# Quality Control and Homogenisation of the Belgian Historical Weather Data

**C. Delvaux, R. Ingels**, V. Vrabel, M. Journée, C. Bertrand

*Royal Meteorological Institute of Belgium*

*Climatological and Meteorological Information Service*

# Introduction

- **Recent digitization project of belgian data**

- Based on monthly climate bulletin
- daily data (temperature and precipitation)
- from 1880 to 1950

- Extend the Belgian daily data already available in our database from 1951 to nowadays

# Introduction

- **Project**

Create high quality climatological long series in Belgium
(*period 1880-2015*)

- **Parameters**

Daily maximum temperature (TX)
Daily minimum temperature (TN)
Precipitation (RR)

- **Main steps to obtain good results**
→ creation of long series
→ quality control of the data
→ monthly homogenization of the data (HOMER - Ongoing work)

# Creation of long series

Almost none of the stations covers the entire period of time ->

**The long series can be a combination of stations**
*maximum distance (10 km)*
*maximum elevation difference (50 m)*

- **Long series** (1880 – 2015)
*27 RR & 16 TT*
- **Short series** (1951 – 2015)
*162 RR & 66 TT*

# Quality control – Part 1

- Specific quality control needed for the new encoded data (1880 – 1950) of the long series (~ 1 million data) because of the bad quality of the data

<div style="border:1px solid;">**Examples of the most frequent errors found**</div>

**1) during encoding** : wrong parameter encoded, duplicated data, confusion between missing/zero values, data attributed to a wrong station, classical typing error

➔ *Automatic and visual tests*



2/1886

201_ANTWERPEN
208_HUY-STATTE
212_OOSTENDE
213_SINT-TRUIDEN
216_LEOPOLDSBURG CAMP DE BEVERLOO
221_GEMBLOUX
223_THIMISTER

easy to detect and correct :

~~121 °C~~ → 12.1 °C

difficult to find, even with accurate test

**9.9 °C -> 5.5 °C**

# Quality control – Part 1

- Specific quality control needed for the new encoded data (1880 – 1949) of the long series (~ 1 million data) because of the bad quality of the data

> **Examples of the most frequent errors found**

2) Observer error : precipitation not measured every day → accumulation

3) Transmission of data :
Bad communication between
the institute and the observer



➔ **About 20 % of data have been modified between 1880 – 1950 !**

# Quality control – Part 2

| Index | Confidence |
|-------|------------|
| v | validated data |
| c | corrected data |
| s | suspicious data |

**Daily temperature data**

Physical Limits Consistency — KO → Index = s1 / OK
Internal Consistecy — KO → Index = s2 / OK
Plausible Value — KO → Index = s3 / OK
Temporal Consistency — KO → Index = s4 / OK
Spatial Consistency — KO → Index = s5 / OK
Already corrected ? — yes → Index = c / no
Index = v

- New quality control procedures
  *Minimum data quality required*

- Applied to all the daily data (1880 – 2015)

- Basic tests to more specific tests

- Apply a quality index for each daily data
  - Validated data *(v)*
  - Suspicious data *(sX)* where *X* explain why the data is suspicious
  - Corrected data *(c)*

- Examples of some tests with TT (made for TN and TX)

# Quality control – Part 2

| Index | Confidence |
|-------|------------|
| v | validated data |
| c | corrected data |
| s | suspicious data |

Daily temperature data

Physical Limits Consistency
KO → Index = s1
OK

Internal Consistecy
KO → Index = s2
OK

Plausible Value
KO → Index = s3
OK

Temporal Consistency
KO → Index = s4
OK

Spatial Consistency
KO → Index = s5
OK

Already corrected ?
yes → Index = c
no

Index = v

**- 50 °C < TT < 50 °C**

**TX > TN**

# Quality control – Part 2



| Index | Confidence |
|-------|-----------|
| v | validated data |
| c | corrected data |
| s | suspicious data |

$T_{min} < TN < T_{max}$

Envelope which assumes that the annual temperature variations follow a sinusoidal wave
-> *Upper and lower bounds by regions*
-> *Based on validated extreme temperature data observed each day*



