



Forecast verification

..... a few comments

Anna Ghelli

anna.ghelli@ecmwf.int



Forecast quality versus forecast value

A forecast has high **QUALITY** if it predicts the observed conditions well according to some objective or subjective criteria.



Quality but no value

A forecast has **VALUE** if it helps the user to make a better decision.



Value but no quality

Verification goals and process

What are our goals with forecast evaluation?

Evaluate usefulness of forecasts

In general?

For specific users?

Improve ensemble and modeling system

Track changes in forecast performance over time

PROCESS

Start by determining

What are the questions we want to answer??

Impact forecasts

**Air
travel**

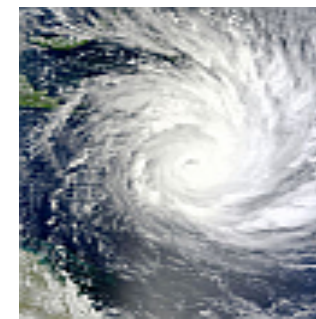


Floods



Courtesy Beth Ebert

Tourism



Energy



Sports



Agriculture



Roads



**Emergency
management**

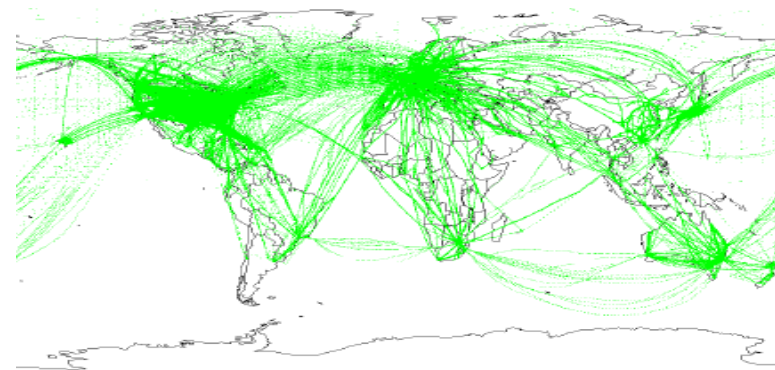
User-relevant verification - Aviation

$$\text{Flight time error (FTE)} = \text{flight_time}_{\text{obs}} - \text{flight_time}_{\text{fcst}}$$

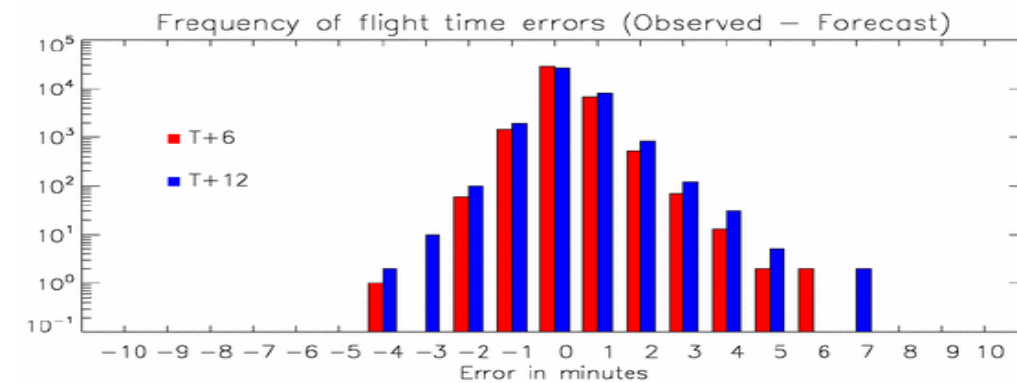
Accurate measure of wind forecast accuracy *directly relevant to airlines*

Calculated using the track that the aircraft actually took

Uses AMDAR observations from real flights rather than model analyses or radiosondes



AMDAR – 1-7 Feb 2010



Courtesy Phil Gill

FMI

dev.hirlam.fmi.fi/Tienpinta_OnlineVerif_Opt/

VARASTOT, SADE, KITKA

Tienpintamalli OnlineVerifointi

EnnVesi EnnLumi EnnKuura EnnJää HavVesi HavLumi HavJää [mm]

