



# Basic Verification Concepts

Barbara Brown  
National Center for Atmospheric Research  
Boulder Colorado USA

[bgb@ucar.edu](mailto:bgb@ucar.edu)

# Basic concepts - outline

---

- What is verification?
- Why verify?
- Identifying verification goals
- Forecast “goodness”
- Designing a verification study
- Types of forecasts and observations
- Matching forecasts and observations
- Statistical basis for verification
- Comparison and inference
- Verification attributes
- Miscellaneous issues
- **Questions to ponder: Who? What? When? Where? Which? Why?**

# What is verification?

---

## Verify: ver·i·fy

Pronunciation: 'ver-&-"fɪ

1 : to confirm or substantiate in law by oath

2 : to establish the **truth**, **accuracy**, or **reality** of <verify the claim>

**synonym** see **CONFIRM**

- Verification is the process of comparing forecasts to relevant observations
  - Verification is one aspect of measuring forecast **goodness**
- Verification measures the **quality** of forecasts (as opposed to their **value**)
- For many purposes a more appropriate term is “**evaluation**”

# Why verify?

---

- Purposes of verification (traditional definition)
  - Administrative
  - Scientific
  - Economic

# Why verify?

---

- Administrative purpose
  - Monitoring performance
  - Choice of model or model configuration (has the model improved?)
- Scientific purpose
  - Identifying and correcting model flaws
  - Forecast improvement
- Economic purpose
  - Improved decision making
  - “Feeding” decision models or decision support systems

# Why verify?

---

- What are some other reasons to verify hydrometeorological forecasts?

# Why verify?

---

- What are some other reasons to verify hydrometeorological forecasts?
  - Help operational forecasters understand model biases and select models for use in different conditions
  - Help “users” interpret forecasts (e.g., “What does a temperature forecast of 0 degrees really mean?”)
  - Identify forecast weaknesses, strengths, differences

# Identifying verification goals

---

- What *questions* do we want to answer?
  - Examples:
    - In what locations does the model have the best performance?
    - Are there regimes in which the forecasts are better or worse?
    - Is the probability forecast well calibrated (i.e., reliable)?
    - Do the forecasts correctly capture the natural variability of the weather?

*Other examples?*



# Identifying verification goals (cont.)

---

- What forecast performance attribute should be measured?
  - Related to the *question* as well as the type of forecast and observation
- Choices of verification statistics/measures/graphics
  - Should match the type of forecast and the attribute of interest
  - Should measure the quantity of interest (i.e., the quantity represented in the question)

# Forecast “goodness”

---

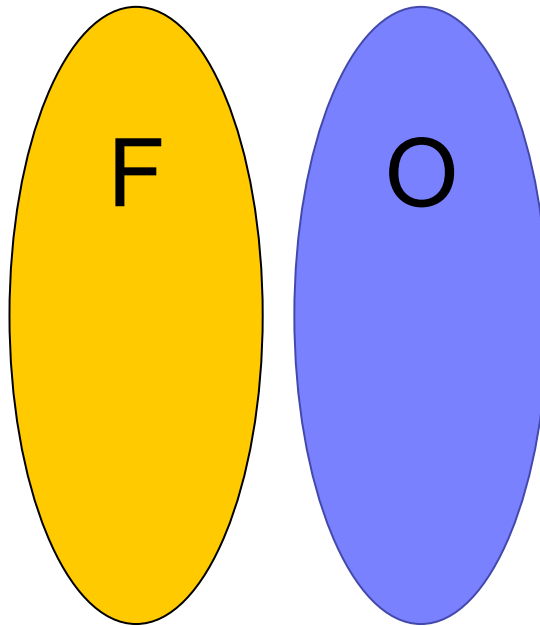
- Depends on the quality of the forecast