*Example for TN – Lemberge (1981)*

# Quality control – Part 2

| Index | Confidence |
|-------|-----------|
| v | validated data |
| c | corrected data |
| s | suspicious data |

**Daily temperature data**

Physical Limits Consistency — KO → Index = s1
— OK →

Internal Consistecy — KO → Index = s2
— OK →

Plausible Value — KO → Index = s3
— OK →

Temporal Consistency — KO → Index = s4
— OK →

Spatial Consistency — KO → Index = s5
— OK →

Already corrected ? — yes → Index = c
— no →

Index = v

**|TN(day)-TN(day-1)| < Ɛ**

- *Ɛ based on the extreme temperature difference observed between two consecutive validated data*

- by month & by regions

| Month | Ɛ (TN) |
|-------|--------|
| Jan | 14.3 |
| Feb | 12.5 |
| Mar | 12 |
| Apr | 10.1 |
| May | 10.4 |
| Jun | 10.3 |
| Jul | 9.6 |
| Aug | 9.4 |
| Sep | 10.6 |
| Oct | 11.7 |
| Nov | 13.2 |
| Dec | 12.6 |

# Quality control – Part 2

| Index | Confidence |
|-------|-----------|
| v | validated data |
| c | corrected data |
| s | suspicious data |



**1) Classic spatial test**

**5 closest neighboring values**
*Based on "distance + 100 * altitude"*

- Inverse Distance Weighting
- Standard Deviation



TX values suspicious (too warm) if :

**TX** > IDW + Standard Deviation + 6°C

*AND*

**TX** > TX *(of the 5 neighbors!)* + 4°C

# Quality control – Part 2

| Index | Confidence |
|-------|-----------|
| v | validated data |
| c | corrected data |
| s | suspicious data |

Daily temperature data

Index = s1 ← KO — Physical Limits Consistency
↓ OK
Index = s2 ← KO — Internal Consistecy
↓ OK
Index = s3 ← KO — Plausible Value
↓ OK
Index = s4 ← KO — Temporal Consistency
↓ OK
Index = s5 ← KO — Spatial Consistency
↓ OK
Index = c ← yes — Already corrected ?
↓ no
Index = v

## 2) Trend test

- Neighbor comparison of daily rise/drop of temperature

| Date | TX (Leuven) | Trend | 5 closest neighbors | Trend Test |
|------|-------------|-------|---------------------|------------|
| 9/07/42 | 26.1 | .. | .. | .. |
| 10/07/42 | 16.3 | ↓↓↓ | ALL ↓↓↓ | OK |
| 11/07/42 | 19.2 | ↑↑↑ | ALL ↓↓↓ | KO |

↑↑↑ = *temperature increase of more than 2 degrees*

↓↓↓ = *temperature decrease of more than 2 degrees*

# Quality control – Part 2

- Quality Control procedures is realized in **two times**

-> First run allows to assign a first quality index to all the data
-> Second run takes only validated data for spatial tests

| Date | TX (Leuven) | Trend | 5 closest neighbors | Quality Index |
|---|---|---|---|---|
| 9/07/1942 | 26.1 | .. | | .. |
| 10/07/1942 | 16.3 (s51!) | ↓↓↓ | ALL ↓↓↓ | s51 |
| 11/07/1942 | 19.2 | ↑↑↑ | ALL ↓↓↓ | ~~S52~~ v |

# Quality control – Part 2

**Results** :

~ 99.5 % of validated temperature data (about 10000 values)

*Can be explained by the basic QC already made for data from 1951*

*About 80 % of the suspicious values are detected by spatial tests*

**Some corrections** when it was possible (especially for *s1* and *s2*)

Only when no doubt

# Homogenisation of temperature short series

# Station locations

- 66 stations (short series TT)

- From 1951 to 2015

- At least 90 % of daily data

- 11 foreign stations

    - 5 FR

    - 3 GE

    - 3 DE

# Metadata

- Station catenations

- Shelter relocation

- Change of shelter type

- Change of instrument

- Automatisation

- Change of observer

- Other things like information on the shelter site

# Methodology

- HOMER

- Trainings

- 3 people working separatly on different cluster composition

- Common breaks (usually big ones)

