o-\mu}{\sigma}) - \frac{1}{\sqrt{\pi}} \right\}$$

$\mu$  = ens. mean,  $\sigma$  = ens. std.,  $o$  = observation,  $\Phi, \phi$  = CDF and PDF of stand. norm. distr.

If the forecast distribution is Gaussian:

$$\longrightarrow f_j^{Cal}(t, \tau) \sim \mathcal{N}(\alpha(t, \tau) + \beta(t, \tau)\mu(t, \tau), \gamma(t, \tau)^2 \sigma^2(t, \tau))$$

...the crps simplifies to:

$$crps(\mathcal{N}(\mu, \sigma^2), o) = \sigma \left\{ \frac{o-\mu}{\sigma} [2\Phi(\frac{o-\mu}{\sigma}) - 1] + 2\phi(\frac{o-\mu}{\sigma}) - \frac{1}{\sqrt{\pi}} \right\}$$

$\mu$  = ens. mean,  $\sigma$  = ens. std.,  $o$  = observation,  $\Phi, \phi$  = CDF and PDF of stand. norm. distr.

The average score over all  $k$  pairs of forecasts and observations is:

$$\Gamma(\mathcal{N}(\alpha + \beta\mu, \gamma^2\sigma^2), o) = \frac{1}{k} \sum_{j=1}^k \sqrt{\gamma^2\sigma_j^2} \{Z_j[2\Phi(Z_j) - 1] + 2\phi(Z_j) - \frac{1}{\sqrt{\pi}}\},$$

$$Z_j = \frac{o_j - (\alpha + \beta\mu_j)}{\sqrt{\gamma^2\sigma_j^2}}$$

If the forecast distribution is Gaussian:

$$\longrightarrow f_i^{Cal}(t, \tau) \sim \mathcal{N}(\alpha(t, \tau) + \beta(t, \tau)\mu(t, \tau), \gamma(t, \tau)^2\sigma^2(t, \tau))$$

...the crps simplifies to:

$$crps(\mathcal{N}(\mu, \sigma^2), o) = \sigma \left\{ \frac{o-\mu}{\sigma} [2\Phi(\frac{o-\mu}{\sigma}) - 1] + 2\phi(\frac{o-\mu}{\sigma}) - \frac{1}{\sqrt{\pi}} \right\}$$

$\mu$  = ens. mean,  $\sigma$  = ens. std.,  $o$  = observation,  $\Phi, \phi$  = CDF and PDF of stand. norm. distr.

The average score over all  $k$  pairs of forecasts and observations is:

$$\Gamma(\mathcal{N}(\alpha + \beta\mu, \gamma^2\sigma^2), o) = \frac{1}{k} \sum_{j=1}^k \sqrt{\gamma^2\sigma_j^2} \{Z_j[2\Phi(Z_j) - 1] + 2\phi(Z_j) - \frac{1}{\sqrt{\pi}}\},$$

$$Z_j = \frac{o_j - (\alpha + \beta\mu_j)}{\sqrt{\gamma^2\sigma_j^2}}$$

$$\alpha = \alpha(t, \tau) = (a_0 + a_1t) + (a_2 + a_3t)\tau + (a_4 + a_5t)\tau^2 + (a_6 + a_7t)\tau^3 + \dots$$

$$\beta = \beta(t, \tau) = (b_0 + b_1t) + (b_2 + b_3t)\tau + (b_4 + b_5t)\tau^2 + (b_6 + b_7t)\tau^3 + \dots$$

$$\gamma = \gamma(t, \tau) = (c_0 + c_1t) + (c_2 + c_3t)\tau + (c_4 + c_5t)\tau^2 + (c_6 + c_7t)\tau^3 + \dots$$

$\longrightarrow$  find an  $a_0, b_0, c_0, \dots, a_7, b_7, c_7$  that minimize  $\Gamma$



## Example: Surface Temperature

## Data:

- ▶ Surface temperature over the North Atlantic region
- ▶ Model: MPI-ESM-LR, Prototype (GECCO2)
- ▶ 15 ensemble members
- ▶ Initialisation years: 1961-2000
- ▶ Annual mean
- ▶ Reference: NCEP 20CR

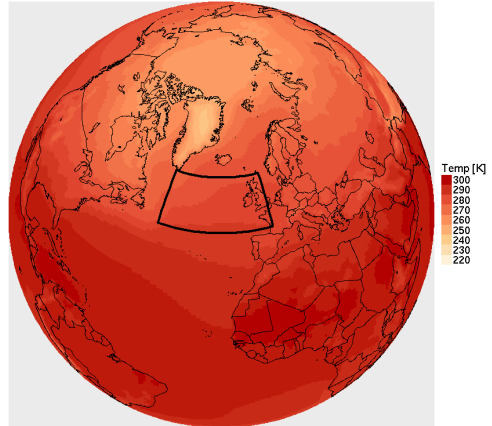
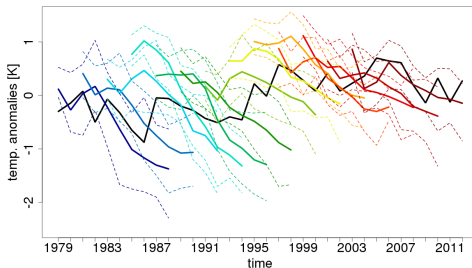


