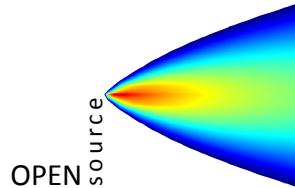

The **openair** manual

open-source tools for analysing air
pollution data

University of York and Ricardo Energy & Environment



David Carslaw

version: 4th May 2019

The **openair** package and the documentation are provided in the hope that it is useful to you. While these tools are free, if you use them, we ask that acknowledgement be given to their source. Please cite:

Carslaw, D.C. and K. Ropkins, (2012). *openair — an R package for air quality data analysis.* Environmental Modelling & Software. Volume 27-28, pp. 52–61.

Carslaw, D.C. (2019). The *openair manual — open-source tools for analysing air pollution data.* Manual for version 2.6-5, University of York.

Every effort has been made to ensure the proper functioning of these tools and their effective and accurate documentation. We would like to be informed of any problems or queries related to their use. All correspondence should be sent to <mailto:david.carslaw@york.ac.uk>.

The development website is <https://github.com/davidcarslaw/openair>. Please use the Github website ‘Issues’ tracker to report any bugs, make suggestions or ask questions.

This document was produced using R version 3.6.0 and **openair** version 2.6-5.

Copyright © 2019 David C. Carslaw.

Contents

I Introduction to R for analysing data	7
1 Introduction	7
1.1 Background	7
1.2 Using this document	7
1.3 The open-source approach	8
1.4 Aims	8
1.5 Can I use these tools for commercial purposes?	9
2 Statistical software R	10
2.1 Introduction to R	10
2.2 Why use R?	11
2.3 Why not use R?	11
2.4 Some useful R resources	12
3 Basic use of R	13
3.1 Introduction	13
3.2 Use it as a calculator	13
3.3 Basic graph plotting	15
3.4 Getting help	16
4 Using R to analyse air pollution monitoring data	17
4.1 Getting data into R	17
4.2 More sophisticated plotting	20
4.3 Plotting time series with different averaging times	20
4.4 Summarising time series data	24
4.5 Relationships between variables	24
5 General use of R — advice and examples	26
5.1 Data input and output	26
5.1.1 Data import	26
5.1.2 Data export	28
5.2 Selecting and replacing parts of vectors and data frames	28
5.3 Combining and cleaning up files	32
5.4 Reshaping data	38
5.5 Example: converting hour-day data to column format	40
5.6 Daily means from hourly means — processing wind direction data	41
5.7 Using an Editor	42
5.7.1 Using the built-in editor	42
5.7.2 Using a dedicated editor	43
5.8 Several plots on one page	44
5.9 Saving and using plots	46
5.10 Graphing lots of data — using level plots	47
5.11 Special symbols for use in plotting air pollution data	47
5.12 Using databases with R	50
6 Multivariate plots — introduction to the Lattice package	50
6.1 Introduction to the Lattice package	50
6.2 Example simple plots	51

7 Functions in R	52
II Dedicated functions for analysing air pollution data	54
8 Introduction	54
8.1 Installing and loading the openair package	54
8.2 Access to source code	55
8.3 Brief introduction to openair functions	55
8.4 Input data requirements	60
8.4.1 Dealing with more than one site	61
8.5 Using colours	61
8.6 Automatic text formatting	62
8.7 Multiple plots on a page	63
8.8 Annotating openair plots	63
8.9 Getting help	69
9 Getting data into openair	71
9.1 Issues related to time zones	71
9.2 The import function	72
9.3 The importAURN function	75
9.4 The importKCL function	76
9.5 Importing data from the CERC ADMS modelling systems	77
9.5.1 An example considering atmospheric stability	78
10 The summaryPlot function	81
11 The cutData function	85
12 The windRose and pollutionRose functions	86
12.1 Purpose	86
12.2 Options available	87
12.3 Example of use	89
13 The percentileRose function	97
13.1 Purpose	97
13.2 Options available	97
13.3 Example of use	99
14 The polarFreq function	102
14.1 Purpose	102
14.2 Options available	102
14.3 Example of use	104
15 The polarPlot and polarCluster functions	108
15.1 Purpose	108
15.2 Options available	110
15.3 Example of use	114
15.3.1 Conditional Probability Function (CPF) plot	118
15.3.2 The polarCluster function for feature identification and extraction	121

16 The <code>polarAnnulus</code> function	126
16.1 Purpose	126
16.2 Options available	126
16.3 Example of use	129
17 The <code>timePlot</code> function	130
17.1 Purpose	130
17.2 Options available	131
17.3 Example of use	134
18 The <code>timeProp</code> function	137
19 The <code>calendarPlot</code> function	140
19.1 Purpose	140
19.2 Options available	140
19.3 Example of use	142
20 The <code>TheilSen</code> function	147
20.1 Purpose	147
20.2 Options available	149
20.3 Example of use	151
20.4 Output	154
21 The <code>smoothTrend</code> function	155
21.1 Purpose	155
21.2 Options available	156
21.3 Example of use	158
22 The <code>timeVariation</code> function	162
22.1 Purpose	162
22.2 Options available	163
22.3 Example of use	165
22.4 Output	171
23 The <code>scatterPlot</code> function	172
23.1 Purpose	172
23.2 Options available	172
23.3 Example of use	176
24 The <code>linearRelation</code> function	182
24.1 Options available	183
24.2 Example of use	184
25 The <code>trendLevel</code> function	185
25.1 Purpose	185
25.2 Options available	185
25.3 Example of use	187
26 openair back trajectory functions	191
26.1 Trajectory gridded frequencies	198
26.2 Trajectory source contribution functions	201
26.2.1 Identifying the contribution of high concentration back trajectories	202
26.2.2 Allocating trajectories to different wind sectors	203

26.2.3 Potential Source Contribution Function (PSCF)	205
26.2.4 Concentration Weighted Trajectory (CWT)	205
26.3 Back trajectory cluster analysis with the <code>trajCluster</code> function	208
27 Model evaluation — the <code>modStats</code> function	213
27.1 Purpose	213
27.2 Options available	215
27.3 Example of use	216
28 Model evaluation — the <code>TaylorDiagram</code> function	218
28.1 Purpose	218
28.2 Options available	219
28.3 Example of use	222
29 Model evaluation — the <code>conditionalQuantile</code> and <code>conditionalEval</code> functions	225
29.1 Purpose	225
29.2 Options available	225
29.3 Example of use	226
30 The <code>calcFno2</code> function—estimating primary NO₂ fractions	234
30.1 Purpose	234
30.2 Options available	235
30.3 Example of use	235
31 Utility functions	237
31.1 Selecting data by date	237
31.2 Selecting run lengths of values above a threshold — pollution episodes .	238
31.3 Calculating rolling means	240
31.4 Aggregating data by different time intervals	241
31.5 Calculating percentiles	246
31.6 The <code>corPlot</code> function — correlation matrices	247
31.7 Preparing data to compare sites, for model evaluation and intervention analysis	249
31.7.1 Intervention analysis	249
31.7.2 Combining lots of sites	251
Acknowledgements	252
Further information and bug reporting	253
A Installing and maintaining R	258
A.1 Downloading and installing R	258
A.2 Maintenance	258
B Bootstrap estimates of uncertainty	259
B.1 The bootstrap	259
B.2 The block bootstrap	260
C A closer look at trends	262
D Production of HYSPLIT trajectory files	267



Part I

Introduction to R for analysing data

1 Introduction

1.1 Background

This document provides information on the use of computer software called ‘R’ to analyse air pollution data. The document supports an initiative to develop and make available a consistent set of tools for analysing and understanding air pollution data in a free, open-source environment.

The amount of monitoring data available is substantial and increasing. In the UK alone there are thought to be over 1000 continuous monitoring sites. Much of the data available is only briefly analysed; perhaps with the aim of comparing pollutant concentrations with national and international air quality limits. However, as it will hopefully be seen, the critical analysis of air pollution data can be highly rewarding, perhaps yielding important information on pollutant sources that was previously unknown or unquantified.

There are however, several barriers that prevent the more effective analysis of air pollution data. These include user knowledge (knowing how to approach analysis and which techniques to use), cost (both in time and money) and access to specialist software that might be necessary to carry out all but the most simple types of analysis. Part of the aim of this document and its associated tools is to overcome some of these barriers.

1.2 Using this document

This document has been specifically written for those with an interest in analysing air quality data, although the techniques also lend themselves to wider atmospheric science problems and other sources of data e.g. traffic data. It is split into two parts.

Part I gives some background information on R and provides examples of using R with an air pollution data set from London. The intention here is to give an overview of how to use R and this is done mostly by considering actual air pollution data.

Part II describes dedicated functions for analysing air pollution data, which are available in the R ‘package’ called **openair** (Carslaw and Ropkins 2012). Even though the capabilities of the functions in Part II are greater than those highlighted in Part I, they are easier to apply.

The document assumes no previous knowledge of R. The document itself contains code which can be used directly in R (just copy and paste it in — see later). All the code used to make plots is shown in this manual. This code is the code that can be typed into R and can be copied directly from this document and pasted into R. The latter makes it easier to become familiar with the R language. The code also uses ‘mark-up’ to highlight functions, options etc. In addition, where a plot is produced, the code immediately precedes it. Users are encouraged to reproduce the plots shown and produce their own variations on them — for example, by plotting different pollutants. The document also contains extensive hypertext links to make it easy to navigate and cross-reference sections, figures etc. The document is meant to be a kind of work book, allowing users

to work through the examples from start to finish. However, it also acts as a reference guide to using R for the specific purposes of analysing monitoring data.

This document was produced entirely using free software. The document itself was produced using the L^AT_EX typesetting programme and all the plots were produced using R.

1.3 The open-source approach

The tools developed that are described here are *open-source*. This means that they are freely available to anybody and all the source code is open to scrutiny. Free software allows users the freedom to run, copy, distribute, study, change and improve the software. This philosophy is espoused at <http://www.gnu.org/philosophy/free-sw.html>.

The open-source approach is fundamental to this initiative. There are many advantages to open-source software tools beyond their zero direct cost. First is the belief that making tools open it will encourage their use and scrutiny. Second, by making tools available in this way it is more likely that others will contribute to them — perhaps identifying or fixing bugs, or maybe developing them further, as is the case for many open-source software projects. As described in Chambers (2007), the open-source approach can help lead to *trustworthy* software, and this is an important component of the aims here. It is a community approach that encourages trust and participation.

There are also difficulties in adopting an open-source approach however, not least the need for organisations to make money from their work. The development of these tools needs to be paid for somehow. However, there are organisations and individuals that see benefit in this way of working and are willing to fund or contribute to this initiative. Those who have contributed so far to this project are listed in the acknowledgements section.

In the environmental field there are maybe even more compelling arguments for using open-source tools. Many believe for example that those affected by environmental decisions reached through using tools/models/data should be able to scrutinise them — and why not? However, the reality is often far from this situation and there is often reliance on ‘black boxes’ where such scrutiny is not possible.

1.4 Aims

The aims of this document are as follows:

- To highlight the importance of looking at data effectively.
- To introduce the statistical software R and provide some background to the language.
- To show how R can be used to look at and understand air pollution monitoring data typically collected by local authorities, on behalf of governments and industry. However, the tools should also be very useful for those involved in academic research.
- To free-up time to think about air quality data rather than spending time preparing and analysing data. Part II in particular is aimed at providing functions that are *easy* to use so that more time can be spent on more productive activities.

The focus is very much on being pragmatic — the basic approach is to learn by doing. There are many books and on-line documents on R and these notes are not meant to

duplicate them. The approach used here is example-based, as it is our experience that this is the best way in which to learn to use R. Also, some of the concepts are easier to digest if applied to a problem or data set that is familiar. This document cannot cover all the aspects of R that may be relevant or useful for the analysis of air pollution, but provides more of an introduction to how R can be used.

It is also important to stress that these functions are most useful when there is a clear purpose to the analysis. While *Exploratory Data Analysis* (EDA) is extremely valuable, there is no substitute for having a clear aim (Tukey 1977). This was perhaps best expressed by the statistician John Tukey who developed the idea of EDA:

The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data.

How this document was produced

One of the aims of this document was to ensure that users are able to reproduce all the analyses *exactly* as intended. This is not a straightforward task for a complex project under continual development. In addition, the large amount of code and functions presented provides many opportunities for things going wrong. It is easy, for example, to show how a function works and provide the results/plot; update the function and then find out that the options have changed and it no longer runs as intended. In other words, the documentation and the tools themselves go out of sync. Even cutting and pasting text can easily go wrong — as we have discovered.

For this reason we have adopted an innovative approach to ensuring that everything works as intended. This document blends text with code in that the whole document must be ‘run’ to produce it. Each time a version of this documentation is produced, all the code is run at the same time to generate all the various outputs e.g. plots. **This means that all users should be able to reproduce exactly the same outputs as shown in this report.¹**

The approach uses a package called *knitr* (Xie 2013a; Xie 2013b). *knitr* mixes a typesetting system (*LATEX*) with R. When a document is produced, blocks of code embedded in the *LATEX* file are recognised and run in R. In some ways it reverses the ‘normal’ way of doing things — rather than document computer code, the documentation is written to contain the code. The document will not compile if the code does not function — it is as simple as that. In our document, most of the outputs are graphics, but increasingly quantitative information will also be produced.

In adopting this approach we found many problems with the manual (and some functions), even though we took care to develop this work. In time we also have ideas for using this approach to automatically carry out analyses. Imagine a report similar to this (but written more as a tutorial) where the data used are your own data. This approach would have the major advantage that all the analyses would be directly relevant to the user in question, and entirely reproducible.

1.5 Can I use these tools for commercial purposes?

In short, the answer is yes. Part of the aim of producing these tools was to allow *anybody* to use them for *any* purpose. Indeed, this is the principal purpose of the Knowledge

Exchange grant. Our work is very much released in the true spirit of the Free Software Foundation <http://www.fsf.org/>. However, there are a few points users should note:

1. If you use these tools in reports, publications etc., we ask that you cite their source (see the preamble at the beginning concerning how to do this).
2. It is not possible to provide a guarantee or warranty for these tools, although we have tried hard to ensure they function as documented and are adopting methods for quality control.
3. We request that should you find these tools useful and enhance them, that you make such enhancements available to us for wider use.
4. For more detailed information on the various licenses under which R and its packages operate, the user is referred to the R website.

2 Statistical software R

2.1 Introduction to R

R is a computer programming language developed specifically for the purposes of analysing data ([R-project](#)). It is variously described as a statistical system, a system for statistical computation and graphics, and an environment for data analysis and statistics. Its origins go back to innovative developments at Bell Laboratories in the USA during the 1970s, where the language developed was called S. Since that time S has become commercial software and is sold as S-Plus by the Insightful Corporation.

Over the past 10 years or so an open-source version of S has been developed called R. Unlike some open-source software R is highly developed, highly capable and well established. It is very robust and works on a wide range of platforms (e.g. Windows, Mac, and Linux). One of its strengths is the large and growing community of leading researchers that contribute to its development. Increasingly, leading statisticians and computational scientists choose to publish their work in R; making their work available worldwide and encouraging the adoption and use of innovative methodologies.

R is available as Free Software under the terms of the [Free Software Foundation's GNU General Public License](#).

Another key strength of R is the package system of its operation. The base software, which is in itself highly capable (e.g. offering for example linear and generalized linear models, nonlinear regression models, time series analysis, classical parametric and nonparametric tests, clustering and smoothing), has been greatly extended by additional functionality. Packages are available to carry out a wide range of analyses including: generalized additive models, linear and non-linear modelling, regression trees, Bayesian approaches etc.

For air pollution purposes, R represents the *ideal* system with which to work. Core features such as effective data manipulation, data/statistical analysis and high quality graphics lend themselves to analysing air pollution data. The ability to develop one's own analyses, invent new plots etc. using R means that advanced tools can be developed for specific purposes. Indeed, Part II of this document is focussed on the use of dedicated tools for air quality analysis. The use of R ensures that analyses and graphics are not constrained to 'off the shelf' tools. These tools will often contain functionalities that are either part of the R base system or that exist through specific packages.

The principal difficulty in using R is the steep learning curve in being able to use it effectively. Many with programming experience (e.g. FORTRAN, C++) tend to find R an

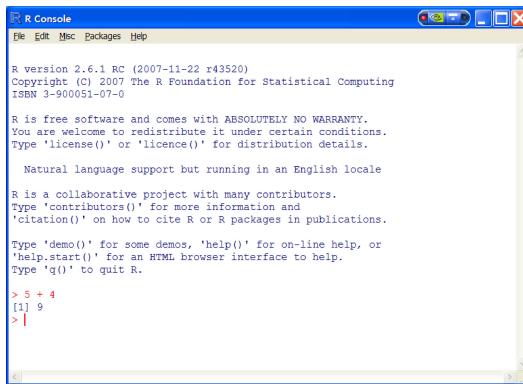


FIGURE 2.1 The basic R console.

easier language to learn than those that have no experience in programming. However, for most users this is not the experience. One of the principal difficulties is the lack of a nice graphical user interface (GUI). However, the more one develops tools in R, the more it is realised that a GUI approach significantly constrains capabilities (try plotting trends for each hour of the day. While it certainly takes considerable effort to learn to use R, the rewards are also high.

2.2 Why use R?

There are numerous reasons why R is a good choice for analysing data. A few are listed below.

- It is **free!** This for many people is the key attraction. For this reason, R has become increasingly popular among a wide range of users including universities and businesses.
- It works on several platforms e.g. Windows, Mac OS, and Linux. This makes it very portable and flexible. It is also extremely robust; it is remarkably bug-free and crashes are very rare.
- It has been designed with data analysis in mind — to carry out analysis quickly, effectively and reliably.
- The base system offers a very wide range of data analysis and statistical abilities.
- Excellent graphics output that will grace any report. Furthermore, all the default plotting options have been carefully thought out unlike Excel, for example, whose default plotting options are very poor. There are over 4000 *packages* that offer all kinds of wonderful analysis techniques not found in any other software package. R continues to rapidly grow in popularity, which has resulted in better information for users e.g. there are now many dedicated books.

2.3 Why not use R?

For all its inherent strengths, R does have drawbacks. Here are a few.

- It is difficult to learn — there is a steep ‘learning curve’.² Many would argue this is not the case; particularly if you are familiar with another programming language

²A steep learning curve means you learn a lot per unit time ...

such as C++ or Fortran. However, our experience is that it is hard work for most people.

- There is no Graphical User Interface; instead one has something that looks like a DOS screen where one types commands. This seems very old-fashioned and at odds with the modern computing experience, but with use you will begin to see this as real advantage. A whole report can be based on a series of ‘scripts’ that can be run to carry out analysis. It is not so easy to record mouse movements and menu choices ...
- There is no help or support — or little that is apparent. This is true, especially compared with commercial software. However, there is extensive on-line help available and many people have written manuals and guides. This document also aims to address the lack of direct, specific support for air pollution analysis.

We have used several systems for data analysis over the years. These have included databases, Visual Basic, GIS, contouring software, statistics software and Excel. We always found it frustrating to do some analysis in Visual Basic, for example, then transfer it to another application for plotting and further analysis. We wrote code to carry out various statistical analyses too. For us, R does all these thing in one place to a very high standard.

Another increasingly important aspect is the ability to run *simulations*. Until relatively recently computers were not powerful enough to routinely run simulations using methods such as randomization, Monte Carlo and bootstrap calculations. All these approaches can greatly enhance many analyses and they are used in many of the functions described in this document. Often, the reason is to obtain a better estimate of uncertainties, which are important to consider when trying to draw inferences from data. Also, many of these methods were essentially inaccessible; or at least beyond consideration for most. This is where R excels — all these methods exist or can be coded building on some of the base functions.

2.4 Some useful R resources

The web is the best place to find information on R. There are many useful documents that people have written (see <http://www.r-project.org/>) under documentation/other. The official user guides can be hard going for many that are completely new to R, but worth a look later.

We also have the *R Book* by Michael Crawley, which is quite useful (Crawley 2007). This is a big book full of examples. It is not well laid-out and has come in for some criticism but we have nevertheless found it to be useful. Another useful book is by John Maindonald and John Braun (Maindonald and Braun 2007). Perhaps the best introductory book is that by Dalgaard (2008), which provides a gentle introduction to R, is well-written and up to date. Dalgaard (2008) also covers basics statistics, with the added benefit of their use in R. It is worth also having a look at the R-project pages as they provide a list of R books—even forthcoming titles not yet published. For those getting into R more seriously, I strongly recommend Spector (2008) for data manipulation and Sarkar (2007) for graph plotting using the *lattice* package.

There are some useful contributed documents on the R web-pages. Under the *Documentation* section on the main page have a look at Other | contributed documentation. We found that ‘An Introduction to R: Software for Statistical Modelling & Computing’ by Petra Kuhnert and Bill Venables was very useful and also the document from David

Rossiter (also check his web site — he has made available lots of high quality information on R and its uses at http://www.itc.nl/personal/rossiter/pubs/list.html#pubs_m_R).

A very good (and free) book is by Helsel and Hirsch (2002) from the US Geological Survey (see <http://pubs.usgs.gov/twri/twri4a3/>). Although the book does not consider R, it provides an excellent grounding in the types of statistics relevant to air pollution monitoring. This book is all about water resources and statistical methods that can be applied to such data. However, there is a great deal of similarity between water pollution and air pollution (e.g. seasonal cycles). One of the authors (Hirsch) is responsible for developing several important methods for analysing trends that are very widely used; including in this document. The document is also very well written and applied in nature — ideal reading for a background in statistical methods relevant to air pollution.

For those that really want to learn R itself, then Matloff (2011) is well worth considering. The book is full of useful examples about how R works and how to use it.

3 Basic use of R

3.1 Introduction

R will of course need to be installed on a computer and this is described in Appendix A. The install procedure is straightforward for Windows systems (the operating system most likely used by those with interests in **openair**). Using R will be a very different experience for most users used to software such as spreadsheets. One of the key differences is that data are stored as *objects*. These objects can take many forms, but one of the most common is called a *data frame*, which is widely used to store spreadsheet-like data. Before working with such data, it is useful to see how simpler data are stored.

3.2 Use it as a calculator

To use R, one should type directly in the console shown in Figure 2.1. Later, it will be shown when more than one line needs to be input, alternative methods can be used to send commands to the console.

R can be used to do simple maths; in this example type in ‘5 + 4’ and press return. The [1] shows that this is the first (and only in this case) result.

```
5 + 4
## [1] 9
```

The output (9) has a ‘[1]’ next to it showing it is the first (and only in this case) result. To assign a value to a variable, in this case x, type

```
x = 5
```

Often it is useful to recall previous commands that have been entered, perhaps modifying them slightly. To recall previous commands the up (\uparrow) and down arrow (\downarrow) keys can be used.

To assign a variable to a value most R users will use the ‘assignment operator’ ($<-$). However, most new users to R find this unnecessarily unusual. In this document we mostly use $<-$ but for most circumstances the $=$ will work the same way.

```
x
## [1] 5
```

Note! R is case sensitive

In the case above using a capital X gives an error:

```
x * 5
## [1] 25
X * 5
## Error in eval(expr, envir, enclos): object 'X' not found
```

It is often necessary to modify a line that has been input previously (this is part of the interactive strength of R). To recall the previous line(s) for editing use the up arrow (\uparrow) on the keyboard.

One of the very important features of R is that it considers ‘vectors’. In the example above, `x` was a single integer 5. However, it can equally represent any sequence of numbers. In the example below, we are going to define two variables `x` and `y` to be samples of 10 random numbers between 0 and 1. In this case there is a special in-built function to generate these numbers called `runif`.

generate
random
numbers

```
x = runif(10)
y = runif(10)
```

To see what `x` looks like, just type it in:

```
x
## [1] 0.49 0.94 0.67 0.82 0.43 0.59 0.23 1.00 0.54 0.48
```

This is one of the most powerful features of R; it also makes it easier to code functions. The ‘traditional’ way of doing this would be to have an array and loops, something like:

```
# for i = 1 to 10
#   x(i) = runif(1)
# next
```

The R code is much neater and can also be extended to matrices and arrays.

Another useful thing to do sometimes is to generate a sequence of numbers. Here the functions `seq` and `rep` are very useful. We first show their basic use and then show a useful example relevant to monitoring data.

To generate a sequence of numbers between 1 and 10:

generate a
sequence
numbers

```
z = seq(1:10)
z
## [1] 1 2 3 4 5 6 7 8 9 10
```

To divide all these number by 10, simply type:

```
z / 10
## [1] 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
```

This code again highlights *vectorised* calculations in R: you need not treat each element of `z` separately. For more information on selecting parts of vectors, and subsets of data in general see (§5.2).

In fact, for this particular example you can just type `z <- 1:10`. Many other possibilities are also available. For example, to generate a sequence of numbers between 1 and 10 using 23 numbers (bit odd, but you get the idea!):

```
z = seq(1, 10, length.out = 23)
z
## [1] 1.0 1.4 1.8 2.2 2.6 3.0 3.5 3.9 4.3 4.7 5.1 5.5 5.9 6.3 6.7
## [16] 7.1 7.5 8.0 8.4 8.8 9.2 9.6 10.0
```

The function `rep` on the other hand repeats numbers a certain number of times. For example, to repeat 1, 10 times:

```
rep(1, 10)
## [1] 1 1 1 1 1 1 1 1 1 1
```

Now, these functions can be combined. A scenario is that you have a data set of concentrations exactly one year long and you want the hour of the day (if you have the date/time there are other ways to do this shown later). A year is 365×24 hours long (8760 hours). What is needed is to repeat the sequence 0–23, 365 times. This is done by:

```
hour = rep(seq(0, 23), 365)
```

Easy! There are loads of variations on this theme too. It's the sort of thing that should be easy but often is a pain to do in other software.

A very common thing to do in R is combine a series of numbers. This is done by *concatenating* them. For example, if one wanted `p` to represent the numbers 1, 3, 5, 10:

```
p = c(1, 3, 5, 10)
p
## [1] 1 3 5 10
```

The use of `c()` is extremely common in R and is used in almost every analysis. One example is when setting the axis limits in plots. A long-hand way of doing this might be to separately set the lower and upper ranges of the limits (e.g. `x.lower = 0, x.upper = 100`). However, with the `c` command, it is more compactly written as `xlim = c(0, 100)`.

3.3 Basic graph plotting

One of R's strengths is the ease with which graphs can be plotted. Almost all graphs in R use an ‘ink on paper’ approach — once something is added it cannot be changed. You will need to plot it again.

To plot `x` against `y` as a scatter `plot`, simply use the `plot` function.

You should end up with something similar to Figure 3.1 (but not exactly the same, because we are plotting random numbers).

repeating a sequence of numbers or characters

```
plot(x, y)
```

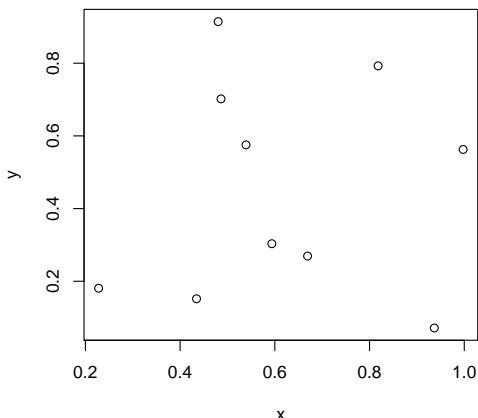


FIGURE 3.1 Plot of 10 random numbers between 0 and 1.

```
z <- rnorm(1000)
hist(z)
```

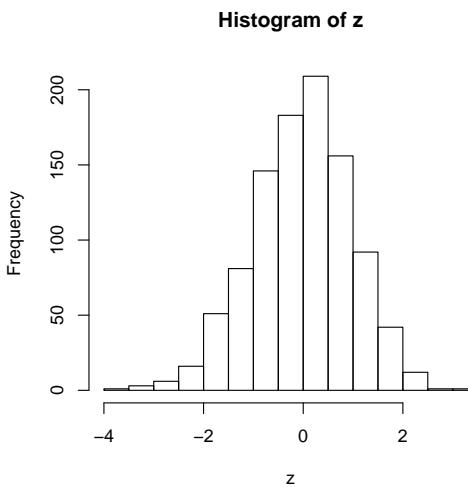


FIGURE 3.2 Histogram of 1000 random numbers drawn from a normal distribution.

Because R is mostly considered as a system for doing statistics, it should come as no surprise that doing some basic statistics and plotting some basic graphs is easy. One useful plot is a histogram. There is an in-built function called `hist`. To make it a bit more interesting, we are first going to draw 1000 random samples from a normal distribution, again the in-built function `rnorm` can do this for you, as shown for Figure 3.2.

The plot should look something like Figure 3.2. This should give you some idea that using R can be very efficient and straightforward for basic plotting. To make a histogram in Excel by contrast is pretty painful. Note also that the bin width size has been automatically chosen — in this case to be 0.05.

3.4 Getting help

There are various ways to get help. For particular functions like `plot`, `hist`, `seq` etc., simply type `help(plot)` or `?plot` (or whatever) to bring up a screen that lists

all the options of a function. Almost always shown are examples of use too. Failing that, try the R-project website, choose ‘Search’ on the first page and then ‘R site search’. These web pages are very useful because they contain questions asked by R-users and answered by R experts. You need to read the ‘posting guide’ before sending question to this site — people supply answers in their own time so make sure you have explored all avenues first.

To search for information in installed packages, `help.search` is useful e.g.

```
help.search("polar plot")
```

4 Using R to analyse air pollution monitoring data

4.1 Getting data into R

So far we have just made up our data using some in-built R functions. A more likely scenario is that you have a data file that you want to analyse. R has many powerful facilities for importing data. However, we are going to keep things simple because it is worth getting your data into a simple format in the first place. Data more often than not is represented as columns on a rectangular grid — and unfortunately is often in Excel with all of its complexity of sheets and formatting. *It is strongly recommended that you work with data in a .csv format.* The files are simple, small and easily read by a wide variety of software.

The importance of data preparation

Experience shows that it is well worth putting a lot of effort into making sure your data are correct before you start analysing them. When you prepare a file carry out some basic checks. For R it is advisable to have variable names with no spaces. It is also more convenient to keep variable names in lower case because it is more difficult to type in capital letters.

It is also convenient to keep your variables names simple because you may need to refer to them lots. Also, we find it is a good idea to keep everything in lower case letters, again because they are quicker to type.

A file containing real monitoring data from Marylebone Road in London has been put together. This is a ‘warts n’ all’ file and contains missing data — typical of that available elsewhere. We have also put in some basic meteorological data (wind speed and direction). The file is called ‘example data long.csv’. This file contains a column representing dates/times in the format dd/mm/yyyy HH:MM, which is very common in data sets of monitoring data. The file is available from <http://www.openair-project.org>.

The easiest way to work with this file is first to set the working directory of R to be the same as that where the file is. This is done using the File|Change dir... command in R.

To read the file in, type (or better still copy and paste this directly into R):

```
## note! - remember to change the directory
mydata <- read.csv("~/My Drive/openair/Data/example data long.csv",
                  header = TRUE)

mydata$date <- as.POSIXct(strptime(mydata$date, format = "%d/%m/%Y %H:%M",
                                    tz = "GMT"))
```

Read in a csv
file

Note that the formatting of date/times can be rather intimidating for new users of R. For this reason it is recommended that the functions in the **lubridate** package are used (Section 9.1 for more information). Using **lubridate** the code above can be more simply written:

```
library(lubridate)
mydata <- read.csv("~/My Drive/openair/Data/example data long.csv",
                      header = TRUE)
mydata$date <- dmy_hm(mydata$date)
```

Note, most of the rest of the manual uses the **lubridate package to format the dates.**

Note that the **openair** package has tools to make this easier — see [Getting data into openair](#). The information presented here shows how to do things using basic R functions.