**AND**

- The user and his/her application of the forecast information

# Good forecast or bad forecast?

---

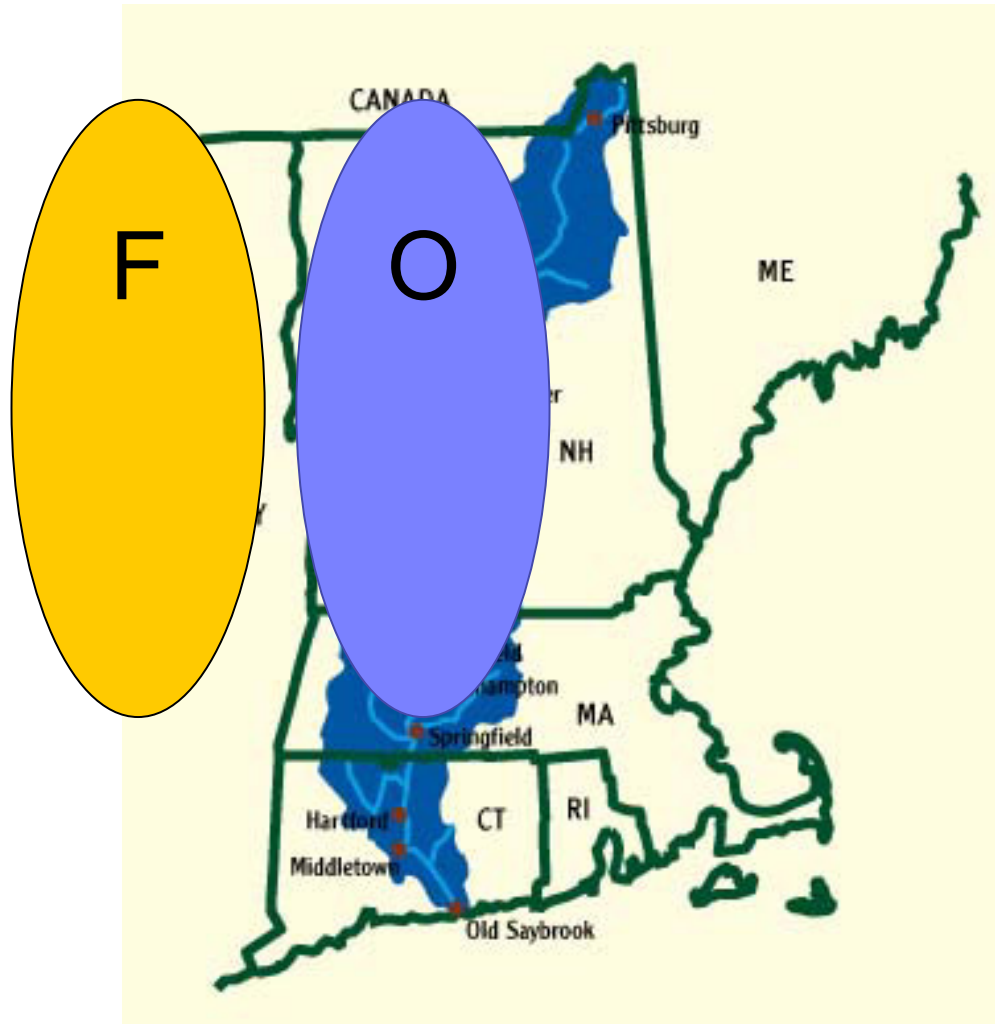


Many verification approaches would say that this forecast has NO skill and is very inaccurate.

# Good forecast or Bad forecast?

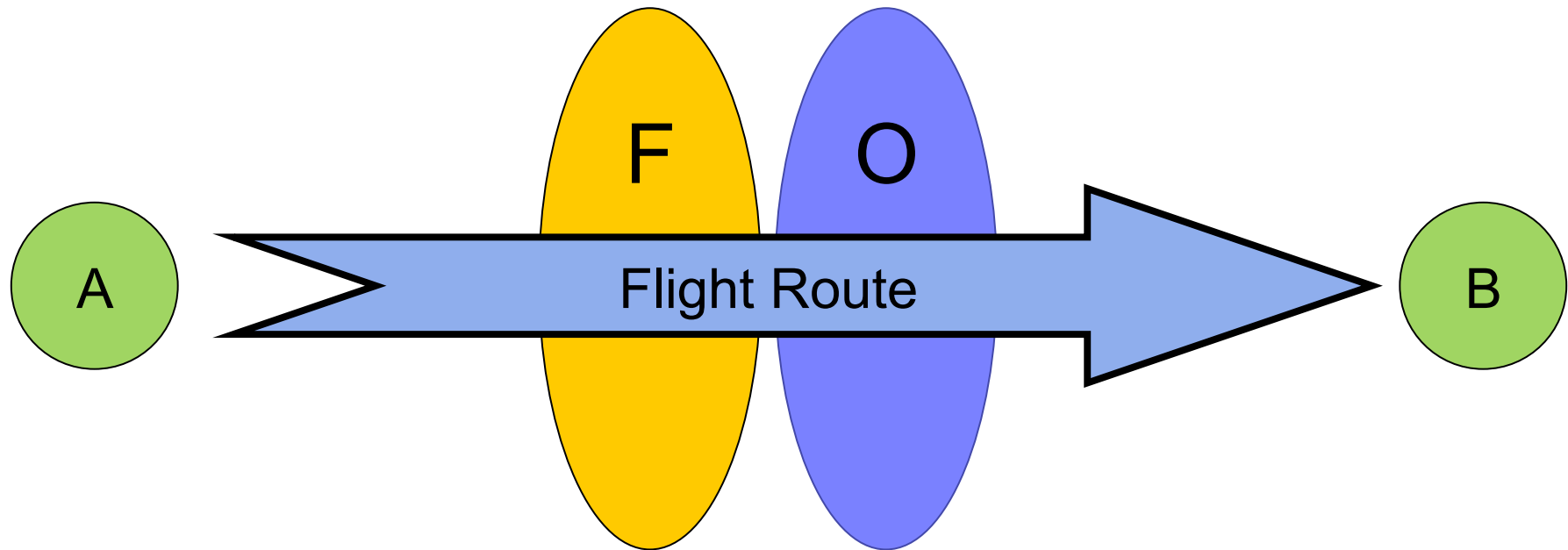
---

If I'm a water manager for this watershed, it's a pretty bad forecast...



# Good forecast or Bad forecast?

---



If I'm an aviation traffic strategic planner...

It might be a pretty good forecast

Different users have different ideas about what makes a forecast good

Different verification approaches can measure different types of "goodness"

# Forecast “goodness”

---

- Forecast quality is only one aspect of forecast “goodness”
- Forecast value is related to forecast quality through complex, non-linear relationships
  - In some cases, *improvements in forecast quality (according to certain measures) may result in a degradation in forecast value for some users!*
- **However** - Some approaches to measuring forecast quality can help understand goodness
  - Examples
    - Diagnostic verification approaches
    - New features-based approaches
    - Use of multiple measures to represent more than one attribute of forecast performance
    - Examination of multiple thresholds

# Basic guide for developing verification studies

---

## **Consider the users...**

- ... of the forecasts
- ... of the verification information
- What aspects of forecast quality are of interest for the user?
  - Typically (always?) need to consider multiple aspects

## **Develop verification questions** to evaluate those aspects/attributes

- Exercise: What verification questions and attributes would be of interest to ...
  - ... operators of an electric utility?
  - ... a city emergency manager?
  - ... a mesoscale model developer?
  - ... aviation planners?

