# Basic Verification Concepts

Barbara Brown
National Center for Atmospheric Research
Boulder Colorado USA

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

# What is verification?

- Verify: **ver·i·fy**
  Pronunciation: 'ver-&-"fI
  **1 :** to confirm or substantiate in law by oath
  **2 :** to establish the truth, accuracy, or reality of
  <*verify* the claim>
  **synonym** see <u>CONFIRM</u>

- 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*)

# Why verify?

- □ Purposes for 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
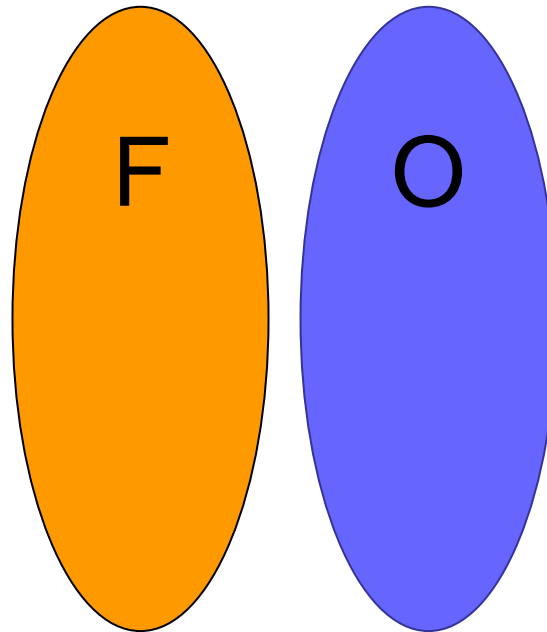    - Should measure the quantity of interest (i.e., the quantity represented in the question)