This reads the file into something called a **data.frame**, which is the most common way in which R stores data (others include vectors, matrices and arrays).

An alternative if you want to bring up a dialog box that allows you to browse the file system (and then format the date) is:

```
mydata <- read.csv(file.choose(), header = TRUE)
mydata$date <- dmy_hm(mydata$date)
```

Another neat way of getting data into R if you are interested in having a quick look at data is to read data directly from the clipboard. If you are using Excel you can copy some data (which may or may not include column headers; better if it does) and read it into R easily by:

```
mydata <- read.delim("clipboard", header = TRUE)
```

Which assumes that the data did have a header field. If not R will provide column names like V1, V2 etc. and the option **header = FALSE** should be used instead.

We will come onto the second line in a minute. It is *always* useful to check to see what has been loaded and using the **summary** command in R is one useful way to check:

format
date-times

view a
summary of
the data

```
summary(mydata)

##      date                  ws                  wd                  nox
##  Min.   :1998-01-01 00:00:00  Min.   : 0   Min.   : 0   Min.   : 0
##  1st Qu.:1999-11-14 15:00:00  1st Qu.: 3   1st Qu.:140  1st Qu.: 82
##  Median :2001-09-27 06:00:00  Median : 4   Median :210  Median :153
##  Mean   :2001-09-27 06:00:00  Mean   : 4   Mean   :200  Mean   :179
##  3rd Qu.:2003-08-10 21:00:00  3rd Qu.: 6   3rd Qu.:270  3rd Qu.:249
##  Max.   :2005-06-23 12:00:00  Max.   :20   Max.   :360  Max.   :1144
##                               NA's   :632  NA's   :219  NA's   :2423
##      no2                  o3                  pm10                 so2                  co
##  Min.   : 0   Min.   : 0   Min.   : 1   Min.   : 0   Min.   : 0
##  1st Qu.: 33  1st Qu.: 2   1st Qu.:22  1st Qu.: 2   1st Qu.: 1
##  Median : 46  Median : 4   Median :31  Median : 4   Median : 1
##  Mean   : 49  Mean   : 7   Mean   :34  Mean   : 5   Mean   : 1
##  3rd Qu.: 61  3rd Qu.:10  3rd Qu.:44  3rd Qu.: 6   3rd Qu.: 2
##  Max.   :206  Max.   :70   Max.   :801  Max.   :63   Max.   :20
##  NA's   :2438 NA's   :2589  NA's   :2162  NA's   :10450  NA's   :1936
##      pm25
##  Min.   : 0
##  1st Qu.: 13
##  Median : 20
##  Mean   : 22
##  3rd Qu.: 28
##  Max.   :398
##  NA's   :8775
```

The **summary** function is extremely useful. It shows that there are 9 variables: date, ws, wd It provides the minimum, maximum, mean, median, and the first and third quantiles. It shows for example that NO_x ranges from 0 to 1144 ppb and the mean is 178.8 ppb. Also shown is something called **NA's**, which are missing data. For NO_x there are 2423 missing values for example. These missing values are very important and it is also important to know how to deal with them. When R read the .csv file, it automatically assigned missing data the value NA.

Also note that the date is read in as a character string and R does not know it is a date. Dealing properly with dates and times in any software can be very difficult and frustrating. There are time zones, leap years, varying lengths of seconds, minutes, hours etc.; all in all it is highly idiosyncratic.³ R does however have a robust way of dealing with dates and times. In this case it is necessary to deal with dates and times and two functions are used to convert the date to something recognised as such by R. The function **strptime** tells R what format the data are in — in this case day/month/year hour:minute i.e. it ‘strips’ the date out of the character string. In addition, R is told that the data are GMT. This allows great flexibility for reading in dates and times in a wide range of formats. The **as.POSIXct** function then converts this to a convenient format to work with. This may appear to be complicated, but it can be applied in the same way to all files and once done, it is possible to proceed without difficulty. The other need for storing dates/times in this way is to deal with GMT/BST (or ‘daylight savings time’). Some of the functions in Part II use time zone information to process the data, because, for example, emissions tend to vary by local time and not GMT. Note that the **openair** package will automatically do these conversions if the data are in a format dd/mm/yyyy HH:MM.⁴ For advice on dealing with other time zones, please see [Section 9.1](#).

If you just want to know what variables have been read in, it is easier to type

show the
variable names

³Just to emphasise this difficulty, note that a ‘leap second’ will be added at the end of 2008 due to the occasional correction needed because of the slowing of the Earth’s rotation.

⁴The **openair** package will in time contain several functions to make it easier to import data and deal with dates and times.

```
names(mydata)
## [1] "date"  "ws"    "wd"    "nox"   "no2"   "o3"    "pm10"  "so2"   "co"    "pm25"
```

To access a particular variable, one has to refer to the data frame name `mydata` and the variable itself. For example, to refer to `nox`, you must type `mydata$nox`. There are other ways to do this such as attaching the data frame and using the `with` command, but this is the basic way of doing it.

Now one can get some summary information on one variable in the data frame e.g.

```
summary(mydata$nox)
##   Min. 1st Qu. Median   Mean 3rd Qu.   Max.   NA's
##       0      82     153     179     249    1144    2423
```

Missing data can have an important effect on various R functions. For example, to find the mean NO_x concentration use the function `mean`:

```
mean(mydata$nox)
## [1] NA
```

The result `NA` is because `nox` contains some missing data. Therefore, it is necessary to exclude them using the `na.rm` command:

```
mean(mydata$nox, na.rm = TRUE)
## [1] 179
```

Often it is sensible (and easier) to remove all the missing data from a data frame:

```
newdata <- na.omit(mydata)
```

This makes a new data frame with *all* the missing data removed. However, we will work with the original data which includes missing data.

4.2 More sophisticated plotting

Let's have a look at a histogram of NO₂ concentrations, shown in Figure 4.1. Looks good, but it can easily be tidied up. It needs a new title and x-axis caption. This is easy using some of the in-built options in the `hist` function (Figure 4.2).

Or maybe you want to shade the bars (Figure 4.3)...

Another very useful plotting function is a *density* plot using the `density` function (Figure 4.4). This has the advantage over the histogram of avoiding trying to select a bin width and for some data can give a much clearer indication of its distribution.

4.3 Plotting time series with different averaging times

This is one of the most useful things to do with air pollution data. It can be a pain in some software to plot time series data with different averaging times — try for example plotting daily or weekly means in Excel. Now that we have a proper date format in R, we can do all sorts of things. First we show the basic plot of hourly data using a line plot. This plots all the data in its native format (hours), shown in Figure 4.5.

missing data
are represented
as `NA` (not
available) in R

```
hist(mydata$no2)
```

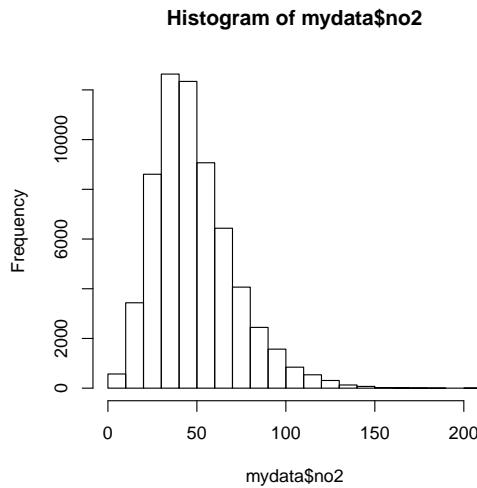


FIGURE 4.1 Histogram of NO₂ concentrations at Marylebone Road.

```
hist(mydata$no2, main = "Histogram of nitrogen dioxide",
     xlab = "Nitrogen dioxide (ppb)")
```

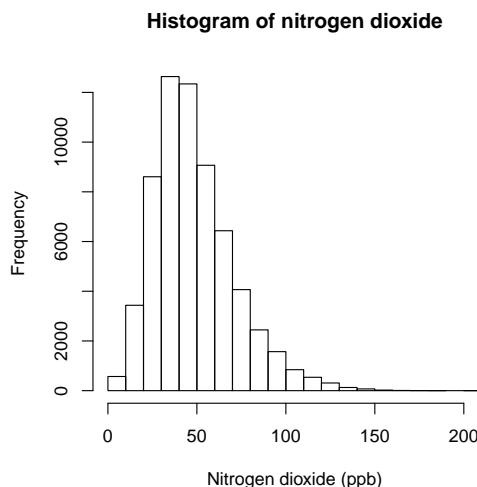


FIGURE 4.2 Histogram of NO₂ concentrations at Marylebone Road — with better labels.

Use functions in `openair` to flexibly aggregate data on different time bases

The `openair` package has functions to calculate and plot data on almost any time-averaging basis. See §(31.4) and §(17) for more details.

Say you just want to plot a section of the data — say the first 500 points. How do you do this? Selecting subsets of data is another of the real strengths of R. Remember that the variables are in vector form. Therefore, to select only the first 500 values of a variable `x` you can type `x[1:500]`, or values from the 300th to 400th data points `x[300:400]`. In the code below we choose the first 500 of date values *and* NO_x. The result is shown in

```
hist(mydata$no2, main = "Histogram of nitrogen dioxide",
     xlab = "Nitrogen dioxide (ppb)", col = "lightblue")
```

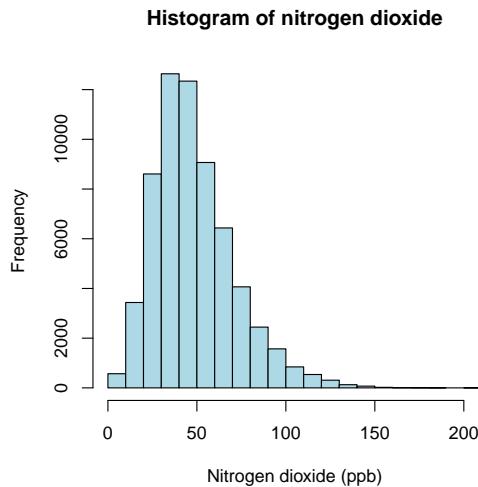


FIGURE 4.3 Histogram of NO₂ concentrations at Marylebone Road – with better labels and some colour.

```
dens <- density(mydata$no2, na.rm = TRUE)
plot(dens, main = "Density plot of nitrogen dioxide",
     xlab = "Nitrogen dioxide (ppb)")
```

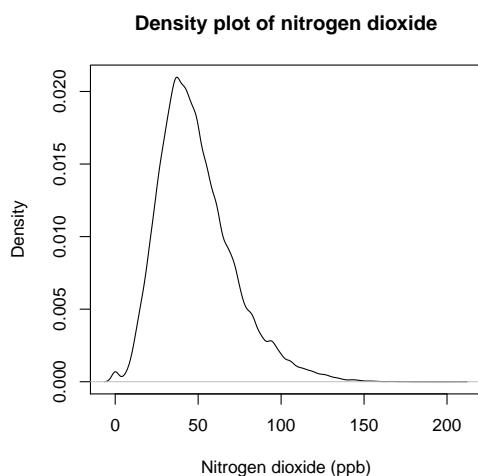


FIGURE 4.4 Density plot of NO₂ concentrations at Marylebone Road.

```
plot(mydata$date, mydata$nox, type = "l", xlab = "year",
     ylab = "Nitrogen oxides (ppb)")
```

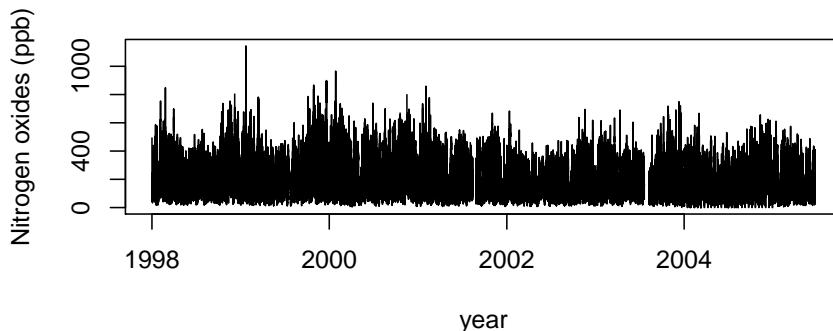


FIGURE 4.5 Hourly time series plot of NO_x at Marylebone Road.

```
plot(mydata$date[1:500], mydata$nox[1:500], type = "l", xlab = "date",
     ylab = "Nitrogen oxides (ppb)")
```

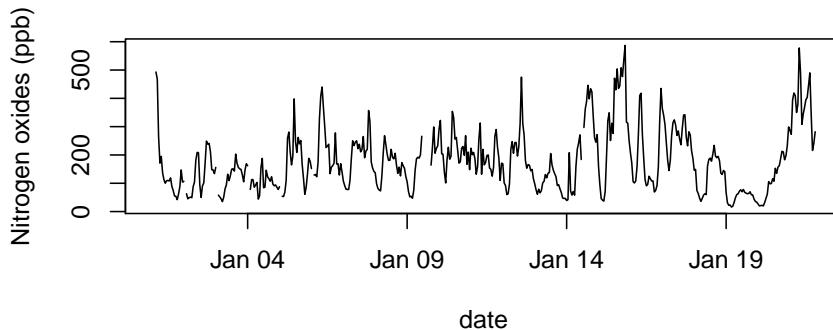


FIGURE 4.6 Hourly time series plot of NO_x at Marylebone Road — first 500 records.

Figure 4.6. Note also that R automatically adjusts the x-axis labels to day of the month, rather than year as before. Also note some missing data before 4 January. This is the advantage of keeping all the data and not removing missing values. If we had removed the values, the gap in data would not have been shown and there would have been a discontinuity. For more information on selecting parts of a data frame and selecting by date, see (§5.2).

openair has a flexible functions called **timeAverage** that makes aggregating data like this much easier so it can be done in a single step for almost any averaging period

TABLE 4.1 Some commonly used date-time formats useful when averaging data.

Code	function
%Y	annual means
%m	monthly means
%Y-%m	monthly averages for whole time series
%Y-%j	daily averages for whole time series
%Y-%W	weekly averages for whole time series
%w-%H	day of week — hour of day

```
plot(as.factor(format(mydata$date, "%m")), mydata$o3)
```

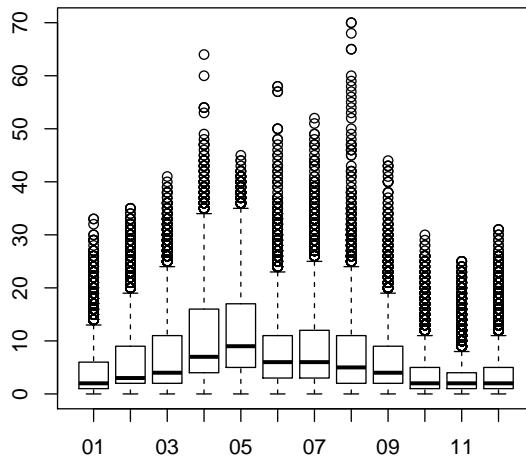


FIGURE 4.7 Monthly box and whisker plot of O₃ at Marylebone Road.

(see Section 31.4).

4.4 Summarising time series data

There are some nice ways of quickly summarising data over different time scales including day of year, month of year and day of week. The way pollutant concentrations vary over different time scales can provide some useful clues as to what sources are important. Again, this is the sort of thing that can be tricky in other software. One of the most useful; plots available is called a *box and whisker* plot, which is a very effective way of summarising large amounts of data. Say for example, we were interested to see how ozone concentrations vary by month of the year as shown in Figure 4.7.

The results shown in Figure 4.7 show several interesting features. The dark lines shows the median concentration, which peaks in May. However, the highest hourly concentrations are observed in August — presumably due to regional-scale photochemical pollution episodes.

It is also possible to view an entire series of data as monthly means, as shown in Figure 4.8. This plot shows quite nicely that the median and the peak concentrations of NO₂ increased in 2003, which is now known to be due to increased emissions of primary NO₂. Other useful summary functions include day of the week (%A) and year (%Y). There are many possibilities for plotting here and it is suggested you try some of your own.

4.5 Relationships between variables

Exploring how variables are related to one another is a very useful thing to do, but often tricky when you have lots of variables and lots of data. At a basic level, a scatter plot of one variable against another is useful. However, when you have more than a few variables it becomes quite an effort to manually plot one against another. Luckily, R has some excellent facilities to help you out; in particular the **pairs** function. We could just use the command **pairs(mydata)** and this would plot each variable all other variables. Our data frame now has 11 columns, so this would be 121 plots! Furthermore, each variable is of length 65,533 long — so the plots would look busy to say the least.

```
plot(as.factor(format(mydata$date, "%Y-%m")), mydata$no2, col = "lightpink")
```

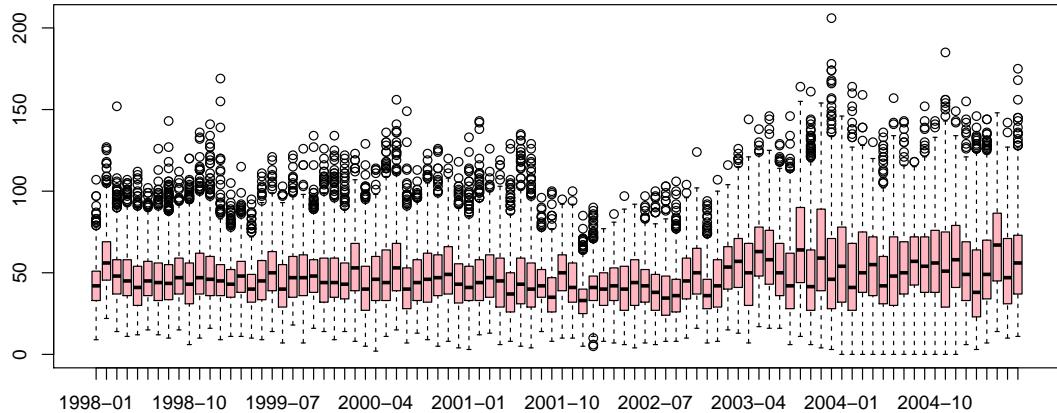


FIGURE 4.8 Yearly-monthly box and whisker plot of NO₂ at Marylebone Road.

We therefore take the opportunity to introduce a new function in R called `sample`. What we want to do is randomly select 500 lines from our data set, which is the function `sample(1:nrow(mydata), 500)`. What this does is randomly select 500 numbers from a sample a long as our data set i.e. `nrow(mydata)`. We also want to limit the columns chosen to plot. We will plot `date`, `ws`, `wd`, `nox`, `no2`; which correspond to columns 11, 2, 3, 4 and 5. You can of course choose others, or more or less than 500 points. The code below is one function but is spread over several lines for clarity. The plot is shown in Figure 4.9.

So what does this plot show? Well, if you look at the first column (date), the plots essentially show the trend in the different variables i.e. no change in wind speed, slight dip in wind direction, a decrease in NO_x in 2001 and an increase in NO₂ from around 2003. They also show that concentrations of NO_x and NO₂ do not change much with wind speed (street canyon effects) and that the highest concentrations are recorded when the wind is westerly. These plots can be extremely useful and are worth exploring in different ways.

```
pairs(mydata[sample(1:nrow(mydata), 500), c(1, 2, 3, 4, 5)],
      lower.panel = panel.smooth,
      upper.panel = NULL,
      col = "skyblue3")
```

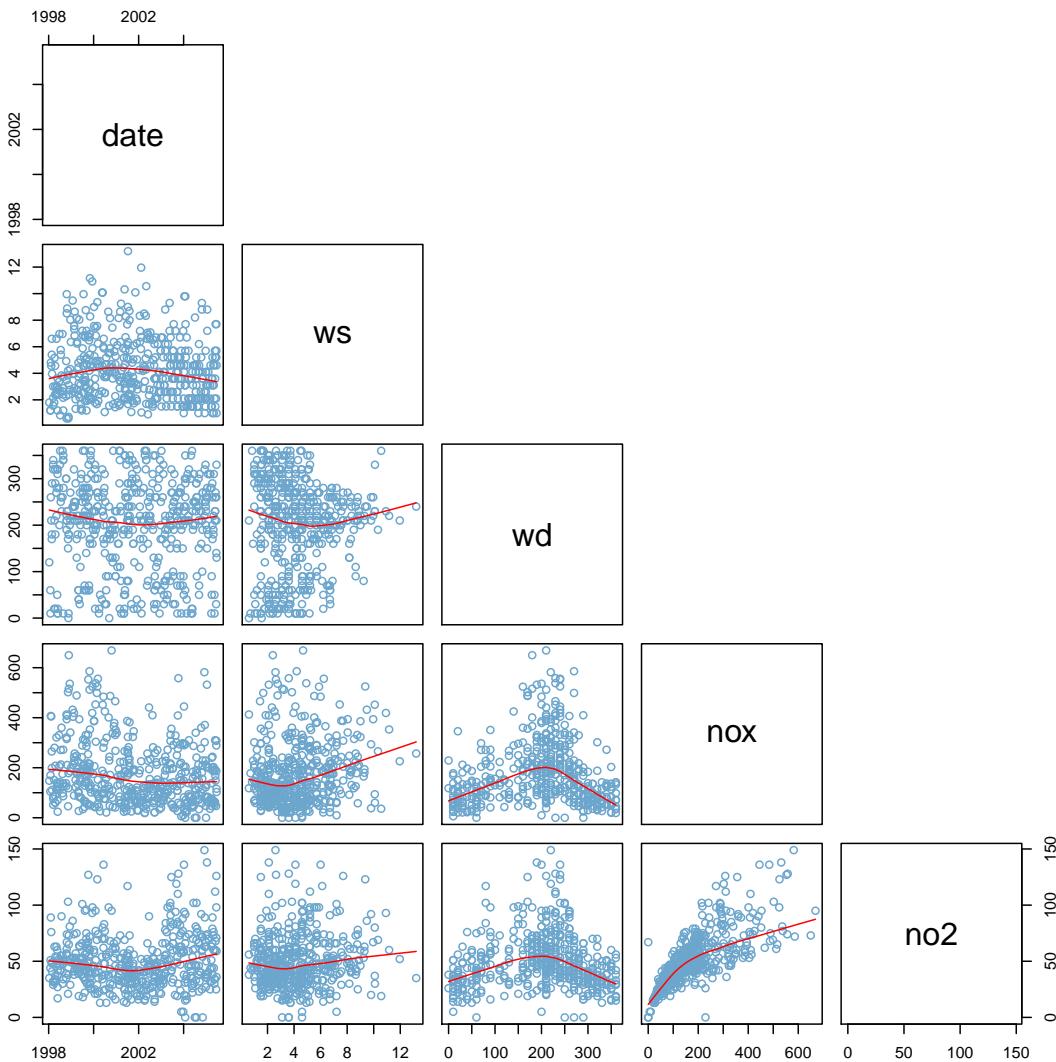


FIGURE 4.9 Pairs plot for 500 randomly selected hours of data from Marylebone Road.

5 General use of R — advice and examples

5.1 Data input and output

Getting data into and out of R is an obvious requirement. There are many options for doing so, but we only cover some of the common methods in this section.

5.1.1 Data import

Sometimes it is useful to input just a few bits of data where it is unnecessary to import from a file. This is simply done using the `c` function:

```
small <- c(1, 5, 10, 16)
small

## [1] 1 5 10 16
```

More commonly it is necessary to import data from a file. As has been discussed before it is best to keep the format of such files as simple as possible e.g. comma-delimited (.csv) or text (.txt). You can use R to list files in a particular directory. Usually a user would have set a ‘working directory’, and all file commands will relate to that (in the Windows version of R, go to File menu and choose Change Dir ...It is also possible to set working directories with a command e.g.

```
setwd("~/My Drive/openair/manual")
```

Or if you want to know what the current working directory is:

```
getwd()

## [1] "G:/My Drive/openair/manual"
```

Rather than using an external program to list files, you can do that in R too (in this case using the ‘pattern’ option to list only csv files:

```
list.files(pattern = ".csv")

## [1] "example data long.csv" "f-no2 input.csv"           "hour-day.csv"
```

Now we have a working directory set it is possible just to refer to the file name and not the path. So, to import the file ‘hour-day.csv’ we can:

```
hour.day <- read.csv("hour-day.csv", header = FALSE)
head(hour.day)

##   V1 V2 V3 V4 V5 V6 V7 V8 V9 V10 V11 V12 V13 V14 V15 V16 V17 V18 V19 V20 V21
## 1 15 27 32 17 21 11 15 32 29 21 25 32 42 34 31 44 63 71 59 34 29
## 2 31 19 25 25 21 32 38 78 109 128 105 88 86 86 94 97 NA 176 208 126 86
## 3 59 31 29 29 31 34 55 65 130 134 113 107 117 115 126 113 103 88 92 50 31
## 4 15 13 10 17 19 31 44 53 126 128 96 71 103 82 90 94 120 97 139 94 74
## 5 53 31 25 31 21 46 53 82 162 174 149 136 113 126 120 157 201 220 174 149 92
## 6 73 52 42 42 29 34 52 59 97 122 204 269 323 285 248 193 264 254 241 183 124
##   V22 V23 V24
## 1   29   25   36
## 2  111   92  113
## 3   36   21   27
## 4   67   65   59
## 5 109  134   86
## 6   97   92   55
```

In this case the data did not have a header line and we told R not to expect one. R automatically labels the fields V1, V2 ...

If the file was tab-delimited then the `read.table` function should be used (`read.csv` is a version of `read.table`).

One potentially useful capability is the direct reading of data over the internet. For example, we have example data on a web server, and this can be read directly just like a file on a computer file system:

```
test.data <- read.csv("http://www.openair-project.org/CSV/OpenAir_example_data_long.csv",
                      header = TRUE)
```

This capability might be useful for making data available for anyone to access easily.

5.1.2 Data export

Exporting data can also take a number of forms. A common issues is exporting a data frame(s) as a csv file.

```
write.csv("exportedData.csv", row.names = FALSE)
```

To save data in an ‘R format’ (.RData file), it is necessary to use the `save` function. If, for example there was a data frame called ‘test.data’, this can be save directly:

```
save(test.data, file = "testData.RData")
```

in fact several objects can be saved at one in this way:

```
save(test.data, more.data, file = "testData.RData")
```

To load this data back into R, use the `load` function:

```
load("testData.RData")
```

5.2 Selecting and replacing parts of vectors and data frames

Selecting parts of a data frame is one of the more useful things that one can learn. This often causes new users lots of difficulties, in part because of the way variables are stored in R. Earlier it was seen that in the data frame `mydata`, one could refer to all `nox` values simply as `mydata$nox`. What if one just wanted to select parts of this large data frame? Here are some examples of the sorts of things you might want to do, and in each case we read data into a new data frame called `subdata`⁵. We first consider various ways of selecting parts of vectors.

Consider the vector `x` defined as integers 1, 4, 5, 18, 22, 3, 10, 33, -2, 0.

```
x = c(1, 4, 5, 18, 22, 3, 10, 33, -2, 0)
```

To select the 4th element, we use the square brackets [] to subsample:

```
x[4]
```

```
## [1] 18
```

To select the 3rd to 6th integers:

```
x[3:6]
```

```
## [1] 5 18 22 3
```

To select everything except the first and second value, elements can be *omitted* using the – sign.

```
x[c(-1, -2)]
```

```
## [1] 5 18 22 3 10 33 -2 0
```

⁵You might not always wish to make a new data frame because it will take up extra memory – many of the examples shown can be done ‘on the fly’.

Values greater than 5:

```
x[x > 5]
## [1] 18 22 10 33
```

The indexes corresponding to integers > 5 can be found using the `which` command. This basically finds the location of numbers in a vector that meet a certain criterion. In this example, the fourth element is 18. This is a very useful function for subsetting.

```
which(x > 5)
## [1] 4 5 7 8
```

To select a specific value it is necessary to use the double = sign i.e.

```
x[x == 18]
## [1] 18
```

It is also easy to *reverse* a sequence of numbers, which is useful on many occasions:

```
rev(x)
## [1] 0 -2 33 10 3 22 18 5 4 1
```

The next thing to do is consider how to replace parts of a vector. This is something that is often necessary when preparing data. To replace the -2 by 0:

```
x[x == -2] = 0
x
## [1] 1 4 5 18 22 3 10 33 0 0
```

Note that all the usual operators such as $>$ can be used here.

Next we are going to select different parts of the data frame `mydata`. This can be more complicated because the data comprise both rows and columns. Select the first 500 rows of data. This selects rows 1 to 500 and the blank space after the comma means ‘select all columns’ i.e. all variables:

```
subdata = mydata[1:500, ]
```

One can check the number of rows selected:

```
nrow(subdata)
## [1] 500
```

Select a few variables from a data frame. Here the function `subset` is easy to use. We select just the `nox` and `no2` data. Note that when using this command, one does not need to use the \$ operator, which makes selecting a bit easier to see.

```
subdata = subset(mydata, select = c(nox, no2))
```

If one wanted to select all `nox`, `no2` and `date` values where `nox` concentrations were greater than 600 ppb:

the `subset`
function is very
useful in R

```
subdata = subset(mydata, nox > 600, select = c(nox, no2, date))
```

Selecting by date is very useful but a bit more complicated. However, once learnt it is extremely flexible and useful. We want to select all `nox`, `no2` and `date` values for 2004, although any start/end time can be used. We start by defining a start and end date, then carry out the actual selection. In this example, we must first convert our date (which is in character format) into a date/time format that R understands. Note that dates/times in R conventionally work in a hierarchical way (biggest to smallest component). Therefore ‘2004-02-03 00:00’ is the 3rd of February and not the 2nd March. In most cases dates would have been read in and converted appropriately anyway, but in this particular case we need to specify a particular date. The conversion from character string to a recognised date/time in R is done using the `as.POSIXct` function. This may seem complicated, but once learnt is both convenient and powerful. The `openair` package makes this much easier — see (§31.1) for more details.

```
start.date <- as.POSIXct("2004-01-01 00:00", tz = "GMT")
end.date <- as.POSIXct("2004-12-31 23:00", tz = "GMT")
subdata <- subset(mydata, date >= start.date & date <= end.date,
                  select = c(date, nox, no2))
```

or, using `lubridate`

```
start.date <- ymd_hm("2004-01-01 00:00")
end.date <- ymd_hm("2004-12-31 23:00")
subdata <- subset(mydata, date >= start.date & date <= end.date,
                  select = c(date, nox, no2))
```

One can easily check what this subset looks like the the functions `head` and `tail`, which give the first and last few lines of a data frame:

```
head(subdata)
## # A tibble: 6 x 3
##   date             nox   no2
##   <dttm>           <int> <int>
## 1 2004-01-01 00:00:00     98    38
## 2 2004-01-01 01:00:00    141    62
## 3 2004-01-01 02:00:00    159    56
## 4 2004-01-01 03:00:00     97    44
## 5 2004-01-01 04:00:00     60    26
## 6 2004-01-01 05:00:00     64    31
```

and,

```
tail(subdata)
## # A tibble: 6 x 3
##   date             nox   no2
##   <dttm>           <int> <int>
## 1 2004-12-31 18:00:00    114    53
## 2 2004-12-31 19:00:00    183    67
## 3 2004-12-31 20:00:00    206    69
## 4 2004-12-31 21:00:00    237    73
## 5 2004-12-31 22:00:00    232    68
## 6 2004-12-31 23:00:00    212    68
```

Another useful way of selecting subsets is using the `%in%` (or `match`) function. Some examples are given below with dates.