EnnSade (int),vrk1:Tutka,vrk2-3:MetEd HavSade (kum)

EnnKitka HavKitka (min=0.10, max=0.82)

LÄMPÖTILA, KELITULKINTA

Tienpintamalli OnlineVerifointi

EnnTtie EnnT2m EnnTD2m HavTtie1 HavTtie2 HavTilma HavTDilma

EnnT2m EnnTD2m: Meteorologin editori

EnnKeli EnnKeli2 HavKeli1 HavKeli2

Enn: 1=kuiva 2=koostea 3=märkä 4=märkä lumi 5=kuura 6=osittain jäinen 7=jäinen 8=luminen
Hav: 1=kuiva 2=koostea 3=märkä 4=märkä/suolattu 5=kuura 6=lumi 7=jää 8=tn.koostea/suolainen

A I K A (UTC)

Tienpintamallin ajoaika

Ens. vrk (21h) ajettu havaintodatalla

TÄNÄÄN

00 12
03 15
06 18
09 21

-1 PÄIVÄ

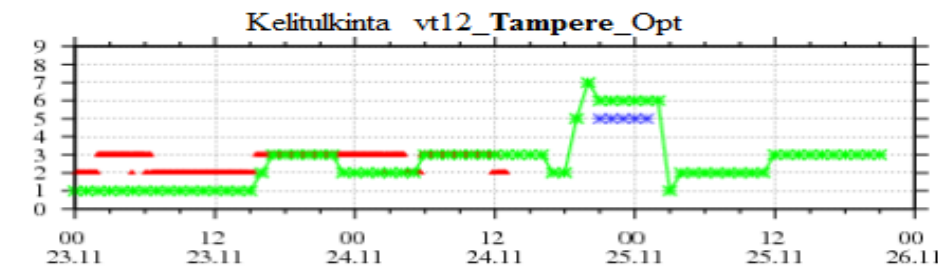
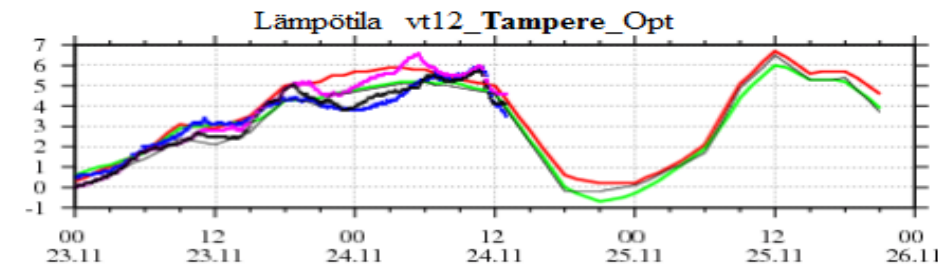
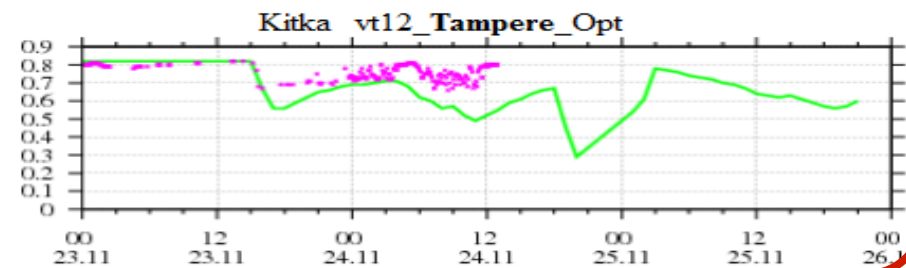
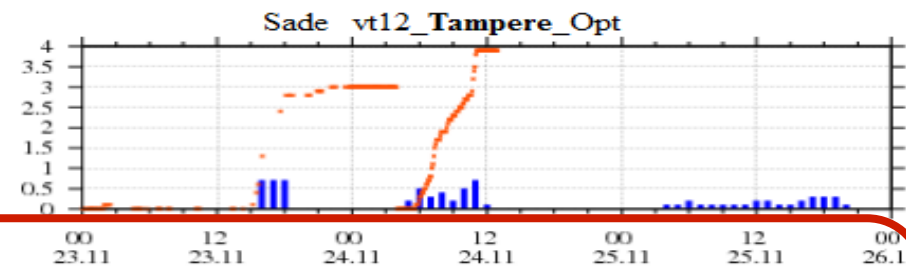
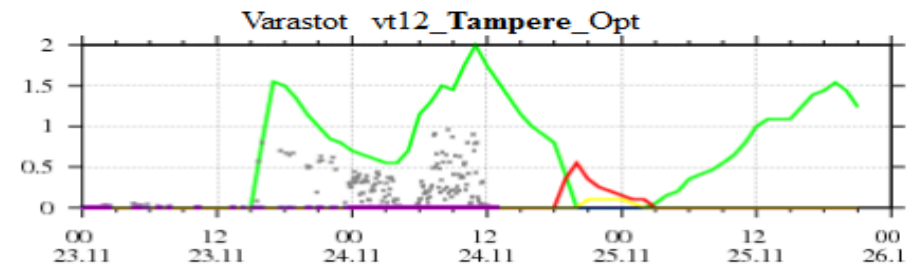
00 12
06 18