# Basic guide for developing verification studies

---

**Identify observations** that represent the event being forecast, including the

- Element (e.g., temperature, precipitation)
- Temporal resolution
- Spatial resolution and representation
- Thresholds, categories, etc.

**Identify multiple verification attributes** that can provide answers to the questions of interest

**Select measures and graphics** that appropriately measure and represent the attributes of interest

**Identify a standard of comparison** that provides a reference level of skill (e.g., persistence, climatology, old model)



# Types of forecasts, observations

---

- **Continuous**

- Temperature
- Rainfall amount
- 500 mb height

- **Categorical**

- **Dichotomous**

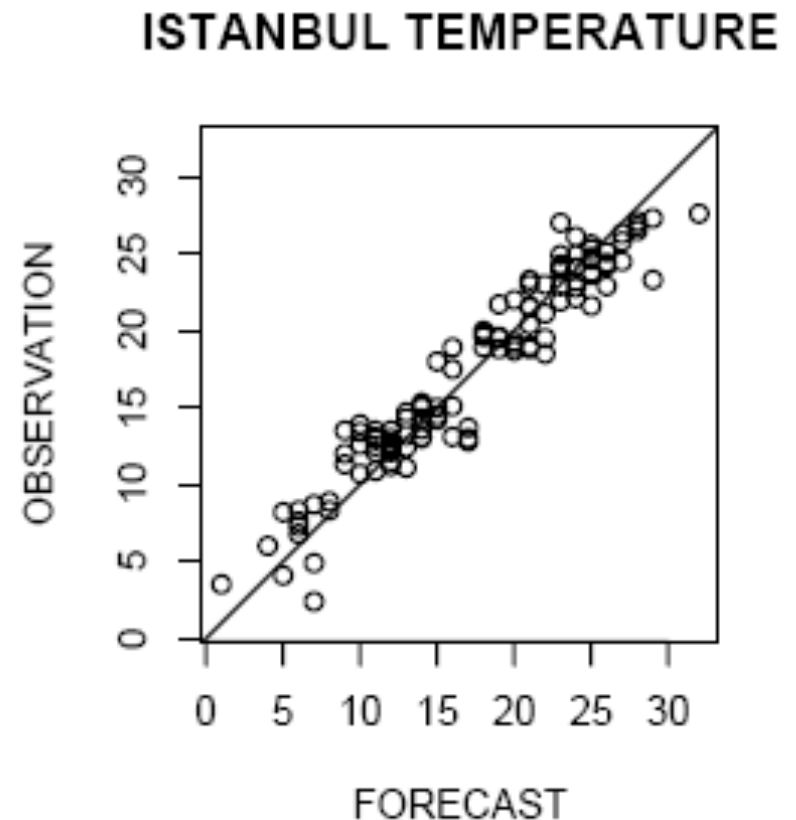
- Rain vs. no rain
- Strong winds vs. no strong wind
- Night frost vs. no frost
- Often formulated as Yes/No

- **Multi-category**

- Cloud amount category
- Precipitation type

- May result from *subsetting* continuous variables into categories

- Ex: *Temperature categories of 0-10, 11-20, 21-30, etc.*

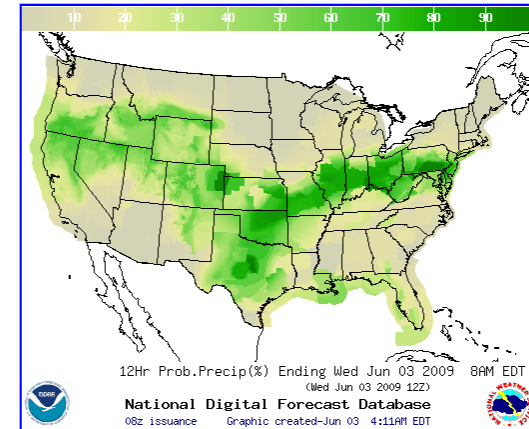


# Types of forecasts, observations

- Probabilistic

- Observation can be **dichotomous**, **category**, or **continuous**
  - Precipitation occurrence – **Dichotomous** (Yes/No)
  - Precipitation type – **Multi-category**
  - Temperature distribution - **Continuous**
- Forecast can be
  - Single probability value (for **dichotomous** events)
  - **Multiple probabilities** (discrete probability distribution for multiple categories)
  - **Continuous** distribution
- For dichotomous or multiple categories, probability values may be limited to certain values (e.g., multiples of 0.1)

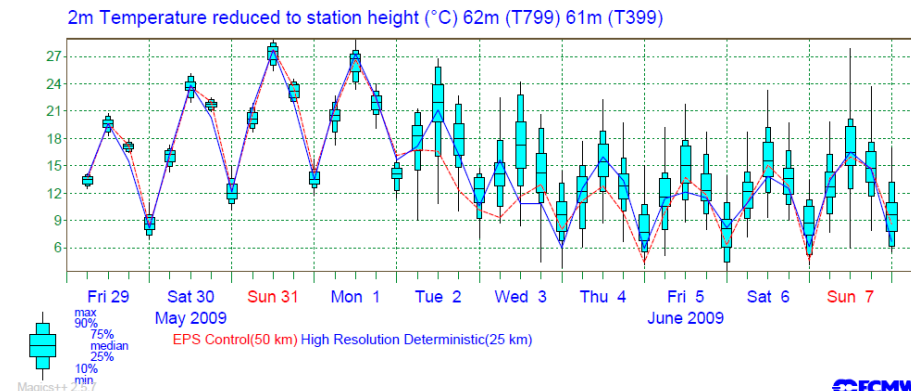
multi-



*2-category precipitation forecast (PoP) for US*

- Ensemble

- Multiple iterations of a **continuous** or **categorical** forecast
  - May be transformed into a probability distribution
- Observations may be **continuous**, **dichotomous** or **multi-category**



