

Forecast evaluation III

Thordis L. Thorarinsdottir

Norwegian Computing Center, Oslo, Norway

www.nr.no/~thordis

CUSO winter school 2021

Outline for this lecture

Assume we have a prediction $p \in \mathcal{P}$ and an observation $o \in \mathcal{O}$ where we wish to measure the skill of the prediction by applying a function

$$s : \mathcal{P} \times \mathcal{O} \longrightarrow \mathbb{R}$$

with a lower function value indicating a better skill.

- What if we only care about a subset of the observations, e.g. the extremes?
- What if we are working in high dimensions?
- What if the observation is also given by a distribution?

Forecast failure: how the Met Office lost touch with reality

Ideology has corrupted a valuable British institution

Rupert Darwall 13 July 2013

118 Comments



Verifying only the extremes erases propriety

Amisano and Giacomini (JBES, 2007) consider the restricted score

$$R^*(F, y) = -\mathbb{1}\{y \geq t\} \log f(y).$$

However, if $g(y) > f(y)$ for all $y \geq t$, then

$$\mathbb{E}R^*(G, y) < \mathbb{E}R^*(F, y)$$

independent of the true sampling density.

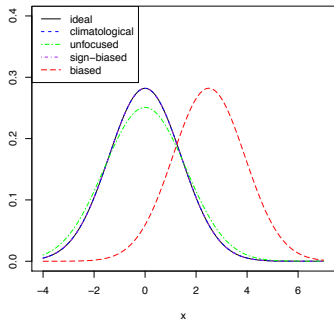
Indeed, if the forecaster's belief is F , his best prediction under R^* is

$$g(y) = \frac{f(y)}{\int_t^\infty f(x) dx} \mathbb{1}\{y \geq t\}$$

(Gneiting and Ranjan, JBES, 2011).

Demonstration by simulation

True data distribution: $G_t = N(\mu_t, 1)$ with $\mu_t \sim N(0, 1)$.

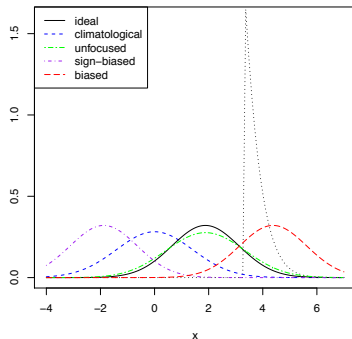


Forecaster	F_t
<i>Ideal</i>	$N(\mu_t, 1)$
<i>Sign-biased</i>	$N(-\mu_t, 1)$
<i>Climatological</i>	$N(0, 2)$
<i>Unfocused</i>	$\frac{1}{2} \{ N(\mu_t, 1) + N(\mu_t + \tau_t, 1) \}$
<i>Biased</i>	$N(\mu_t + 2.5, 1)$

Here, $\tau_t = \pm 1$ with probability 1/2.

(Lerch et al., SS, 2015)

Results for $y > 4.65$ (99th percentile)



Forecaster	CRPS*	LogS*
<i>Ideal</i>	1.36	8.47
<i>Sign-biased</i>	5.01	16.87
<i>Climatological</i>	2.92	4.75
<i>Unfocused</i>	1.34	2.69
<i>Biased</i>	0.55	1.38

Better: Use threshold-weighted scoring rules

Diks *et al.* (JE, 2011) propose the **conditional likelihood score**

$$R(F, y) = -\omega(y) \log \left(\frac{f(y)}{\int \omega(x) f(x) dx} \right)$$

and the **censored likelihood score**

$$R(F, y) = -\left[\omega(y) \log f(y) + (1 - \omega(y)) \log \left(1 - \int \omega(x) f(x) dx \right) \right].$$

Better: Use threshold-weighted scoring rules

Gneiting and Ranjan (JBES, 2011) propose the **threshold weighted CRPS**

$$\begin{aligned} R(F, y) &= \int (F(x) - \mathbb{1}\{y \leq x\})^2 \omega(x) dx \\ &= \int_0^1 (F^{-1}(\tau) - y) (\mathbb{1}\{y \leq F^{-1}(\tau)\} - \tau) \omega(\tau) d\tau. \end{aligned}$$