-2 PÄIVÄÄ

00 12
06 18

(Shift+) Reload/Refresh-näppäimellä saat uusimman havainnon

LINKKEJÄ



Road surface friction = slipperiness

Forecast (green) vs. Observed (magenta)

Courtesy FMI

Uncertainty in observations

As models improve, can no longer ignore observation error!

Remove observation *bias* errors where possible

Effects of *random* obs error on verification

“Noise” leads to poorer scores for deterministic forecasts

Ensemble forecasts have poorer reliability & ROC

What can we do?

Error bars in scatter plots

Quantitative reference to “gold standard”

Correct for systematic error in observations

RMSE – Ciach & Krajewski (*Adv. Water Res.*, 1999)

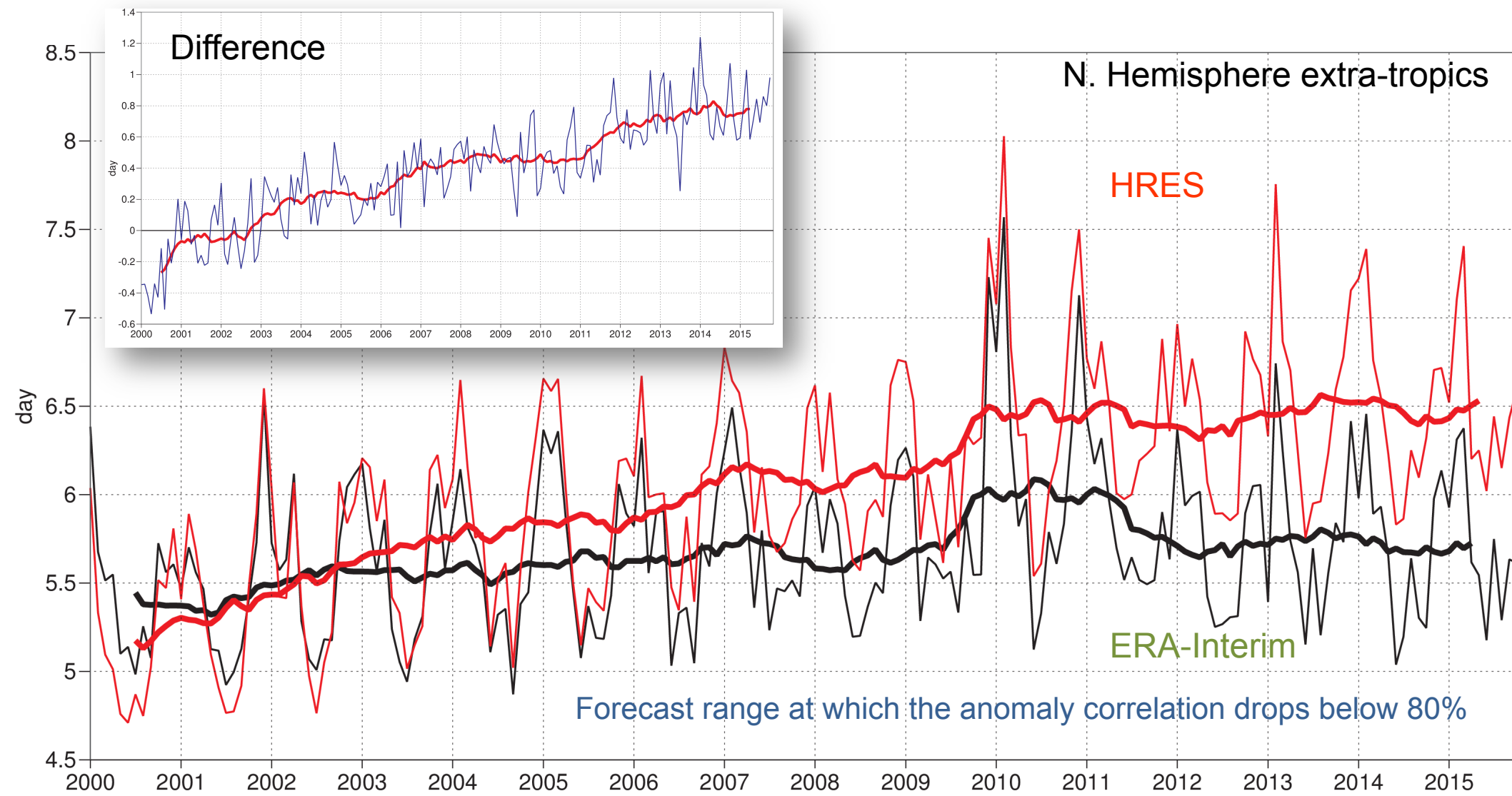
Categorical scores – Briggs et al. (*MWR*, 2005), Bowler (*MWR*, 2006)