View the first
or last few
lines of a data
frame

selecting
different time
periods

```

subdata <- subset(mydata, format(date, "%Y") %in% 1998)

# select 1998 and 2005
subdata <- subset(mydata, format(date, "%Y") %in% c(1998, 2005))

# select weekends
subdata <- subset(mydata, format(date, "%A") %in% c("Saturday", "Sunday"))

```

This function is very useful for selecting subsets of data where there are multiple search criteria. For example, if a data frame had a field such as site name and the aim was to select data from several sites, this would be a good way to do it.

selecting
columns based
on characters
in them

It is sometimes useful to select columns (or rows) of a data frame based on their names. One extremely powerful command in R is `grep`. `grep` does character matching. It is potentially useful in numerous circumstances, but we only consider a simple case here. Say, for example we had a very large data frame with 50 column names, but we only want to extract those with the characters ‘nox’ in. We could search through and find those columns by number and refer to them in that way — but that requires a lot of manual work and has lots of potential to go wrong. In the example below we create a simple dummy data frame as an example.

```

test.dat <- data.frame(lond.nox = 1, lond.no2 = 3, nox.back = 4, no2.back = 1)
test.dat

##   lond.nox lond.no2 nox.back no2.back
## 1         1       3        4       1

```

First, for information we can print the names of the data frame:

```

names(test.dat)

## [1] "lond.nox" "lond.no2" "nox.back" "no2.back"

```

To find those names that contain the character string ‘nox’ we use `grep`:

```

grep(pattern = "nox", names(test.dat))

## [1] 1 3

```

So, columns 1 and 3 contain the character string ‘nox’. We can put this altogether and do it in one line to select those columns in the data frame that contain ‘nox’:

```

sub.dat <- test.dat[ , grep(pattern = "nox", names(test.dat))]
sub.dat

##   lond.nox nox.back
## 1         1        4

```

The `grep` command is potentially useful for selecting pollutants to plot in `openair` plots e.g. to choose any column with ‘pm’ (PM₁₀ and PM_{2.5}) in it:

```

timePlot(mydata, pollutant = names(mydata)[grep(pattern = "pm",
names(mydata))])

```

5.3 Combining and cleaning up files

So far the emphasis has been on manually importing a single .csv file to a data frame. Often with monitoring data there are numerous files all in the same format that somehow need to be read and merged. R has some very powerful and convenient ways of dealing with this situation and only the simplest case is shown here. The scenario is that you have loads of .csv files in a directory, all the same headings (although not necessarily so) and the aim is to read and combine them all. This can be done using the code below.

reading in lots
of files

```
path.files <- "D:\\temp\\\" # directory containing files
test.data <- lapply(list.files(path = path.files, pattern = ".csv"),
                    function(.file) read.csv(paste(path.files, .file, sep = ""),
                                              header = TRUE))
test.data <- do.call(rbind, test.data)
```

There are a few things to note here. In R for Windows, file paths are shown using '\\'. The function `list.files` will search for files (in this case .csv) in the D:\Temp. In the code above it is assumed a header is also present. For more refined searching see `help(list.files)`. The `lapply` function is extremely useful in R and can help avoid looping through data. In this case the function `function(.file)` is applied to the list of file names/paths supplied by `list.files`. This is a neat way of applying a function without knowing beforehand how many files there are. The traditional way of doing it would be to have a loop such as `for i = 1 to n` where `n` would be the number of files.

Note, different numbers of columns can also be dealt with using the `rbind.fill` function from the `reshape2` package as described below. In this case, the `do.call(rbind, test.data)` would be modified to `do.call(rbind.fill, test.data)`.

A common task is combining different files into one for processing. First we consider the scenario of a file with air pollution measurements and another with meteorological data. The aim is to combine them into one data frame. Rather than import data, we generate it instead. The first is a data frame called `airpol` with 1 day of data at the beginning of 2007 with pollutants NO_x and SO₂. The other is a meteorological data set, with the same dates but with wind speed and direction.

combine two
data frames

```
library(lubridate) # if not loaded already...
airpol <- data.frame(date = seq(ymd_hm("2007-01-01 00:00"),
                             by = "hours", length = 24),
                      nox = 1:24, so2 = 1:24)

met <- data.frame(date = seq(ymd_hm("2007-01-01 00:00"),
                            by = "hours", length = 24),
                   ws = rep(1, 24), wd = rep(270, 24))
```

You can check the contents of these data frames:

```

head(airpol)

##           date nox so2
## 1 2007-01-01 00:00:00 1 1
## 2 2007-01-01 01:00:00 2 2
## 3 2007-01-01 02:00:00 3 3
## 4 2007-01-01 03:00:00 4 4
## 5 2007-01-01 04:00:00 5 5
## 6 2007-01-01 05:00:00 6 6

head(met)

##           date ws  wd
## 1 2007-01-01 00:00:00 1 270
## 2 2007-01-01 01:00:00 1 270
## 3 2007-01-01 02:00:00 1 270
## 4 2007-01-01 03:00:00 1 270
## 5 2007-01-01 04:00:00 1 270
## 6 2007-01-01 05:00:00 1 270

```

To combine them, use the `merge` function:

```

test.data <- merge(airpol, met)
head(test.data)

##           date nox so2 ws  wd
## 1 2007-01-01 00:00:00 1 1 1 270
## 2 2007-01-01 01:00:00 2 2 1 270
## 3 2007-01-01 02:00:00 3 3 1 270
## 4 2007-01-01 03:00:00 4 4 1 270
## 5 2007-01-01 04:00:00 5 5 1 270
## 6 2007-01-01 05:00:00 6 6 1 270

```

When called like this the `merge` function combines data frames only where both had data. So, for example, if the `met` data frame only had the first 12 hours of 2007, merging would produce a file with only 12 hours i.e. where they match (a *natural join* in database terminology). The behaviour can be changed by selecting various options in `merge`. Following on from the previous example, the option `all` could have been set to `TRUE`, thus ensuring *all* records from each data frame would be combined—with the missing 12 hours in the `met` data frame included as `NA`. Type `help(merge)` to see the details. Functions of this type can save lots of time aligning various time series in spreadsheets.

Note, that given a data frame with multiple air pollution sites and a column called ‘site’ (i.e. values for the field ‘date’ are repeated the same number of times there are numbers of sites) it is easy to merge a *single* meteorological data set. This is the type of analysis where several air quality sites in a region are associated with a single meteorological data set. Given a data frame `aq` with multiple sites in a format like ‘date’, ‘nox’, ‘site’ and a meteorological data set `met` in the form something like ‘date’, ‘ws’, ‘wd’ then the merging is done by:

```
all.data <- merge(aq, met, by = "date", all = TRUE)
```

This code ensures that for each site for a particular date/time there are associated meteorological values. In other words, it is not necessary to think about separately joining meteorological and air quality data for each individual air quality site. See (§9) for scenarios where doing this may be useful, such as importing data for multiple AURN sites from the UK air quality archive.

Sometimes it is necessary to combine data frames that have the *same* field names. For example, data from two monitoring sites that measure the same pollutants. In the example below, we make two copies of the data frame `airpol` and name them `site1` and `site2`, respectively. Normally, of course, the data frames would contain different data, perhaps spanning different time periods. A new data frame is made using the `merge` function but with additional options set. Now, we explicitly state that we want to merge on the date field (`by = "date"`). In order to tell the NO_x and SO_2 fields apart, suffixes are used. The resulting data frame has now been merged and the NO_x from `site1` is called `nox.st1` etc.

```
site1 <- airpol
site2 <- airpol

both <- merge(site1, site2, by = "date",
              suffixes = c(".st1", ".st2"), all = TRUE)

head(both)

##           date nox.st1 so2.st1 nox.st2 so2.st2
## 1 2007-01-01 00:00:00      1      1      1      1
## 2 2007-01-01 01:00:00      2      2      2      2
## 3 2007-01-01 02:00:00      3      3      3      3
## 4 2007-01-01 03:00:00      4      4      4      4
## 5 2007-01-01 04:00:00      5      5      5      5
## 6 2007-01-01 05:00:00      6      6      6      6
```

A problem that is often encountered is combining files for different years, perhaps with different numbers of columns. We consider the slightly more difficult latter situation; although the former one is tackled in the same straightforward way. This situation can arise frequently with monitoring data. For example, in year 1, two pollutants are measured (say NO_x and NO_2), then in year 2 another pollutant is added as the monitoring is expanded. In year 2 data are available for NO_x , NO_2 and PM_{10} . This is a straightforward enough problem to deal with but can be surprisingly frustrating and time consuming to do in spreadsheets (particularly if the column order changes). However, help is at hand with the `merge` function. Given the situation mentioned, `merge` will deal with this:

```
# make some data
year1 <- data.frame(date = seq(ymd_hm("2007-01-01 00:00"),
                            by = "hours", length = 24),
                     nox = 1:24, so2 = 1:24)

year2 <- data.frame(date = seq(ymd_hm("2008-01-01 00:00"),
                            by = "hours", length = 24), nox = 1:24,
                     so2 = 1:24, pm10 = 1:24)

test.data <- merge(year1, year2, all = TRUE)
head(test.data)

##           date nox so2 pm10
## 1 2007-01-01 00:00:00    1    1   NA
## 2 2007-01-01 01:00:00    2    2   NA
## 3 2007-01-01 02:00:00    3    3   NA
## 4 2007-01-01 03:00:00    4    4   NA
## 5 2007-01-01 04:00:00    5    5   NA
## 6 2007-01-01 05:00:00    6    6   NA
```

In this example, `year1` contains hourly data for all of 2007 for NO_x and NO_2 , and

`year2` contains hourly data for all of 2008 for NO_x, NO₂ and PM₁₀. The data frame `test.data` then contains two years of data and has all variables present. For year 1 where there are no PM₁₀ data, these data are shown as missing i.e. `NA`.

Another useful application of the `merge` function is to fill in gaps due to missing data. The scenario is that you have a file (say a year long of hourly data), but some lines are missing. This sort of situation arises frequently and can be time consuming to sort out. What is needed is to ‘pad out’ the file and fill in the gaps with the missing dates and set the other fields to missing (`NA`, in R-speak). To show this, we first deliberately remove 2 of the hours from the `airpol` data frame. We then create a data frame with all the hours (note that only 24 are used here, but it can of course be any length), then the data are merged:

padding-out
missing hours

```
airpol <- airpol[-c(2, 3), ] # select everything except record 2 and 3
# create all the dates that should exist
all.dates = data.frame(date = seq(ymd_hm("2007-01-01 00:00"),
                                   by = "hours", length = 24))
# merge the two
test.data <- merge(all.dates, airpol, all = TRUE)
head(test.data)

##           date nox so2
## 1 2007-01-01 00:00:00   1   1
## 2 2007-01-01 01:00:00  NA  NA
## 3 2007-01-01 02:00:00  NA  NA
## 4 2007-01-01 03:00:00   4   4
## 5 2007-01-01 04:00:00   5   5
## 6 2007-01-01 05:00:00   6   6
```

The missing hours are thus inserted, but the variables themselves are set to missing.

Finally (and not surprisingly) a package already exists that does this for you. The `reshape2` package can manipulate data in very flexible ways. However, the function `rbind.fill` is particularly useful if you have lots of different data frames to combine because `merge` can only merge two data frames at once. Note you can download this package from CRAN.

```
library(reshape2)
test.data <- rbind.fill(year1, year2)
```

interpolate
missing data

Sometimes it is useful to *fill in missing data* rather than ignore it. Here, we show two options from the `zoo` (zero-ordered observations) package. The first function `na.locf` will fill missing data with the value of the last non-missing data point. To do a linear interpolation between points, the `na.locf` function should be used:

```

library(zoo)
# make some data with missing values
a = c(1, NA, NA, NA, 3, NA, 18, NA, NA, 20)
# show data
a

## [1] 1 NA NA NA 3 NA 18 NA NA 20

# fill with last non-missing point
na.locf(a)

## [1] 1 1 1 1 3 3 18 18 18 20

# interpolate missing points
na.approx(a)

## [1] 1.0 1.5 2.0 2.5 3.0 10.5 18.0 18.7 19.3 20.0

```

There are various other options that can be used with these functions, which can be considered by typing `help(zoo)`. These functions can also be applied to data frames (or columns of). Say we want to interpolate all the missing NO_x concentrations:

```
mydata$nox <- na.approx(mydata$nox, na.rm = FALSE)
```

Note that the `na.rm = FALSE` option ensures that trailing NAs are not removed, making the number of records the same as the original data. The code above would replace the NO_x concentrations. If preferred, a new column could be made, in this case called `nox.all`:

```
mydata$nox.all <- na.approx(mydata$nox, na.rm = FALSE)
```

Once data are imported into R — say by loading a .csv file into a data frame, there are often tasks that need to be carried out to alter the data before processing it. Some of these common tasks are considered in this section.

One of the most immediate tasks with air pollution data is to convert the date/time field into something understood by R and this has already been discussed.

Next, it may be necessary to change the name of a variable. To list the existing variables names:

```

names(mydata)

## [1] "date"      "ws"        "wd"        "nox"       "no2"       "o3"        "pm10"
## [8] "so2"       "co"        "pm25"     "nox.all"

```

You can also refer to a single column name. In the code below, we show an example of how to change one of the names (in this case `nox`) to `nitrogen.oxides`.

changing
variable names

```

names(mydata)[4] # display name of 4th column (nox)

## [1] "nox"

names(mydata)[4] = "nitrogen.oxides" # change the name
names(mydata) # show new names

## [1] "date"           "ws"            "wd"            "nitrogen.oxides"
## [5] "no2"            "o3"            "pm10"          "so2"
## [9] "co"             "pm25"          "nox.all"        ""

# change it back again
names(mydata)[4] = "nox"

```

To change more than one name at a time (say the 4th and 5th column names):

```
names(mydata)[c(4, 5)] = c("new1", "new2")
```

If you have imported data that has lots of upper case names and you want them all in lower (because they are easier to refer to), use the **tolower** function e.g.

```

names <- c("NOx", "PM10") #make some upper case names
tolower(names)

## [1] "nox"   "pm10"

```

If you import data that has rather verbose descriptions, which become a pain to refer to, you can abbreviate them. In this example, we have two site names and the aim is to abbreviate them using only two letters.

```

names = c("North Kensington", "Marylebone Road")
abbreviate(names, 2)

## North Kensington Marylebone Road
##           "NK"           "MR"

```

There is potential for this going wrong, if, for example two of the names were very similar:

```

names <- c("North Kensington roadside", "North Kensington background",
               "Marylebone Road")
abbreviate(names, 2)

##   North Kensington roadside North Kensington background
##                   "NKr"           "NKh"
##           Marylebone Road
##                   "MR"

```

However, R is clever enough to work this out, and uses an extra letter as necessary. The **abbreviate** function can be very effective at simplifying files for processing and generally makes logical simplifications. Note that in the examples above, one could have chosen to abbreviate the names to any length.

Data can easily be ordered and this might be necessary if for example, the date field was not sequential in time. An example is:

```
mydata <- mydata[order(mydata$date), ]
```

which keeps data in the same data frame **mydata** but ensures that *all* the data are

ordered by date.

5.4 Reshaping data

Data are stored in a wide variety of ways and it is often necessary to do some data manipulation in order to analyse or plot data. This section distinguishes between two main storage options: stacked or column format (narrow or wide). By way of an example, consider the simple case of two sites each measuring NO_x. One way of storing all this data in a single data frame would be to have columns: ‘date’, ‘site1.nox’, ‘site2.nox’. An alternative would be to stack the data and have columns ‘date’, ‘nox’, ‘site’. For such a simple example there isn’t much difference between the two options. But what if there were 10, 20 or 100 sites? Having columns ‘site1.nox’, ‘site2.nox’ ...would get rather tedious, whereas the stacked data would still only have three columns.

For **openair** functions there is a *big* advantage in stacking data like this, and all the **openair** import functions do this. This is because it then becomes easy to plot any number of quantities **without referring to them individually and without knowing how many there are**. This will become clearer as **openair** functions are used, but imagine trying to plot NO_x at 10 sites using the two different approaches using the **openair timePlot** function:

For column format:

```
timePlot(mydata, pollutant = c("site1.nox", "site2.nox", "site3.nox", ...,
                                "site10.nox"))
```

And stacked data:

```
timePlot(mydata, pollutant = "nox", type = "site")
```

The latter example works for any number of sites without having to know the number.

This section provides some examples of using two recent packages for manipulating typical air quality data. Use is made of two packages: ‘tidyverse’ for cleaning and reorganising data and ‘dplyr’ for carrying out various manipulations.

The main interest here is ‘reshaping’ data, which is often necessary to aid analysis, plotting etc. and is a very common task in data analysis.

Load the packages we need:

```
library(openair)
library(tidyverse)
library(dplyr)
```

Single site with several pollutants

In this case we have a typical situation where we have a file where each column is a variable and a column representing the site name (or site code) and date. This is the way data from the **openair importAURN** data are stored.

```
single <- importAURN(site = "kc1", year = 2014)
head(single)

##           date code          site      o3     no2     co
## 1 2014-01-01 00:00:00 KC1 London N. Kensington 68.85165 11.27738 0.188600
## 2 2014-01-01 01:00:00 KC1 London N. Kensington 67.20520 14.70330 0.213049
## 3 2014-01-01 02:00:00 KC1 London N. Kensington 60.86885 17.34159 0.218870
## 4 2014-01-01 03:00:00 KC1 London N. Kensington 56.02928 19.58830 0.218870
## 5 2014-01-01 04:00:00 KC1 London N. Kensington 58.12476 17.98180 0.204317
## 6 2014-01-01 05:00:00 KC1 London N. Kensington 62.36563 12.89790 0.221780
##       so2 pm10    nox   pm2.5 nv2.5 v2.5 nv10 v10   ws   wd
## 1 0.85592 14.5 12.12716 0.55422  9.1   5.7  3.1   8.6 5.4 6.7 213.6
## 2 0.54371 13.1 15.76601 0.69308  7.7   2.0  5.4   9.8 2.9 6.6 216.6
## 3 0.69982 10.2 18.19122 0.55411  3.1   1.4  1.5   6.6 3.3 6.3 220.0
## 4 0.58340 14.9 20.43793 0.55411  2.3   2.0  0.2   9.5 4.9 6.0 219.2
## 5 0.29137 16.5 18.61867 0.41535  5.6   4.0  1.4  13.3 2.7 6.0 216.5
## 6 0.64128 18.7 13.37555 0.31151  7.5   4.7  2.5  15.2 2.9 6.2 209.1
```

The first question is, how do we go from this ‘wide’ format to a ‘narrow’ format i.e. have a column for the pollutant name and stack pollutants on top of one another. This is easy enough with the `tidyR gather` function.

```
narrow <- gather(single, key = pollutant, value = value, o3:wd)
head(narrow)

##           date code          site pollutant      value
## 1 2014-01-01 00:00:00 KC1 London N. Kensington      o3 68.85165
## 2 2014-01-01 01:00:00 KC1 London N. Kensington      o3 67.20520
## 3 2014-01-01 02:00:00 KC1 London N. Kensington      o3 60.86885
## 4 2014-01-01 03:00:00 KC1 London N. Kensington      o3 56.02928
## 5 2014-01-01 04:00:00 KC1 London N. Kensington      o3 58.12476
## 6 2014-01-01 05:00:00 KC1 London N. Kensington      o3 62.36563
```

When used like this the `gather` function first takes the name of the data frame to process (`single`), the ‘key’ is the name you want to give to the gathered columns (‘pollutant’ in this case), ‘value’ is the name for the new column of values of the gathered variables and finally the ‘o3:wd’ means gather the columns from ‘o3’ to ‘wd’ (basically all the variables). The other way of selecting the variable columns is to say what you *don’t want*. In this case rather than ‘o3:wd’ it could be `-c(date, site, code)`.

What if you want to go the other way and put the data frame back to columns? This time you must provide a ‘key’ (which you can think of as the column of variables that will be separated into individual columns) and the associated value of those columns, which is the ‘value’ column in this case.

```
wide <- spread(narrow, key = pollutant, value = value)
```

Multiple sites with multiple pollutants

A more complicated example is if we import two sites, which will stack sites on top of one another. Have a look at:

```
single <- importAURN(site = c("kc1", "my1"), year = 2014)
```

This is the same format as before. But what if we wanted all the measurement data by site as columns — so for example, Marylebone Road NO_x, North Kensington NO_x etc? For this, we need a key that combines the site name and pollutant.

```
narrow <- gather(single, key = pollutant, value = value, o3:wd)
```

To make things a bit easier we will drop the site name and use the site code to identify the site.

```
narrow <- select(narrow, -site)
```

To make a column that combines site code and pollutant we use the ‘unite’ function.

```
narrow <- unite(narrow, site_pol, c(code, pollutant))
```

Now we can spread the data by site and pollutant:

```
wide <- spread(narrow, key = site_pol, value = value)
```

Now we have columns:

```
names(wide)
```

```
## [1] "date"      "KC1_co"     "KC1_no"     "KC1_no2"    "KC1_nox"    "KC1_nv10"
## [7] "KC1_nv2.5" "KC1_o3"     "KC1_pm10"   "KC1_pm2.5" "KC1_so2"    "KC1_v10"
## [13] "KC1_v2.5"  "KC1_wd"     "KC1_ws"     "MY1_co"     "MY1_no"     "MY1_no2"
## [19] "MY1_nox"   "MY1_nv10"   "MY1_nv2.5"  "MY1_o3"    "MY1_pm10"   "MY1_pm2.5"
## [25] "MY1_so2"   "MY1_v10"    "MY1_v2.5"   "MY1_wd"    "MY1_ws"
```

These functions are very useful for getting data into the right shape for analysis.

5.5 Example: converting hour-day data to column format

In this example we deal with a common problem in manipulating data – reformatting data in one format to another. Often data are stored as rows representing days and columns representing hours. This is often a format used by AEA for AURN data. However, this example highlights a more general requirement to reformat data.

The aim is to convert the 24×365 ‘matrix’ of data to a single column of data. An example file is provided called ‘hour-day.csv’. To make it simple the file only contains hourly data with no column headings (hour of day) or row names (days). Fortunately there are some built-in functions available in R that make the reformatting of these data easy:

```
nox <- read.csv("~/My Drive/openair/Data/hour-day.csv",
                    header = FALSE)
nox <- as.data.frame(t(nox))
nox <- stack(nox)
nox <- nox$values
```

In the code above, the data are first read in — note the option `header = FALSE` in this case. Next, the data are transposed using the `t` (transpose) function, which produces a matrix of data and transposes the rows/columns. Transposing the data ensures that the hours are now in columns. Note that this operation instead of representing hours in rows, puts them into 365 columns. The columns can now be stacked on top of each other. We convert the matrix back to a data frame with the `as.data.frame` function. Next we use the `stack` function that literally *stacks* columns of data, working from column 1 to column 365. Finally, we extract the values resulting from applying the `stack` function. This sort of data manipulation is straightforward in R but would be much trickier in Excel. The code can actually be written in two lines, but becomes less easy to understand:

```
nox <- read.csv("~/My Drive/openair/Data/hour-day.csv",
                    header = FALSE)
nox <- stack(as.data.frame(t(nox)))[, "values"]
```

Even if you don't understand this code, this example should provide sufficient information on how to apply it.

In fact, in this case, there is an easier way to do this:

```
nox <- read.csv("~/My Drive/openair/Data/hour-day.csv",
                    header = FALSE)
nox <- t(nox)
nox <- as.vector(nox)
```

In this code, the data are transformed as before, producing a matrix, and the matrix is converted to a vector. When converting a matrix to a vector, it works on columns rather than by rows. The **stack** function is therefore a better choice if we were interested in groups of data for further processing, as in the example below:

```
test.data = data.frame(grp1 = 1:3, grp2 = 10:12, grp3 = 20:22)
test.data

##   grp1 grp2 grp3
## 1     1    10    20
## 2     2    11    21
## 3     3    12    22

stacked = stack(test.data)
stacked

##   values ind
## 1      1 grp1
## 2      2 grp1
## 3      3 grp1
## 4     10 grp2
## 5     11 grp2
## 6     12 grp2
## 7     20 grp3
## 8     21 grp3
## 9     22 grp3
```

This then makes it much easier to work with the different groups e.g. calculate means, or plotting the data.

5.6 Daily means from hourly means — processing wind direction data

Sometimes it is necessary or useful to calculate daily means from hourly data. Many particle measurements, for example, are measured as daily means and not hourly means. If we want to analyse such particle data for example, by considering how it varies with meteorological data, it is necessary to express the meteorological (and maybe other data) as daily means. It is of course straightforward to calculate daily means of concentrations and wind speeds, as shown elsewhere in this document. However, this is not the case for wind directions. For example the average of 10° and 350° is 0° (or 360°) and not 180° .

The way to deal with this is to average with u and v wind components. A function has been written to do this:

```

dailymean <- function(mydata) {
  ## for wind direction, calculate the components
  mydata$u = sin(2 * pi * mydata$wd / 360)
  mydata$v = cos(2 * pi * mydata$wd / 360)
  dailymet = aggregate(mydata, list(Date = as.Date(mydata$date)), mean,
                        na.rm = TRUE)
  ## mean wd
  dailymet = within(dailymet, wd <- atan2(u, v) * 360 / 2 / pi)
  ## correct for negative wind directions
  ids = which(dailymet$wd < 0) # ids where wd < 0
  dailymet$wd[ids] = dailymet$wd[ids] + 360
  dailymet = subset(dailymet, select = c(-u, -v, -date))
  dailymet
}

```

In this function a data frame is supplied containing hourly data and the returned data frame contains correctly formatted daily data. Note that very similar functions can be used to calculate means over other time periods e.g. hourly means from 15-minute data. The code below shows the use of this function.

```

mydaily = dailymean(mydata) # calculate daily means
# show top of data frame
head(mydaily)

##           Date   ws   wd  nox  no2   o3  pm10  so2    co  pm25 nox.all
## 1 1998-01-01  6.8 190 174  39 6.9   18 3.2 2.7   NaN   174
## 2 1998-01-02  7.1 226 130  39 6.5   28 3.9 1.8   NaN   130
## 3 1998-01-03 11.0 221 120  38 8.4   20 3.2 1.7   NaN   120
## 4 1998-01-04 11.5 219 106  35 9.6   21 3.0 1.6   NaN   106
## 5 1998-01-05  6.6 238 171  46 5.0   24 4.5 2.1   NaN   171
## 6 1998-01-06  4.4 196 211  45 1.3   35 5.7 2.5   NaN   211

```

5.7 Using an Editor

5.7.1 Using the built-in editor

As you begin to use R, you will quickly realise there are more efficient ways to do things other than just typing in commands. Often, you will want access to a series of commands for say, plotting a graph with defaults of your choosing. What is needed is an *Editor*. R has an in-built editor — just select File|New script... and the editor window will open. This is a bit like Notepad in Windows. While working it can be useful to put together lines of code in the editor, select the code, right-click and run it.

Advice

As you increase the amount of coding you do, it becomes increasingly difficult to remember what the code actually does. It is always a good idea to liberally comment your code with lines starting with a `#`. This is especially important if you intend making your code available to others. See [Figure 5.1](#) for an example of how commented lines are coloured differently in a dedicated editor, making it easy to distinguish between the code and code comments.

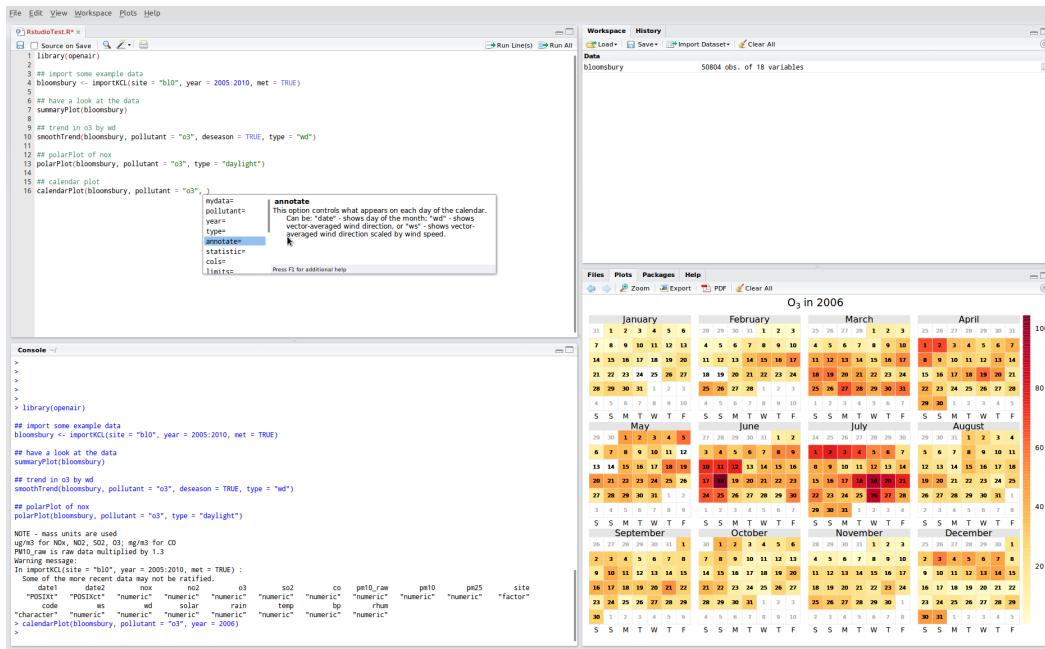


FIGURE 5.1 RStudio is one of the best R editors around.

5.7.2 Using a dedicated editor

The built-in editor is useful for small amounts of work. However, with use you will find a ‘dedicated’ editor easier to use. We recommend something called *RStudio*, a screen shot is shown in Figure 5.1. At the time of writing RStudio is a beta version of the software, but is already very good — and likely to get much better in time. More details can be found at <http://rstudio.org/>. There are also a lot of well-written ‘knowledge base’ articles, see <http://support.rstudio.org/help/kb>. Below are a few features that makes RStudio useful for working with **openair** and other R projects.

- It works on Windows, linux and Apple Mac.
- It has been developed by people that clearly use R.
- In the top left pane of Figure 5.1 is where you can work on your R script e.g. use it to develop and save a series of analyses for a report. Note that if you type the name of a function (or part of) R Studio will offer *completions* if you press TAB. This feature also works in the R console, shown in the bottom left pane.
- It is easy to send a selection or line from the script to the R console by selecting ‘Run Line(s)’.
- In the top right pane you can view the objects in your workspace. If you click on one you can view all or most of an object as appropriate. The ‘history’ tab gives a summary of all the commands you have input, which you can search through.
- At the bottom right there are several tabs. In Figure 5.1 the plot tab is shown, which shows the most recent plot made using the console. The ‘Packages’ tab provides a neat way of loading a package — just select the one you want and it will load it and all dependent packages.

5.8 Several plots on one page

Often it is useful or necessary to plot more than one plot on a page. This is the sort of task that can be fiddly to carry out if the plots are produced separately and then need to be combined. Problems include alignment and sizing, which as a minimum can be frustrating to get right. R makes it easy to plot any number of plots on a page in a neat and consistent way. The key is to use the `par` function to set up the page as you want it *before* you plot your graphs. The `par` function can control and fine-tune a vast number of plot options — see `help(par)` for specific information.⁶

To set up a page to plot several plots in a regular grid, the esoterically named `mfrow` or `mfcol` option is set. For example, to plot two graphs side-by-side one types in:

```
par(mfrow = c(1, 2))
```

This sets up the plot window for 1 row and 2 columns. And `par(mfrow = c(2, 3))` therefore would allow for six plots in 2 rows and 3 columns etc.