# Forecast "goodness"

- Depends on the quality of the forecast

## AND

- The user and his/her application of the forecast information
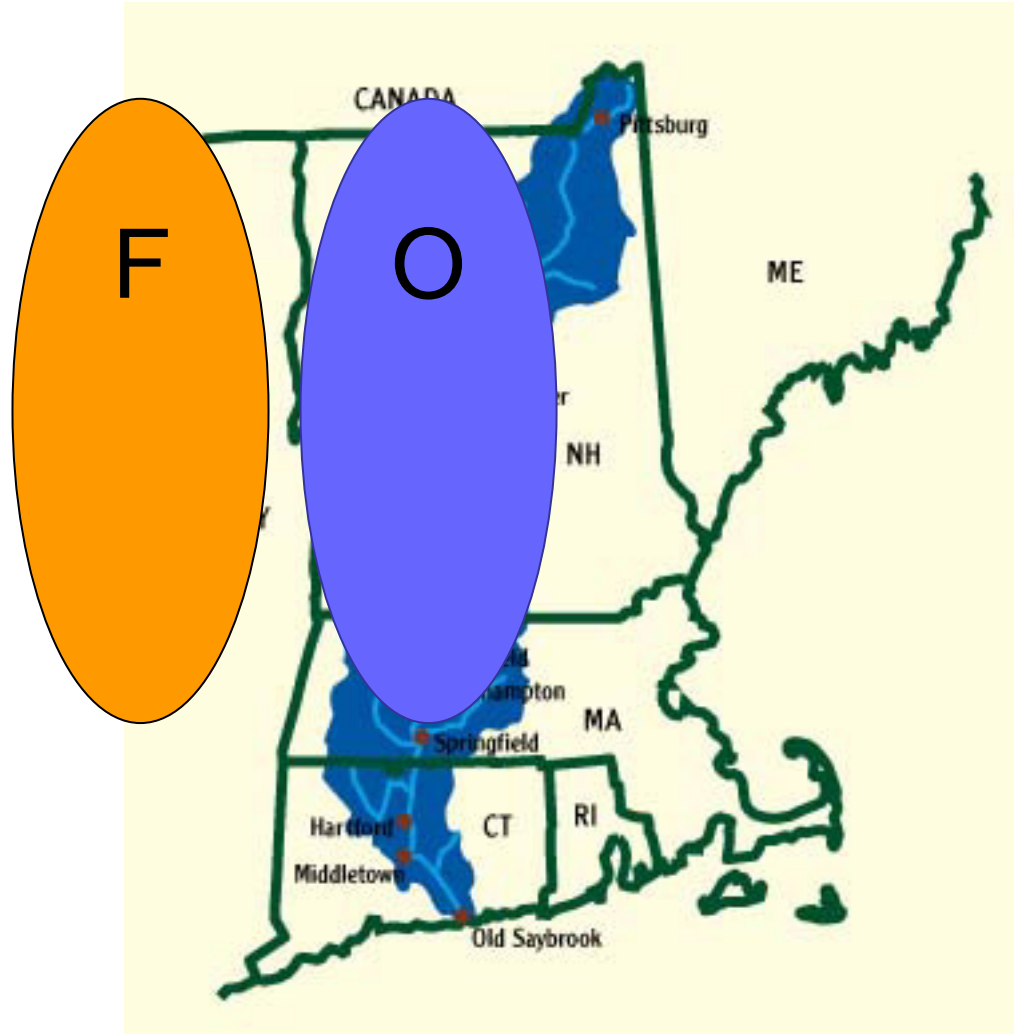
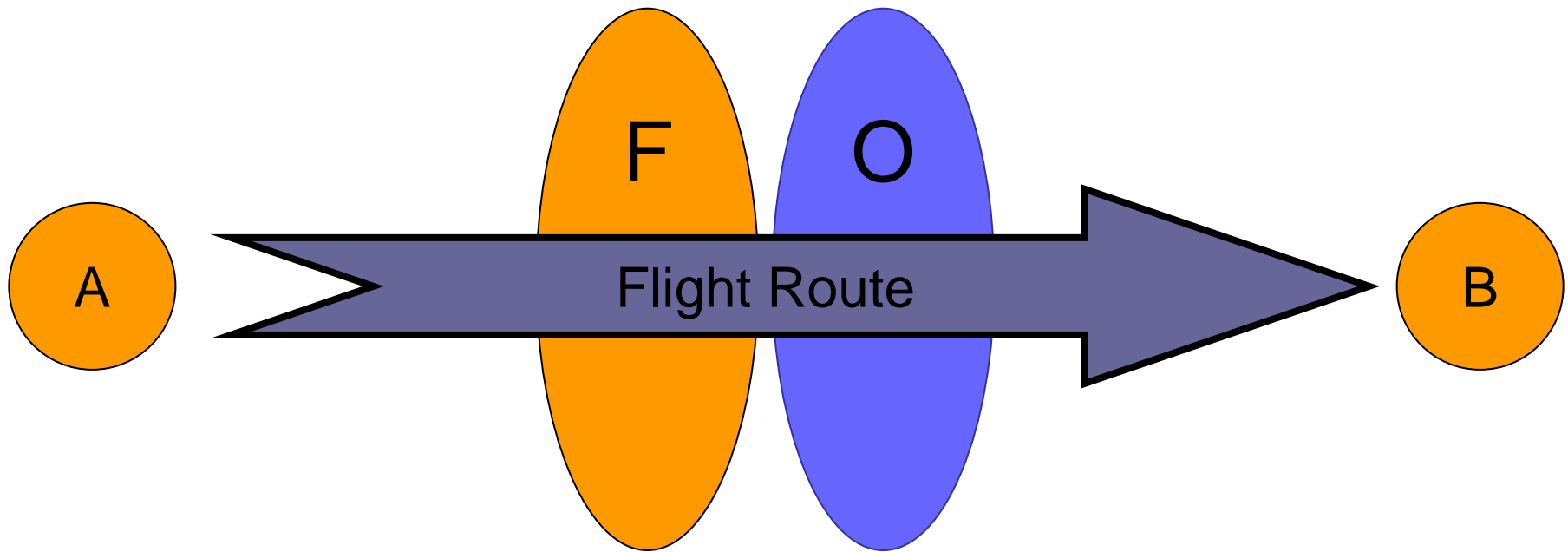# Good forecast or bad forecast?