*ECMWF 2-m temperature meteogram for Helsinki*

# Matching forecasts and observations

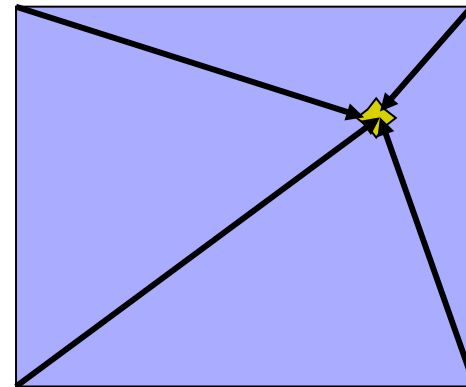
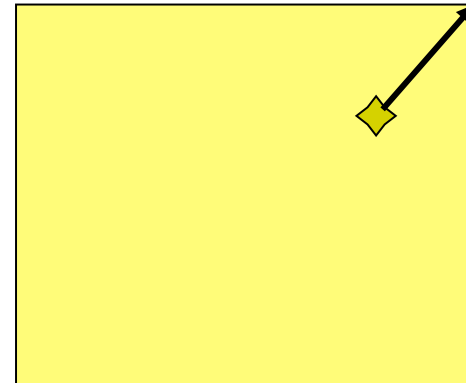
---

- May be the *most difficult* part of the verification process!
- Many factors need to be taken into account
  - Identifying observations that represent the forecast event
    - Example: Precipitation accumulation over an hour at a point
  - For a gridded forecast there are many options for the matching process
    - Point-to-grid
      - Match obs to closest gridpoint
    - Grid-to-point
      - Interpolate?
      - Take largest value?

# Matching forecasts and observations

---

- Point-to-Grid and Grid-to-Point
- Matching approach can impact the results of the verification



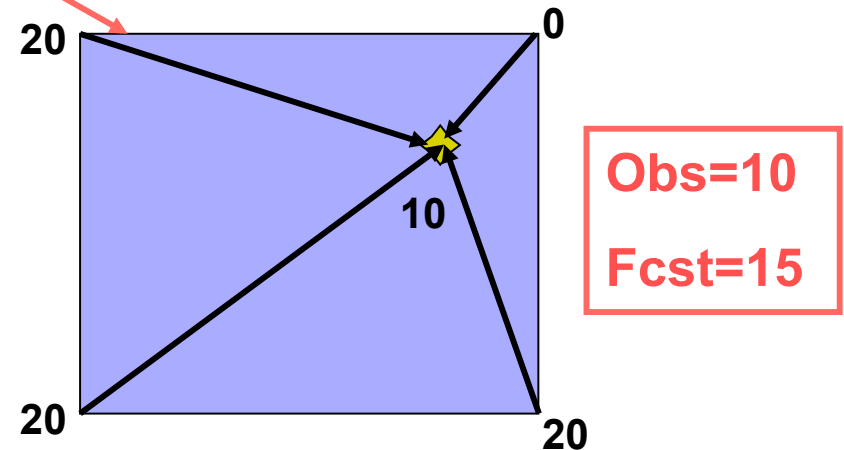
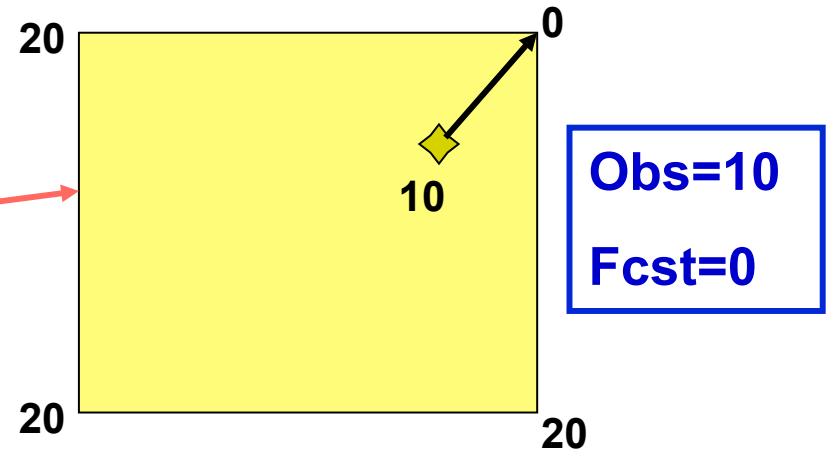
# Matching forecasts and observations

## Example:

- Two approaches:
  - Match rain gauge to nearest gridpoint **or**
  - Interpolate grid values to rain gauge location
    - Crude assumption: equal weight to each gridpoint
- Differences in results associated with matching:

*“Representativeness”  
difference*

*Will impact most  
verification scores*



# Matching forecasts and observations

---

## Final point:

- It is not advisable to use the model analysis as the verification “observation”
- Why not??

# Matching forecasts and observations

---

## Final point:

- It is not advisable to use the model analysis as the verification “observation”
- Why not??
- Issue: Non-independence!!