Therefore, the code below produces two plots, side-by-side of NO_x and NO_2 , of monthly mean concentrations as shown in Figure 5.2.

Note, however that all `openair` functions use `lattice` graphics, where a slightly different approach is required. See Section 8.7 for more details on how to plot several `openair` plots on one page.

⁶Note that this will only work with base graphics and not `lattice`.

```

par(mfrow = c(1, 2))
# first plot
means <- tapply(mydata$nox, format(mydata$date, "%Y-%m"), mean, na.rm = TRUE)
dates <- seq(mydata$date[1], mydata$date[nrow(mydata)], length = nrow(means))
plot(dates, means,
      type = "l",
      col = "darkgreen",
      xlab = "year",
      ylab = "nitrogen oxides (ppb)")
means <- tapply(mydata$no2, format(mydata$date, "%Y-%m"), mean, na.rm = TRUE)
dates <- seq(mydata$date[1], mydata$date[nrow(mydata)], length = nrow(means))

# second plot
plot(dates, means,
      type = "l",
      col = "skyblue3",
      xlab = "year",
      ylab = "nitrogen dioxide (ppb)")

```

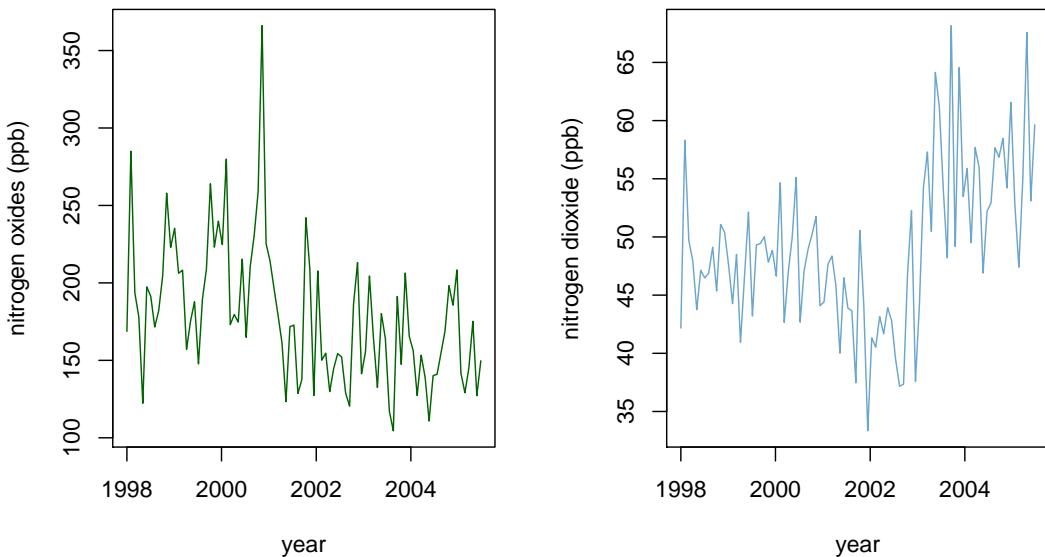


FIGURE 5.2 Plotting two plots side-by-side using the `par` setting.

5.9 Saving and using plots

There are several ways of saving and using plots and some general advice is given in this section. The simplest way to use a plot from the R GUI is to size to the required size, right-click on it and choose to copy as a metafile or a bitmap; this being on Windows. A *metafile* is a Windows-specific format, whereas a *bitmap* is a more general system independent format. An important difference between the two formats is that a metafile is in a *vector* graphics format whereas a bitmap is a *raster* format made up of a regular grid of pixels. Vector-based graphics are made up of lines and shapes that remain sharp at any resolution. Zoom in on a metafile plot and it will remain sharp. By contrast zooming in on a bitmap image will reveal fuzziness — and eventually the individual pixels. Other vector formats include pdf and encapsulated postscript (eps) and other bitmap images include png (portable network graphics) and jpeg (the format most often used for photographs).

Most plots in **openair** work best using vector graphics. This manual for example uses LaTeX and the plots are in pdf format i.e. they are sharp when you zoom in.

However, most users of **openair** probably use applications such as Microsoft Word. Word does accept a wide range of graphic formats, but unfortunately even recent versions cannot use pdf graphics. In this case it will generally be best to copy plots as metafiles. Unfortunately, Windows metafiles in R cannot handle alpha transparency. What this means is for certain plots that use semi-transparent shading (e.g. `timeVariation`, Section 22), copying the plot as a metafile will not work for those semi-transparent areas in the plot. The other problem with metafiles is that a Word document containing many of them will be slow to render and frustrating to work with, as well as being potentially very large in file size.

So what's best? A good compromise is the png format. These images capture **openair** graphics well, are generally small file sizes and can be used by almost all applications. They also capture alpha transparency. It is possible to save as a png format using the R GUI and RStudio from the file menu, or the export button, respectively. By default the png will be saved at 72 dots per inch, which might not be sufficiently detailed, especially if it is resized later. A better way to work (in my view) is to always use scripts to make the plots. This has two main advantages. First, all analyses *including plots* will entirely and easily reproducible, which is a major advantage. Second, it is possible to have more control over how the plots are saved.

The way to save a plot using code is as follows.

```
png("myPlot.png")
polarPlot(mydata, pollutant = "so2")
dev.off()
```

The code above first prepares a *device* for printing on (in this case a png device) and save *myPlot.png* in your working directory. Second, a plot is produced and third (importantly!) the device is closed. Note that if you fail to add the final line the device will remain open and the file will not be written. By default, `png` will save a file 480×480 pixels; see `?png` for details. This resolution might be fine for simple graphics, but a much better approach is to be sure of the final resolution. Most commonly it is recommended that bitmap graphics in printed documents are 300 dpi (dots per inch). So, an approach that guarantees this resolution is as follows. Think about the size of the plot needed in inches (in this case 5 inches wide by 4.5 inches high) and use this in the call to `png`:

```
png("myPlot.png", width = 5 * 300, height = 4.5 * 300, res = 300)
polarPlot(mydata, pollutant = "so2")
dev.off()
```

This will produce a sharp plot which at 5×4.5 inches will be 300 dpi. In some cases e.g. in quality publications, it might be necessary to save at an even higher resolution. For example to save at 600 dpi just replace the 300 above with 600.

Users often ask how to change font size in `openair` plots, which can be controlled through the `fontsize` option. However, it should also be noted that saving plots in the way described above can usefully control the relative size of the font in a plot. This method also has the advantage of making all outputs reproducible. To see the effect try comparing the output of the plot above (5×4.5 inches) with the one below (3.5×3 inches). You will see that the latter has font that is relatively larger compared with the rest of the plot.

```
png("myPlot.png", width = 3.5 * 300,
    height = 3 * 300, res = 300)
polarPlot(mydata, pollutant = "so2")
dev.off()
```

Another point worth mentioning is that in some cases it is useful to overlay plots on a map. In these cases it is useful to set the background to be transparent rather than white (the default). To do this just supply the option `bg = "transparent"`.

5.10 Graphing lots of data — using level plots

It is often the case when plotting monitoring data that there are so many data points it gets hard to see relationships. Consider the scatter plot for NO_x and NO_2 shown in Figure 5.3 — it is very difficult to see how NO_x and NO_2 are related because of the large number of points. An alternative way of plotting such data is to ‘bin’ it first and count the number of points in each bin and plot it as a *level* plot.

The first part of the code in above will produce a basic plot using the base graphics function `image`, shown in Figure 5.3. Now it is possible to see the relationship between NO_x and NO_2 much more clearly. It also has the benefit of showing where most of the data are, which is not very apparent in Figure 5.3.

A better looking plot can be produced with a bit more work using *lattice* graphics, as shown in Figure 5.4. This plot has the advantage of also showing a scale, which in this case in the number of points in each bin. Lattice graphic are considerd more in (§6).

5.11 Special symbols for use in plotting air pollution data

Air pollution concentrations are expressed in many ways, perhaps most commonly in $\mu\text{g m}^{-3}$. It is always preferable to display these units properly rather than, for example as ug/m³. The same is also true for subscripts in pollutant names such as NO_x and $\text{PM}_{2.5}$. R has its own way of dealing with specialist symbols, which is similar to L^AT_EX. This section provides code for commonly used expressions. For more information type `help(plotmath)`. We illustrate the use of these symbols through examples shown in Table 5.1, by setting the y-axis label. However, these labels can be used elsewhere too, such as in titles or to annotate specific parts of a plot.

To demonstrate what these symbols look like, Figure 5.5 provides an example. The code is shown below.

```

x <- mydata$nox
y <- mydata$no2
# find maximum values
x.max <- max(x, na.rm = TRUE)
y.max <- max(y, na.rm = TRUE)
# set the bin interval
x.int <- 5
y.int <- 2
# bin the data
x.bin <- cut(x, seq(0, x.max, x.int))
y.bin <- cut(y, seq(0, y.max, y.int))
# make a frequency table
freq <- table(x.bin, y.bin)
# define x and y intervals for plotting
x.range <- seq(0 + x.int/2, x.max - x.int/2, x.int)
y.range <- seq(0 + y.int/2, y.max - y.int/2, y.int)
# plot the data
image(x = x.range, y = y.range, freq, col = rev(heat.colors(20)))

```

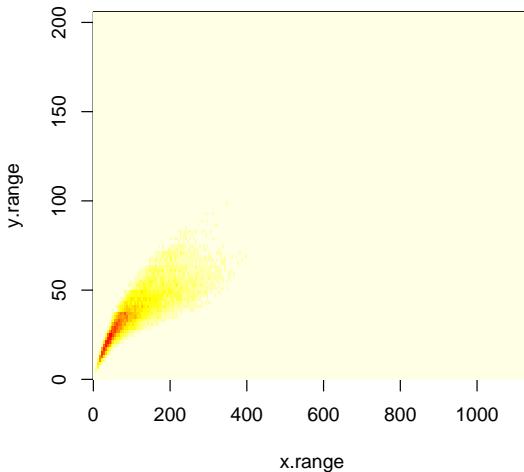


FIGURE 5.3 Using the `image` function to plot NO_x against NO₂.

TABLE 5.1 Examples of commonly used text formats for air pollution.

Text required	Expression
NO _x	<code>ylab = expression("NO"[X])</code>
PM _{2.5}	<code>ylab = expression("PM"[2.5])</code>
($\mu\text{g m}^{-3}$)	<code>ylab = expression("(" * mu * "g m" ^-3 * ")")</code>
PM ₁₀ ($\mu\text{g m}^{-3}$)	<code>ylab = expression("PM"[10] * " (" * mu * "g m" ^-3 * ")")</code>
Temperature (°C)	<code>xlab = expression("Temperature (" * degree * "C)")</code>

```
library(lattice)
grid <- expand.grid(x = x.range, y = y.range)
z <- as.vector(freq)
grid <- cbind(grid, z)
levelplot(z ~ x * y, grid, col.regions = rev(heat.colors(20)))
```

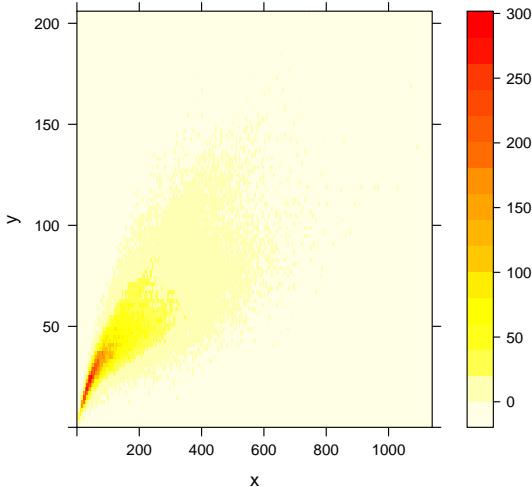


FIGURE 5.4 Using the lattice graphics `levelplot` function to plot NO_x against NO_2 .

```
plot(1, 1,
      xlab = expression("Temperature (" * degree * "C"),
      ylab = expression("PM[10]~(" * mu * "g m"^-3 * ")"),
      main = expression("PM[2.5] * " and NO[x] * " at Marylebone Road"),
      type = "n")
text(1, 1, expression(bar(x) == sum(frac(x[i], n), i==1, n)))
text(.8, .8, expression(paste(frac(1, sigma*sqrt(2*pi)), " ",
plain(e)^{frac(-(x-mu)^2, 2*sigma^2)}))), cex = 1.2)
text(1.2, 1.2, expression(hat(beta) == (X^t * X)^{-1} * X^t * y))
```

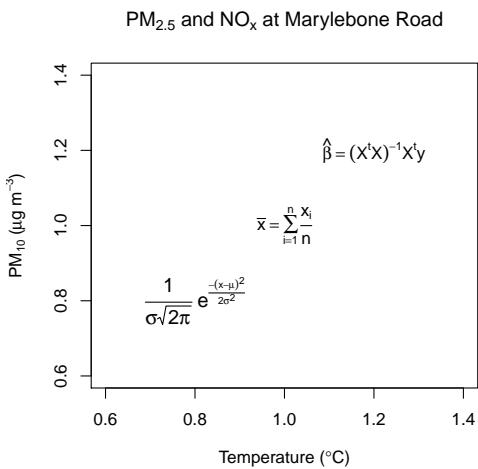


FIGURE 5.5 Examples of different symbols that can be used in R plots.

5.12 Using databases with R

So far the discussion has focussed on working with .csv files, which might be adequate for most purposes. However, as the amount of data increases, the storage of it in this way is not to be recommended. A much better approach is to store data in a database. There are several advantages in doing so. First, it forces a more disciplined approach to storage (e.g. variables names and formats). Second, it is possible to store a lot more information in this way. Finally, for very large amounts of information R can run out of memory because all the calculations are done in RAM.⁷ In the latter case it can be much better to use the SQL database language to do some of the work first and then bring smaller data sets into R.

There are several database types that R can work with; perhaps the most common being Microsoft Access. For those interested in open-source databases MySQL is highly recommended (we use this with many large data sets). However, a discussion of MySQL goes beyond the aims of this document. It should also be noted that you do not actually need to have Microsoft Access to read or write data to it.

The **RODBC** package will automatically try and preserve data formats, including those for date/time. We have experienced a few difficulties here to do with British Summertime and GMT. **RODBC** will bring data in in a format consistent with what the operating system is set to, which can be either BST or GMT (or other time zones). The best option in our view is that before the data are imported, set the system to GMT as above. This will avoid all sorts of potential problems.

6 Multivariate plots — introduction to the Lattice package

6.1 Introduction to the Lattice package

In (§2.1) one of the benefits highlighted in using R was the extensive number of *packages* available that extend the core features of R. One package called **lattice** is particularly useful for plotting and analysing monitoring data. The **lattice** package is based on the original S (S-Plus) *Trellis* package that provides excellent multivariate plotting capabilities.⁸ This is one of the stronger capabilities that R has and greatly enhances the possibilities for plotting monitoring data. The original trellis graphics were designed by William Cleveland at Bell laboratories and were based on research into how best to visualise graphics. For those interested in this there are a couple of books available (Cleveland 1985; Cleveland 1993). Lattice graphics can be used to produce similar plots to those shown elsewhere in this document, but in some cases can produce better looking plots with better-chosen defaults. However, the real strength of Lattice is the ability to deal with multivariate data and to plot several plots on one page.

Installing and loading a package in R

The capabilities of R are greatly enhanced by a number of optional *packages*. To use different packages, they must first be installed. Many packages such as **lattice** are installed as part of the R installation itself. However, they need to be loaded to use them. This can be done in two principal ways: use the menu and choose Packages | Load package... and choose from the available packages listed; or in code you can issue a command **library(package name)**.

⁷Now that there is a 64-bit version of R for Windows, this is less of a problem and the limitation is more to do with the amount of RAM the computer has.

⁸The names trellis and lattice are meant to reflect the idea of multiple plots i.e. like a garden trellis.

```
xyplot(nox ~ date, data = mydata, type = "l")
```

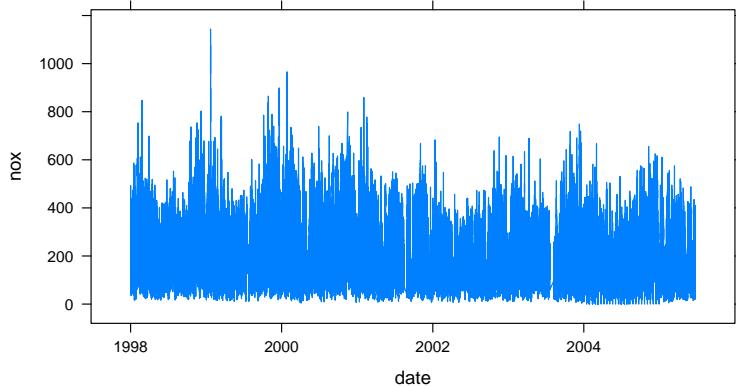


FIGURE 6.1 Example plot using the `lattice` package to plot hourly concentrations of NO_x at Marylebone Road

In many cases the package you want may not be installed on your system. In this case you can choose Packages | Install package(s)..., where you are then prompted for a location to install from (choose one in the UK). It is possible that this option will not work due to firewalls etc. information is being downloaded from a remote server. An alternative way of doing this is to go the main R web pages, select CRAN (Comprehensive R Archive Network), choose the appropriate web site address (again UK) and under the heading ‘Software’ choose ‘packages’. Choose the package you want (for Windows, choose Windows binary zip file). Download the file to your hard disk and then choose Packages | Install package(s) from local zip files....

The drawback of using the `lattice` package is that most new users (and some experienced ones too) find it difficult to use. The emphasis is very much on the use of code to make plots. The focus of this section therefore is to provide some examples of the use of Lattice graphics on our data set that can be readily applied to other data sets with minimum or no alteration.

6.2 Example simple plots

We start with plotting the basic time series of NO_x as shown in Figure 4.5. The code is shown below. For basic plotting like this, the terminology is straightforward. One of the first things to note is the use of a formula to represent the plot `nox ~ date`. You can think of this as an equation for plotting data i.e. $y = f(x)$. Therefore, what appears on the y-axis is given first (in this case `nox`), then the x-axis data (in this case `date`). Also given is the argument `data = mydata`, and the type of plot `type = "l"` as before.

The plot generated is shown in Figure 6.1, which can be compared with Figure 4.5. There are a few differences to note: the default colour is blue, the y-axis labels are horizontal (for easier reading) and their tick marks are shown on all four sides (outside the plot so that they do not clutter up the data actually shown). The lattice plot can be annotated in much the same way as the base plots, with options such as `ylim`, `ylab` etc.

Much of the power of lattice graphics lies in the ability to plot one variable against another *dependent* on a third. To get an idea of what is meant here, consider how NO_2

```
# weekday names
weekdays <- c("Monday", "Tuesday", "Wednesday", "Thursday", "Friday",
           "Saturday", "Sunday")

# make data frame
test.data = data.frame(nox = 5 * runif(70), no2 = runif(70),
                       weekday = as.factor(rep(weekdays, each = 10)))

# order the days i.e. do not want the default alphabetical
test.data$weekday = ordered(test.data$weekday, levels = weekdays)

# plot
xyplot(no2 ~ nox | weekday, data = test.data, as.table = TRUE)
```

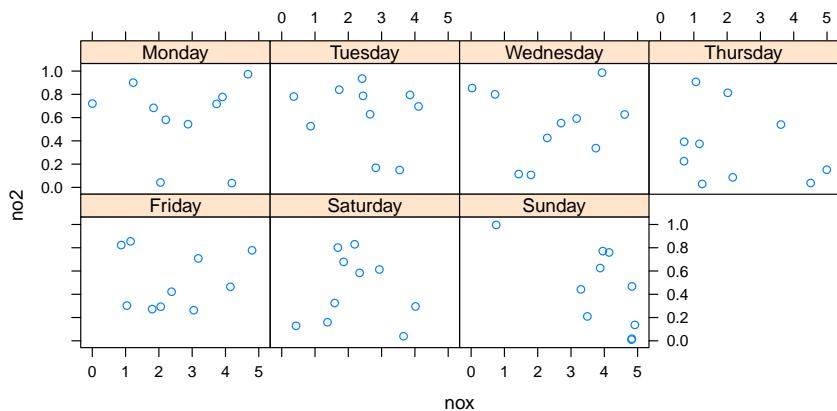


FIGURE 6.2 Example plot using the `lattice` package to plot two variables against each other (`nox` and `no2`), but dependent on a third (`weekday`).

varies by NO_x by day of the week. Now, day of the week is a *categorical* variable, which in R is referred to as a *factor*. We illustrate the use of this type of plotting by making some simple artificial data. In the code below we first define the days of the week. We then make a data frame where the NO_x concentrations are 70 random numbers between 0 to 5 (`nox = 5 * runif(70)`) i.e. 10 for each day of the week. The NO_2 concentrations are similarly assumed to be 70 random numbers, which are between 0 to 1 in this case and the days of the week `weekday` are each of the days repeated 10 times. When lattice plots a factor, it does so alphabetically. However, this makes little sense for the days of the week and therefore we force the ordering of the days with the code shown. Finally, a plot is produced. Note that in the plot command the formula `no2 ~ nox | weekday` is used. This means in simple terms ‘plot no2 against nox, but show how it varies by day of the week’. For some reason, lattice always fills plots from the bottom left corner to the top right. The `as.table = TRUE` command forces it to plot it from the top left to the bottom right, which more most applications seems like a more logical way of plotting. The result of the plotting is shown in Figure 6.2.

7 Functions in R

This section highlights the importance of functions in R to carry out specific tasks. Most of the rest of Part II of this report considers the use of dedicated functions written to analyse air pollution data. This section gives an overview of why functions are important and how they work. Functions are useful in many ways:

- For repetitive tasks they help take the effort out of processing data. An example would include a function to import and clean-up data.
- Functions provide a much more structured way of working. They help to break down big problems into smaller bits that are easier to work with.
- For air pollution analysis, dedicated functions can (and have) been written to offer unique analysis capabilities that are not offered in any other software. R offers excellent capabilities here for creating new analyses and plots.

Let's consider a very simple function that adds two numbers (although one would not actually ever need to write such a simple function!):

```
add.two <- function(a = 1, b = 2)
{a + b}
```

The function name is **add.two**, it accepts two arguments **a** and **b**. The body of the function is written within the braces **{}**. This function can be read into R simply by pasting it in — and it then becomes available for use. Let's use the function to add 10 and 5:

```
add.two(10, 5)
## [1] 15

add.two(c(1, 2, 3), 2)
## [1] 3 4 5

add.two(c(1, 2, 3), c(5, 6, 7))
## [1] 6 8 10
```

Easy! Note that in the definition of the function itself, we provided default values for **a** and **b**. If one called the function without supplying **a** or **b** it would return 3:

```
add.two()
## [1] 3
```

This functionality is useful for testing purposes, but also for providing defaults for some of the variables. If the user does not want to alter a default value, there is no need to supply it.



Part II

Dedicated functions for analysing air pollution data

8 Introduction

Part II of this report focuses on the development and use of dedicated functions written to process air pollution data. These functions greatly extend the capabilities outlined in Part I, where the focus was on developing an understanding of R.

While many of the options in these functions allow quite a sophisticated analysis to be undertaken, the defaults generally use the simplest (and fastest) assumptions. A more detailed analysis can refine these assumptions e.g. by accounting for autocorrelation, or fine-control over the appearance of a plot.

It should be noted that while the aim is to keep this documentation up to date, the primary source of information related to the different functions is contained within the package itself. Once loaded, type `?openair` to see all the help pages associated with the package. The website for `openair` is <http://www.openair-project.org>.

This section contains important information on loading the `openair` package for the first time and the input data requirements. Users will need to consider the advice in this section to ensure that `openair` can be used without problems.

8.1 Installing and loading the `openair` package

The dedicated functions for the analysis of air pollution data have been made available in the `openair` package. As of September 2010, `openair` is available on CRAN. This means it should be very straightforward to install. In Windows, choose the packages menu in R and then choose ‘Install package(s)’. You will be prompted for a location from which to download — so scroll down to the appropriate country. Once selected you will then be shown a large list of available packages — choose `openair`.

Second approach to installing `openair`

The second approach to installing `openair` is slightly more involved but should still be easy. This situation arises if you are using a computer that does not let R communicate externally. These are the steps:

1. Download the **Windows binary** (zip) files from <http://cran.r-project.org/web/packages/> for the following packages and store them somewhere convenient on your computer: ‘`openair`’, ‘`reshape2`’, ‘`plyr`’ and ‘`RColorBrewer`’. This is done by clicking on the link for each package, which will show a page with a downloads section where the zip file is shown.
2. Install `openair` and all the other dependent packages by choosing ‘Packages’ from the R menu, then ‘Install packages(s) from local zip files...’, and choose all the zip files that were downloaded.

The package can be tested by going into the ‘Packages’ menu in R and choosing ‘Load package’ and then choose `openair`. The package comes with some test data — several years of data from the Marylebone Road site in London (‘`mydata`’). Test it by trying a function e.g.

```
summaryPlot(mydata)
```

Note that it is only necessary to install packages once — unless a package has been updated or a new version of R is installed. Occasionally it is useful to update the packages that have been installed through the ‘Update packages’ option under the Packages menu. Because the **openair** package (and R itself) are continually updated, it will be useful to know this document was produced using R version 3.6.0 and **openair** version 2.6-5.

8.2 Access to source code

All R code is accessible. On CRAN, you will see there are various versions of packages: Package source, MacOS X binary and Windows binary. The source code is contained in the Package source, which is a tar.gz (compressed file).

For **openair** all development is carried out using Github for version control. Users can access all code used in openair at <https://github.com/davidcarslaw/openair>.

8.3 Brief introduction to **openair** functions

This section gives a brief overview of the functions in **openair**. The core functions are summarised in [Table 8.3](#), which shows the input variables required, the main purpose of the function, whether multiple pollutants can be considered and a summary of the **type** option. The **type** option given in [Table 8.3](#) gives the maximum number of conditioning variables allowed in each function — more on this later.

Having read some data into a data frame it is then straightforward to run any function. Almost all functions are run as:

```
functionName(thedata, options, ...)
```

The usage is best illustrated through a specific example, in this case the **polarPlot** function. The details of the function are shown in [Section 15](#) and through the help pages (type `?polarPlot`). As it can be seen there are a large number of options associated with **polarPlot** — and most other functions and each of these has a default. For example, the default pollutant considered in **polarPlot** is ‘nox’. If the user has a data frame called **theData** then **polarPlot** could minimally be called by:

```
polarPlot(theData)
```

which would plot a ‘nox’ polar plot if ‘nox’ was available in the data frame **theData**.

Note that the options do not need to be specified in order nor is it always necessary to write the whole word. For example, it is possible to write:

```
polarPlot(theData, type = "year", poll = "so2")
```

In this case writing **poll** is sufficient enough to uniquely identify that the option is **pollutant**.

Also there are many common options available in functions that are not explicitly documented, but are part of lattice graphics. Some of the common ones are summarised in [Table 8.1](#). The **layout** option allows the user to control the layout of multi-panel plots e.g. `layout = c(4, 1)` would ensure a four-panel plot is 4 columns by 1 row.

TABLE 8.1 Common options used in **openair** plots that can be set by the user but are generally not explicitly documented.

option	description
xlab	x-axis label
ylab	y-axis label
main	title of the plot
pch	plotting symbol used for points
cex	size of symbol plotted
lty	line type used
lwd	line width used
layout	the plot layout e.g. c(2, 2)

Controlling font size

All **openair** plot functions have an option **fontsize**. Users can easily vary the size of the font for each plot e.g.

```
polarPlot(mydata, fontsize = 20)
```

The font size will be reset to the default sizes once the plot is complete. Finer control of individual font sizes is currently not easily possible. See also [Section 5.9](#), which also describes a way of controlling font size when saving graphics.

The **openair** ‘type’ option

One of the central themes in **openair** is the idea of *conditioning* plots. Rather than plot x against y, considerably more information can usually be gained by considering a third variable, z. In this case, x is plotted against y for many different intervals of z. This idea can be further extended. For example, a trend of NO_x against time can be *conditioned* in many ways: NO_x vs. time split by wind sector, day of the week, wind speed, temperature, hour of the day...and so on. This type of analysis is rarely carried out when analysing air pollution data, in part because it is time consuming to do. However, thanks to the capabilities of R and packages such as **lattice**, it becomes easier to work in this way.

In most **openair** functions conditioning is controlled using the **type** option. **type** can be any other variable available in a data frame (numeric, character or factor). A simple example of **type** would be a variable representing a ‘rbefore’ and ‘after’ situation, say a variable called **period** i.e. the option **type = "period"** is supplied. In this case a plot or analysis would be separately shown for ‘before’ and ‘after’. When **type** is a numeric variable then the data will be split into four *quantiles* and labelled accordingly. Note however the user can set the quantile intervals to other values using the option **n.levels**. For example, the user could choose to plot a variable by different levels of temperature. If **n.levels = 3** then the data could be split by ‘low’, ‘medium’ and ‘high’ temperatures, and so on. Some variables are treated in a special way. For example if **type = "wd"** then the data are split into 8 wind sectors (N, NE, E, ...) and plots are organised by points of the compass.

There are a series of pre-defined values that **type** can take related to the temporal components of the data. To use these there *must* be a **date** field so that the can be calculated. These pre-defined values of **type** are shown in [Table 8.2](#) are both useful and convenient. Given a data frame containing several years of data it is easy to analyse the data e.g. plot it, by year by supplying the option **type = "year"**. Other useful and

straightforward values are "hour" and "month". When `type = "season"` **openair** will split the data by the four seasons (winter = Dec/Jan/Feb etc.). Note for southern hemisphere users that the option `hemisphere = "southern"` can be given. When `type = "daylight"` is used the data are split between nighttime and daylight hours. In this case the user can also supply the options `latitude` and `longitude` for their location (the default is London).

TABLE 8.2 Pre-defined time-based values for the **openair** `type` option.

option	splits data by ...
"year"	year
"hour"	hour of the day (0 to 23)
"month"	Month of the year
"season"	spring, summer, autumn, winter
"weekday"	Monday, Tuesday, ...
"weekend"	Saturday, Sunday, weekday
"monthlyear"	every month-year combination
"gmtbst"	separately considers daylight saving time periods
"daylight"	nighttime and daylight

Table 8.3 summarises the functions that accept the option 'type' and the number of types that can be set. Numerous examples of conditioning are given throughout this document.

TABLE 8.3 Summary of main **openair** analysis functions. Click on function name to be taken to the section on that function.