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?



**F** **O**

A Flight Route B

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

# Simple 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?
- ❑ 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?

# Simple 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)
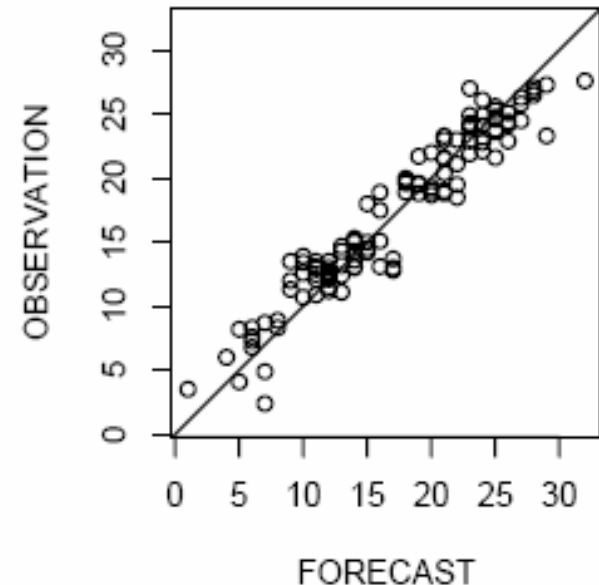
# Types of forecasts, observations

- ❑ **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
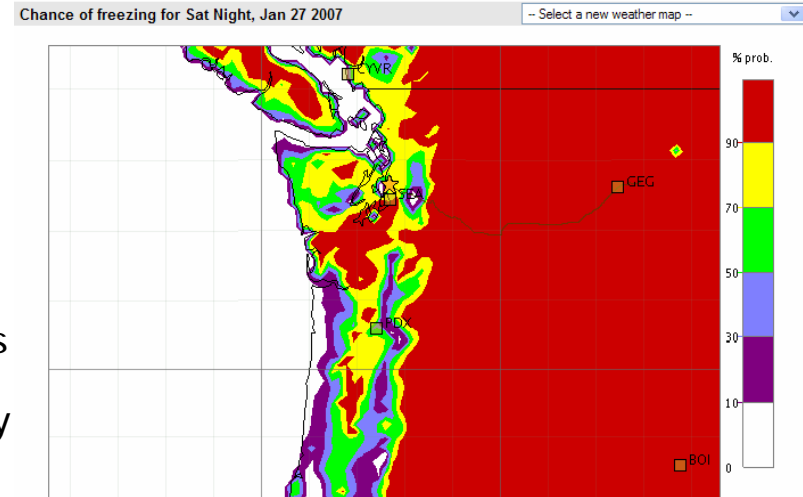
**ISTANBUL TEMPERATURE**

# Types of forecasts, observations

- Probabilistic
  - Observation can be dichotomous, multi-category, or continuous
    - Precipitation occurrence – Yes/No
    - Precipitation type
    - Temperature distribution
  - 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)

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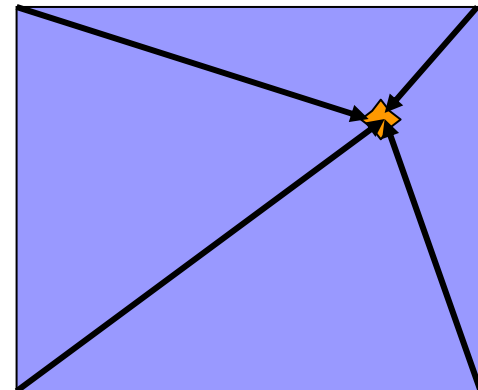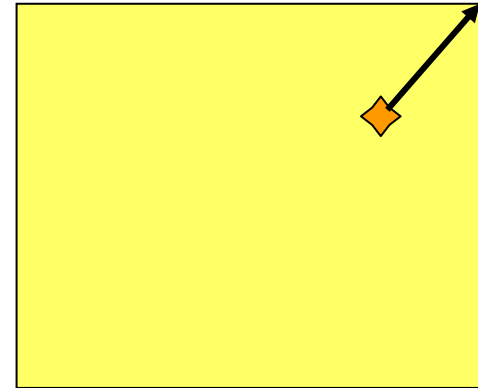  - Observations may be continuous, dichotomous or multi-category