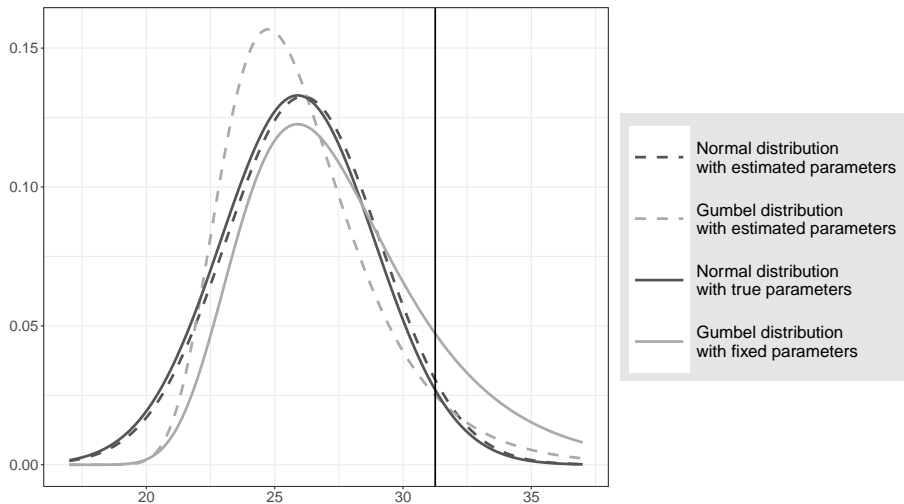
Here, we may e.g. set

$$w_1(x) = \mathbb{1}\{x \geq u\}$$

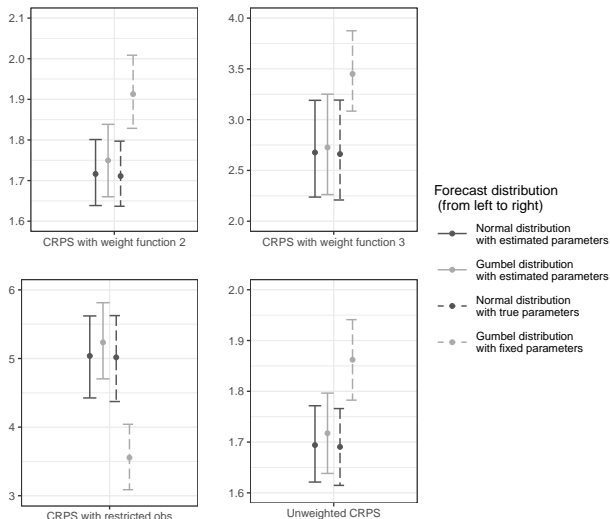
$$w_2(x) = 1 + \mathbb{1}\{x \geq u\}$$

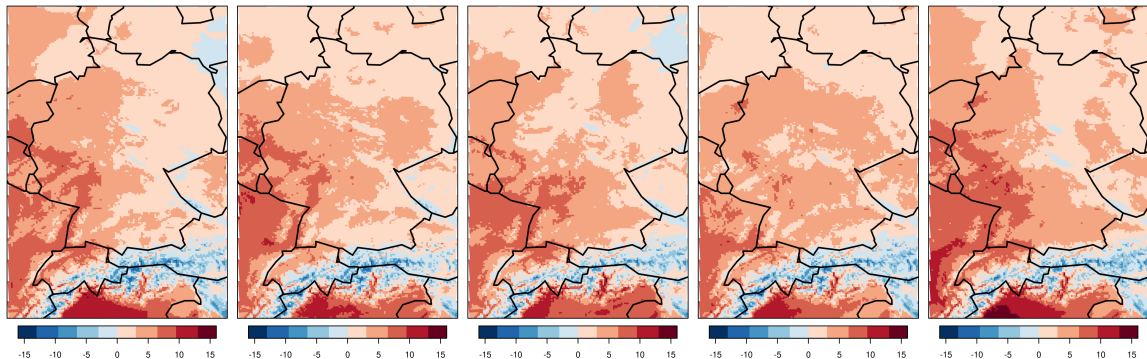
$$w_3(x) = 1 + \mathbb{1}\{x \geq u\} u$$