# Statistical basis for verification

---

- **Joint**, **marginal**, and **conditional** distributions are useful for understanding the statistical basis for forecast verification
  - These distributions can be related to specific summary and performance measures used in verification
  - Specific attributes of interest for verification are measured by these distributions



# Statistical basis for verification

---

Basic (**marginal**) probability

$$p_x = \Pr(X = x)$$

is the probability that a random variable,  $X$ , will take on the value  $x$

Example:

- $X =$  gender of tutorial participant (students + teachers)
- What is an estimate of  $\Pr(X=female)$  ?

# Statistical basis for verification

---

Basic (**marginal**) probability

$$p_x = \Pr(X = x)$$

is the probability that a random variable,  $X$ , will take on the value  $x$

Example:

- $X =$  gender of tutorial participant (students + teachers)
- What is an estimate of  $\Pr(X=female)$  ?

Answer:

# Female participants: 13 (36%)

# Male participants: 23 (64%)

**$\Pr(X=female)$  is  $13/36 = 0.36$**

# Basic probability

---

## Joint probability

$$p_{x,y} = \Pr(X = x, Y = y)$$

= probability that **both** events  $x$  and  $y$  occur

Example: What is the probability that a participant is female and is from the Northern Hemisphere?

# Basic probability

---

## Joint probability

$$p_{x,y} = \Pr(X = x, Y = y)$$

= probability that **both** events  $x$  and  $y$  occur

Example: What is the probability that a participant is female and is from the Northern Hemisphere?

11 participants (of 36) are Female and are from the Northern Hemisphere

$$\Pr(X=Female, Y=Northern Hemisphere) = 11/36 = 0.31$$

# Basic probability

---

## Conditional probability

$$p_{x,y} = \Pr(X = x \mid Y = y)$$

= probability that event  $x$  is true (or occurs) given that event  $y$  is true (or occurs)

Example: If a participant is from the Northern Hemisphere, what is the likelihood that he/she is female?

# Basic probability

---

## Conditional probability

$$p_{x,y} = \Pr(X = x | Y = y)$$

= probability that event  $x$  is true (or occurs) given that event  $y$  is true (or occurs)

Example: If a participant is from the Northern Hemisphere, what is the likelihood that he/she is female?

Answer: 26 participants are from the Northern Hemisphere. Of these, 11 are female.

$$\Pr(X=Female | Y=Northern Hemisphere) = 11/26 = 0.42$$

[Note: This prob is somewhat larger than  $\Pr(X=Female) = 0.36$ ]

# What does this have to do with verification?

---

Verification can be represented as the process of evaluating the **joint** distribution of forecasts and observations,  $p(f, x)$

- All of the information regarding the forecast, observations, and their relationship is represented by this distribution
- Furthermore, the joint distribution can be factored into two pairs of **conditional** and **marginal** distributions:

$$p(f, x) = p(F = f | X = x)p(X = x)$$

$$p(f, x) = p(X = x | F = f)p(F = f)$$

# Decompositions of the joint distribution

---

- Many forecast verification attributes can be derived from the conditional and marginal distributions
- Likelihood-base rate decomposition

$$p(f, x) = \underbrace{p(F = f | X = x)}_{\text{Likelihood}} \underbrace{p(X = x)}_{\text{Base rate}}$$

- Calibration-refinement decomposition

$$p(f, x) = \underbrace{p(X = x | F = f)}_{\text{Calibration}} \underbrace{p(F = f)}_{\text{Refinement}}$$



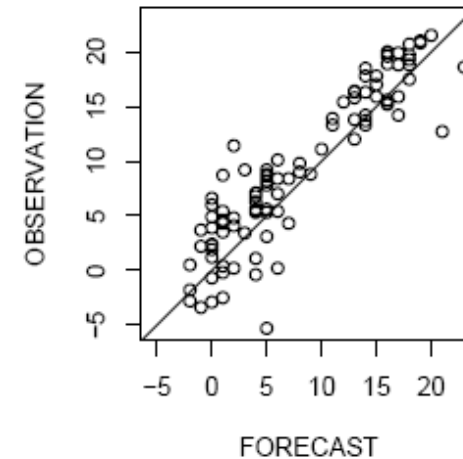
# Graphical representation of distributions