Chance of freezing for Sat Night, Jan 27 2007          -- Select a new weather map --

*Probability of freezing temperatures; from U. Washington*

# 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
    - ❑ <u>Example</u>: 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

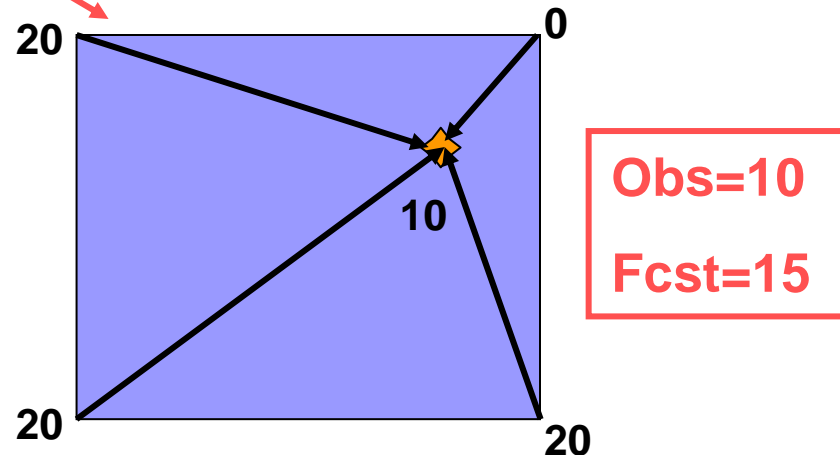- Matching approach can impact the results of the verification

# 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" error*



**Obs=10**

**Fcst=0**



**Obs=10**

**Fcst=15**

# 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 lecturer
  - 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 lecturer*
    - *What is an estimate of $\Pr(X{=}female)$ ?*
    - Answer:
        - Female lecturers (3): B. Brown, B. Casati, E. Ebert
        - Male lecturers (4): S. Mason, P. Nurmi, M. Pocernich, L. Wilson
        - Total number of lecturers: 7

      $\Pr(X{=}female)$ is 3/7 = 0.43

# 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 lecturer is female and has dark hair?

# 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 lecturer is female  and has dark hair?

  □ Answer:

  3 of 7 lecturers are female

  2 of the 3 female lecturers have dark hair

  Thus, the probability that a lecturer is female and has dark hair is

$$\Pr(X = female, Y = dark\ hair) = \frac{2}{7}$$

where $X$ is gender and $Y$ is hair color

# Basic probability

- Conditional probability

$$\Pr(X = x \,|\, Y = y) \quad \text{or} \quad \Pr(Y = y \,|\, X = x)$$

- Probability associated with one attribute – given the value of a second attribute
- Example
  - What is the probability that hair color is dark, given a lecturer is female?

    $$\Pr(Y = dark\ hair \,|\, X = female)$$

  - What is the probability that a lecturer is male, given hair color is dark?

    $$\Pr(X = male \,|\, Y = dark\ hair)$$

    (Note: "Gray" counts as "light" ☺ )

# 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:

$$p(f, x) = p(F = f \mid X = x) p(X = x)$$

$$p(f, x) = p(X = x \mid F = f) p(F = f)$$

# Decompositions of the joint distribution

- □ Many forecast verification attributes can be derived from the conditional and marginal distributions