An extreme version of the example from last lecture



Results for 1 000 forecast-observation pairs





Three scores for multivariate forecasts

- ① The **Dawid-Sebastiani (DS) score**

$$R(F, y) = \log \det \Sigma_F + (y - \mu_F)^\top \Sigma_F^{-1} (y - \mu_F)$$

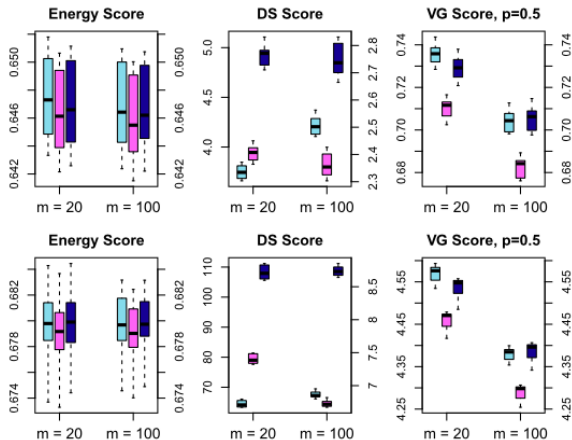
- ② The **energy score (ES)**

$$R(F, y) = \mathbb{E}_F \|X - y\| - \frac{1}{2} \mathbb{E}_F \|X - X'\|$$

- ③ The **variogram score**

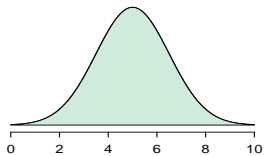
$$\mathbb{R}_p(F, y) = \sum_{i=1}^d \sum_{j=1}^d \omega_{ij} (|y_i - y_j|^p - \mathbb{E}_F |X_i - X_j|^p)^2$$

ES lacks discrimination; DS hard to estimate

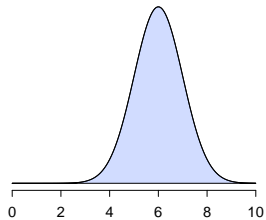


Too weak (light blue), adequate (violet) and too strong (dark blue) correlation in 5 (top) and 15 (bottom) dimensions (Scheuerer and Hamill, MWR, 2015)

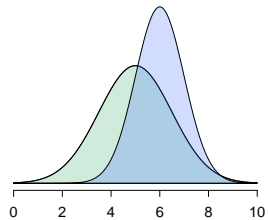
(a) Forecast



(b) Observation



(c) Comparison



Propriety condition for divergences

Two distributions may be compared using a **divergence function**,

$$d : \mathcal{F} \times \mathcal{F} \rightarrow [0, \infty], \quad d(F, F) = 0 \quad \forall F \in \mathcal{F}.$$

Definition (Thorarinsdottir, Gneiting and Gissibl, 2013)

Let $Y_1, \dots, Y_k \sim G$ and G_k be the corresponding empirical CDF. A divergence function d is *k-proper* if

$$\mathbb{E} d(G, G_k) \leq \mathbb{E} d(F, G_k).$$

Similarly, d is *asymptotically proper* if

$$\lim_{k \rightarrow \infty} \mathbb{E} d(G, G_k) \leq \lim_{k \rightarrow \infty} \mathbb{E} d(F, G_k),$$

for all $F, G \in \mathcal{F}$.

Many well known distances don't fulfill this condition

The area validation metric is given by

$$d(F, G) = \int |F(t) - G(t)| dt$$

Let $G \sim \mathcal{U}([0, 1])$ and F_k discrete with probability mass $1/k$ in $x = i/(k+1)$ for $i = 1, \dots, k$.
Then

$$\frac{1}{4} = \mathbb{E}_G d(F_1, \hat{G}_1) < \mathbb{E} d(G, \hat{G}_1) = \frac{1}{3}.$$

Similar example can be constructed for the Kolmogorov-Smirnov distance

$$d(F, G) = \sup_{t \in \mathbb{R}} |F(t) - G(t)|.$$

Every proper scoring rule defines a k -proper divergence function

Theorem (Thorarinsdottir, Gneiting and Gissibl, 2013)

Assume that $R(G, G) \neq +\infty$ and let

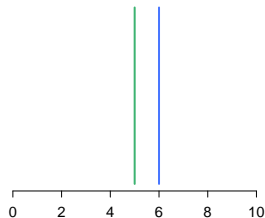
$$d(F, G) = R(F, G) - R(G, G),$$

where R is a proper scoring rule. Then d is k -proper for all $k = 1, 2, \dots$

Note that

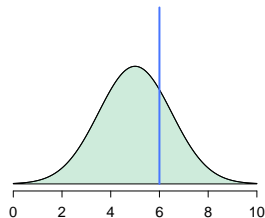
- $d(F_m, G_k)$ and $\frac{1}{k} \sum_i^k R(F_m, y_i)$ will result in the same ranking of F_1, \dots, F_M .
- it holds that $d(G, G) = 0$, while $R(G, G)$ might depend on G .