Function	Mandatory variables	Purpose	Multiple pollutants	type option
calcFno2	see §(30) for details	estimate primary NO ₂ emissions ratio from monitoring data	no	no
calendarPlot	date, one numeric field	Calendar-type view of mean values	no	no
conditionalEval	observed and modelled values and other variables(s)	extensions to conditionalQuantile	no	yes [1]
conditionalQuantile	observed and modelled values	quantile comparisons for model evaluation	no	yes [2]
kernelExceed	date, ws, wd, one other numeric field	bivariate kernel density estimates for exceedance statistics	no	Yes [1]
linearRelation	date, two numeric fields	explore linear relationships between variables in time	no	limited
TheilSen	date, one numeric field	Calculate Theil-Sen slope estimates and uncertainties	no	Yes [2]
modStats	observed and modelled values	calculate a range of model evaluation statistics	no	yes [≥ 1]
percentileRose	wd, one other numeric field	percentiles by wind direction	no	Yes [2]
polarAnnulus	date, ws, wd, one other numeric field	polar annulus plot for temporal variations by wind direction	yes	Yes [2]
polarCluster	ws, wd, one other numeric field	cluster analysis of bi-variate polar plots for feature extraction	No	No
polarFreq	ws, wd	alternative to wind rose/pollution rose	no	Yes [2]
polarPlot	ws, wd, one other numeric field	bi-variate polar plot	yes	Yes [2]
pollutionRose	ws, wd, one other numeric field	pollution rose	no	Yes [2]
scatterPlot	x and y values to plot	traditional scatter plots with enhanced options	no	Yes [2]
smoothTrend	date, one numeric field	smooth trend estimates	yes	Yes [2]
summaryPlot	date, one numeric field	summary view of a data frame	yes	no
TaylorDiagram	two numeric fields	model evaluation plot	no	Yes [2]
timePlot	date, one numeric field	Time-series plotting	yes	Yes [1]
timeProp	date, one numeric, one category field	Time-series plotting with categories as stacked bar chart	yes	Yes [1]
timeVariation	date, one numeric field	diurnal, day of week and monthly variations	yes	Yes [1]
trajCluster	data from importTraj	HYSPLIT back trajectory cluster analysis	no	Yes [2]
trajPlot	data from importTraj	HYSPLIT back trajectory plots — points of lines	no	Yes [2]
trajLevel	data from importTraj	HYSPLIT back trajectory plots — binned or smoothed	no	Yes [2]
trendLevel	date, one other numeric field	flexible level plots or ‘heat maps’	no	Yes [2]
windRose	date, ws, wd	traditional wind rose	no	Yes [2]

TABLE 8.4 Summary of **openair** utility functions. Click on function name to be taken to the section on that function.

Function	Mandatory variables	Purpose	Multiple pollutants	type option
calcPercentile	date, one numeric variable	calculate percentiles for numeric variables in a data frame	NA	NA
corPlot	a data frame	correlation matrix with conditioning	yes	yes [1]
cutData	a data frame	partition data into groups for conditioning plots and analysis	yes	yes [≥ 1]
importADMS	an ADMS output file e.g. .pst, .met, .mop, .bgd	import outputs from the ADMS suite of dispersion models (McHugh et al. 1997)	NA	NA
importAURN	site code and year	import hourly data from the UK air quality data archive (http://www.airquality.co.uk/data_and_statistics.php)	NA	NA
importKCL	site code and year	import hourly data from the London Air data archive (http://www.londonair.org.uk/LondonAir/Default.aspx)	NA	NA
importTraj	site code and year	import HYSPLIT back trajectory data from KCL servers	NA	NA
quickText	a string	properly format common pollutant names and units	NA	NA
selectByDate	date and one other variable	flexibly select date periods by year, day of week etc.	NA	NA
selectRunning	date and one other variable	select contiguous periods of a certain run-length above a specified threshold	NA	NA
splitByDate	date and one other variable	partition and label a data frame by time periods	NA	NA
timeAverage	date, one numeric variable	calculate statistics over flexible time periods account for data capture rates etc.	NA	NA

8.4 Input data requirements

The **openair** package applies certain constraints on input data requirements. **It is important to adhere to these requirements to ensure that data are correctly formatted for use in **openair**.** The principal reason for insisting on specific input data format is that there will be less that can go wrong and it is easier to write code for a more limited set of conditions.

The **openair** package requires as an input a data frame, which generally consists of hourly date/time, pollution and meteorological data. As shown elsewhere in this document, the recommended way of inputting data into R is through reading a .csv file. This in itself avoids potential issues with ‘awkward’ file formats e.g. with varying header lines. Of course, anyone familiar with R will know how to do this and may choose to import their data from a range of sources such as databases. A few important requirements and advice are given below.

Use **openair** functions to help import data!

There are several functions in (§9) that make the process of importing data into **openair** much simpler. Where possible, these functions should be used. (§9) also contains some useful functions for manipulating data.

1. Data should be in a ‘rectangular’ format i.e. columns of data with a header on the first line. The file ‘example data long.csv’ provides a template for the format and users should refer to that file if in doubt. The best approach is to use the **import** function that is part of **openair**, described in (§9).
2. Where fields should have numeric data e.g. concentrations of NO_x, then the user should ensure that no other characters are present in the column, accept maybe something that represents missing data e.g. ‘no data’. Even here, it is essential to tell the **import** function how missing data are represented; see (§9).
3. The date/time field should be given the heading **date** — note the lower case. No other name is acceptable.
4. The wind speed and wind direction should be named **ws** and **wd**, respectively (note again, lower case). There is an implicit assumption that wind speed data are in units of m s⁻¹. Most functions have been written assuming reasonable ranges in wind speed in m s⁻¹. However, the functions will work if the units were in knots, for example and several functions allow the user to annotate the plots with the correct units. Wind directions follow the UK Met Office format and are represented as degrees from north e.g. 90° is east. North is taken to be 360°.
5. Other variables names can be upper/lower case *but should not start with a number*. If column names do have white spaces, R will automatically replace them with a full-stop. While ‘PM2.5’ as a field name is perfectly acceptable, it is a pain to type it in—better just to use ‘pm25’ (**openair** will recognise pollutant names like this and automatically format them as PM_{2.5} in plots).

Note if users wish to assume non-zero wind speeds to be calm e.g. any wind speed below 0.5 m s⁻¹, then these can be set directly e.g.

```
mydata$ws[mydata$ws < 0.5] <- 0
```

It should be mentioned again that any reasonably large amount of data should be kept in a database and not Excel sheets and the like. Much less will go wrong if this is the case; see §(5.12) for some information on Access databases.

8.4.1 Dealing with more than one site

In many situations users will have more than one site available and most **openair** functions can deal with this situation. However, it does require that the data are in a certain format. If the data are available via the AURN archive or via the KCL LAQN then it is possible to use the **importAURN** or **importKCL** functions to select multiple sites at once and the data will be correctly formatted for use by the functions.

If it is not possible to import the data in this way, it is necessary to format the data in such a way that can be used. The format is very similar to that described above for several pollutants at a single site. With more than one site it is necessary to have another column (with name **site**) with the site name in. Data are therefore ‘stacked’.

Sometimes data will not be in this format and site data will be in separate columns. (§31.7.2) shows the approach that can be used to format such data.

If users need help with formatting their data, please contact us for advice.

8.5 Using colours

Type **colors()**
or **colours()**
into R to see
full list of
named colours

Many of the functions described require that colour scales are used; particularly for plots showing surfaces. It is only necessary to consider using other colours if the user does not wish to use the default scheme, shown at the top of Figure 8.1. The choice of colours does seem to be a vexing issue as well as something that depends on what one is trying to show in the first place. For this reason, the colour schemes used in **openair** are very flexible: if you don’t like them, you can change them easily. R itself can handle colours in many sophisticated ways; see for example the **RColorBrewer** package.

Several pre-defined colour schemes are available to make it easy to plot data. In fact, for most situations the default colour schemes should be adequate. The choice of colours can easily be set; either by using one of the pre-defined schemes or through a user-defined scheme. More details can be found in the **openair openColours** function. Some of the defined colours are shown in Figure 8.1, together with an example of a user defined scale that provides a smooth transition from yellow to blue. The code that produced this plot is shown for Figure 8.1:⁹

The user-defined scheme is very flexible and the following provides examples of its use. In the examples shown next, the **polarPlot** function is used as a demonstration of their use.

```
# use default colours - no need to specify
polarPlot(mydata)

# use pre-defined "jet" colours
polarPlot(mydata, cols = "jet")

# define own colours going from yellow to green
polarPlot(mydata, cols = c("yellow", "green"))

# define own colours going from red to white to blue
polarPlot(mydata, cols = c("red", "white", "blue"))
```

⁹This is given for interest, the user does not need to know this to use the colours.

```

library(openair)
## small function for plotting
printCols <- function(col, y) {
  rect((0:200) / 200, y, (1:201) / 200, y + 0.1, col = openColours(col, n = 201),
    border = NA)
  text(0.5, y + 0.15, deparse(substitute(col)))
}

## plot an empty plot
plot(1, xlim = c(0, 1), ylim = c(0, 1.6), type = "n", xlab = "", ylab = "",
  axes = FALSE)
printCols("default", 0)
printCols("increment", 0.2)
printCols("heat", 0.4)
printCols("jet", 0.6)
printCols("hue", 0.8)
printCols("brewer1", 1.0)
printCols("greyscale", 1.2)
printCols(c("tomato", "white", "forestgreen" ), 1.4)

```



FIGURE 8.1 Pre-defined colour scales in **openair**. The top colour scheme is a user-defined one.

8.6 Automatic text formatting

openair will increasingly try to automate the process of annotating plots. It can be time consuming (and tricky) to repetitively type in text to represent $\mu\text{g m}^{-3}$ or PM_{10} ($\mu\text{g m}^{-3}$) etc. in R. For this reason, an attempt is made to automatically detect strings such as 'nox' or 'NOx' and format them correctly. Where a user needs a y-axis label such as NO_x ($\mu\text{g m}^{-3}$) it will only be necessary to type `ylab = "nox (ug/m3)"`. The same is also true for plot titles.

Over time we will add to the number of text strings that could be automatically formatted. It is suggested that users get in touch if they have a specific request that is not yet covered. Most functions have an option called `auto.text` that is set to TRUE by default. Users can override this option by setting it to FALSE.

Note that there will be occasions when the user will want to format the text themselves, as shown by the examples in Table 5.1. In this case the option `auto.text = FALSE` should be set when using a function and the user should supply their own expression.

8.7 Multiple plots on a page

We often get asked how to combine multiple plots on one page. Recent changes to **openair** makes this a bit easier. Note that because **openair** uses **lattice** graphics the base graphics **par** settings will not work.

It is possible to arrange plots based on a column \times row layout. Let's put two plots side by side (2 columns, 1 row). First it is necessary to assign the plots to a variable:

```
a <- windRose(mydata)
b <- polarPlot(mydata)
```

Now we can plot them using the split option:

```
print(a, split = c(1, 1, 2, 1))
print(b, split = c(2, 1, 2, 1), newpage = FALSE)
```

In the code above for the ‘split’ option, the last two numbers give the overall layout (2, 1) — 2 columns, 1 row. The first two numbers give the column/row index for that particular plot. The last two numbers remain constant across the series of plots being plotted.

There is one difficulty with plots that already contain sub-plots such as **timeVariation** where it is necessary to identify the particular plot of interest (see the **timeVariation** help for details). However, say we want a polar plot (**b** above) and a diurnal plot:

```
c <- timeVariation(mydata)
print(b, split = c(1, 1, 2, 1))
print(c, split = c(2, 1, 2, 1), subset = "hour", newpage = FALSE)
```

For more control it is possible to use the **position** argument. **position** is a vector of 4 numbers, $c(xmin, ymin, xmax, ymax)$ that give the lower-left and upper-right corners of a rectangle in which the plot is to be positioned. The coordinate system for this rectangle is [0–1] in both the x and y directions.

As an example, consider plotting the first plot in the lower left quadrant and the second plot in the upper right quadrant:

```
print(a, position = c(0, 0, 0.5, 0.5), more = TRUE)
print(b, position = c(0.5, 0.5, 1, 1))
```

The position argument gives more fine control over the plot location.

8.8 Annotating **openair** plots

A frequently asked question about **openair** and requested feature is how to annotate plots. While all **openair** functions could have options to allow annotations to be made, this would make the functions cumbersome and reduce flexibility. Nevertheless it is useful to be able to annotate plots in lots of different ways. Fortunately there are existing functions in packages such as **lattice** and **latticeExtra** that allow for plots to be updated. An example of the sorts of annotation that are possible is shown in Figure 8.2, which is an enhanced version of Figure 17.1. These annotations have been subsequently added to Figure 17.1 and built up in layers. This section considers how to annotate **openair** plots more generally and uses Figure 8.2 as an example of the types of annotation possible. Also considered specifically is the annotation of plots that are in polar coordinates, as these can sometimes benefit from different types of annotation.

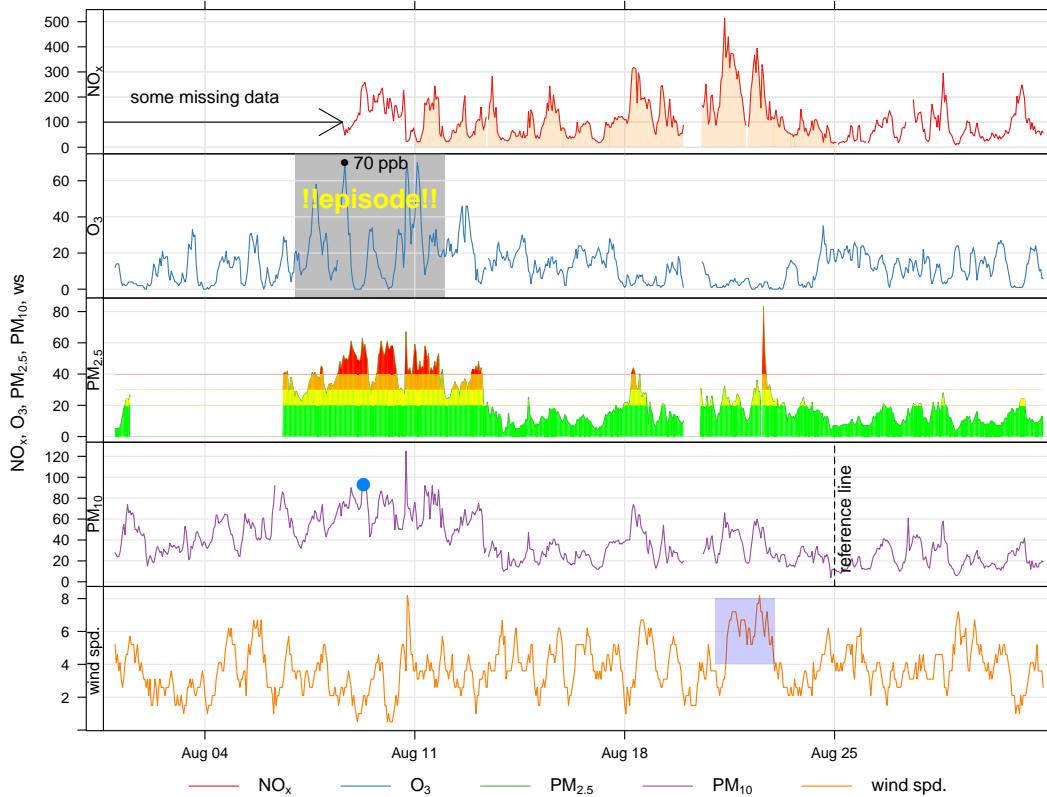


FIGURE 8.2 Examples of different ways of annotating a plot in **openair**.

There are several different types of objects that can be useful to add to plots including text, shapes, lines and other shading. Given that many **openair** plots can consist of multiple panels, it is also useful to think about how to annotate specific panels. The examples given in this section will apply to all **openair** plot, the only difference being the coordinate system used in each case.

The basis of **openair** annotations is through the use of the **latticeExtra** package, which should already be installed as part of **openair**. In that package there is a function called **layer** that effectively allows annotations to be built up ‘layer by layer’.

Adding text

To add text (or other annotations) it is necessary to know the coordinates on a plot for where the text will go, which will depend on the data plotted. In this extended example using the **timePlot** function, the y-axis will be in ordinary numerical units, whereas the x-axis will be in a date-time format (**POSIXct**).

There are various ways that annotations can be added, but the method used here is to add to the previous plot using a function called **trellis.last.object()** to which we want to add a later. This may seem complicated, but once a few examples are considered, the method becomes very powerful, flexible and straightforward. In a multi-panel plot such as Figure 8.2 it is also useful to specify which rows/columns should be added to. If they are not specified then the annotation will appear in all panels.

First, a plot should be produced to which we wish to add some text.

```
## make sure latticeExtra is loaded
library(latticeExtra)
timePlot(selectByDate(mydata, year = 2003, month = "aug"),
         pollutant = c("nox", "o3", "pm25", "pm10", "ws"))
```

So, considering Figure 8.2, this is how the text ‘some missing data’ was added to the top panel.

```
trellis.last.object() +
  layer(ltext(x = ymd("2003-08-04"), y = 200,
              labels = "some missing data"), rows = 1)
```

So what does this do? First, the `trellis.last.object()` is simply the last plot that was plotted. Next the `layer` function is used to add some text. The text itself is added using the `ltext` (lattice) function. It is worth having a look at the help for `ltext` as that gives an overview of all the common annotations and other options. We have chosen to plot the text at position $x = '2003-08-04'$ and $y = 200$ and the label itself. A useful option to `ltext` is `pos`. Values can be 1, 2, 3 and 4, and indicate positions below (the default), to the left of, above and to the right of the specified coordinates

Adding text and a shaded area

This time we will highlight an interval in row 2 (O_3) and write some text on top. Note that this time we use the `lpolygon` function and choose to put it under everything else on the plot. For the text, we have chosen a colour (yellow) font type 2 (bold) and made it a bit bigger (`cex = 1.5`). Note also the ‘y’ values extend beyond the actual limits shown on the plot — just to make sure they cover the whole region.

The polygon could of course be horizontal and more than one producing a series of ‘band’ e.g. air quality indexes. A more sophisticated approach is shown later for $PM_{2.5}$.

```
## add shaded polygon
trellis.last.object() +
  layer(lpolygon(x = c(ymd("2003-08-07"),
                     ymd("2003-08-07"), ymd("2003-08-12"),
                     ymd("2003-08-12")), y = c(-20, 600, 600, -20),
                     col = "grey", border = NA), under = TRUE, rows = 2)

## add text
trellis.last.object() +
  layer(ltext(x = ymd_hm("2003-08-09 12:00"), y = 50,
              labels = "!!episode!!", col = "yellow",
              font = 2, cex = 1.5), rows = 2)
```

The small shaded, semi-transparent area shown in the bottom panel was added as follows:

```
## add shaded polygon
plt <- plt +
  layer(lpolygon(x = c(ymd("2003-08-21"), ymd("2003-08-21"),
                     ymd("2003-08-23"), ymd("2003-08-23")),
                     y = c(4, 8, 8, 4), col = "blue", border = NA,
                     alpha = 0.2), rows = 5)
```

Adding an arrow

The arrow shown on the first panel of Figure 8.2 was added as follows. Note the `code = 3` placed arrows at both ends. Note that `angle` is the angle from the shaft of the

arrow to the edge of the arrow head.

```
trellis.last.object() +
  layer(larrows(ymd("2003-08-01"), 100,
                ymd_hm("2003-08-08 14:00"),
                100, code = 3, angle = 30), rows = 1)
```

Adding a reference line and text

This code adds a vertical dashed reference line shown in the 4th panel (PM_{10}) along with some text aligned at 90 degrees using the `srt` option of `ltext`.

```
trellis.last.object() +
  layer(panel.abline(v = ymd("2003-08-25"), lty = 5),
        rows = 4)
trellis.last.object() +
  layer(ltext(x = ymd_hm("2003-08-25 08:00"), y = 60,
              labels = "reference line", srt = 90), rows = 4)
```

Highlight a specific point

Up until now annotations have been added using arbitrary coordinates in each panel. What if we wanted to highlight a particular point, or more generally work with the actual data that are plotted. Knowing how to refer to existing data greatly extends the power of these functions.

It is possible to refer to a specific point in a panel simply by indexing the point of interest i.e. `x`, `y`. For example, to mark the 200th PM_{10} concentration (without knowing the actual date or value):

```
## add a specific point
trellis.last.object() +
  layer(lpoints(x[200], y[200], pch = 16, cex = 1.5),
        rows = 4)
```

What if we wanted to highlight the maximum O_3 concentration? It is possible to work out the index first and then use that to refer to that point. Note the ‘;’ to allow for the code to span multiple commands.

```
## add a point to the max O3 concentration
trellis.last.object() +
  layer({maxy <- which.max(y);
         lpoints(x[maxy], y[maxy], col = "black", pch = 16)},
        rows = 2)

## label max ozone
trellis.last.object() +
  layer({maxy <- which.max(y);
         ltext(x[maxy], y[maxy], paste(y[maxy], "ppb"),
               pos = 4)}, rows = 2)
```

Add a filled polygon

It can be seen in the top panel of [Figure 8.2](#) that some of the data are highlighted by filling the area below the line. This approach can be useful more generally in plotting. While it is possible to draw polygons easily and refer to the data itself, there needs to

be a way for dealing with gaps in data, otherwise these gaps could be filled in perhaps unpredictable ways. A function has been written to draw a polygon taking into account gaps (`poly.na`).

```
poly.na <- function(x1, y1, x2, y2, col = "black", alpha = 0.2) {
  for(i in seq(2, length(x1)))
    if (!any(is.na(y2[c(i - 1, i)])))
      lpolygon(c(x1[i - 1], x1[i], x2[i], x2[i - 1]),
                c(y1[i - 1], y1[i], y2[i], y2[i - 1]),
                col = col, border = NA, alpha = alpha)
}
```

This time we work out the ids of the data spanning an area of interest. Then the `poly.na` function is used. Note that the alpha transparency is by default 0.2 but another value can easily be supplied, as shown in the air quality ‘bands’ example.

```
trellis.last.object() +
  layer({id <- which(x >= ymd("2003-08-11") &
                      x <= ymd("2003-08-25"));
        poly.na(x[id], y[id], x[id], rep(0, length(id)),
                col = "darkorange")}, rows = 1)
```

Add air quality bands as polygons

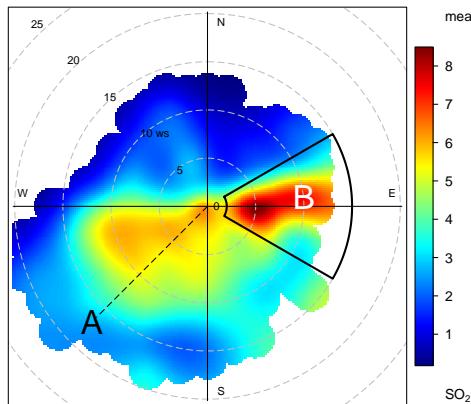
It is a simple extension to go from using a polygon below the data to polygons at certain intervals e.g. air quality indexes. These are shown for PM_{2.5} and the bands considered are 0–20, 20–30, 30–40 and >40.

```
trellis.last.object() +
  layer(poly.na(x, y, x, rep(0, length(x)),
                col = "green", alpha = 1), rows = 3)
trellis.last.object() +
  layer(poly.na(x, ifelse(y <20, NA, y), x,
                rep(20, length(x)), col = "yellow", alpha = 1),
        rows = 3)
trellis.last.object() +
  layer(poly.na(x, ifelse(y <30, NA, y),
                x, rep(30, length(x)),
                col = "orange", alpha = 1), rows = 3)
trellis.last.object() +
  layer(poly.na(x, ifelse(y <40, NA, y),
                x, rep(40, length(x)),
                col = "red", alpha = 1), rows = 3)
```

Polar plot examples

Many of the examples considered above are relevant to all other functions e.g. how to add text, choosing rows and columns to plot in. Polar coordinate plots are different because of the coordinate system used and this section considers a few examples.

One useful approach is to be able to draw an arc, perhaps highlighting an area of interest. A simple, but flexible function has been written to do this. It takes arguments `theta1` and `theta2` that define the angular area of interest and `lower` and `upper` to set the lower and upper wind speed, respectively. It also has additional arguments `theta3` and `theta4` which optionally set the angles for the ‘upper’ wind speed.

FIGURE 8.3 `polarPlot` for SO₂ with annotations.

```

arc <- function(theta1 = 30, theta2 = 60, theta3 = theta1, theta4 = theta2,
                 lower = 1, upper = 10){
  ## function to work out coordinates for an arc sector
  if (theta2 < theta1) {
    ang1 <- seq(theta1, 360, length = abs(theta2 - theta1))
    ang2 <- seq(0, theta2, length = abs(theta2 - theta1))
    angles.low <- c(ang1, ang2)

    ## for upper angles
    ang1 <- seq(theta1, 360, length = abs(theta4 - theta3))
    ang2 <- seq(0, theta2, length = abs(theta4 - theta3))
    angles.high <- c(ang1, ang2)
  } else {
    angles.low <- seq(theta1, theta2, length = abs(theta2 - theta1))
    angles.high <- seq(theta3, theta4, length = abs(theta4 - theta3))
  }
  x1 <- lower * sin(pi * angles.low / 180)
  y1 <- lower * cos(pi * angles.low / 180)
  x2 <- rev(upper * sin(pi * angles.high / 180))
  y2 <- rev(upper * cos(pi * angles.high / 180))
  data.frame(x = c(x1, x2), y = c(y1, y2))
}

}

```

Following on from the previous examples, some annotations have been added to the basic polar plot for SO₂ as shown in Figure 8.3. Note that in these plots (0, 0) is the middle of the plot and the radial distance will be determined by the wind speed — or whatever the radial variable is. This way of plotting arcs can also be applied to other functions that show directional data.

```

polarPlot(mydata, pollutant = "so2", col = "jet")
trellis.last.object() + layer(ltext(-12, -12, "A", cex = 2))
trellis.last.object() + layer(ltext(10, 2, "B", cex = 2, col = "white"))
trellis.last.object() + layer(lsegments(0, 0, -11.5, -11.5, lty = 5))
## add an arc to highlight area of interest
trellis.last.object() +
  layer(lpolygon(x = arc(theta1 = 60, theta2 = 120, lower = 2,
                         upper = 15)$x, y = arc(theta1 = 60,
                         theta2 = 120, lower = 2,
                         upper = 15)$y, lty = 1, lwd = 2))

```

Using grid graphics — identify locations interactively

The examples above provide a precise way of annotating plots for single or multi-pans **openair** displays. However, these methods won't work for plots that consist of completely separate plots such as the four plots in **timeVariation**. There are however other methods that can be used to annotate such plots using the package **grid**, which forms the basis of **lattice** graphics. There is enormous capability for annotating plots using the **grid** package and only a few simple examples are given here.

Given a plot such as [Figure 22.1](#), how could texts be added at any location — say in the middle monthly plot? One very useful function for this type of annotation that allows the user to interactively choose a location is the **grid.locator()** function in the **grid** package. That function can be called with different coordinate systems — but the one we want defines the bottom-left corner as (0, 0) and the top right as (1, 1).

First of all, make a **timeVariation** plot like [Figure 22.1](#).

```
timeVariation(mydata)
```

Now let's choose a location on the plot interactively using the mouse and selecting somewhere in the middle of the monthly plot.

```
library(grid)
## bring up the interactive location chooser
grid.locator(unit="npc")
```

What should happen is that in the R console the coordinates are given for that point. In my case these were $x = 0.503$ and $y = 0.338$. These coordinates can now be used as the basis of adding some text or other annotation. In the example below, the **grid.text** function is used to add some text for these coordinates making the font bigger (**cex** = 2), bold (**font** = 2) and blue (**col** = "blue").

```
grid.text(x = 0.503, y = 0.338, label = "here!",
          gp = gpar(cex = 2, font = 2, col = "blue"))
```

Even with this basic approach, some sophisticated annotation is possible with *any* **openair** plot. There are many other functions that can be used from the **grid** package that would allow for polygons, segments and other features to be drawn in a similar way to the examples earlier in this section. Continuing with the same example, here is how to add an arrow pointing to the maximum NO_x concentration shown on the top plot for Saturday (again using the **grid.locator** function).

```
grid.lines(x = c(0.736, 0.760), y = c(0.560, 0.778),
            arrow = arrow())
grid.text(x = 0.736, y = 0.560, label = "maximum", just = "left")
```

8.9 Getting help

The principal place for seeking help with **openair** functions is through the software itself. The document you are reading will increasingly give the background to the ideas and wider information. Also, the package itself will always contain the most up to date help. Furthermore, the process of building and checking packages is strict. For example, it is checked to see if all the options in a function match with descriptions in the help files, and all examples given in the help (and there are many) are run to ensure they all work. Nevertheless, the options shown for each function in this document are

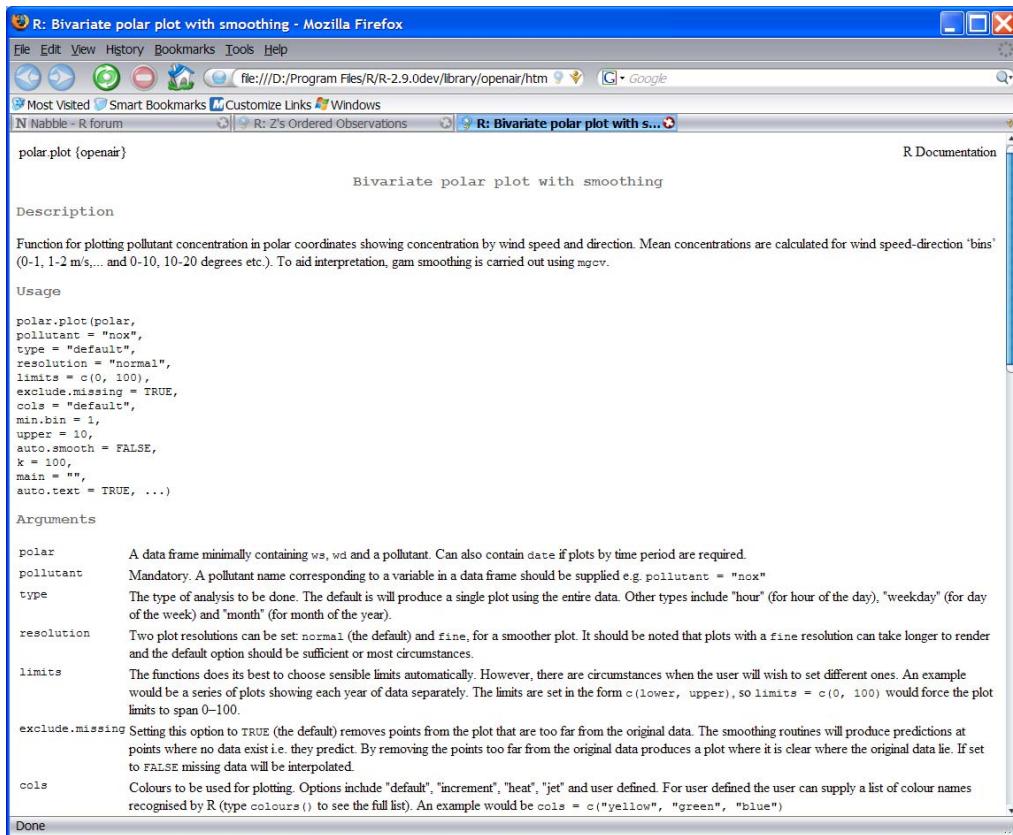


FIGURE 8.4 The help screen for the function `polarPlot`.

parsed directly from the **openair** package ensuring consistency between this document and the package help. To bring up the general help page (assuming you have loaded **openair**), type `?openair`, which will bring up the main **openair** page, from which there are links to all functions. Similarly, if you want help with a specific function more directly, type something like `?polarPlot`.

The help screen will provide the most up to date information on the function including: a short description, a description of all the options, a more detailed description and links to other similar functions. Importantly, each function help will have several examples given of its use, which are easily reproducible; just copy them into R. These examples use the data set ‘example data long.csv’ mentioned previously.