- Improving breaks list

- Re-do homogenisation with final clusters

# Clusters

Creation based on :
- proximity
- correlation
- climatic area



- 5 clusters
- Around 15 stations

# First results with HOMER

TN

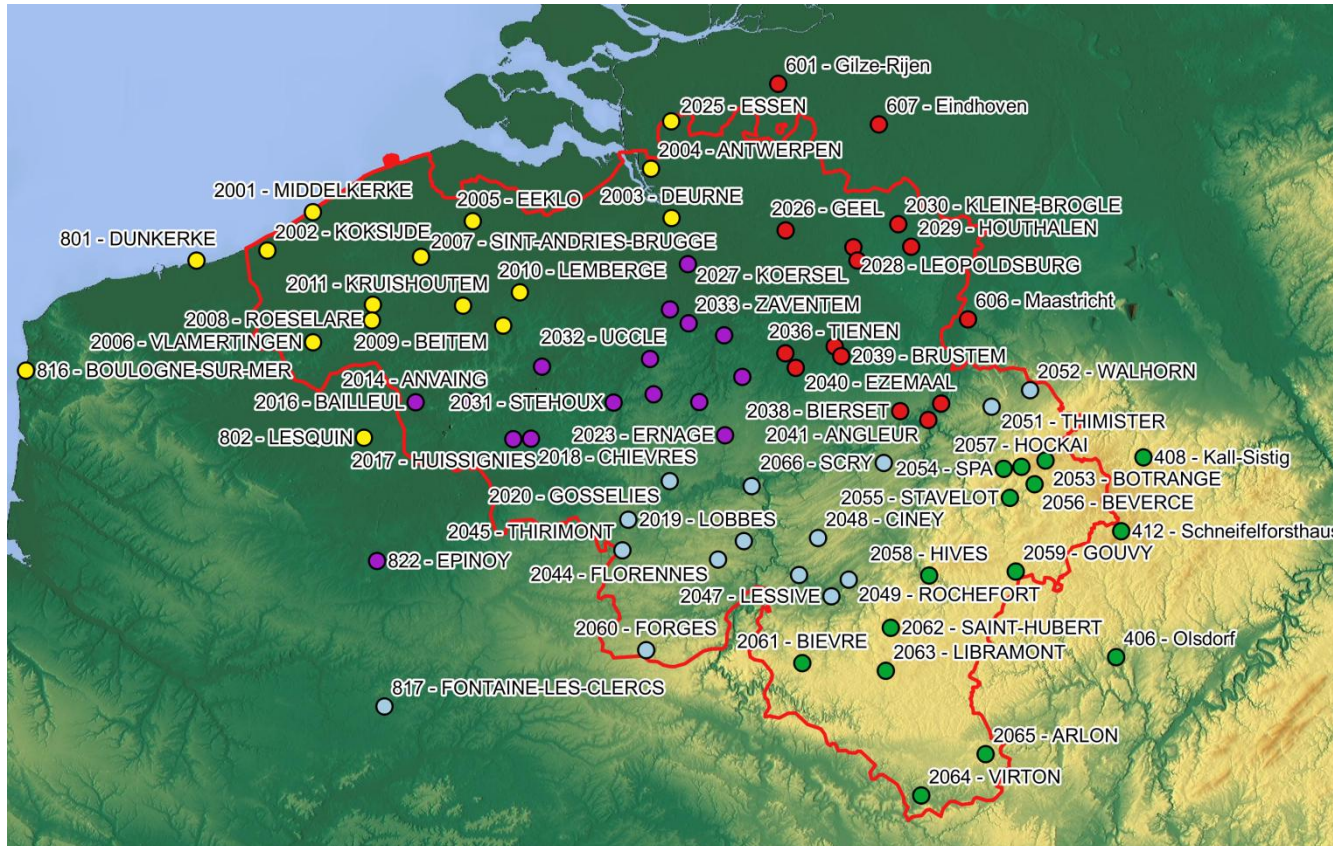| TT code | TT name | date | amplitude | MMD | metadata |
|---|---|---|---|---|---|
| 817 | FONTAINE-LES-CLERCS | 12/1967 | -0.52 | | instrument change |
| 2019 | LOBBES | 1/1984 | 0.72 | | station catenation |
| 2020 | GOSSELIES | 8/1974 | 0.20 | | station catenation |
| 2043 | DENEE-MAREDSOUS | 12/1981 | -1.33 | | station catenation |
| 2044 | FLORENNES | 12/1987 | 0.31 | | |
| 2045 | THIRIMONT | 4/1974 | 0.46 | | station catenation |
| 2045 | THIRIMONT | 12/1993 | -0.63 | | |
| 2046 | MALONNE | 10/1997 | -0.67 | | station catenation |
| 2047 | LESSIVE | | no breaks | | |
| 2048 | CINEY | 12/2007 | -0.34 | | |
| 2049 | ROCHEFORT | 12/1966 | -0.22 | | |
| 2049 | ROCHEFORT | 12/2008 | -0.25 | | |
| 2050 | HOUYET | | no breaks | | |
| 2051 | THIMISTER | | no breaks | | |
| 2052 | WALHORN | | no breaks | | |
| 2060 | FORGES | | no breaks | | |
| 2066 | SCRY | 12/1962 | -1.21 | | station catenation |
| 2066 | SCRY | 12/2005 | 0.39 | | |