Examples



$$(x - y)^2$$

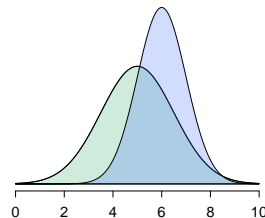
$$|x - y|$$



$$(\text{mean}(F) - y)^2$$

$$\int [F(t) - \mathbb{1}\{t \geq y\}]^2 dt$$

$$-\log(f(y))$$



$$(\text{mean}(F) - \text{mean}(G))^2$$

$$\int [F(t) - G(t)]^2 dt$$

$$\int g(u) \log \frac{g(u)}{f(y)} d\lambda(u)$$



A practical example: Climate services

NORSK KLIMASERVICESENTER – NEDLASTING AV GRIDDATA

Nedlastning av griddata

- Velg hva du vil laste ned under
- Marker ønsket område i kartet eller velg fra listen til høyre
- Legg utvalget i nedlastingskurven

Utslippsscenario og modell

Utslippsscenario

- ☐ RCP8.5
- ☐ RCP4.5

Klimamodell

- ☐ CNRM, CCLM, 1971-2100
- ☐ CNRM, RCA, 1971-2100
- ☐ EC-EARTH, CCLM, 1971-2100
- ☐ EC-EARTH, HIRHAM, 1971-2100
- ☐ EC-EARTH, RACMO, 1971-2100
- ☐ EC-EARTH, RCA, 1971-2100
- ☐ HADGEM, RCA, 1971-2100
- ☐ IPSL, RCA, 1971-2100

Før du kan laste ned må du:

- velge utslippsscenario
- velge klimamodell
- angi en tidsperiode
- velge klima/hydrologisk variabel
- velge område

Velg fylke:

- ☐ Akershus
- ☐ Aust-Agder
- ☐ Buskerud
- ☐ Finnmark
- ☐ Hedmark
- ☐ Hordaland
- ☐ Mere og Romsdal
- ☐ Nord-Trøndelag
- ☐ Nordland
- ☐ Oppland
- ☐ Oslo
- ☐ Rogaland
- ☐ Sogn og Fjordane
- ☐ Svalbard
- ☐ Sør-Trøndelag
- ☐ Telemark
- ☐ Troms
- ☐ Vest-Agder
- ☐ Vestfold
- ☐ Østfold

Topografisk gråtonekart

Meteorologisk institutt

BIERKNES CENTRE for Climate Research

Om griddata

Kontakt oss

©2017 NVE

Climate models and climate projections

Schematic for Global Atmospheric Model

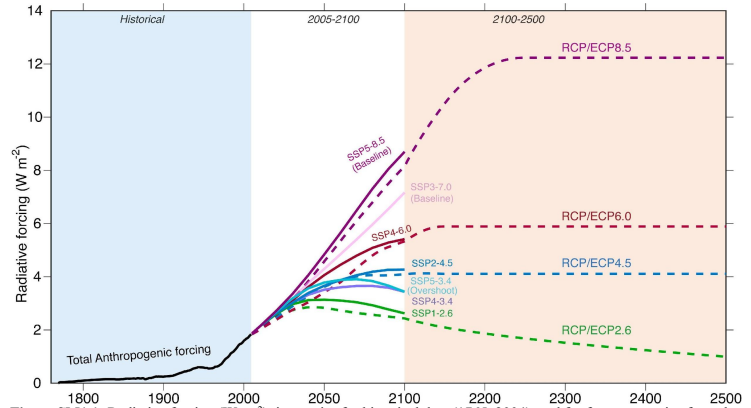
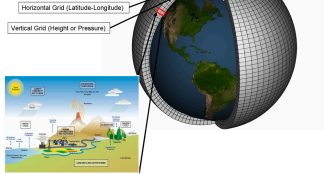


Figure on the right from IPCC.