Multiple observation sources



Courtesy Beth Ebert

Model performance: HRES relative to ERA-I



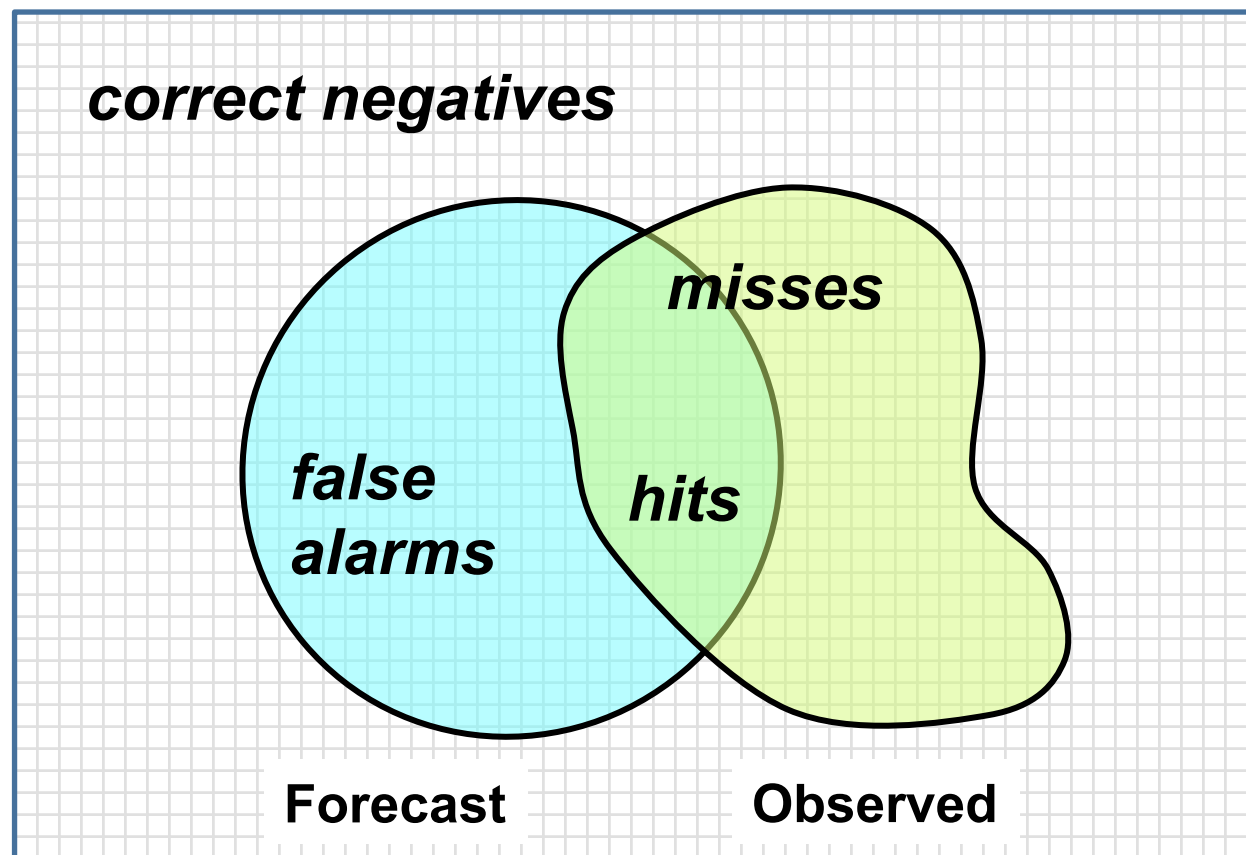


EXTRA

 **ECMWF**

Measure	Attribute evaluated	Comments
Probability forecasts		
Brier score	Accuracy	Based on squared error
Resolution	Resolution (resolving different categories)	Compares forecast category climatologies to overall climatology
Reliability	Calibration	
Skill score	Skill	Skill involves <i>comparison</i> of forecasts
Sharpness measure	Sharpness	Only considers distribution of forecasts
ROC	Discrimination	Ignores calibration
C/L Value	Value	Ignores calibration
Ensemble distribution		
Rank histogram	Calibration	Can be misleading
Spread-skill	Calibration	Difficult to achieve
CRPS	Accuracy	Squared difference between forecast and observed distributions Analogous to MAE in limit
IGN score	Accuracy	Local score, rewards for correct category; infinite if observed category has 0 density

Traditional spatial verification



$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (F_i - O_i)^2}$$

		Observed	
		yes	no
Predicted	yes	<i>hits</i>	<i>false alarms</i>
	no	<i>misses</i>	<i>correct negatives</i>

$$POD = \frac{hits}{hits + misses}$$

$$FAR = \frac{false\ alarms}{hits + false\ alarms}$$

$$TS = \frac{hits}{hits + misses + false\ alarms}$$

Verifying rare extreme values