---

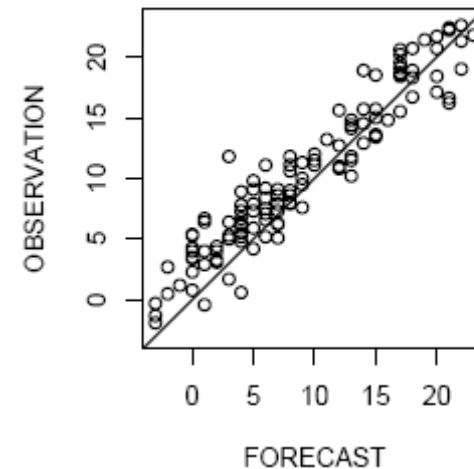
## Joint distributions

- Scatter plots
- Density plots
- 3-D histograms
- Contour plots

OSLO TEMPERATURE



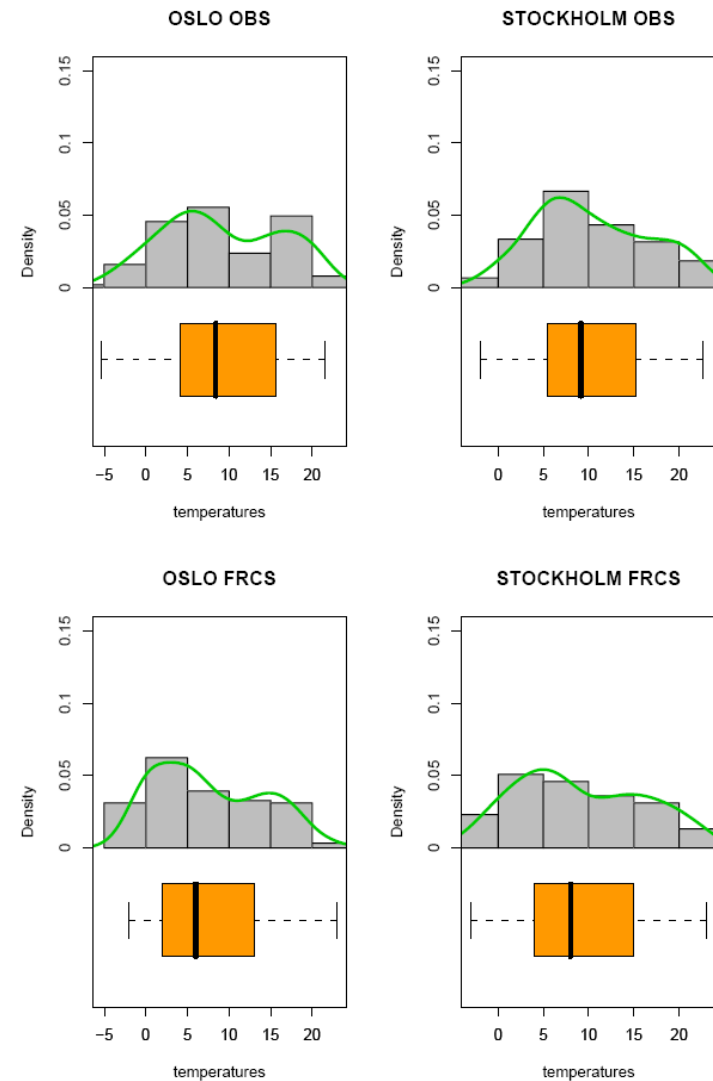
STOCKHOLM TEMPERATURE



# Graphical representation of distributions

## Marginal distributions

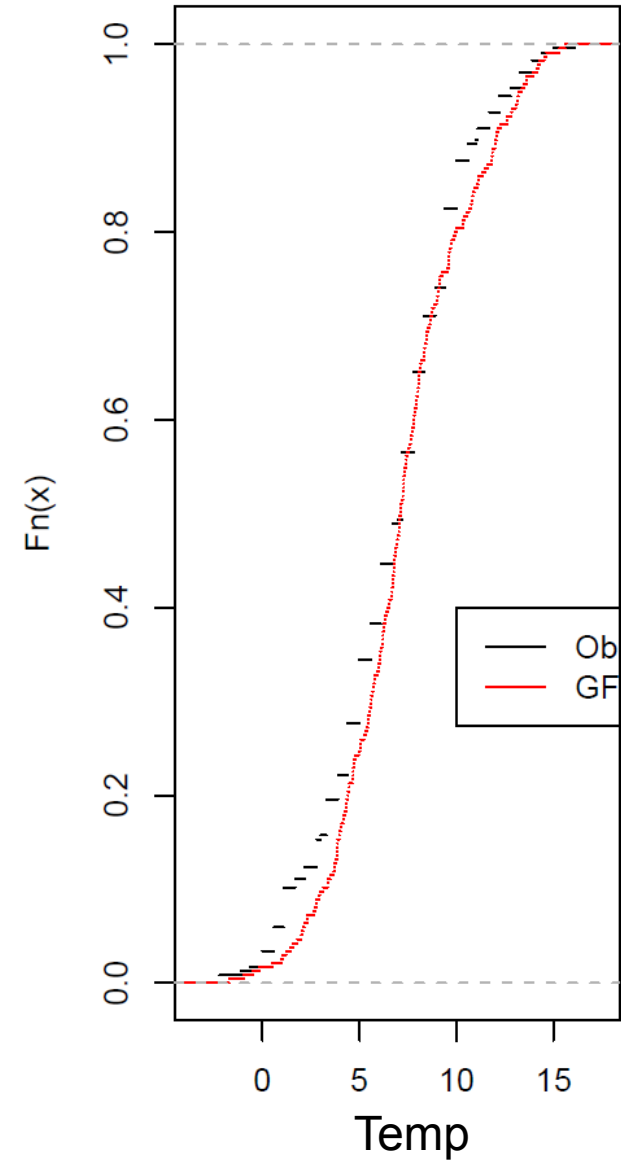
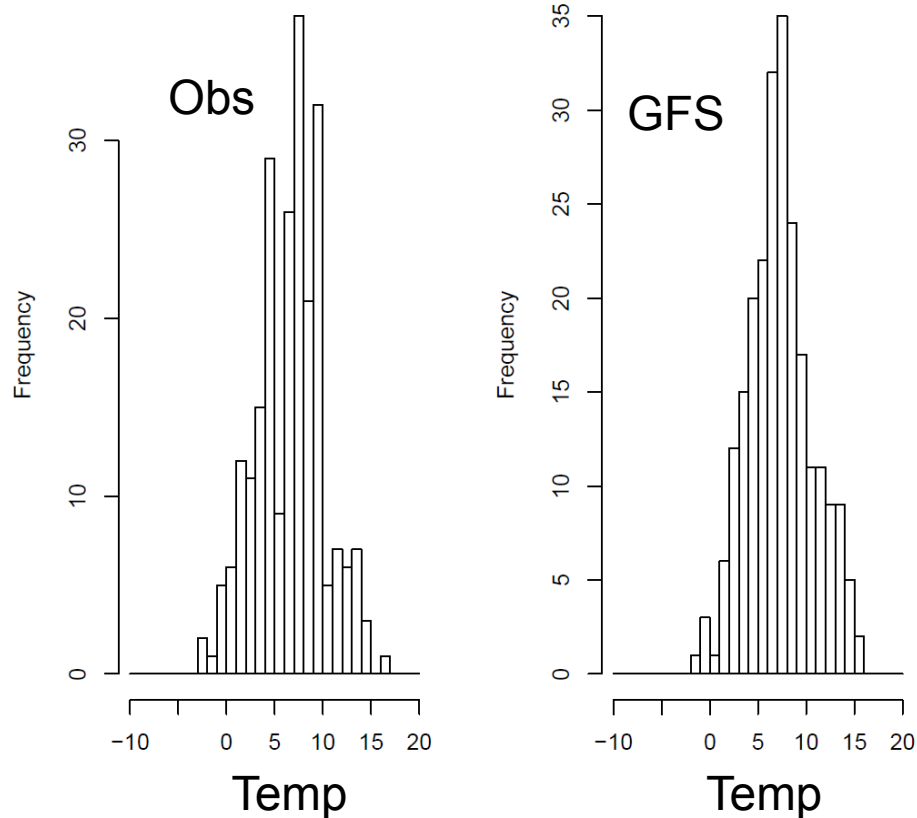
- Stem and leaf plots
- Histograms
- Box plots
- Cumulative distributions
- Quantile-Quantile plots