How to evaluate climate predictions/projections?

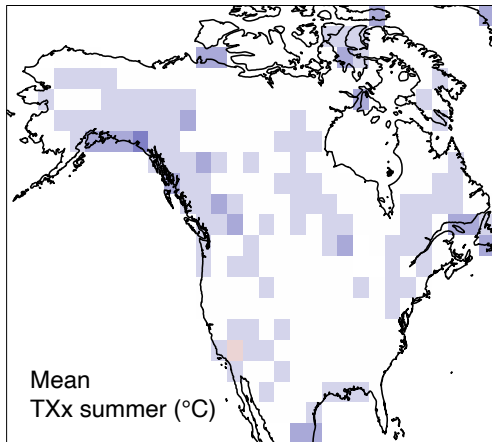
Climate models are difficult to compare to data. Often climatologists compute some summary statistic (...) and compare climate models using observed (or rather estimated) forcings to the observed (or rather estimated) temperatures.

(...) it seems more appropriate to compare the distribution (over time and space) of climate model output to the corresponding distribution of observed data.

Guttorp (E, 2011)

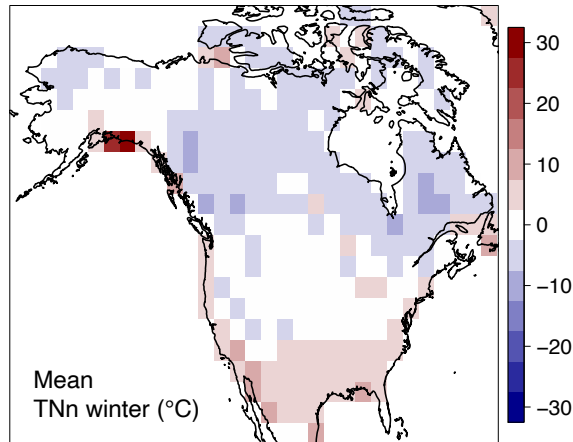
Which is the better truth, model or data?

ERA5 minus HadEX2



min. = -14.96, max. = 3.23,
mean = -2.42, MAE = 2.65

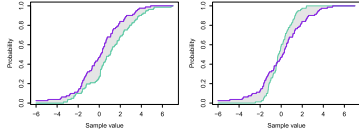
ERA5 minus HadEX2



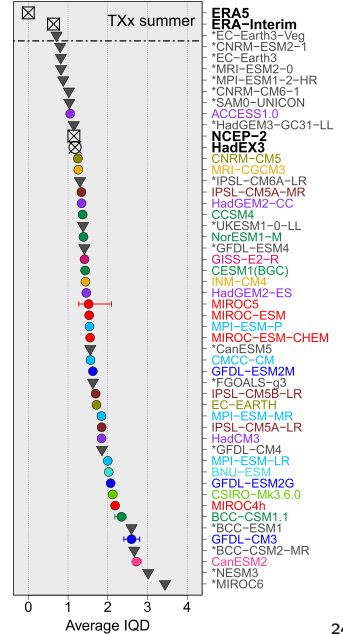
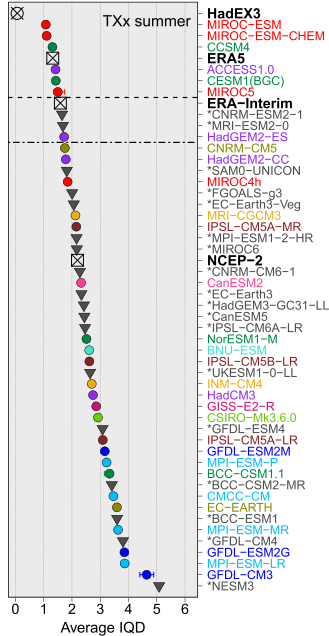
min. = -9.45, max. = 27.87,
mean = -0.85, MAE = 3.91

We use the IQD

$$d(F, G) = \int (F(t) - G(t))^2 dt$$

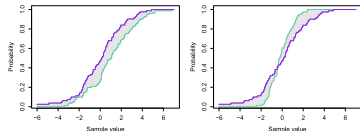


(T et al., ERL, 2020)