Categorical scores

Metrics should reward hits, penalise misses and false alarms

For rare events, traditional categorical scores like TS $\rightarrow 0$

New extremal dependence scores:

Event forecast	Event observed		Marginal total
	Yes	No	
Yes	a	b	a + b
No	c	d	c + d
Marginal total	a + c	b + d	a + b + c + d = n

$H = a / (a+c)$, hit rate

$F = b / (b+d)$, false alarm rate

$p = (a+c) / n$, base rate

$q = (a+b) / n$, relative frequency of forecasted events

$$EDS = \frac{\log p - \log H}{\log p + \log H}$$

$$SEDS = \frac{\log q - \log H}{\log p + \log H}$$

$$EDI = \frac{\log F - \log H}{\log F + \log H}$$

$$SEDI = \frac{\log F - \log H - \log(1 - F) + \log(1 - H)}{\log F + \log H + \log(1 - F) + \log(1 - H)}$$

Ferro & Stephenson, *Weather & Forecasting*, 2011

Scores

Root Mean Square Error:

$$E = \sqrt{(fc - an)^2}$$

Measures accuracy
Range: 0 to infinity perfect score = 0

Bias:

$$BIAS = \overline{FC - OBS}$$

Measures bias
Range: -infinity to +infinity
perfect score = 0

Mean Absolute Error :

$$MAE = \overline{|FC - OBS|}$$

Measures accuracy
Range: 0 to infinity perfect score = 0

Anomaly Correlation:

$$ACC = \frac{(fc - c)(an - c)}{\sqrt{A_{fc} A_{an}}}$$

$$A_{fc} = \overline{(fc - c)^2}$$

$$A_{an} = \overline{(an - c)^2}$$

Measures accuracy
Range: -100% to 100%
perfect score = 100%