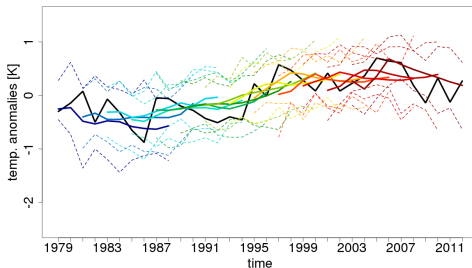
Figure: ST time mean for NCEP 20CR

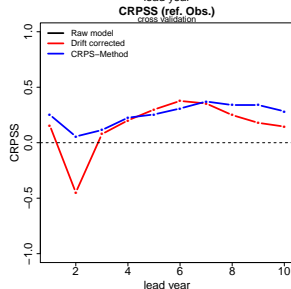
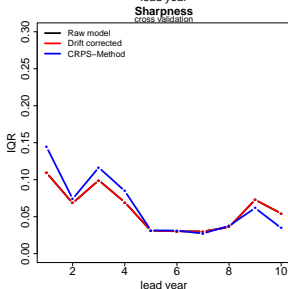
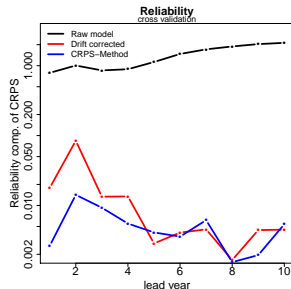
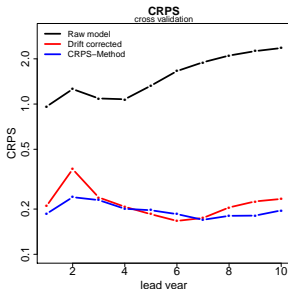
# Apply re-calibration method to decadal prediction

**Before:**



**After:**





## What is a good probabilistic forecast?

*„... an important goal is to maximize sharpness without sacrificing reliability.“ (Wilks, 2011; Gneiting, 2007; Murphy and Winkler, 1987)*

- ▶ CRPS minimization method by *Gneiting et al. 2005* addresses to reliability and sharpness for seasonal prediction.
- ▶ The developed extension to decadal predictions also includes a lead time dependent drift correction.

## Validation

- ▶ CRPS method is mostly superior to drift correction and climatology w.r.t. predictive skill.
- ▶ Sharpness will be decreased to obtain good reliability.

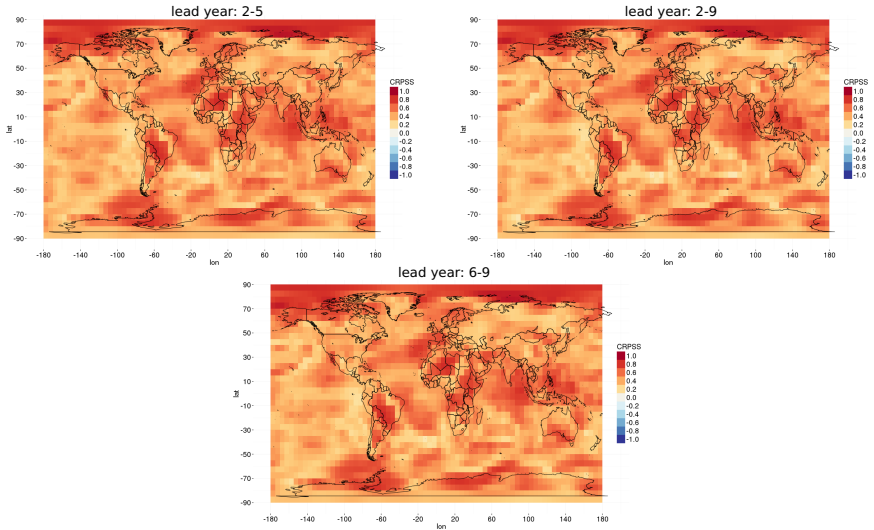
**Thank you for your attention!**

- ▶ Gneiting, T., A. E. Raftery, A. H. Westveld, and T. Goldman, 2005: Calibrated probabilistic forecasting using ensemble model output statistics and minimum crps estimation. *Monthly Weather Review*, 133, 1098–1118.
- ▶ Gneiting, T., and A. E. Raftery, 2007: Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102 (477), 359–378.
- ▶ Weigel, A., M. A. Liniger, and C. Appenzeller., 2009: Seasonal ensemble forecasts: Are recalibrated single models better than multimodels? *Mon. Weather Rev.*, 137(4), 1460–1479.
- ▶ Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. Academic Press.

# Appendix



## CRPSS: Re-calibrated vs. Climatology (cross validation)



## CRPSS: Re-calibrated vs. Drift Corrected (cross validation)