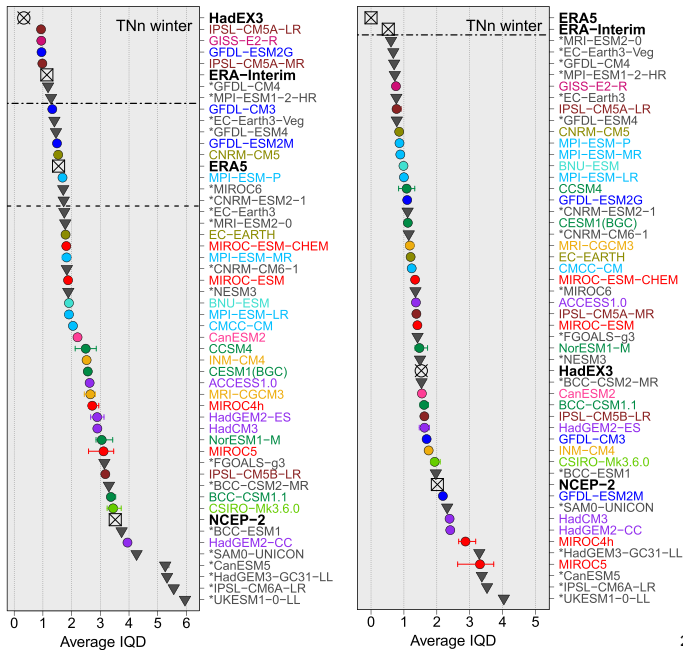


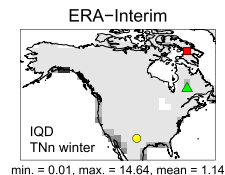
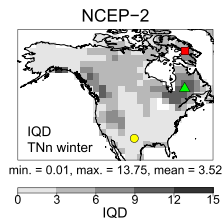
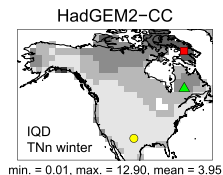
We use the IQD

$$d(F, G) = \int (F(t) - G(t))^2 dt$$

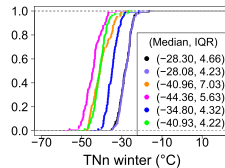
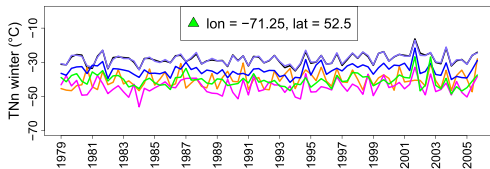
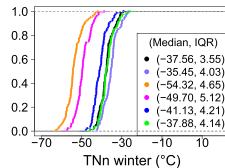
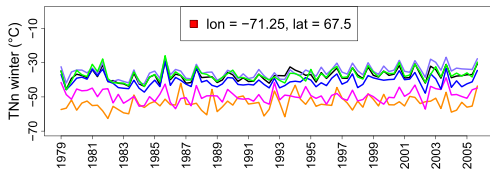


(T et al., ERL, 2020)





HadEX2 HadEX3 HadGEM2-CC NCEP-2 ERA-Interim ANUSPLIN+Livneh



On the climate scale, we generally work with anomalies rather than absolute values