- □ Likelihood-base rate decomposition

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

**Likelihood**    **Base rate**
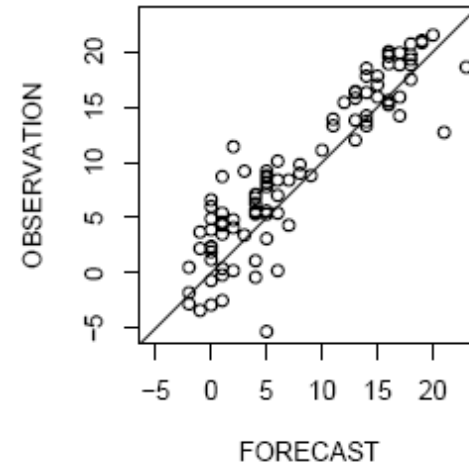
- □ Calibration-refinement decomposition

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

**Calibration**    **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
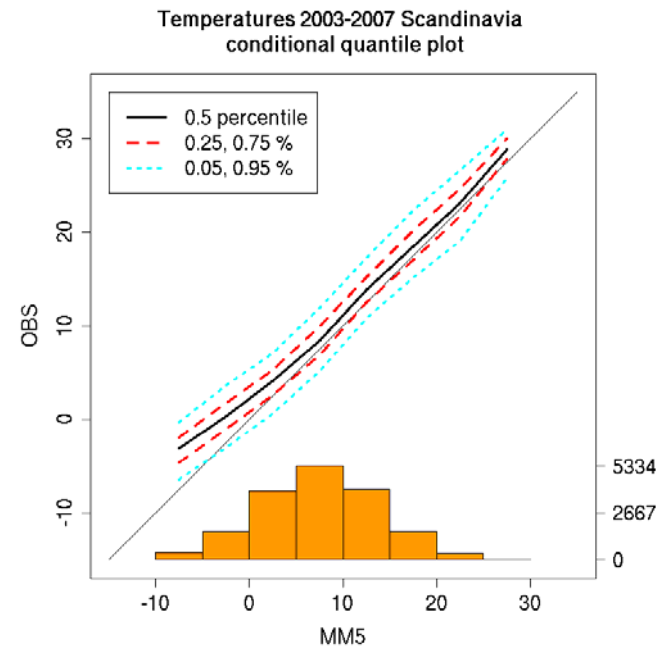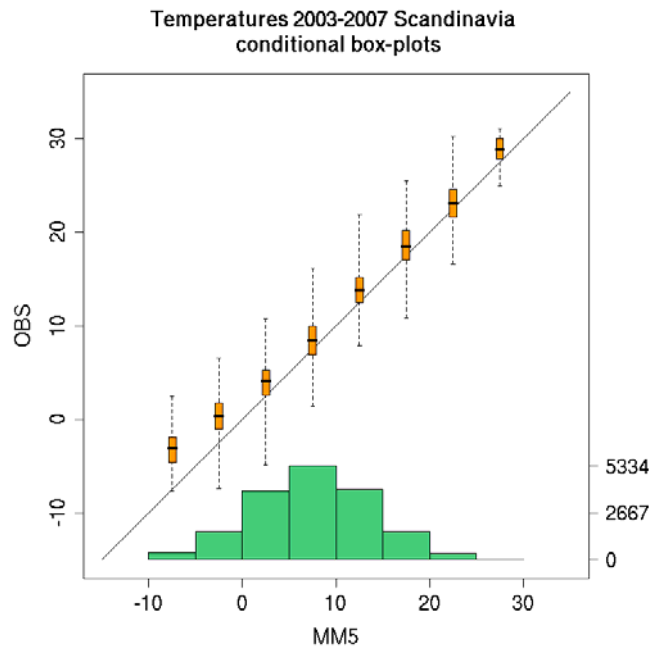
# Graphical representation of distributions