Handy tip: use TAB for word completion

If you are typing directly into R you do not always need to type the whole word of a function or option. Taking the **calendarPlot** function as an example, type ‘calen’ then press TAB and it will complete the whole string ‘calendarPlot’. Similarly, when typing the function options such as ‘pollutant’, just type the first few lines ‘poll’, press TAB and it will complete as ‘pollutant=’. This makes R much quicker to work with. It takes a bit of experimentation to get a feel for how many letters are required before a unique function name or option can be completed.

9 Getting data into **openair**

Importing data is usually the first step involved in data analysis using **openair**. As has been stressed before, the key issue is ensuring the data are in a simple format avoiding any unnecessary formatting. For this reason data are best stored either in a database or a .csv file. R itself has lots of capabilities for importing data and these will be useful in many situations e.g. `read.table` and `read.csv`. However, **openair** has several dedicated functions to make data import easier for users, as well as some more specific functions for particular data types. These are described below.

9.1 Issues related to time zones

Time zones can very difficult to work with in R and other software. Issues related to Daylight Savings Time (DST) can be particularly challenging (see http://en.wikipedia.org/wiki/List_of_zoneinfo_time_zones for information in World time zones). In summary, **openair** assumes times that do not include DST. To do this assumptions are made about times when data are imported to R using the **openair import** function and users should take care if importing data in other ways.

Ensuring certainty in time zones used

It is recommended that when importing data into **openair** using either the `import` function or by other means that the time zone of the original data are set to GMT or a fixed offset from GMT as outlined in this section. This way of importing data and setting the time zone should avoid any complexities related to Daylight Savings Time.

If the **openair import** function is used then the user can set the time zone for the original data. By default ‘GMT’ is assumed. For non-GMT time zones a GMT offset can be assumed. The option `tzone` (see below) can be used to control the time zone. For example, to set the time zone of the data to the time zone in New York (EST, 5 hours behind GMT) set `tzone = "Etc/GMT+5"`. To set the time zone of the data to Central European Time (CET, 1 hour ahead of GMT) set `tzone = "Etc/GMT-1"`. Note that *the positive and negative offsets are opposite to what most users expect*.

When **openair** started the **lubridate** package did not exist. **lubridate** makes formatting dates and times much easier and avoids much of the complexity that users often face with these issues. While there are many useful functions in **lubridate**, the convenience functions such as `dmy`, `ymd`, `dmy_hm` etc. are very useful and make it much easier to format dates and times. Note that these functions will by default return data in the UTC time zone in POSIXct format, which is required in **openair** and other R uses. It is worth taking a look at the vignette for **lubridate** for an overview of capabilities. However, a few examples of formatting dates are shown below:

```

library(lubridate)
dmy("1/5/2015")

## [1] "2015-05-01"

dmy("1-5-2015")

## [1] "2015-05-01"

dmy_hm("1/5/2015 12:00")

## [1] "2015-05-01 12:00:00 UTC"

mdy_hms("1/5/2015 12:00:00")

## [1] "2015-01-05 12:00:00 UTC"

```

Similarly the time zone should be set in the same way if the data are imported manually. For example:

```

dat <- read.csv("~/My Drive/openair/Data/example data long.csv",
               header = TRUE)

## set time zone to equivalent of CET but with no DST:
dat$date <- dmy_hm(dat$date, tz = "Etc/GMT-1")

```

For most functions these issues will not matter. However, if users are including data from other sources or are working with data that are GMT/UTC (Hysplit back trajectories for example), then it will be important to understand how to deal with these issues.

It may also be easier to set your R session to a non-DST time zone in a similar way to make sure everything is displayed in a consistent format e.g.

```
Sys.setenv(TZ = "Etc/GMT-1")
```

9.2 The **import** function

A flexible function **import** has been written to import .csv or .txt file data and format the date/time correctly. The main purpose of this function is to help format dates etc. for use in **openair** and R. This is the principal means by which most users should import data unless the data are from UK networks. It is simple to use with its default assumptions e.g. header on the first line and a column ‘date’ in the format dd/mm/yyyy HH:MM:

```
mydata <- import()
```

Typing this into R will bring up an ‘open file’ dialog box, from which you can choose a .csv file. Try importing the ‘example data long.csv’ file in this was to see how it works. Used without any options like this, it assumes that the **date** field is in the format dd/mm/yyyy HH:MM and is called ‘date’. By default the time zone is set to be GMT. However, users can apply a GMT offset as described in [Section 9.1](#).

Often it is better to supply the file path because this makes the analysis more reproducible e.g.

```
mydata <- import("d:/temp/my interesting data.csv")
```

The **import** function is actually very flexible and can take account of different date formats, header lines etc. See the options below. For most users, few if any of these options will need to be used, but for ‘difficult’ data, the flexibility should be helpful. One option that is often useful is to tell R how missing data are represented. If the fields are left blank, they will automatically be set to **NA**. However, it may be that the file identifies missing data by ‘NoData’, or ‘-999’. In this case, **import** should be called like:

```
mydata <- import(na.strings = "NoData")
```

or

```
mydata <- import(na.strings = "-999")
```

In the case of missing data being represented by several strings e.g. ‘-99.0’ and ‘-999’, it should be called like

```
mydata <- import(na.strings = c("-99.0", "-999"))
```

It is essential to supply the **import function with details of how missing data are represented if they are not represented by either a blank cell or **NA**.** This is because if text is present in a column that should be numeric, then R will consider that the column is a character and not numeric. When using the **import** function, details of the format of each field are printed in R. The user can check that fields that should be numeric appear as either ‘numeric’ or ‘integer’ and not ‘character’ or ‘factor’.

Another example is a file that has separate date and time fields e.g. a column called ‘mydate’ and a separate column called ‘mytime’. Further, assume that date is in the format dd.mm.YYYY e.g. 25.12.2010, and time is in the format HH:MM. Then the file could be imported as:

```
import("c:/temp/test.csv", date = "mydate", date.format = "%d.%m.%Y",
       time = "mytime", time.format = "%H:%M")
```

What if the date was in the format mm.dd.YYYY?:

```
import("c:/temp/test.csv", date = "mydate", date.format = "%m.%d.%Y",
       time = "mytime", time.format = "%H:%M")
```

...and the time was just the hour as an integer (0–23):

```
import("c:/temp/test.csv", date = "mydate", date.format = "%m.%d.%Y",
       time = "mytime", time.format = "%H")
```

Another common situation is that hour is represented as 1–24 in a date-time field. In this case it is necessary to correct for this. R stores **POSIXct** format as seconds, so 3600 need to be subtracted to ensure the time is correct. Note that if there is a separate column for hour then **import** will correct that automatically. So, for the date-time situation:

```
import("c:/temp/test.csv", date = "mydate", date.format = "%m.%d.%Y %H",
       correct.time = -3600)
```

Note if time was expressed as HH:MM:ss, then the option **time.format = "%H:%M:%S"** should be used.

There are other options for ignoring the first n lines i.e. due to header information and so on. The user can specify the header line row (**header.at**) and the row the data

starts at (**data.at**).

Note also that **import** assumes there are no daylight saving time (DST) issues in the original data i.e. a missing hour in spring and a duplicate hour in autumn. Dealing with these issues in R rapidly gets too complicated ... users should therefore ensure the original data do not consider DST.

The options for the **import** function are:

file	The name of the file to be imported. Default, file = file.choose() , opens browser. Alternatively, the use of read.table (in utils) also allows this to be a character vector of a file path, connection or url.
file.type	The file format, defaults to common ‘csv’ (comma delimited) format, but also allows ‘txt’ (tab delimited).
sep	Allows user to specify a delimiter if not ‘,’ (csv) or TAB (txt). For example ‘;’ is sometimes used to delineate separate columns.
header.at	The file row holding header information or NULL if no header to be used.
data.at	The file row to start reading data from. When generating the data frame, the function will ignore all information before this row, and attempt to include all data from this row onwards.
date	Name of the field containing the date. This can be a date e.g. 10/12/2012 or a date-time format e.g. 10/12/2012 01:00.
date.format	The format of the date. This is given in ‘R’ format according to strptime . For example, a date format such as 1/11/2000 12:00 (day/month/year hour:minutes) is given the format “%d/%m/%Y %H:%M”. See examples below and strptime for more details.
time	The name of the column containing a time — if there is one. This is used when a time is given in a separate column and date contains no information about time.
time.format	If there is a column for time then the time format must be supplied. Common examples include “%H:%M” (like 07:00) or an integer giving the hour, in which case the format is “%H”. Again, see examples below.
tzone	The time zone for the data. In order to avoid the complexities of DST (daylight savings time), openair assumes the data are in GMT (UTC) or a constant offset from GMT. Users can set a positive or negative offset in hours from GMT. For example, to set the time zone of the data to the time zone in New York (EST, 5 hours behind GMT) set tzzone = "Etc/GMT+5" . To set the time zone of the data to Central European Time (CET, 1 hour ahead of GMT) set tzzone = "Etc/GMT-1" . Note that the positive and negative offsets are opposite to what most users expect.
na.strings	Strings of any terms that are to be interpreted as missing (NA). For example, this might be “-999”, or “n/a” and can be of several items.
quote	String of characters (or character equivalents) the imported file may use to represent a character field.
ws	Name of wind speed field if present if different from “ws” e.g. ws = "WSPD" .

wd	Name of wind direction field if present if different from “wd” e.g. <code>wd = "WDIR"</code> .
correct.time	Numerical correction (in seconds) for imported date. Default <code>NULL</code> turns this option off. This can be useful if the hour is represented as 1 to 24 (rather than 0 to 23 assumed by R). In which case <code>correct.time = -3600</code> will correct the hour.
...	Other arguments passed to <code>read.table</code> .

9.3 The **importAURN** function

While `import` is a useful function for ad-hoc data import, much of the data stored in the UK and beyond resides on central repositories that are available over the Internet. The UK AURN archive and King’s College London’s London Air Quality Network (LAQN) are two important and large databases of information that allow free public access. Storing and managing data in this way has many advantages including consistent data format, and underlying high quality methods to process and store the data. We are working with AEA and KCL to make things easier to link with `openair` functions.

Many users download hourly data from the air quality archive at <http://www.airquality.co.uk>. Most commonly, the data are emailed to the user as .csv files and have a fixed format as shown below. This is a useful facility but does have some limitations and frustrations, many of which have been overcome using a new way of storing and downloading the data described below.

The `importAURN` function has been written to make it easy to import data from the UK AURN. AEA have provided .RData files (R workspaces) of all individual sites and years for the AURN. These files are updated on a daily basis. This approach requires a link to the Internet to work.

There are several advantages over the web portal approach where .csv files are downloaded. First, it is quick to select a range of sites, pollutants and periods (see examples below). Second, storing the data as .RData objects is very efficient as they are about four times smaller than .csv files (which are already small) — which means the data downloads quickly and saves bandwidth. Third, the function completely avoids any need for data manipulation or setting time formats, time zones etc. Finally, it is easy to import many years of data beyond the current limit of about 64,000 lines. The final point makes it possible to download several long time series in one go.

The site codes and pollutant names can be upper or lower case. The function will issue a warning when data less than six months old is downloaded, which may not be ratified. Type `?importAURN` for a full listing of sites and their codes.

Note that currently there is no meteorological data associated with the archive. To use the full flexibility of the functions it is recommended that the AURN data are combined with appropriate local meteorological data. See (§5.3) for more details on how to combine separate air pollution and meteorological files.

The function has the following options.

site	Site code of the AURN site to import e.g. “my1” is Marylebone Road. Several sites can be imported with <code>site = c("my1", "nott")</code> — to import Marylebone Road and Nottingham for example.
year	Year or years to import. To import a sequence of years from 1990 to 2000 use <code>year = 1990:2000</code> . To import several specific years use <code>year = c(1990, 1995, 2000)</code> for example.

pollutant Pollutants to import. If omitted will import all pollutants from a site. To import only NOx and NO2 for example use `pollutant = c("nox", "no2")`.

- hc** A few sites have hydrocarbon measurements available and setting `hc = TRUE` will ensure hydrocarbon data are imported. The default is however not to as most users will not be interested in using hydrocarbon data and the resulting data frames are considerably larger.
- meta** Should meta data be returned? If `TRUE` the site type, latitude and longitude are returned.
- verbose** Should the function give messages when downloading files? Default is `FALSE`.

Some examples of usage are shown below.

```
## import all pollutants from Marylebone Rd from 1990:2009
mary <- importAURN(site = "my1", year = 2000:2009)

## import nox, no2, o3 from Marylebone Road and Nottingham Centre for 2000
thedata <- importAURN(site = c("my1", "nott"), year = 2000,
                        pollutant = c("nox", "no2", "o3"))

## import over 20 years of Mace Head O3 data!
o3 <- importAURN(site = "mh", year = 1987:2009)
## import hydrocarbon data from Marylebone Road
hc <- importAURN(site = "my1", year = 2008, hc = TRUE)
```

In future, **openair** functions will recognise AURN data and capture units, thus enabling plots to be automatically annotated. Furthermore, there is the potential to include lots of other ‘meta data’ such as site location, site type etc., which will be added to the function in due course.

9.4 The **importKCL** function

King’s College London also make available their data in a similar way to the **importAURN**. One difference compared to the **importAURN** is the availability of meteorological data in London, which is accessible through the option **met**. We have provided a ‘typical’ meteorological data set representing London, which is a composite of data from several instruments co-located with air pollution monitoring sites. Access to reliable meteorological data can be difficult and expensive, and this is an issue we hope to improve in time.

The options for **importKCL** are:

- site** Site code of the network site to import e.g. “my1” is Marylebone Road. Several sites can be imported with `site = c("my1", "kc1")` — to import Marylebone Road and North Kensington for example.
- year** Year or years to import. To import a sequence of years from 1990 to 2000 use `year = 1990:2000`. To import several specific years use `year = c(1990, 1995, 2000)` for example.
- pollutant** Pollutants to import. If omitted will import all pollutants from a site. To import only NOx and NO2 for example use `pollutant = c("nox", "no2")`.

met	Should meteorological data be added to the import data? The default is FALSE . If TRUE wind speed (m/s), wind direction (degrees), solar radiation and rain amount are available. See details below.
	Access to reliable and free meteorological data is problematic.
units	By default the returned data frame expresses the units in mass terms (ug/m ³ for NOx, NO ₂ , O ₃ , SO ₂ ; mg/m ³ for CO). Use units = "volume" to use ppb etc. PM10_raw TEOM data are multiplied by 1.3 and PM2.5 have no correction applied. See details below concerning PM10 concentrations.
extra	Not currently used.
meta	Should meta data be returned? If TRUE the site type, latitude and longitude are returned.

Examples of importing data are given in the help files as part of **openair**. However, the example below shows how to import data from the Bexley 1 site (code BX1), together with the meteorological data.

Like the **importAURN** function the selection is only possible by site code (all the site codes and full site descriptions are shown in the help file). In time we will include other information such as site location and type and develop the functions to make it easy to use these other fields.

Below is an example of importing data from the Bexley 1 site.

```
bx1 <- importKCL(site = "bx1", year = 2000:2009, met = TRUE)
```

in which case a dialog box will appear prompting the user for a file location. The data frame **mydata** is now ready for use in **openair**.

9.5 Importing data from the CERC ADMS modelling systems

The ADMS suite of models is widely used in the UK and beyond. These models are used for a wide range of purposes and one of the benefits of **openair** is that many of the functions are potentially useful for model evaluation. One of the principal benefits of linking **openair** with the ADMS models is the access to meteorological data that is possible. In the UK, the Met Office provides meteorological data in a specific format for use in ADMS models.¹⁰ It is useful to be able to easily import the meteorological data into **openair** because analyses are often limited by the availability of representative meteorological data. However, the use of directly measured input data is only one possibility. When ADMS models run they use a sophisticated meteorological pre-processor to calculate other quantities that are not directly measured, but are important to dispersion modelling. Examples of these other variables are boundary layer height and surface sensible heat flux. These and many other quantities are calculated by the met pre-processor and output to a .MOP file. Access to these other quantities greatly increases the potential for model evaluation and in general provides a much richer source of information for analysis.

Many users may have meteorological data in the ADMS format. This is the format provided by the UK Met Office for the ADMS model. A an example of the format is shown in Figure 9.1, which is a simple text file. The **importADMSMet** function imports such data into R in a format suitable for **openair**.

This can be done, for example by:

¹⁰Specifically hourly sequential data and not statistical summaries of data.

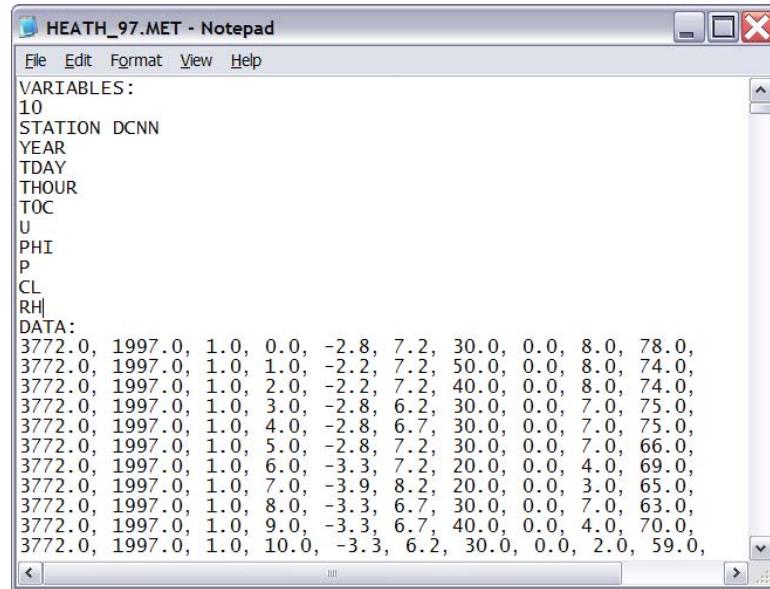


FIGURE 9.1 Typical format of an hourly ADMS met file.

```
met <- importADMS("d:/temp/heathrow01.met")
```

If no file name is supplied, the user will be prompted for one.

Sometimes it may be necessary to import several years. Here's one approach for doing so assuming the files are in a folder d:/metdata and all have a file extension.met:

```
all.met <- lapply(list.files(path = "d:/metdata", pattern = ".met",
                             full.names = TRUE),
                  function(.file) importADMS(.file))
all.met <- do.call(rbind.fill, all.met)
```

all.met will then contain met data for all years in one data frame.

9.5.1 An example considering atmospheric stability

One of the significant benefits of working with ADMS output files is having access to the outputs from the meteorological pre-processor. ADMS uses readily available meteorological variables such as wind speed, temperature and cloud cover and calculates parameters that are used in the dispersion algorithms. When ADMS is run it produces a .MOP file with all these inputs and processed quantities in. Access to parameters such as boundary layer height, Monin-Obukov length and so on can *greatly* increase the opportunities for insightful data analysis using existing **openair** functions. This is almost certainly an area we will cover in more depth later; but for now, here are a few examples.

We are going to use a .MOP file from 2001 following some dispersion modelling of stacks in London. The interest here is to use the results from the met pre-processor to better understand sources in the east of London at the Thurrock background site. First, we can import the Thurrock data (type **?importKCL** for site code listing) using the **importKCL** function:

```
tk1 <- importKCL(site = "tk1", year = 2001)
```

```
## show first few lines of tk1
head(tk1)

##                      date  nox  no2 o3      so2    co pm10_raw pm10
## 41862 2001-01-01 00:00:00   NA   NA NA      NA    NA     5.2  5.2
## 41863 2001-01-01 01:00:00  7.68 5.76 50 46.40133 0.232   11.7 11.7
## 41864 2001-01-01 02:00:00  5.76 3.84 52 57.55550 0.232    7.8  7.8
## 41865 2001-01-01 03:00:00  5.76 3.84 56 14.72350 0.232    5.2  5.2
## 41866 2001-01-01 04:00:00  1.92 1.92 54 10.70800 0.232   10.4 10.4
## 41867 2001-01-01 05:00:00  3.84 1.92 54 10.70800 0.232   14.3 14.3
##                      site code
## 41862 Thurrock - London Road (Grays)  TK1
## 41863 Thurrock - London Road (Grays)  TK1
## 41864 Thurrock - London Road (Grays)  TK1
## 41865 Thurrock - London Road (Grays)  TK1
## 41866 Thurrock - London Road (Grays)  TK1
## 41867 Thurrock - London Road (Grays)  TK1
```

Next we will import the .MOP file. The function automatically lists all the variables imported:

```
met <- importADMS("~/My Drive/openair/Data/met01.MOP")

##          date1        date2       line      run
## "POSIXct" "POSIXt" "integer" "factor"
##         fr           ws      ws.gstar      wd
## "numeric" "numeric" "numeric" "numeric"
##        delta.wd      ftheta0        k
## "numeric" "numeric" "numeric" "numeric"
##          h            nu      delta.theta      temp
## "numeric" "numeric" "numeric" "numeric"
##          p            cl      albedo.met albedo.disp
## "numeric" "numeric" "numeric" "numeric"
## alpha.met      alpha.disp      tsea      delta.t
## "numeric" "numeric" "numeric" "numeric"
## sigma.theta      rhu        q0
## "numeric" "numeric" "numeric" "numeric"
##        rhu.1      drhdzu process.ws.star process.ws.g
## "numeric" "numeric" "numeric" "numeric"
## process.ws.gstar      process.wd.0 process.wd.g process.delta.wd
## "numeric" "numeric" "numeric" "numeric"
## process.wd.sec      process.wstar process.ftheta0 process.k
## "numeric" "numeric" "numeric" "numeric"
## process.recip.lmo      process.h      process.nu process.delta.theta
## "numeric" "numeric" "numeric" "numeric"
## process.temp      process.p      process.delta.t process.sigma.theta
## "numeric" "numeric" "numeric" "numeric"
## process.q0      process.lambdae      process.rhu process.drhdzu
## "numeric" "numeric" "numeric" "numeric"
## process.z0.met      process.z0.disp      process.z0
## "numeric" "numeric"
```

Now we need to merge these two files using ‘date’ as the common field using the **merge** function, which is part of the base R system:

```
tk1 <- merge(tk1, met, by = "date")
```

Now we have a data frame with all the pollution measurements and meteorological variables matched up. A nice first example is to make use of variables that are not readily available. In particular, those representing atmospheric stability are very useful. So, let’s see what a polar plot looks like split by different levels of the atmospheric

```
polarPlot(tk1, pollutant = "so2", type = "process.recip.lmo",
min.bin = 2)
```

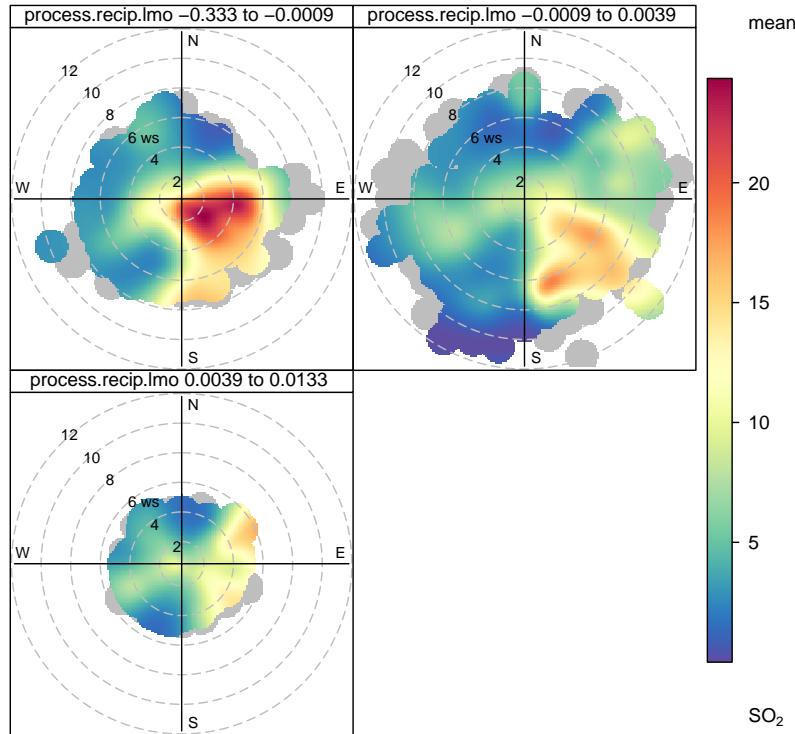


FIGURE 9.2 Use of the `importADMS` function to access atmospheric stability parameters for use in a polar plot. In this case $\frac{1}{L_{MO}}$ is split by three different levels, approximately corresponding to unstable, neutral and stable atmospheric conditions.

stability parameter the reciprocal of the Monin-Obukov length, $\frac{1}{L_{MO}}$. This has the name `process.recip.lmo`. Note we also set the option `min.bin = 2`, to ensure the output is not overly affected by a single high concentration.

The results are shown in Figure 9.2. So what does this tell us? Well, first $\frac{1}{L_{MO}}$ has been split into three different levels (broadly speaking the more negative the value of $\frac{1}{L_{MO}}$ the more unstable the atmosphere is and the more positive $\frac{1}{L_{MO}}$ is, the more stable the atmosphere is). In Figure 9.2 the plot shows what we might think of unstable, neutral and stable atmospheric conditions.

The first thing to note from Figure 9.2 is that lower wind speeds are associated with stable and unstable atmospheric conditions — shown by the smaller plot areas for these conditions (neutral conditions have a larger ‘blob’ extending to higher wind speeds). This is entirely expected. Starting with the unstable conditions (top left panel), SO₂ concentrations are dominated by easterly and south-easterly winds. These concentrations are likely dominated by tall stack emissions from those wind directions. For stable conditions (plot at the bottom), three sources seem to be important. There is a source to the north-east, the south-east and higher concentrations for very low wind speeds. The latter is likely due to road vehicle emissions of SO₂. The neutral conditions are perhaps revealing two sources to the south-east. Taken together, plotting the data in this way is beginning to reveal a potentially large number of sources in the area. Combined with the results from a dispersion model, or knowledge of local stacks, there is a good chance

that these sources can be identified.

A polar plot on its own does not reveal such detailed information. Try it:

```
polarPlot(tk1, pollutant = "so2")
```

Of course care does need to be exercised when interpreting these outputs, but the availability of wider range of meteorological data can only improve inference.

Here are some other analyses (not plotted, but easily run). For NO_x:

```
## dominated by stable conditions and low
## wind speeds (traffic sources)
polarPlot(tk1, pollutant = "nox", type = "process.recip.lmo",
min.bin = 2)
```

PM₁₀:

```
## complex, but dominated by stable/unstable easterly conditions
polarPlot(tk1, pollutant = "pm10", type = "process.recip.lmo",
min.bin = 2)
```

How about the ratio of two pollutants, say the ratio of SO₂/NO_x? First calculate the ratio:

```
tk1 <- transform(tk1, ratio = so2 / nox)
## evidence of a source with high so2/nox ratio ro the SSE
polarPlot(tk1, pollutant = "ratio", type = "process.recip.lmo",
min.bin = 2)
```

And don't forget all the other parameters available such as boundary layer height etc. — and all the other functions in **openair** that can be used.

10 The **summaryPlot** function

The **summaryPlot** function is a way of rapidly summarising important aspects of data. While many statistical summaries are possible to calculate with R, the **summaryPlot** function has been written specifically for monitoring data. The function provides key graphical and statistical summaries. **summaryPlot** has the following options:

- mydata** A data frame to be summarised. Must contain a **date** field and at least one other parameter.
- na.len** Missing data are only shown with at least **na.len** contiguous missing values. The purpose of setting **na.len** is for clarity: with long time series it is difficult to see where individual missing hours are. Furthermore, setting **na.len = 96**, for example would show where there are at least 4 days of continuous missing data.
- clip** When data contain outliers, the histogram or density plot can fail to show the distribution of the main body of data. Setting **clip = TRUE**, will remove the top 1 yield what is often a better display of the overall distribution of the data. The amount of clipping can be set with **percentile**.
- percentile** This is used to clip the data. For example, **percentile = 0.99** (the default) will remove the top 1 percentile of values i.e. values greater than the 99th percentile will not be used.

<code>type</code>	<code>type</code> is used to determine whether a histogram (the default) or a density plot is used to show the distribution of the data.
<code>pollutant</code>	<code>pollutant</code> is used when there is a field <code>site</code> and there is more than one site in the data frame.
<code>period</code>	<code>period</code> is either <code>years</code> (the default) or <code>months</code> . Statistics are calculated depending on the <code>period</code> chosen.
<code>avg.time</code>	This defines the time period to average the time series plots. Can be "sec", "min", "hour", "day" (the default), "week", "month", "quarter" or "year". For much increased flexibility a number can precede these options followed by a space. For example, a timeAverage of 2 months would be <code>avg.time = "2 month"</code> .
<code>print.datacap</code>	Should the data capture % be shown for each period?
<code>breaks</code>	Number of histogram bins. Sometime useful but not easy to set a single value for a range of very different variables.
<code>col.trend</code>	Colour to be used to show the monthly trend of the data, shown as a shaded region. Type <code>colors()</code> into R to see the full range of colour names.
<code>col.data</code>	Colour to be used to show the <i>presence</i> of data. Type <code>colors()</code> into R to see the full range of colour names.
<code>col.mis</code>	Colour to be used to show missing data.
<code>col.hist</code>	Colour for the histogram or density plot.
<code>cols</code>	Predefined colour scheme, currently only enabled for " <code>greyscale</code> ".
<code>date.breaks</code>	Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of <code>date.breaks</code> up or down.
<code>auto.text</code>	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO ₂ .
<code>...</code>	Other graphical parameters. Commonly used examples include the axis and title labelling options (such as <code>xlab</code> , <code>ylab</code> and <code>main</code>), which are all passed to the plot via <code>quickText</code> to handle routine formatting. As <code>summaryPlot</code> has two components, the axis labels may be a vector. For example, the default case (<code>type = "histogram"</code>) sets y labels equivalent to <code>ylab = c("", "Percent of Total")</code> .

the
`summaryPlot`
 function
 should be used
 for checking
 input data
 before
 applying other
 functions

It is called in a very simple way:¹¹

An example of using `summaryPlot` shown in Figure 10.1. For each numerical variable in a data frame, a plot is made, shown in the left panel, showing where data exist (blue) and missing data (red). For clarity, only running sequences of ≥ 24 hours of missing data are shown. It is easy to see therefore that the beginning part of the time

¹¹Note that a data frame `mydata` is automatically loaded when loading the `openair` package. The data set consists of several years of pollutant data from Marylebone Road in London.