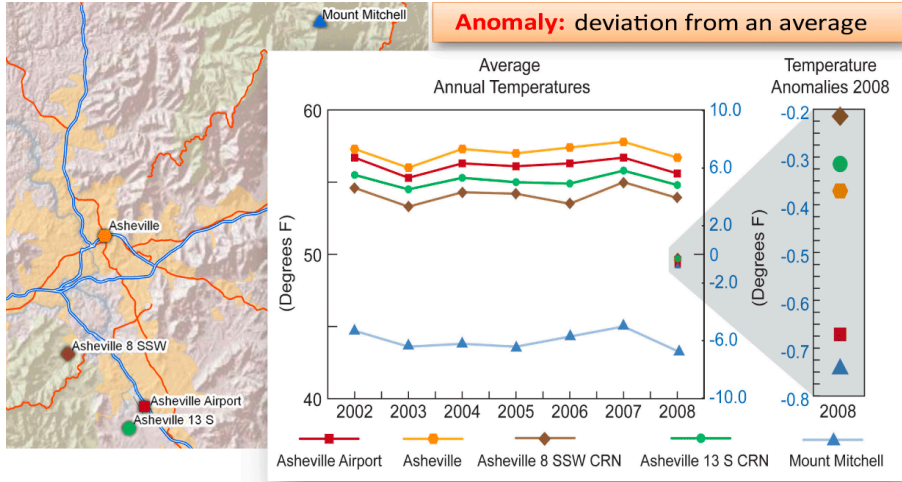
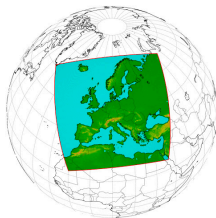
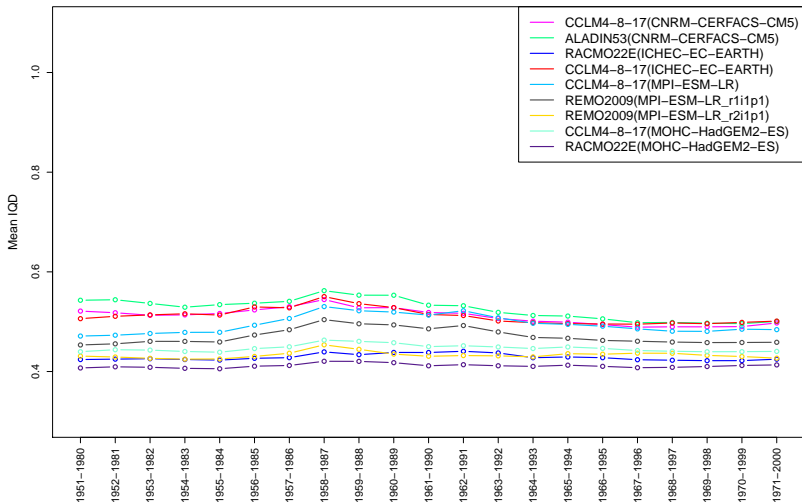


Figure from ncdc.noaa.gov.