# Graphical representation of distributions

## Marginal distributions

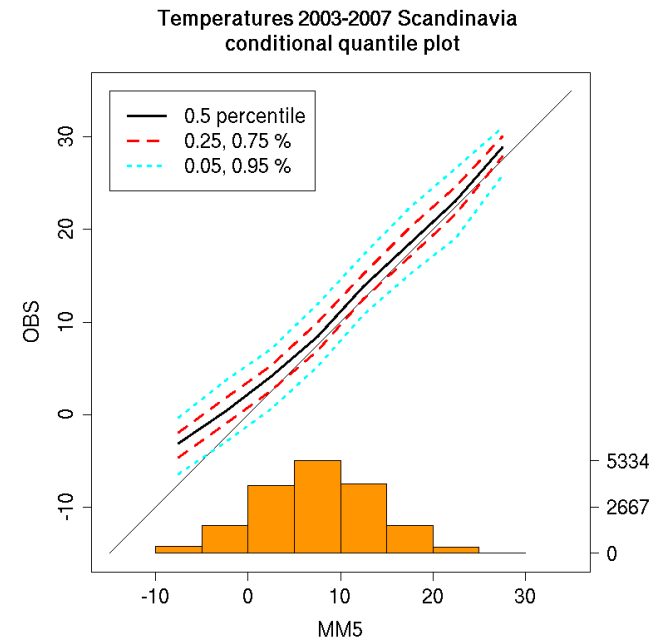
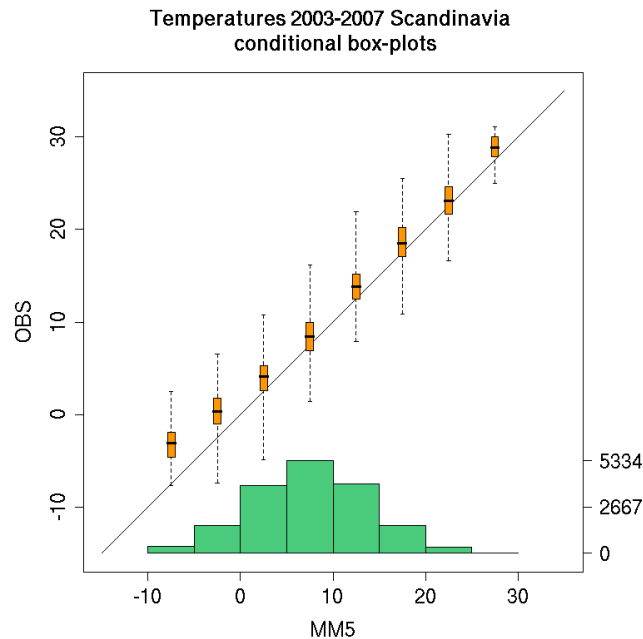
- Density functions
- Cumulative distributions



# Graphical representation of distributions

## Conditional distributions

- Conditional quantile plots
- Conditional boxplots
- Stem and leaf plots



# Stem and leaf plots: Marginal and conditional distributions

---

**Marginal distribution of Tampere probability forecasts**

	Forecast probability			
0.0				
0.1	X	X	X	
0.2	X	X	X	X
0.3	X			
0.4	X			
0.5				
0.6				
0.7	X	X	X	
0.8				
0.9				
1.0	X			

**Conditional distributions of Tampere probability forecasts**

Obs precip = No					Obs precip = Yes			
				0.0				
X	X	X		0.1				
X	X	X		0.2	X			
		X		0.3				
				0.4	X			
				0.5				
				0.6				
				0.7	X	X	X	
				0.8				
				0.9				
				1.0	X			