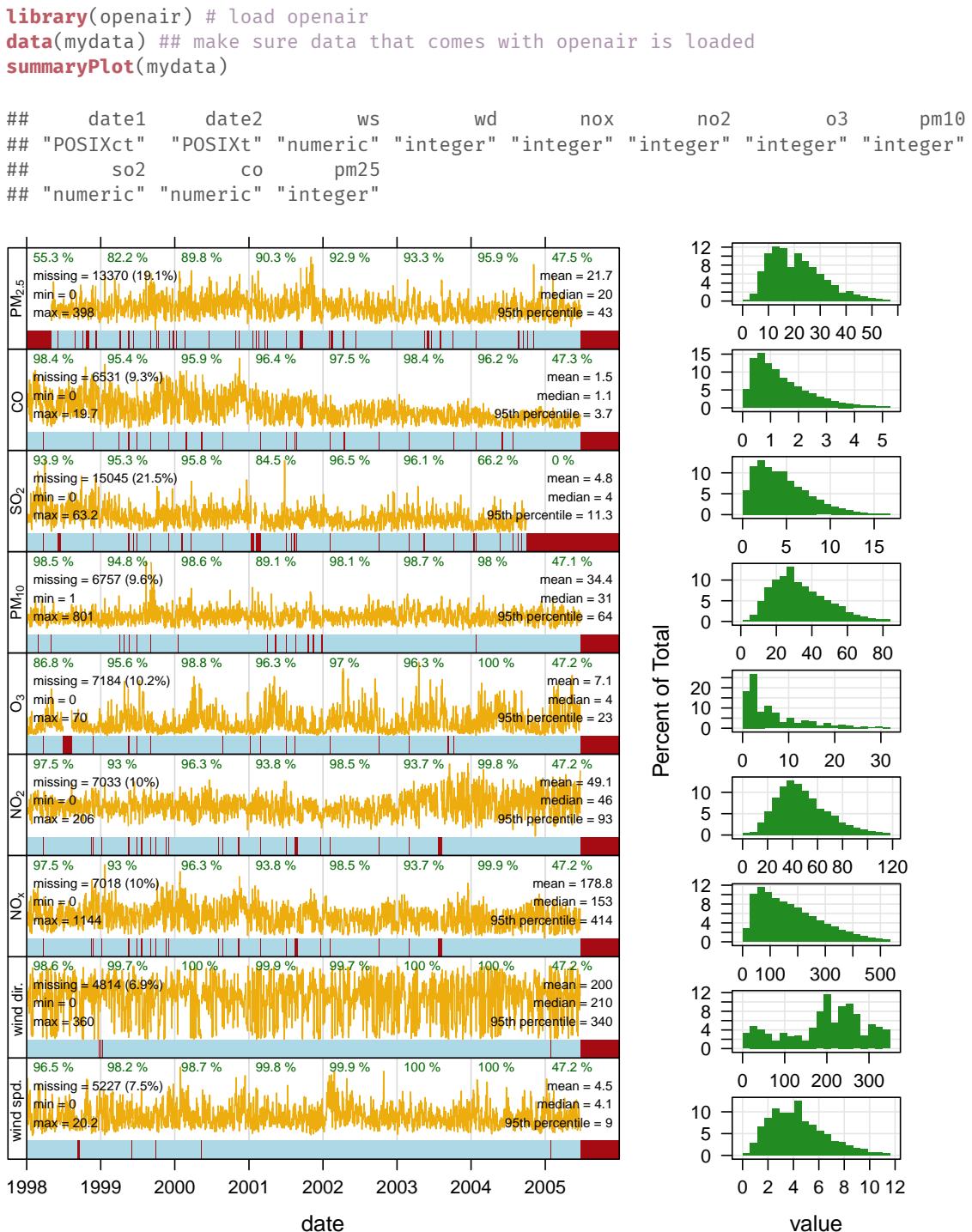


FIGURE 10.1 Use of **summaryPlot** function applied to the **mydata** data frame. The plots in the left panel show the time series data, where blue shows the presence of data and red missing data. The daily mean values are also shown in pale yellow scaled to cover the range in the data from zero to the maximum daily value. As such, the daily values are indicative of an overall trend rather than conveying quantitative information. For each pollutant, the overall summary statistics are given. For each year the percentage data capture is shown in green font. The panels on the right show the distribution of each species using a histogram plot.

series for PM_{2.5} is missing and the end part of SO₂. It is also clear that the time series stops half way through 2005. Also shown in each panel are statistical summaries, which include: number of missing points (with percentage shown in parentheses), minimum, maximum, mean, median and the 95th percentile. For each year, the data capture (%) is shown in green font. So, for example, the data capture for NO_x in 2000 was 96.3 %.

The pale yellow line gives an indication of the variation in values over time expressed as a daily mean. It is in indication because no numerical scale is given. The data are formatted so that 0 is placed at the lower part of the scale (top of the data indicator strip) and the maximum value at the top of the graph. The intention is to give the user a feel for how the data vary over the length of the time series.

The plots shown in the right panel are histograms. It is also possible to show *density plots*. A density plot is somewhat similar to a histogram but avoids having to arbitrarily select a ‘bin’ size. The choice of bin size in histograms can often lead to a misleading impression of how data are distributed — simply because of the bin size chosen. The default behaviour of this function ‘clips’ the data and excludes the highest 1 % of values. This is done to help highlight the shape of the bulk of the data and has the effect of removing the long tail, typical of air pollution concentration distributions.

It is possible, however, not to clip the histogram or density plot data and select various other options:

```
summaryPlot(mydata, clip = FALSE) # do not clip density plot data

summaryPlot(mydata, percentile = 0.95) # exclude highest 5 % of data etc.

# show missing data where there are at least 10 continuous missing values
summaryPlot(mydata, na.len = 10)

summaryPlot(mydata, col.data = "green") # show data in green

summaryPlot(mydata, col.mis = "yellow") # show missing data in yellow

summaryPlot(mydata, col.dens = "black") # show density plot line in black
```

Depending on the data available, there may be too many plots shown on one page, making it difficult to see the detail. Currently, the simplest way to reduce what is shown is to limit the data to be plotted. In the code below, for example, only columns 2 and 5 to 7 are plotted (column 1 in this case is the date and must always be supplied). Alternatively, the **subset** function could be used:

```
summaryPlot(mydata[, c(1, 2, 5:7)]) # only plot columns 2 and 5-7
summaryPlot(subset(mydata, select = c(date, nox, no2, co))) # alternative selecting
```

So far the **summaryPlot** function has been described and used in terms of plotting many variables from a single site. What happens if there is more than one site? Because the plot already produces a dense amount of information it does not seem sensible to plot several variables across several sites at the same time. Therefore, if there is a **site** field, **summaryPlot** will provide summary data for *a single pollutant* across all sites. See **?summaryPlot** for more details.

Use **summaryPlot** first

It is recommended that the **summaryPlot** function is used before moving on to using the other functions detailed below. One of the reasons (apart from getting to know your data) is that it also acts as a way of ensuring that other functions should

work. For example, if wind speed was missing, or was formatted as a character rather than a number, it will not show up in the summary plot. In time we intend to use this function to carry out many data checks and issue warnings if problems are detected.

11 The `cutData` function

The `cutData` function is a utility function that is called by most other functions but is useful in its own right. It's main use is to partition data in many ways, many of which are built-in to `openair`.

Note that all the date-based types e.g. month/year are derived from a column `date`. If a user already has a column with a name of one of the date-based types it will not be used.

For example, to cut data into seasons:

```
mydata <- cutData(mydata, type = "season")
head(mydata)

## # A tibble: 6 x 11
##   date                  ws    wd   nox   no2    o3  pm10    so2    co  pm25
##   <dttm>              <dbl> <int> <int> <int> <int> <dbl> <dbl> <int>
## 1 1998-01-01 00:00:00  0.6   280   285   39     1    29  4.72  3.37  NA
## 2 1998-01-01 01:00:00  2.16   230    NA    NA    NA    37  NA    NA    NA
## 3 1998-01-01 02:00:00  2.76   190    NA    NA    3    34  6.83  9.60  NA
## 4 1998-01-01 03:00:00  2.16   170   493    52     3    35  7.66 10.2   NA
## 5 1998-01-01 04:00:00  2.4    180   468    78     2    34  8.07  8.91  NA
## 6 1998-01-01 05:00:00  3     190   264    42     0    16  5.50  3.05  NA
## # ... with 1 more variable: season <ord>
```

This adds a new field ‘season’ that is split into four seasons. There is an option `hemisphere` that can be used to use southern hemisphere seasons when set as `hemisphere = "southern"`.

The `type` can also be another field in a data frame e.g.

```
mydata <- cutData(mydata, type = "pm10")
head(mydata)

## # A tibble: 6 x 10
##   date                  ws    wd   nox   no2    o3  pm10      so2    co  pm25
##   <dttm>              <dbl> <int> <int> <int> <int> <fct>    <dbl> <dbl> <int>
## 1 1998-01-01 00:00:00  0.6   280   285   39     1 pm10 22 t~  4.72  3.37  NA
## 2 1998-01-01 01:00:00  2.16   230    NA    NA    NA pm10 31 t~  NA    NA    NA
## 3 1998-01-01 02:00:00  2.76   190    NA    NA    3 pm10 31 t~  6.83  9.60  NA
## 4 1998-01-01 03:00:00  2.16   170   493    52    3 pm10 31 t~  7.66 10.2   NA
## 5 1998-01-01 04:00:00  2.4    180   468    78    2 pm10 31 t~  8.07  8.91  NA
## 6 1998-01-01 05:00:00  3     190   264    42    0 pm10 1 to~  5.50  3.05  NA

data(mydata) ## re-load mydata fresh
```

This divides PM₁₀ concentrations into four *quantiles* — roughly equal numbers of PM₁₀ concentrations in four levels.

Most of the time users do not have to call `cutData` directly because most functions have a `type` option that is used to call `cutData` directly e.g.

```
polarPlot(mydata, pollutant = "so2", type = "season")
```

However, it can be useful to call `cutData` before supplying the data to a function in a few cases. First, if one wants to set seasons to the southern hemisphere as above. Second, it is possible to override the division of a numeric variable into four quantiles by using the option `n.levels`. More details can be found in the `cutData` help file.

The `cutData` function has the following options:

- `x` A data frame containing a field `date`.
- `type` A string giving the way in which the data frame should be split. Pre-defined values are: "default", "year", "hour", "month", "season", "weekday", "site", "weekend", "monthlyear", "daylight", "dst" (daylight saving time).
`type` can also be the name of a numeric or factor. If a numeric column name is supplied `cutData` will split the data into four quantiles. Factors levels will be used to split the data without any adjustment.
- `hemisphere` Can be "northern" or "southern", used to split data into seasons.
- `n.levels` Number of quantiles to split numeric data into.
- `start.day` What day of the week should the `type = "weekday"` start on? The user can change the start day by supplying an integer between 0 and 6. Sunday = 0, Monday = 1, ...For example to start the weekday plots on a Saturday, choose `start.day = 6`.
- `is.axis` A logical (`TRUE/FALSE`), used to request shortened cut labels for axes.
- `local.tz` Used for identifying whether a date has daylight savings time (DST) applied or not. Examples include `local.tz = "Europe/London"`, `local.tz = "America/New_York"` i.e. time zones that assume DST. http://en.wikipedia.org/wiki/List_of_zoneinfo_time_zones shows time zones that should be valid for most systems. It is important that the original data are in GMT (UTC) or a fixed offset from GMT. See `import` and the openair manual for information on how to import data and ensure no DST is applied.
- `latitude` The decimal latitude used in `type = "daylight"`.
- `longitude` The decimal longitude. Note that locations west of Greenwich are negative.
- ... All additional parameters are passed on to next function(s).

12 The `windRose` and `pollutionRose` functions

12.1 Purpose

see also
`polarFreq`
`percentileRose`

The wind rose is a very useful way of summarising meteorological data. It is particularly useful for showing how wind speed and wind direction conditions vary by year. The `windRose` function can plot wind roses in a variety of ways: summarising all available wind speed and wind direction data, plotting individual wind roses by year, and also by month. The latter is useful for considering how meteorological conditions vary by season.

Data are summarised by direction, typically by 45 or 30° and by different wind speed categories. Typically, wind speeds are represented by different width 'paddles'. The plots show the proportion (here represented as a percentage) of time that the wind is from a certain angle and wind speed range.

The `windRose` function also calculates the percentage of ‘calms’ i.e. when the wind speed is zero. UK Met Office data assigns these periods to 0 degrees wind direction with valid northerly winds being assigned to 360 degrees.

The `windRose` function will also correct for bias when wind directions are rounded to the nearest 10 degrees but are displayed at angles that 10 degrees is not exactly divisible into e.g. 22.5 degrees. When such data are binned, some angles i.e. N, E, S, W will comprise of three intervals whereas others will comprise of two, which can lead to significant bias. This issue and its solution is discussed by Droppo and Napier (2008) and Applequist (2012).¹² `openair` uses a simple method to correct for the bias by globally rescaling the count in each wind direction bin by the number of directions it represents relative to the average. Thus, the primary four directions are each reduced by a factor of 0.75 and the remaining 12 directions are multiplied by 1.125.

12.2 Options available

The `windRose` function has the following options:

<code>mydata</code>	A data frame containing fields <code>ws</code> and <code>wd</code>
<code>ws</code>	Name of the column representing wind speed.
<code>wd</code>	Name of the column representing wind direction.
<code>ws2</code>	The user can supply a second set of wind speed and wind direction values with which the first can be compared. See details below for full explanation.
<code>wd2</code>	see <code>ws2</code> .
<code>ws.int</code>	The Wind speed interval. Default is 2 m/s but for low met masts with low mean wind speeds a value of 1 or 0.5 m/s may be better. Note, this argument is superseded in <code>pollutionRose</code> . See <code>breaks</code> below.
<code>angle</code>	Default angle of “spokes” is 30. Other potentially useful angles are 45 and 10. Note that the width of the wind speed interval may need adjusting using <code>width</code> .
<code>type</code>	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. “season”, “year”, “weekday” and so on. For example, <code>type = "season"</code> will produce four plots — one for each season. It is also possible to choose <code>type</code> as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another. Type can be up length two e.g. <code>type = c("season", "weekday")</code> will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

¹²Thanks to Philippe Barnéoud of Environment Canada for pointing this issue out.

bias.corr When `angle` does not divide exactly into 360 a bias is introduced in the frequencies when the wind direction is already supplied rounded to the nearest 10 degrees, as is often the case. For example, if `angle = 22.5`, N, E, S, W will include 3 wind sectors and all other angles will be two. A bias correction can made to correct for this problem. A simple method according to Applequist (2012) is used to adjust the frequencies.

cols	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet”, “hue” and user defined. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue", "black")</code> .
grid.line	Grid line interval to use. If <code>NULL</code> , as in default, this is assigned by <code>windRose</code> based on the available data range. However, it can also be forced to a specific value, e.g. <code>grid.line = 10</code> . <code>grid.line</code> can also be a list to control the interval, line type and colour. For example <code>grid.line = list(value = 10, lty = 5, col = "purple")</code> .
width	For <code>paddle = TRUE</code> , the adjustment factor for width of wind speed intervals. For example, <code>width = 1.5</code> will make the paddle width 1.5 times wider.
seg	For <code>pollutionRose seg</code> determines with width of the segments. For example, <code>seg = 0.5</code> will produce segments $0.5 * \text{angle}$.
auto.text	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO ₂ .
breaks	Most commonly, the number of break points for wind speed in <code>windRose</code> or pollutant in <code>pollutionRose</code> . For <code>windRose</code> and the <code>ws.int</code> default of 2 m/s, the default, 4, generates the break points 2, 4, 6, 8 m/s. For <code>pollutionRose</code> , the default, 6, attempts to breaks the supplied data at approximately 6 sensible break points. However, <code>breaks</code> can also be used to set specific break points. For example, the argument <code>breaks = c(0, 1, 10, 100)</code> breaks the data into segments <1, 1-10, 10-100, >100.
offset	The size of the ‘hole’ in the middle of the plot, expressed as a percentage of the polar axis scale, default 10.
normalise	If <code>TRUE</code> each wind direction segment of a pollution rose is normalised to equal one. This is useful for showing how the concentrations (or other parameters) contribute to each wind sector when the proportion of time the wind is from that direction is low. A line showing the probability that the wind directions is from a particular wind sector is also shown.
max.freq	Controls the scaling used by setting the maximum value for the radial limits. This is useful to ensure several plots use the same radial limits.
paddle	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> plots rose using ‘paddle’ style spokes. If <code>FALSE</code> plots rose using ‘wedge’ style spokes.
key.header	Adds additional text/labels above and/or below the scale key, respectively. For example, passing <code>windRose(mydata, key.header = "ws")</code> adds the addition text as a scale header. Note: This argument is passed to

`drawOpenKey` via `quickText`, applying the `auto.text` argument, to handle formatting.

`key.footer` see `key.footer`.

`key.position` Location where the scale key is to plotted. Allowed arguments currently include “top”, “right”, “bottom” and “left”.

`key` Fine control of the scale key via `drawOpenKey`. See `drawOpenKey` for further details.

`dig.lab` The number of significant figures at which scientific number formatting is used in break point and key labelling. Default 5.

`statistic` The `statistic` to be applied to each data bin in the plot. Options currently include “prop.count”, “prop.mean” and “abs.count”. The default “prop.count” sizes bins according to the proportion of the frequency of measurements. Similarly, “prop.mean” sizes bins according to their relative contribution to the mean. “abs.count” provides the absolute count of measurements in each bin.

`pollutant` Alternative data series to be sampled instead of wind speed. The `windRose` default NULL is equivalent to `pollutant = "ws"`.

`annotate` If `TRUE` then the percentage calm and mean values are printed in each panel together with a description of the statistic below the plot. If “” then only the static is below the plot. Custom annotations may be added by setting value to `c("annotation 1", "annotation 2")`.

`angle.scale` The wind speed scale is by default shown at a 315 degree angle. Sometimes the placement of the scale may interfere with an interesting feature. The user can therefore set `angle.scale` to another value (between 0 and 360 degrees) to mitigate such problems. For example `angle.scale = 45` will draw the scale heading in a NE direction.

`border` Border colour for shaded areas. Default is no border.

... For `pollutionRose` other parameters that are passed on to `windRose`. For `windRose` other parameters that are passed on to `drawOpenKey`, `lattice:xyplot` and `cutData`. Axis and title labelling options (`xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

12.3 Example of use

The function is very simply called as shown for Figure 12.1.

Figure 12.2 highlights some interesting differences between the years. In 2000, for example, there were a large number of occasions when the wind was from the SSW and 2003 clearly had more occasions when the wind was easterly. It can also be useful to use `type = "month"` to get an idea of how wind speed and direction vary seasonally.

The `type` option is very flexible in `openair` and can be used to quickly consider the dependencies between variables. Section 11 describes the basis of this option in `openair` plot. As an example, consider the question: what are the meteorological conditions that control high and low concentrations of PM₁₀? By setting `type = "pm10"`, `openair` will split the PM₁₀ concentrations into four *quantiles* i.e. roughly equal numbers of

```
windRose(mydata)
```

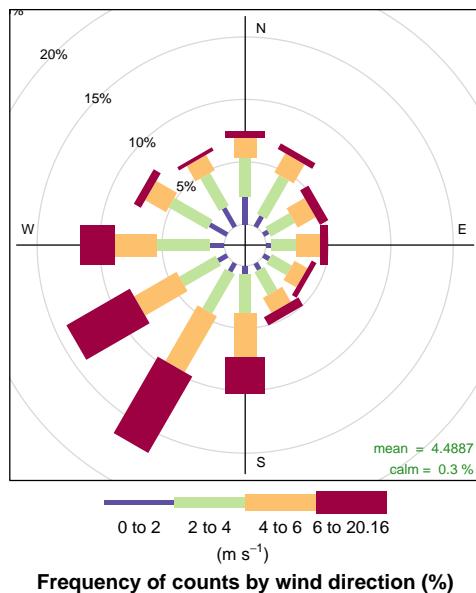


FIGURE 12.1 Use of `windRose` function to plot wind speed/direction frequencies. Wind speeds are split into the intervals shown by the scale in each panel. The grey circles show the % frequencies.

```
windRose(mydata, type = "year", layout = c(4, 2))
```

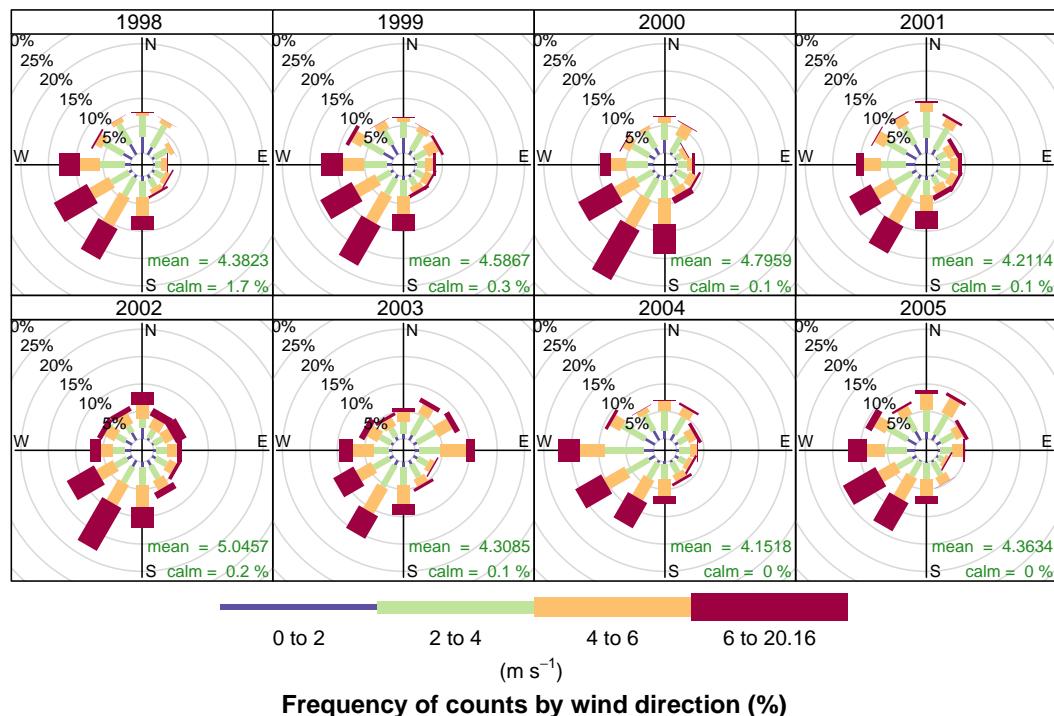


FIGURE 12.2 Use of `windRose` function to plot wind speed/direction frequencies by year. Wind speeds are split into the intervals shown by the scale in each panel. The grey circles show the 10 and 20 % frequencies.

```
windRose(mydata, type = "pm10", layout = c(4, 1))
```

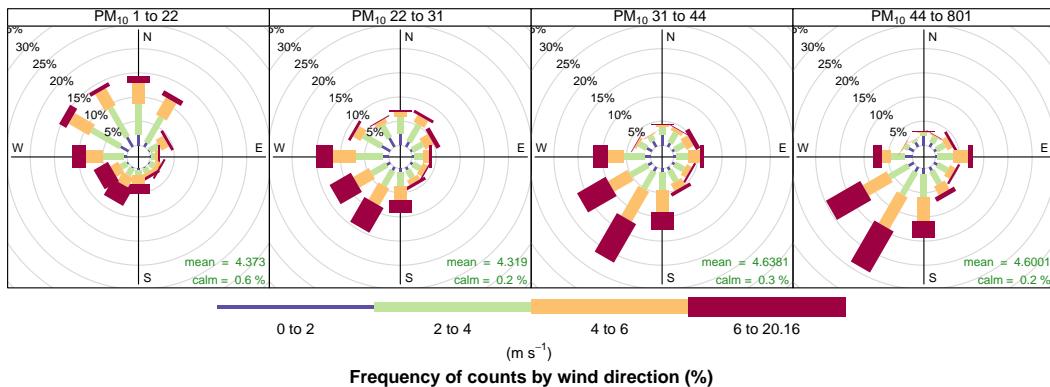


FIGURE 12.3 Wind rose for four different levels of PM₁₀ concentration. The levels are defined as the four *quantiles* of PM₁₀ concentration and the ranges are shown on each of the plot labels.

points in each level. The plot will then show four different wind roses for each quantile level, although the default number of levels can be set by the user — see `?cutData` for more details. Figure 12.3 shows the results of setting `type = "pm10"`. For the lowest concentrations of PM₁₀ the wind direction is dominated by northerly winds, and relatively low wind speeds. By contrast, the highest concentrations (plot furthest right) are dominated by relatively strong winds from the south-west. It is therefore very easy to obtain a good idea about the conditions that tend to lead to high (or low) concentrations of a pollutant. Furthermore, the `type` option is available in almost all `openair` functions.

A comparison of the effect that bias has can be seen by plotting the following. Note the prominent frequencies for W, E and N in particular that are due to the bias issue discussed by Applequist (2012).

```
## no bias correction
windRose(mydata, angle = 22.5, bias.cor = FALSE)

## bias correction (the default)
windRose(mydata, angle = 22.5)
```

`pollutionRose` is a variant of `windRose` that is useful for considering pollutant concentrations by wind direction, or more specifically the percentage time the concentration is in a particular range. This type of approach can be very informative for air pollutant species, as demonstrated by Ronald Henry and co-authors in a recent paper (Henry et al. 2009).

You can produce similar pollution roses using the `pollutionRose` function in recent versions of `openair`, e.g. as in Figure 12.4:

`pollutionRose` is *wrapper* for `windRose`. It simply replaces the wind speed data series in the supplied data set with another variable using the argument `pollutant` before passing that on to `windRose`. It also modifies `breaks` to estimate a sensible set of break points for that pollutant and uses a slightly different set of default options (key to right, wedge style plot) but otherwise handles arguments just like the parent `windRose` function.

While Figure 12.4 indicates that most higher NO_x concentrations are also associated with the SW, conditioning allows you to be much informative. For example, conditioning

```
pollutionRose(mydata, pollutant = "nox")
```

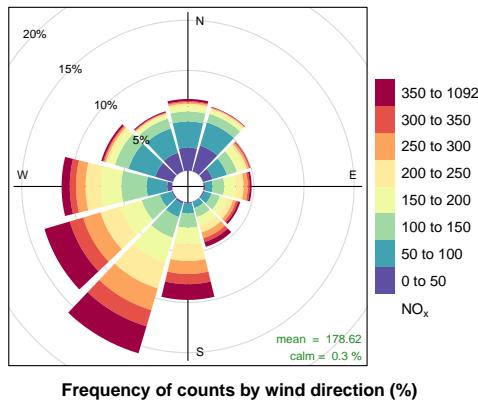


FIGURE 12.4 NO_x pollution rose produced using `pollutionRose` and default `pollutionRose` settings.

```
pollutionRose(mydata, pollutant = "nox", type = "so2", layout = c(4, 1))
```

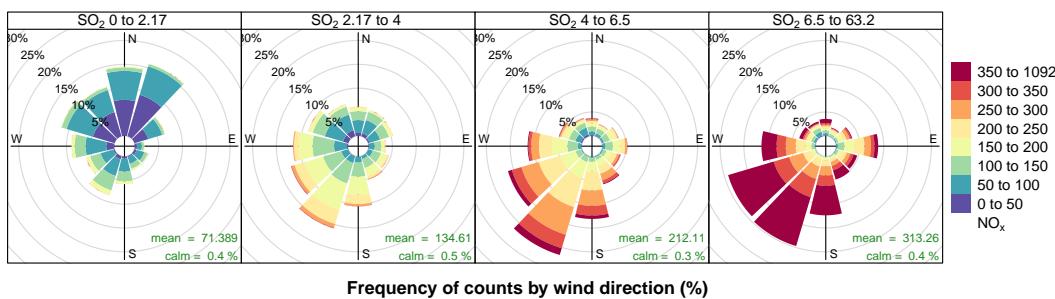


FIGURE 12.5 NO_x pollution rose conditioned by SO₂ concentration.

by SO₂ (Figure 12.5) demonstrates that higher NO_x concentrations are associated with the SW and much of the higher SO₂ concentrations. However, it also highlights a notable NO_x contribution from the E, most apparent at highest SO₂ concentrations that is obscured in Figure 12.4 by a relatively high NO_x background (Figure 12.5).

`pollutionRose` can also usefully be used to show which wind directions dominate the overall concentrations. By supplying the option `statistic = "prop.mean"` (proportion contribution to the mean), a good idea can be gained as to which wind directions contribute most to overall concentrations, as well as providing information on the different concentration levels. A simple plot is shown in Figure 12.6, which clearly shows the dominance of south-westerly winds controlling the overall mean NO_x concentrations at this site. Indeed, almost half the overall NO_x concentration is contributed by two wind sectors to the south-west. The `polarFreq` function can also show this sort of information, but the pollution rose is more effective because both length and colour are used to show the contribution. These plots are very useful for understanding which wind directions control the overall mean concentrations.

It is sometimes useful to more clearly understand the contributions from wind directions that have low frequencies. For example, for a pollution rose of SO₂ there are few occurrences of easterly winds making it difficult to see how the concentration intervals are made up. Try:

```
pollutionRose(mydata, pollutant = "nox", statistic = "prop.mean")
```

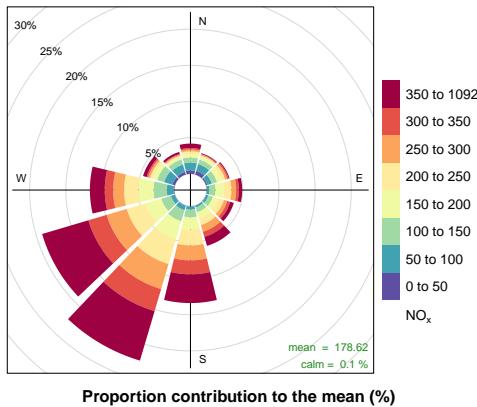


FIGURE 12.6 Pollution rose showing which wind directions contribute most to overall mean concentrations.

```
pollutionRose(mydata, pollutant = "so2", normalise = TRUE, seg = 1)
```

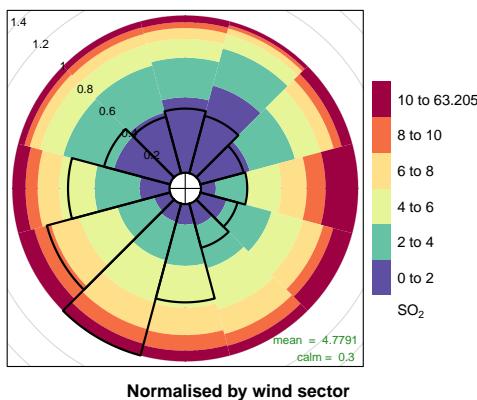


FIGURE 12.7 SO₂ pollution rose produced using `pollutionRose` normalised by each wind sector.

```
pollutionRose(mydata, pollutant = "so2", seg = 1)
```

However, each wind sector can be normalised to give a probability between 0 and 1 to help show the variation within each wind sector more clearly. An example is shown in Figure 12.7 where for easterly winds it is now clearer that a greater proportion of the time the concentration is made up of high SO₂ concentrations. In this plot each wind sector is scaled between 0 and 1. Also shown with a black line is an indication of the wind direction frequency to remind us that winds from the east occur with a low frequency.

Comparing two meteorological data sets

The `pollutionRose` function is also useful for comparing two meteorological data sets. In this case a ‘reference’ data set is compared with a second data set. There are many reasons for doing so e.g. to see how one site compares with another or for meteorological

model evaluation (more on that in later sections). In this case, `ws` and `wd` are considered to the the reference data sets with which a second set of wind speed and wind directions are to be compared (`ws2` and `wd2`). The first set of values is subtracted from the second and the differences compared. If for example, `wd2` was biased positive compared with `wd` then `pollutionRose` will show the bias in polar coordinates. In its default use, wind direction bias is colour-coded to show negative bias in one colour and positive bias in another.

Note that this plot is mostly aimed at showing wind direction biases. It does also show the wind speed bias *but only if there is a wind direction bias also*. However, in most practical situations the plot should show both wind speed and direction biases together. An example of a situation where no wind speed bias would be shown would be for westerly winds where there was absolutely no bias between two data sets in terms of westerly wind direction but there was a difference in wind speed. Users should be aware of this limitation.