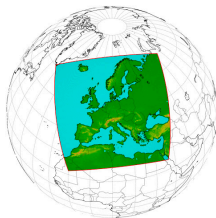
Standard reference periods are 30 years



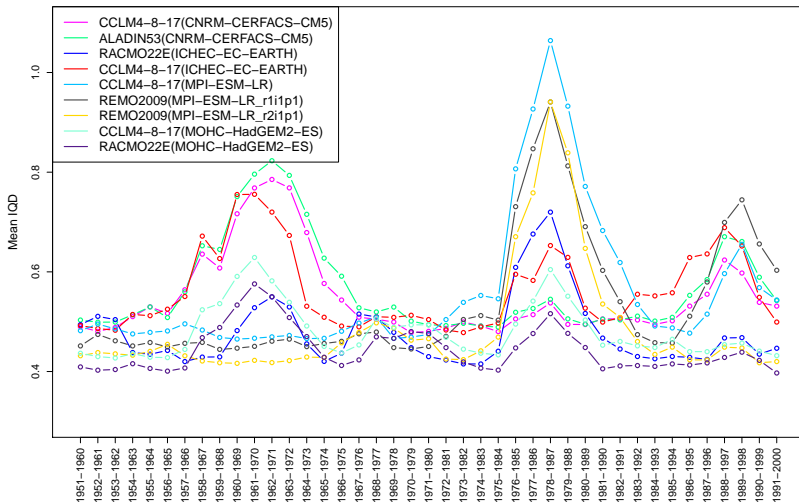
Anomalies winter (DJF), 1950–2005



10 year reference periods result in unstable rankings

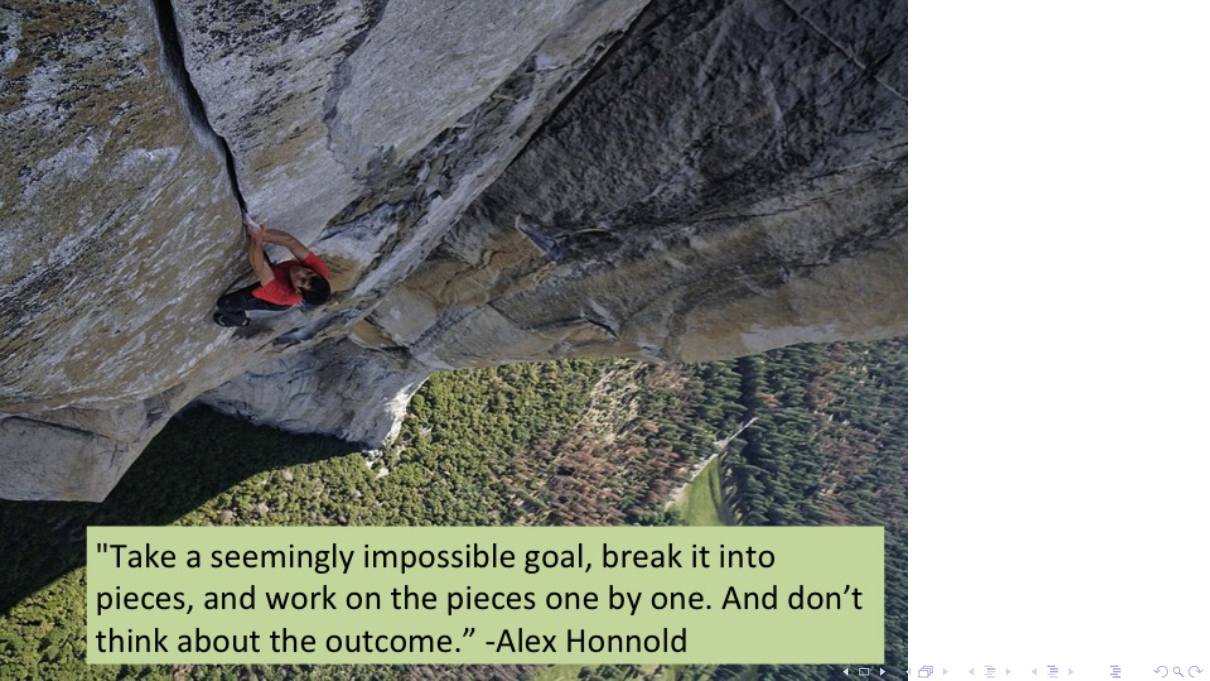


Anomalies winter (DJF), 1950–2005



Conclusions

- From this morning: Performance measure should be selected with care, preferably used in groups.
- Forecaster's dilemma: Verification on extreme events only is bound to discredit skillful forecasters. The only remedy is to consider all available cases when evaluating the models.
- Careful application of weight functions can help interpreting prediction skill in certain regions of interest. In particular, the weighted versions of the CRPS share (almost all of) the desirable properties of the unweighted CRPS.
- Overall: The framework presented here provides a unified setting for comparing two values, a value and a distribution, or two distributions.

A photograph of Alex Honnold climbing the sheer, light-colored rock face of El Capitan. He is wearing a red shirt and black pants, and is positioned on the left side of the frame, reaching up. The background shows a dense green forest below the cliff. A green text box is overlaid at the bottom of the image.

"Take a seemingly impossible goal, break it into pieces, and work on the pieces one by one. And don't think about the outcome." -Alex Honnold

References

- 1 A Amisano and R Giacomini (2007): Comparing density forecasts via weighted likelihood ratio tests. *Journal of Business and Economic Statistics*, 25(2), 177-190.
- 2 T Gneiting and R Ranjan (2011): Comparing density forecasts using threshold-and quantile-weighted scoring rules. *Journal of Business and Economic Statistics*, 29(3), 411-422.
- 3 S Lerch, T Thorarinsdottir, F Ravazzolo and T Gneiting (2017): Forecaster's dilemma: Extreme events and forecast evaluation. *Statistical Science*, 32(1), 106-127.
- 4 C Diks, V Panchenko and D van Dijk (2011): Likelihood-based scoring rules for comparing density forecasts in tails. *Journal of Econometrics*, 163(2), 215-230.
- 5 M Scheuerer and T Hamill (2015): Variogram-based proper scoring rules for probabilistic forecasts of multivariate quantities. *Monthly Weather Review*, 143(4), 1321-1334.
- 6 T Thorarinsdottir, T Gneiting and N Gissibl (2013): Using proper divergence functions to evaluate climate models. *SIAM/ASA Journal on Uncertainty Quantification*, 1(1), 522-534.
- 7 T Thorarinsdottir, J Sillmann, M Haugen, N Gissibl and M Sandstad (2020): Evaluation of CMIP5 and CMIP6 simulations of historical surface air temperature extremes using proper evaluation methods. *Environmental Research Letters*, 15(12), 124041.