# Comparison and inference

---

## Skill scores

- A skill score is a measure of *relative performance*
  - **Ex:** *How much more accurate are my temperature predictions than climatology? How much more accurate are they than the model's temperature predictions?*
  - *Provides a comparison to a **standard***
- Generic skill score definition: 
$$\frac{M - M_{ref}}{M_{perf} - M_{ref}}$$

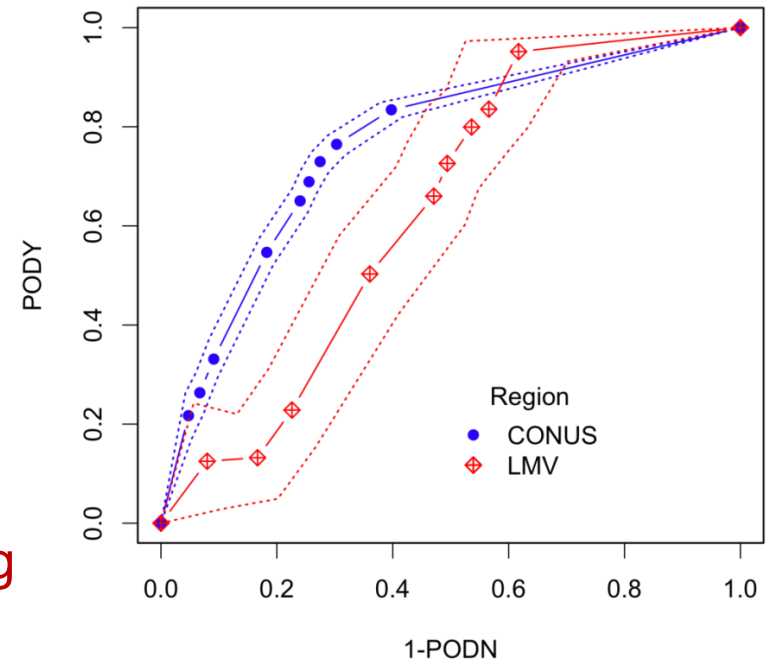
Where  $M$  is the verification measure for the forecasts,  $M_{ref}$  is the measure for the reference forecasts, and  $M_{perf}$  is the measure for perfect forecasts

- Positively oriented (larger is better)
- Choice of the standard matters (*a lot!*)

# Comparison and inference

Uncertainty in scores and measures should be estimated whenever possible!

- Uncertainty arises from
  - Sampling variability
  - Observation error
  - Representativeness differences
  - Others?
- **Erroneous conclusions can be drawn regarding improvements in forecasting systems and models**
- Methods for *confidence intervals* and *hypothesis tests*
  - Parametric (i.e., depending on a statistical model)
  - Non-parametric (e.g., derived from re-sampling procedures, often called “bootstrapping”)



More on this topic to be presented by Ian Jolliffe

# Verification attributes

---

- Verification attributes measure different aspects of forecast quality
  - Represent a range of characteristics that should be considered
  - Many can be related to joint, conditional, and marginal distributions of forecasts and observations



# Verification attribute examples

---

- Bias
  - (Marginal distributions)
- Correlation
  - Overall association (Joint distribution)
- Accuracy
  - Differences (Joint distribution)
- Calibration
  - Measures conditional bias (Conditional distributions)
- Discrimination
  - Degree to which forecasts discriminate between different observations (Conditional distribution)

# Desirable characteristics of verification measures

---

- Statistical validity
- Properness (probability forecasts)
  - “Best” score is achieved when forecast is consistent with forecaster’s best judgments
  - “Hedging” is penalized
  - Example: Brier score
- Equitability
  - Constant and random forecasts should receive the same score
  - Example: Gilbert skill score (2x2 case); Gerrity score
  - No scores achieve this in a more rigorous sense
    - Ex: Most scores are sensitive to bias, event frequency

# Miscellaneous issues

---

- In order to be *verified*, forecasts must be formulated so that they are *verifiable*!
  - Corollary: All forecast should be verified – if something is worth forecasting, it is worth verifying
- Stratification and aggregation
  - Aggregation can help increase sample sizes and statistical robustness but can also hide important aspects of performance
    - Most common regime may dominate results, mask variations in performance
  - Thus it is very important to *stratify results into meaningful, homogeneous sub-groups*

# Verification issues cont.

---

- Observations
  - No such thing as “truth”!!
  - Observations generally are more “true” than a model analysis (at least they are relatively more independent)
  - Observational uncertainty should be taken into account in whatever way possible
    - e.g., how well do adjacent observations match each other?

# Some key things to think about ...

---

## Who...

- ...wants to know?

## What...

- ... does the user care about?
- ... kind of parameter are we evaluating? What are its characteristics (e.g., continuous, probabilistic)?
- ... thresholds are important (if any)?
- ... forecast resolution is relevant (e.g., site-specific, area-average)?
- ... are the characteristics of the obs (e.g., quality, uncertainty)?
- ... are appropriate methods?

## Why...

- ...do we need to verify it?

# Some key things to think about...

---

## How...

- ...do you need/want to present results (e.g., stratification/aggregation)?

## Which...

- ...methods and metrics are appropriate?
- ... methods are required (e.g., bias, event frequency, sample size)

# Suggested exercise

---

**This exercise will show you some different ways of looking at distributions of data**

- Open `brown.R.txt` using WordPad
- In R, open the “File” menu
  - Select “Change dir”
  - Select the “Brown” directory
- In R, open the “File” menu
  - Select “Open script”
  - Under “Files of type” select “All files”
  - Select the text file “`brown.R`”
- Highlight each section of “`brown.R`” individually and copy into the “R console” window using Ctrl-R