TX

| TT code | TT name | date | amplitude | MMD | metadata |
|---|---|---|---|---|---|
| 817 | FONTAINE-LES-CLERCS | | no breaks | | |
| 2019 | LOBBES | 1/1984 | -0.56 | | station catenation |
| 2020 | GOSSELIES | 8/1974 | 0.28 | | station catenation |
| 2020 | GOSSELIES | 12/1995 | -0.27 | | |
| 2043 | DENEE-MAREDSOUS | 12/1966 | 0.30 | | |
| 2044 | FLORENNES | 12/1987 | 0.46 | | |
| 2044 | FLORENNES | 12/1996 | -0.44 | | |
| 2045 | THIRIMONT | 4/1974 | -0.93 | | station catenation |
| 2046 | MALONNE | 10/1997 | -0.58 | | station catenation |
| 2047 | LESSIVE | 5/1974 | 2.79 | | observer change |
| 2047 | LESSIVE | 9/1999 | 0.53 | | station catenation |
| 2048 | CINEY | 12/1981 | -0.23 | | |
| 2049 | ROCHEFORT | 12/2009 | 0.25 | | |
| 2050 | HOUYET | | no breaks | | |
| 2051 | THIMISTER | 12/1975 | 0.28 | | |
| 2051 | THIMISTER | 4/1989 | -0.52 | | relocation |
| 2051 | THIMISTER | 12/2004 | 0.39 | | station catenation |
| 2052 | WALHORN | | no breaks | | |
| 2060 | FORGES | 12/1963 | -0.29 | | |
| 2066 | SCRY | 12/1962 | -1.42 | | station catenation |

# First results with HOMER

TX

Lessive station

| TT code | TT name | date | amplitude | MMD | metadata |
|---------|---------|------|-----------|-----|----------|
| 817 | FONTAINE-LES-CLERCS | | no breaks | | |
| 2019 | LOBBES | 1/1984 | -0.56 | | station catenation |
| 2020 | GOSSELIES | 8/1974 | 0.28 | | station catenation |
| 2020 | GOSSELIES | 12/1995 | -0.27 | | |
| 2043 | DENEE-MAREDSOUS | 12/1966 | 0.30 | | |
| 2044 | FLORENNES | 12/1987 | 0.46 | | |
| 2044 | FLORENNES | 12/1996 | -0.44 | | |
| 2045 | THIRIMONT | 4/1974 | -0.93 | | station catenation |
| 2046 | MALONNE | 10/1997 | 0.58 | | station catenation |
| 2047 | LESSIVE | 5/1974 | 2.79 | | observer change |
| 2047 | LESSIVE | 9/1999 | 0.53 | | station catenation |
| 2048 | CINEY | 12/1981 | -0.23 | | |
| 2049 | ROCHEFORT | 12/2009 | 0.25 | | |
| 2050 | HOUYET | | no breaks | | |
| 2051 | THIMISTER | 12/1975 | 0.28 | | |
| 2051 | THIMISTER | 4/1989 | -0.52 | | relocation |
| 2051 | THIMISTER | 12/2004 | 0.39 | | station catenation |
| 2052 | WALHORN | | no breaks | | |
| 2060 | FORGES | 12/1963 | -0.29 | | |
| 2066 | SCRY | 12/1962 | -1.42 | | station catenation |

