1901	θT	01	-99.9	-3.1	-0.0
1901	01	02	-99.9	-1.3	-3.6
1901	01	03	-99.9	-0.5	-7.9
1901	01	04	-99.9	-1.0	-9.1
1901	01	05	-99.9	-1.8	-8.4
1901	01	06	-99.9	-7.8	-11.5
1901	01	07	-99.9	-6.6	-12.2
1901	01	08	-99.9	-0.6	-9.4
1901	01	09	-99.9	4.2	-2.7
1901	01	10	-99.9	5.9	-1.4
1901	01	11	-99.9	4.9	-7.8
1901	01	12	-99.9	-2.6	-9.0
1901	01	13	-99.9	-1.8	-8.2
1901	01	14	-99.9	3.0	-7.2
1901	01	15	-99.9	5.2	-3.8
1901	01	16	-99.9	4.0	-4.5
1901	01	17	-99.9	6.4	-2.5
1901	01	18	-99.9	4.1	-4.5
1901	01	19	-99.9	6.3	-0.4
1901	01	20	-99.9	7.1	4.3
1901	01	21	-99.9	10.1	6.1
1901	01	22	-99.9	8.3	6.8
1901	01	23	-99.9	8.3	-1.9
1901	01	24	-99.9	5.5	-1.2
1901	01	25	-99.9	7.4	4.8
1901	01	26	-99.9	6.8	3.0
1901	01	27	-99.9	9.1	3.8
1901	01	28	-99.9	5.2	0.4
1901	01	29	-99.9	2.7	-3.7
1901	01	30	-99.9	3.8	-2.0
1901	01	31	-99.9	2.1	-0.2
1901	02	01	-99.9	2.4	-1.8
1901	02	02	-99.9	-0.1	-3.5
1901	02	03	-99.9	2.7	-6.0
1901	02	04	-99.9	-1.6	-8.6
1901	02	05	-99.9	1.8	-3.6
1901	02	06	-99.9	0.3	-4.4
1901	02	07	-99.9	2.8	-10.5
1901	<u>й2</u>	Ω8	_99_9	6.0	_0 8

Data Quality

Thomas C. Peterson President, WMO CCl Principal Scientist, NOAA's National Climatic Data Center Asheville, North Carolina, USA

With material from Malcolm Haylock, Climatic Research Unit, UK Xuebin Zhang, Environment Canada



and Enric Aguilar Universitat Rovira i Virgili Tarragona, Spain



Outline

- Importance of data quality
- Errors
 - Metadata
 - Data
 - Sources of error
 - Outliers
- Identifying Problems
 - Visualisation
 - RClimDex
 - Lucie will address this part
- Correcting Problems

Importance

- Errors associated with measurement can be much larger than the signals that we are trying to measure in climate change detection studies
- Some data problems are unique to a particular country but most are not. Your experiences are very helpful to others!

Metadata errors

- Metadata
 - data about data
 - station history
 - details instrument changes, site moves
- Biggest problem usually lack of accessible metadata
- Correct location (including elevation)
- Units (decimal degrees or deg-min-sec)

Data errors: biases

- Changes in instruments or microscale environment around the station can create biases in the data
- These subtle problems will be addressed in the homogeneity analyses later

Data errors: sources

• Observing

- Was the thermometer or rain gauge read correctly?
- Recording / Digitisation
 - Was the data properly recorded and digitized?
 - "Fat Finger Errors" typing the wrong key
 - Missing minus signs
- Systematic
 - Are the data in the right units
 - Are the data what they are suppose to be
 - E.g., humidity data typed where precipitation should be
 - Are missing set to zero?
 - Are three day accumulations of precip appearing as daily precip?

Data errors: observing

- precision (rounding)
 - Noticeable but probably not important

Australia 003015 Daily Rainfall Pre-1974



Quality Control tests

- Errors (tmax<tmin, out of bounds)
- Data plots (time series, stem plots, histograms)
- Thresholds excedance (4 SD) for temperature
- Comparison with neighbours
- Checks of previous days in precipitation in search of missing values
- Validation with external sources (workshop participants, hurricane papers, other databases)

Missing Weekend Observations



From Viney, N.R. and Bates, B.C., 2004. It never rains on Sunday: the prevalence and implications of untagged multi-day rainfall accumulations in the australian high quality data set. *Int. J. Climatol.* 24: 1171–1192

Data errors: systematic

• units

- change of units
- correct units for analysis
- time of observation



Histogram for Station:elpapalon.editado of PRCP>=1mm



Outliers

- Unreasonable extreme values
- Indices sensitive to extreme values
- Throwing out valid extreme values can cause errors as easily as keeping erroneous extreme values

Detecting Outliers (general concept)

- Measure each observation with reference to distribution of observations
- Temperature
 - gaussian (normal) distribution
 - use standard deviation as reference
- Rainfall
 - positively skewed
 - quasi-normalise
 - use cube-root rainfall
 - use other distribution
 - gamma distribution
 - generalised pareto distribution

When is an outlier an outlier?

- Err on side of caution
 - If in doubt, evaluate it
 - but don't necessarily set it to missing
- Look for collaborative evidence
 - Surrounding stations
 - Other data e.g. low Tmax for high rainfall, synoptic pattern
 - Very hot, cold or wet days usually occur as part of a spell

Correcting errors

- Very difficult to correct daily data
- tmin > tmax
 - swap? Minus sign missing? One or both in error?
- outlier: typographic (Fat Finger) error
 - e.g. 31.0 should be 13.0 or 21.0
 - Check with day before and day after
 - check with neighbouring stations

Conclusion

- Data quality of prime importance

 Accuracy and consistency
- Problems may be identified but more difficult to correct
- Undetected problems always possible
- Exactly how we will do QC will be shown next