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



Temperatures 2003-2007 Scandinavia conditional box-plots

Temperatures 2003-2007 Scandinavia conditional quantile plot

# Exercise: Stem and leaf plots

**Probability forecasts (Tampere)**

| Date 2003 | Observed rain?? | Forecast (probability) |
|-----------|-----------------|------------------------|
| Jan 1     | No              | 0.3                    |
| Jan 2     | No              | 0.1                    |
| Jan 3     | No              | 0.1                    |
| Jan 4     | No              | 0.2                    |
| Jan 5     | No              | 0.2                    |
| Jan 6     | No              | 0.1                    |
| Jan 7     | Yes             | 0.4                    |
| Jan 8     | Yes             | 0.7                    |
| Jan9      | Yes             | 0.7                    |
| Jan 12    | No              | 0.2                    |
| Jan 13    | Yes             | 0.2                    |
| Jan 14    | Yes             | 1.0                    |
| Jan 15    | Yes             | 0.7                    |

# Stem and leaf plots: Marginal and conditional

**Marginal distribution of Tampere probability forecasts**

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

**Conditional distributions of Tampere probability forecasts**

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

**Instructions**: Mark X's in the appropriate cells, representing the forecast probability values for Tampere.

The resulting plots are one simple way to look at marginal and conditional distributions.

# R exercise: Representation of distributions

# Comparison and inference

- ❑ Skill scores
  - A skill score is a measure of *relative performance*
    - ❑ *Ex: How much more accurate are Laurie Wilson's temperature predictions than climatology (or my temperature predictions)?*
    - ❑ *Provides a comparison to a **standard***
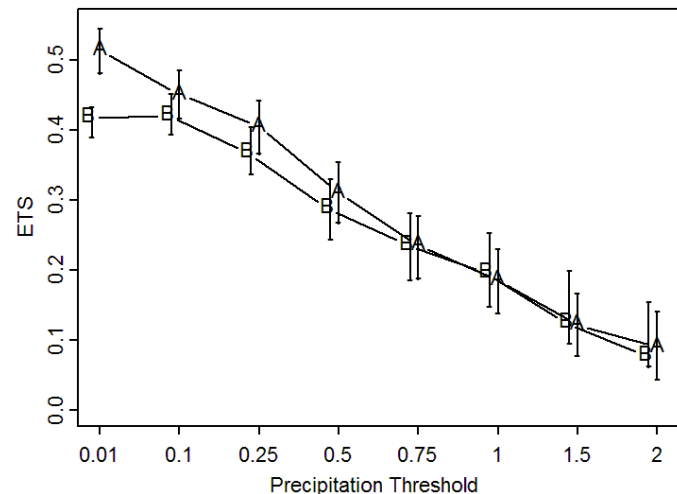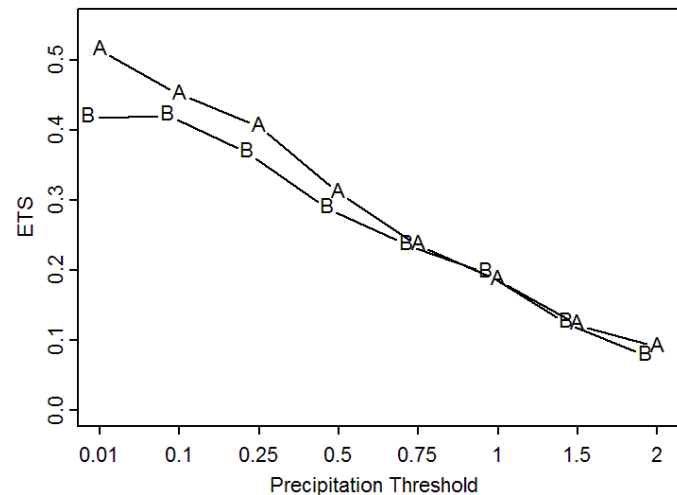  - Generic skill score definition

$$\frac{M - M_{ref}}{M_{perf} - M_{ref}}$$

  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)

# Comparison and inference

- Uncertainty in scores and measures should be estimated whenever possible!
  - Uncertainty arises from
    - Sampling variability
    - Observation error
    - Representativeness error
    - Others?
  - Erroneous conclusions can be drawn regarding improvements in forecasting systems and models
  - 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")

# 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*!
  - <u>Corollary</u>: 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 <u>but</u> 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?

# Stem and leaf plots: Marginal and conditional

**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 | | |