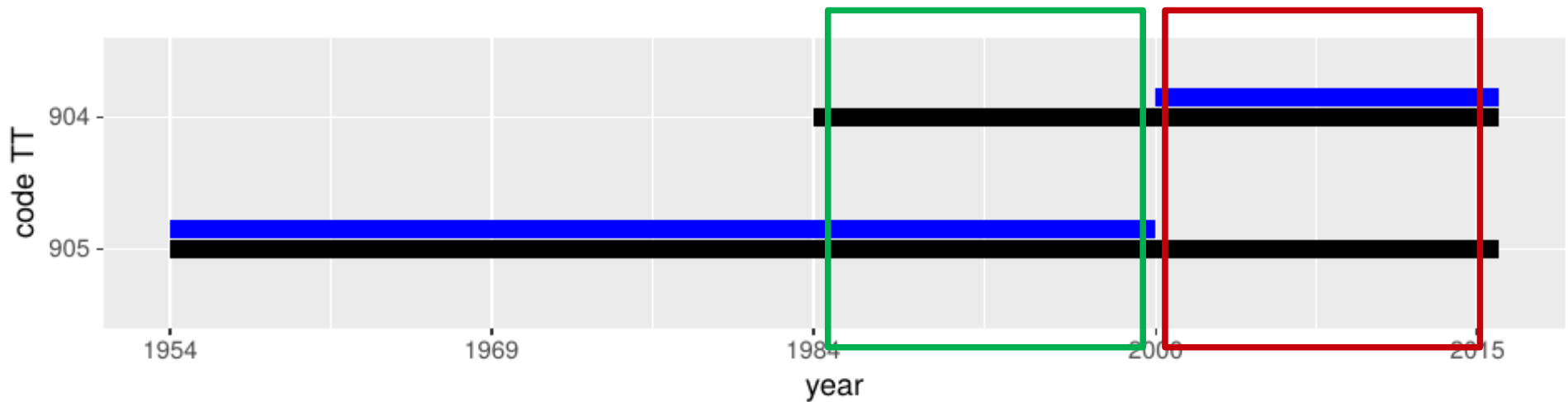# Compare HOMER with something else ?

Parallel data on Uccle site→ extrapolation by linear regression

# Linear regression

Sensitivity study between HOMER and linear regression

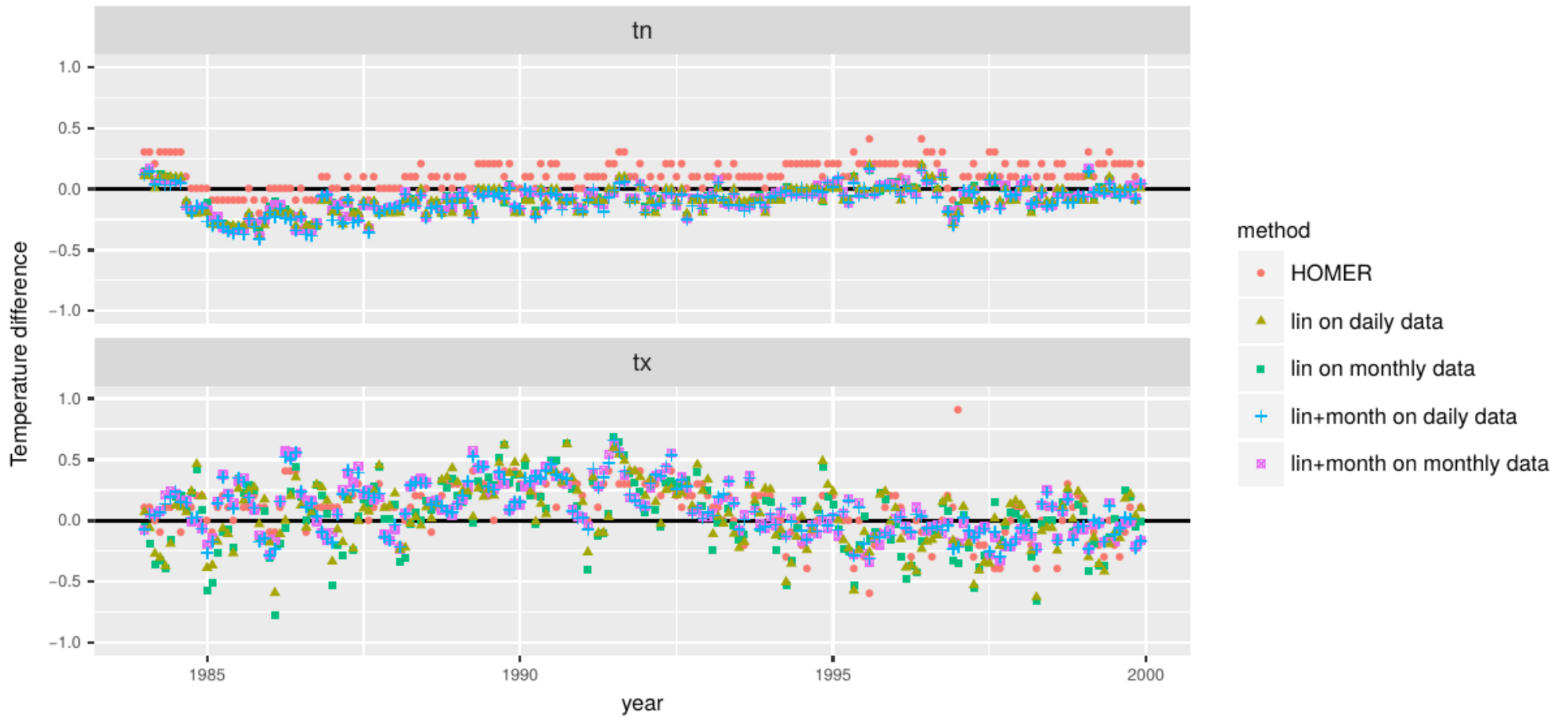Two regression models where computed for TN and TX



| code TT | type | 1954-1983 | 1984-1999 | 2000-2015 |
|---------|------|-----------|-----------|-----------|
| 904 | tn | | 6.7 | 7.3 |
| 905 | tn | 6.0 | 6.6 | 7.1 |
| 904 | tx | | 14.0 | 14.7 |
| 905 | tx | 14.4 | 14.9 | 15.7 |

| msr | linear model trained on daily data | linear model + monthly corr. trained on daily data | linear model trained on monthly data | linear model + monthly corr. trained on monthly data | homer |
|-----|-----|-----|-----|-----|-----|
| tn | 0.14 | 0.15 | 0.13 | 0.14 | 0.19 |
| tx | 0.26 | 0.23 | 0.27 | 0.23 | 0.23 |

# Linear regression



Residuals for monthly data during 1984–1999

# HOMER oddities

- Begining and ending of series (interpolation and homogenisation)

- Climatic events

- ACMANT ?  Didn't use

- Order of break implantation consequence

- Difference TX/TN ?

# Conclusion & perspectives

Sometimes very difficult to get all metadata, especially for older stations before 50'

Homer → Human factor very important, be very careful
   → Sensitivity test ok

What's next ? → Spatial interpolation of homogenised temperature, long series and precipitation, daily homogenisation