In the next example, some artificial wind direction data are generated by adding a positive bias of 30 degrees with some normally distributed scatter. Also, the wind speed data are given a positive bias. The results are shown in [Figure 12.8](#). The Figure clearly shows the mean positive bias in wind direction i.e. the direction is displaced from north (no bias). The colour scale also shows the extent to which wind speeds are biased i.e. there is a higher proportion of positively biased wind speeds shown by the red colour compared with the negatively biased shown in blue. Also shown in [Figure 12.8](#) is the mean wind speed and direction bias as numerical values.

Note that the `type` option can be used in [Figure 12.8](#) e.g. `type = "month"` to split the analysis in useful ways. This is useful if one wanted to see whether a site or the output from a model was biased for different periods. For example, `type = "daylight"` would show whether there are biases between nighttime and daytime conditions.

An example of using user-supplied breaks is shown in [Figure 12.9](#). In this case six intervals are chosen including one that spans -0.5 to $+0.5$ that is useful to show wind speeds that do not change.

```
## $example of comparing 2 met sites
## first we will make some new ws/wd data with a positive bias
mydata <- transform(mydata,
  ws2 = ws + 2 * rnorm(nrow(mydata)) + 1,
  wd2 = wd + 30 * rnorm(nrow(mydata)) + 30)

## need to correct negative wd
id <- which(mydata$wd2 < 0)
mydata$wd2[id] <- mydata$wd2[id] + 360

## results show positive bias in wd and ws
pollutionRose(mydata, ws = "ws", wd = "wd", ws2 = "ws2", wd2 = "wd2", grid.line = 5)
```

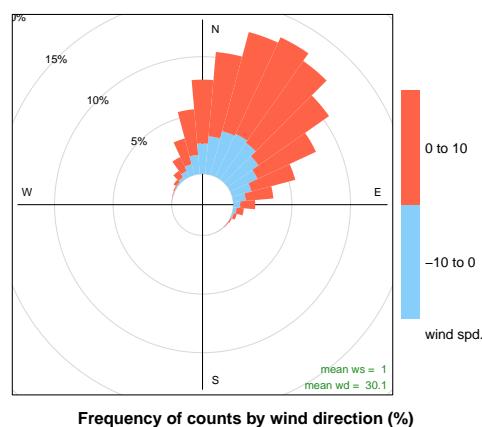


FIGURE 12.8 Pollution rose showing the difference between two meteorological data sets. The colours are used to show whether data tend to be positively or negatively biased with respect to the reference data set.

```

## add some wd bias to some nighttime hours
id <- which(as.numeric(format(mydata$date, "%H")) %in% c(23, 1, 2, 3, 4, 5))
mydata$wd2[id] <- mydata$wd[id] + 30 * rnorm(length(id)) + 120
id <- which(mydata$wd2 < 0)
mydata$wd2[id] <- mydata$wd2[id] + 360

pollutionRose(mydata, ws = "ws", wd = "wd", ws2 = "ws2", wd2 = "wd2",
               breaks = c(-11, -2, -1, -0.5, 0.5, 1, 2, 11),
               cols = c("dodgerblue4", "white", "firebrick"),
               grid.line = 5, type = "daylight")

```

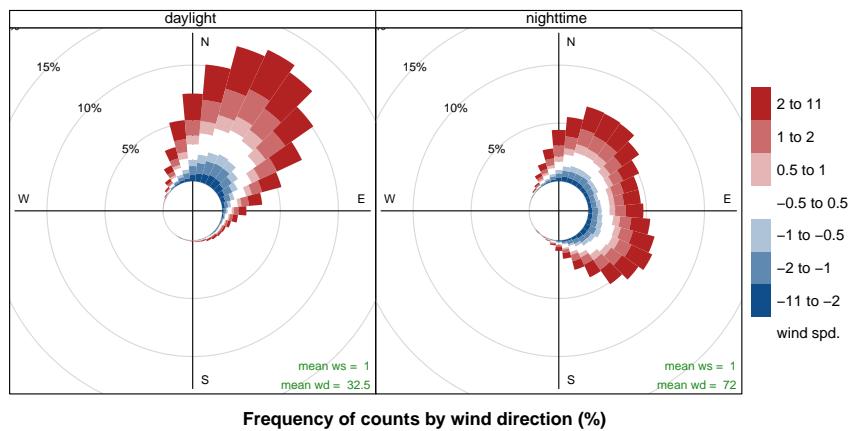


FIGURE 12.9 Pollution rose showing the difference between two meteorological data sets. The colours are used to show whether data tend to be positively or negatively biased with respect to the reference data set. In this case the example shows how to use user-defined breaks and split the data by day/night for a latitude assumed to be London.

13 The percentileRose function

13.1 Purpose

see also
[windRose](#),
[polarPlot](#)
[pollutionRose](#)
[polarAnnulus](#)

`percentileRose` calculates percentile levels of a pollutant and plots them by wind direction. One or more percentile levels can be calculated and these are displayed as either filled areas or as lines.

By default the function plots percentile concentrations in 10 degree segments. Alternatively, the levels by wind direction are calculated using a cyclic smooth cubic spline. The wind directions are rounded to the nearest 10 degrees, consistent with surface data from the UK Met Office before a smooth is fitted.

The `percentileRose` function complements other similar functions including `windRose`, `pollutionRose`, `polarFreq` or `polarPlot`. It is most useful for showing the distribution of concentrations by wind direction and often can reveal different sources e.g. those that only affect high percentile concentrations such as a chimney stack.

Similar to other functions, flexible conditioning is available through the `type` option. It is easy for example to consider multiple percentile values for a pollutant by season, year and so on. See examples below.

13.2 Options available

The `percentileRose` function has the following options:

mydata A data frame minimally containing `wd` and a numeric field to plot — `pollutant`.

pollutant Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. `pollutant = "nox"`. More than one pollutant can be supplied e.g. `pollutant = c("no2", "o3")` provided there is only one `type`.

wd Name of the wind direction field.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

percentile The percentile value(s) to plot. Must be between 0–100. If `percentile = NA` then only a mean line will be shown.

smooth Should the wind direction data be smoothed using a cyclic spline?

<code>method</code>	When <code>method = "default"</code> the supplied percentiles by wind direction are calculated. When <code>method = "cpf"</code> the conditional probability function (CPF) is plotted and a single (usually high) percentile level is supplied. The CPF is defined as $CPF = my/ny$, where my is the number of samples in the wind sector y with mixing ratios greater than the <i>overall</i> percentile concentration, and ny is the total number of samples in the same wind sector (see Ashbaugh et al., 1985).
<code>cols</code>	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours — see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
<code>angle</code>	Default angle of “spokes” is when <code>smooth = FALSE</code> .
<code>mean</code>	Show the mean by wind direction as a line?
<code>mean.lty</code>	Line type for mean line.
<code>mean.lwd</code>	Line width for mean line.
<code>mean.col</code>	Line colour for mean line.
<code>fill</code>	Should the percentile intervals be filled (default) or should lines be drawn (<code>fill = FALSE</code>).
<code>intervals</code>	User-supplied intervals for the scale e.g. <code>intervals = c(0, 10, 30, 50)</code>
<code>angle.scale</code>	The pollutant scale is by default shown at a 45 degree angle. Sometimes the placement of the scale may interfere with an interesting feature. The user can therefore set <code>angle.scale</code> to another value (between 0 and 360 degrees) to mitigate such problems. For example <code>angle.scale = 315</code> will draw the scale heading in a NW direction.
<code>auto.text</code>	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO2.
<code>key.header</code>	Adds additional text/labels to the scale key. For example, passing options <code>key.header = "header"</code> , <code>key.footer = "footer"</code> adds addition text above and below the scale key. These arguments are passed to <code>drawOpenKey</code> via <code>quickText</code> , applying the <code>auto.text</code> argument, to handle formatting.
<code>key.footer</code>	<code>key.header</code> .
<code>key.position</code>	Location where the scale key is to plotted. Allowed arguments currently include <code>"top"</code> , <code>"right"</code> , <code>"bottom"</code> and <code>"left"</code> .
<code>key</code>	Fine control of the scale key via <code>drawOpenKey</code> . See <code>drawOpenKey</code> for further details.
<code>...</code>	Other graphical parameters are passed onto <code>cutData</code> and <code>lattice:xyplot</code> . For example, <code>percentileRose</code> passes the option <code>hemisphere = "southern"</code>

```
percentileRose(mydata, pollutant = "o3")
```

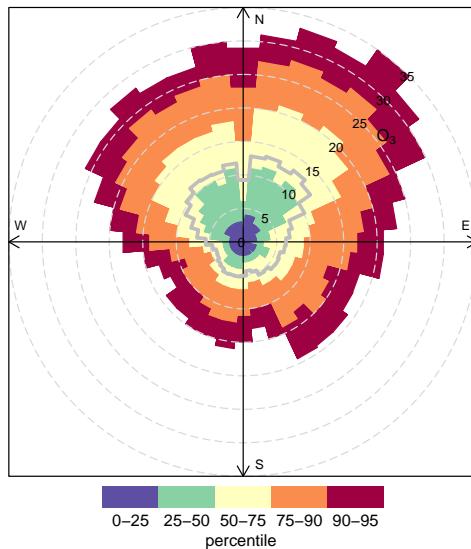


FIGURE 13.1 A `percentileRose` plot of O₃ concentrations at Marylebone Road. The percentile intervals are shaded and are shown by wind direction. It shows for example that higher concentrations occur for northerly winds, as expected at this location. However, it also shows, for example the actual value of the 95th percentile O₃ concentration.

on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common graphical arguments, such as `xlim` and `ylim` for plotting ranges and `lwd` for line thickness when using `fill = FALSE`, are passed on `xyplot`, although some local modifications may be applied by `openair`. For example, axis and title labelling options (such as `xlab`, `ylab` and `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

13.3 Example of use

The first example is a basic plot of percentiles of O₃ shown in Figure 13.1.

A slightly more interesting plot is shown in Figure 13.2 for SO₂ concentrations. We also take the opportunity of changing some of the default options. In this case it can be clearly seen that the highest concentrations of SO₂ are dominated by east and south-easterly winds; likely reflecting the influence of stack emissions in those directions.

A smoothed version of Figure 13.2 can be plotted by:

```
percentileRose(mydata, pollutant = "so2",
               percentile = c(25, 50, 75, 90, 95, 99, 99.9),
               col = "brewer1", key.position = "right",
               smooth = TRUE)
```

Lots more insight can be gained by considering how percentile values vary by other factors i.e. conditioning. For example, what do O₃ concentrations look like split by season and whether it is daylight or nighttime hours? We can set the type to consider season *and* whether it is daylight or nighttime.¹³ This Figure reveals some interesting

¹³In choosing `type = "daylight"` the default is to consider a latitude of central London (or close to). Users can set the latitude in the function call if working in other parts of the world.

```
percentileRose(mydata, pollutant = "so2",
               percentile = c(25, 50, 75, 90, 95, 99, 99.9),
               col = "brewer1", key.position = "right", smooth = TRUE)
```

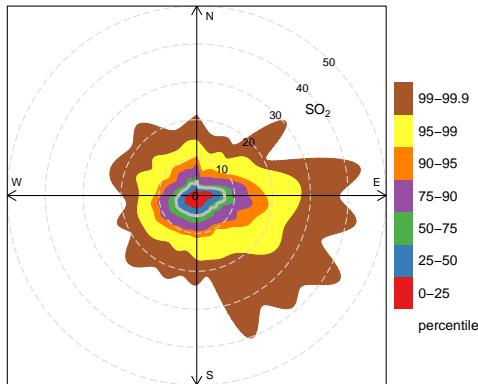


FIGURE 13.2 A `percentileRose` plot of SO₂ concentrations at Marylebone Road. The percentile intervals are shaded and are shown by wind direction. This plot sets some user-defined percentile levels to consider the higher SO₂ concentrations, moves the key to the right and uses an alternative colour scheme.

features. First, O₃ concentrations are higher in the spring and summer and when the wind is from the north. O₃ concentrations are higher on average at this site in spring due to the peak of northern hemispheric O₃ and to some extent local production. This may also explain why O₃ concentrations are somewhat higher at nighttime in spring compared with summer. Second, peak O₃ concentrations are higher during daylight hours in summertime when the wind is from the south-east. This will be due to more local (UK/European) production that is photochemically driven — and hence more important during daylight hours.

The `percentileRose` function can also plot conditional probability functions (CPF) (Ashbaugh et al. 1985). The CPF is defined as $CPF = m_\theta / n_\theta$, where m_θ is the number of samples in the wind sector θ with mixing ratios greater than some ‘high’ concentration, and n_θ is the total number of samples in the same wind sector. CPF analysis is very useful for showing which wind directions are dominated by high concentrations and give the probability of doing so. In `openair`, a CPF plot can be produced as shown in Figure 13.4. Note that in these plots only one percentile is provided and the method must be supplied. In Figure 13.4 it is clear that the high concentrations (greater than the 95th percentile of *all* observations) is dominated by easterly wind directions. There are very low conditional probabilities of these concentrations being experienced for other wind directions.

It is easy to plot several species on the same plot and this works well because they all have the same probability scale (i.e. 0 to 1). In the example below (not shown) it is easy to see for each pollutant the wind directions that dominate the contributions to the highest (95th percentile) concentrations. For example, the highest CO and NO_x concentrations are totally dominated by south/south-westerly winds and the probability of their being such high concentrations from other wind directions is effectively zero.

```
percentileRose(mydata, poll=c("nox", "so2", "o3", "co", "pm10", "pm25"),
               percentile = 95, method = "cpf", col = "darkorange",
               layout = c(2, 3))
```

```
percentileRose(mydata, type = c("season", "daylight"), pollutant = "o3",
               col = "Set3", mean.col = "black")
```

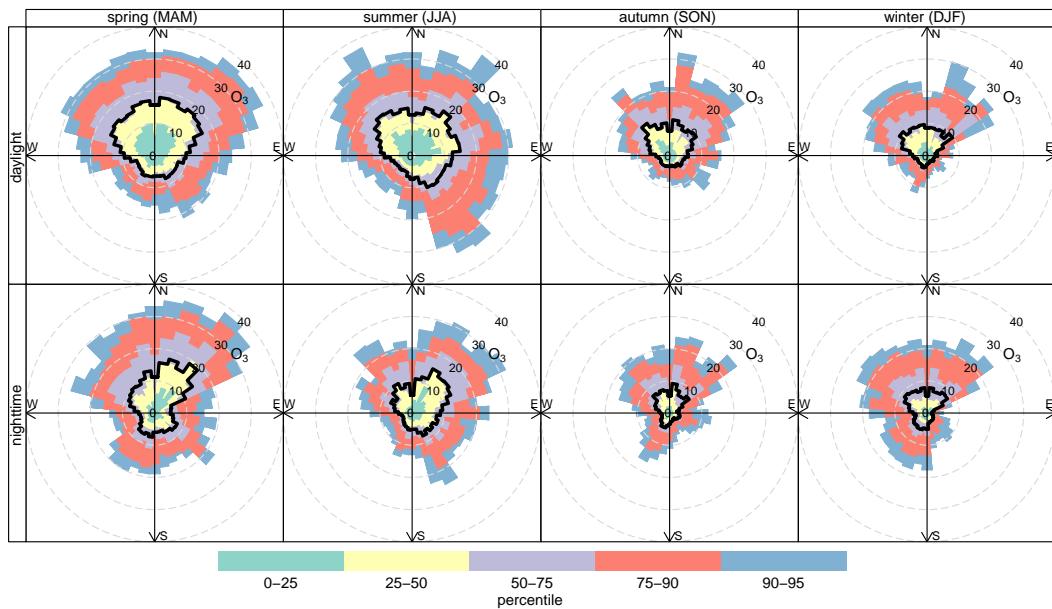


FIGURE 13.3 A `percentileRose` plot of O₃ concentrations at Marylebone Road. The percentile intervals are shaded and are shown by wind direction. The plot show the variation by season and whether it is nighttime or daylight hours.

```
percentileRose(mydata, poll="so2", percentile = 95, method = "cpf",
               col = "darkorange", smooth = TRUE)
```

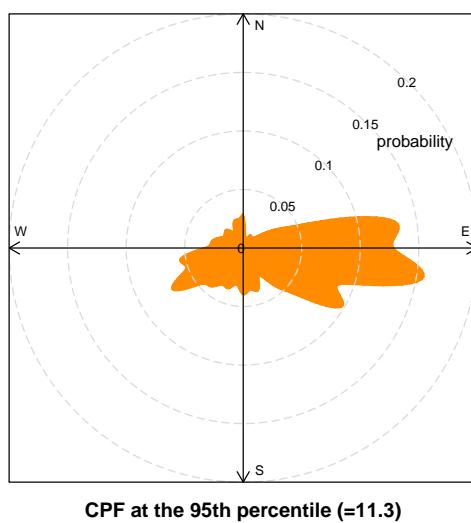


FIGURE 13.4 A CPF plot of SO₂ concentrations at Marylebone Road.

14 The `polarFreq` function

14.1 Purpose

see also
`windRose`
`percentileRose`
`polarPlot`

This is a custom-made plot to compactly show the distribution of wind speeds and directions from meteorological measurements. It is similar to the traditional wind rose, but includes a number of enhancements to also show how concentrations of pollutants and other variables vary. It can summarise all available data, or show it by different time periods e.g. by year, month, day of the week. It can also consider a wide range of statistics.

14.2 Options available

The `polarFreq` function has the following options:

- `mydata`** A data frame minimally containing `ws`, `wd` and `date`.
- `pollutant`** Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. `pollutant = "nox"`
- `statistic`** The statistic that should be applied to each wind speed/direction bin. Can be “frequency”, “mean”, “median”, “max” (maximum), “stdev” (standard deviation) or “weighted.mean”. The option “frequency” (the default) is the simplest and plots the frequency of wind speed/direction in different bins. The scale therefore shows the counts in each bin. The option “mean” will plot the mean concentration of a pollutant (see next point) in wind speed/direction bins, and so on. Finally, “weighted.mean” will plot the concentration of a pollutant weighted by wind speed/direction. Each segment therefore provides the percentage overall contribution to the total concentration. More information is given in the examples. Note that for options other than “frequency”, it is necessary to also provide the name of a pollutant. See function `cutData` for further details.
- `ws.int`** Wind speed interval assumed. In some cases e.g. a low met mast, an interval of 0.5 may be more appropriate.
- `wd.nint`** Number of intervals of wind direction.
- `grid.line`** Radial spacing of grid lines.
- `breaks`** The user can provide their own scale. `breaks` expects a sequence of numbers that define the range of the scale. The sequence could represent one with equal spacing e.g. `breaks = seq(0, 100, 10)` - a scale from 0-10 in intervals of 10, or a more flexible sequence e.g. `breaks = c(0, 1, 5, 7, 10)`, which may be useful for some situations.
- `cols`** Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and `RColorBrewer` colours— see the `openair openColours` function for more details. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`
- `trans`** Should a transformation be applied? Sometimes when producing plots of this kind they can be dominated by a few high points. The default therefore is `TRUE` and a square-root transform is applied. This results in a

non-linear scale and (usually) a better representation of the distribution. If set to `FALSE` a linear scale is used.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

min.bin The minimum number of points allowed in a wind speed/wind direction bin. The default is 1. A value of two requires at least 2 valid records in each bin and so on; bins with less than 2 valid records are set to NA. Care should be taken when using a value > 1 because of the risk of removing real data points. It is recommended to consider your data with care. Also, the `polarPlot` function can be of use in such circumstances.

ws.upper A user-defined upper wind speed to use. This is useful for ensuring a consistent scale between different plots. For example, to always ensure that wind speeds are displayed between 1-10, set `ws.int = 10`.

offset `offset` controls the size of the ‘hole’ in the middle and is expressed as a percentage of the maximum wind speed. Setting a higher `offset` e.g. 50 is useful for `statistic = "weighted.mean"` when `ws.int` is greater than the maximum wind speed. See example below.

border.col The colour of the boundary of each wind speed/direction bin. The default is transparent. Another useful choice sometimes is “white”.

key.header, key.footer Adds additional text/labels to the scale key. For example, passing options `key.header = "header", key.footer = "footer"` adds addition text above and below the scale key. These arguments are passed to `drawOpenKey` via `quickText`, applying the `auto.text` argument, to handle formatting.

key.position Location where the scale key is to plotted. Allowed arguments currently include `"top", "right", "bottom" and "left"`.

key Fine control of the scale key via `drawOpenKey`. See `drawOpenKey` for further details.

auto.text Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO₂.

```
polarFreq(mydata)
```

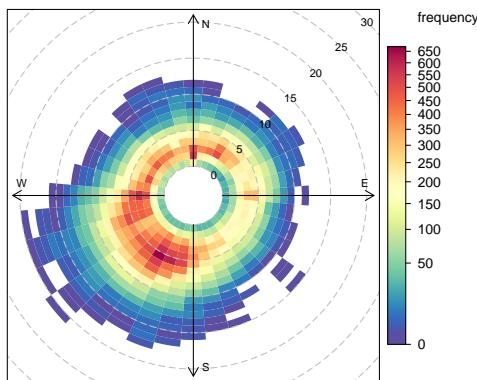


FIGURE 14.1 Use of `polarFreq` function to plot wind speed/directions. Each cell gives the total number of hours the wind was from that wind speed/direction in a particular year. The number of hours is coded as a colour scale shown to the right. The scale itself is non-linear to help show the overall distribution. The dashed circular grey lines show the wind speed scale. The date range covered by the data is shown in the strip.

...

Other graphical parameters passed onto `lattice:xyplot` and `cutData`. For example, `polarFreq` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

For `type = "site"`, it is necessary to format the input data into columns `date`, `ws`, `wd`, `site` (and maybe `pollutant`). This means that `date`, for example is repeated a number of times equal to the number of sites.

14.3 Example of use

This section shows an example output and use, using our data frame `mydata`. The function is very simply run as shown in Figure 14.1. This produces the plot shown in Figure 14.1.

By setting `type = "year"`, the frequencies are shown separately by year as shown in Figure 14.2, which shows that most of the time the wind is from a south-westerly direction with wind speeds most commonly between $2\text{--}6 \text{ m s}^{-1}$. In 2000 there seemed to be a lot of conditions where the wind was from the south-west (leading to high pollutant concentrations at this location). The data for 2003 also stand out due to the relatively large number of occasions where the wind was from the east. Note the default colour scale, which has had a square-root transform applied, is used to provide a better visual distribution of the data.

The `polarFreq` function can also usefully consider pollutant concentrations. Figure 14.3 shows the mean concentration of SO_2 by wind speed and wind direction and clearly highlights that SO_2 concentrations tend to be highest for easterly winds and for 1998 in particular.

By weighting the concentrations by the frequency of occasions the wind is from a certain direction and has a certain speed, gives a better indication of the conditions that dominate the overall mean concentrations. Figure 14.4 shows the weighted mean

```
polarFreq(mydata, type = "year")
```

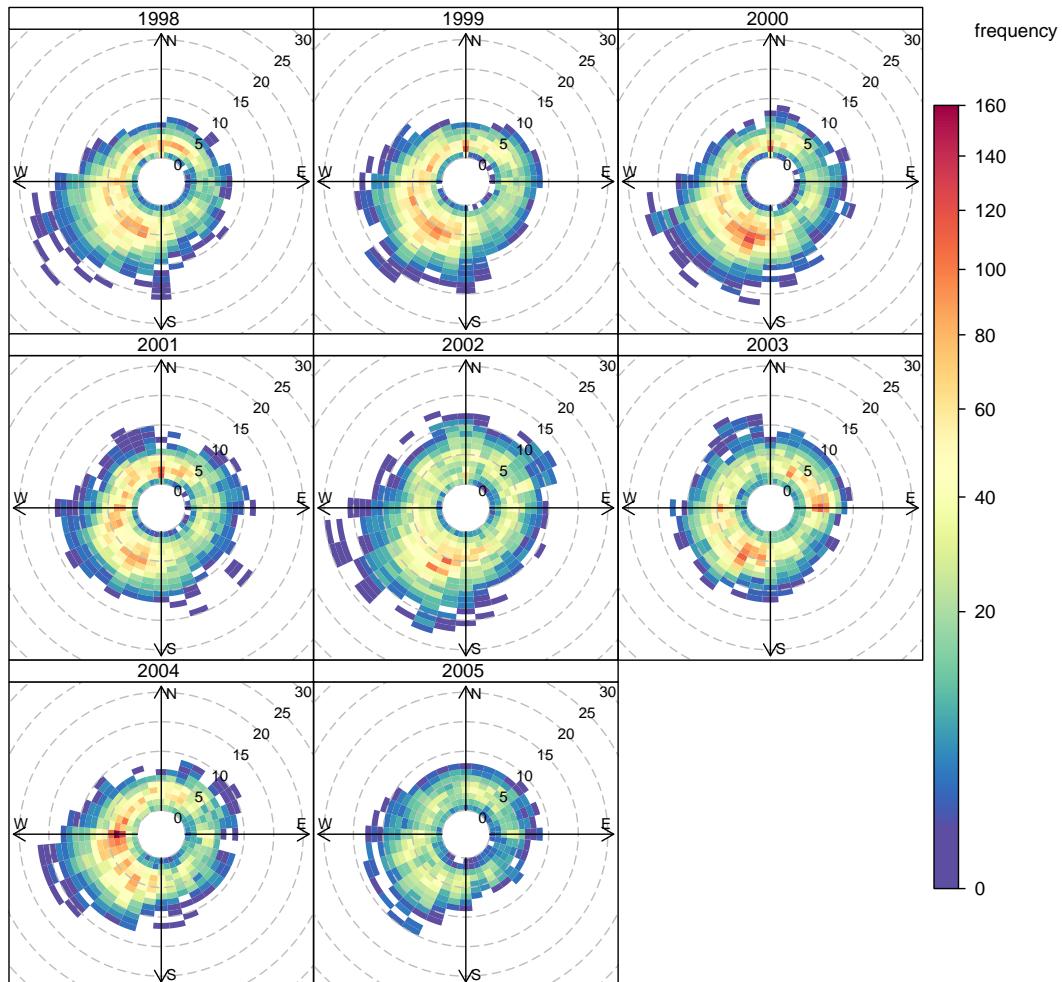


FIGURE 14.2 Use of **polarFreq** function to plot wind speed/directions by year. Each cell gives the total number of hours the wind was from that wind speed/direction in a particular year. The number of hours is coded as a colour scale shown to the right. The scale itself is non-linear to help show the overall distribution. The dashed circular grey lines show the wind speed scale.

concentration of SO₂ and highlights that annual mean concentrations are dominated by south-westerly winds i.e. contributions from the road itself — and not by the fewer higher hours of concentrations when the wind is easterly. However, 2003 looks interesting because for that year, significant contributions to the overall mean were due to easterly wind conditions.

These plots when applied to other locations can reveal some useful features about different sources. For example, it may be that the highest concentrations measured only occur infrequently, and the weighted mean plot can help show this.

The code required to make Figure 14.3 and 14.4 is shown below.

Users are encouraged to try out other plot options. However, one potentially useful plot is to select a few specific years of user interest. For example, what if you just wanted to compare two years e.g. 2000 and 2003? This is easy to do by sending a subset of data to the function. Use here can be made of the **openair selectByDate** function.

```
polarFreq(mydata, pollutant = "so2", type = "year",
statistic = "mean", min.bin = 2)
```

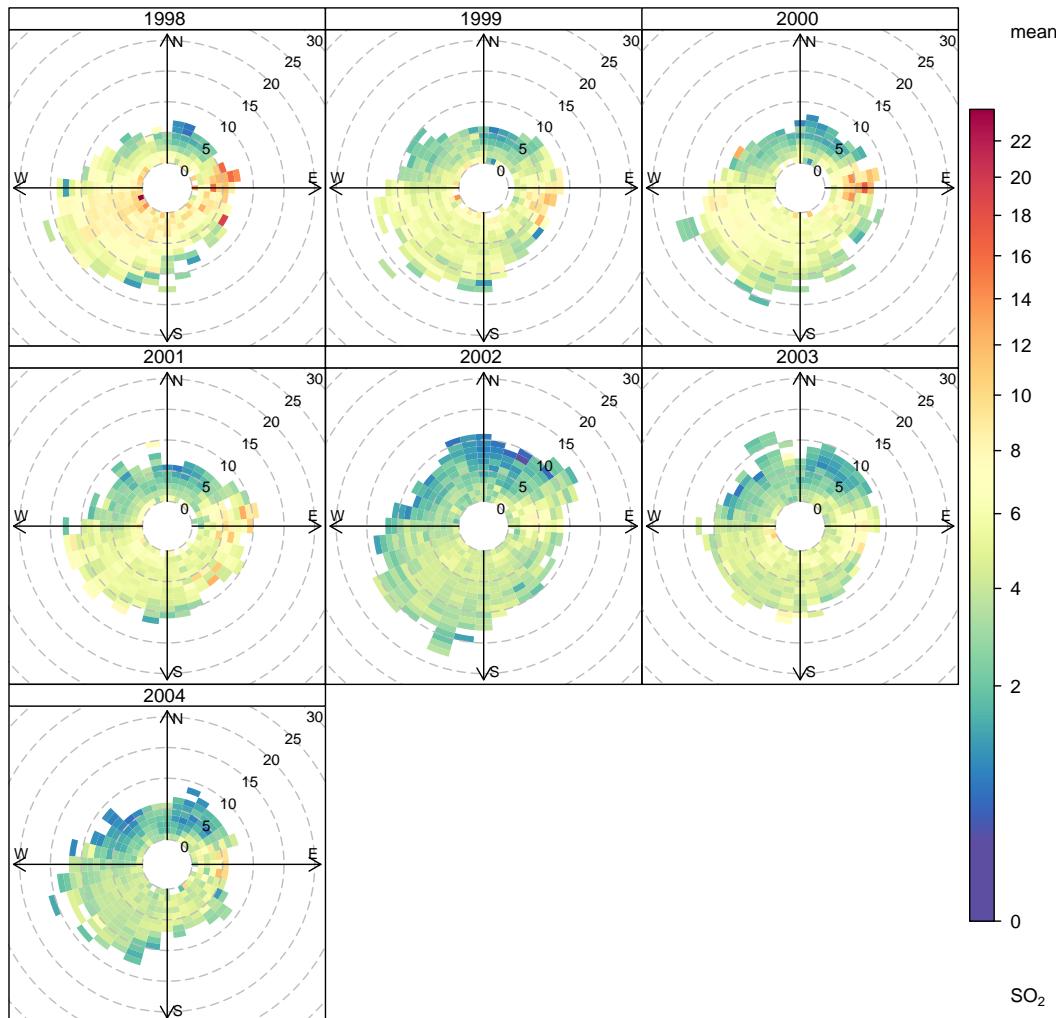


FIGURE 14.3 Use of **polarFreq** function to plot mean SO_2 concentrations (ppb) by wind speed/directions and year.

```
# wind rose for just 2000 and 2003
polarFreq(selectByDate(mydata, year = c(2000, 2003)), cols = "jet",
type = "year")
```

The **polarFreq** function can also be used to gain an idea about the wind directions that contribute most to the overall mean concentrations. As already shown, use of the option **statistic = "weighted.mean"** will show the percentage contribution by wind direction and wind speed bin. However, often it is unnecessary to consider different wind speed intervals. To make the plot more effective, a few options are set as shown in Figure 14.5. First, the **statistic = "weighted.mean"** is chosen to ensure that the plot shows concentrations weighted by their frequency of occurrence of wind direction. For this plot, we are mostly interested in just the contribution by wind direction and not wind speed. Setting the **ws.int** to be above the maximum wind speed in the data set ensures that all data are shown in one interval. Rather than having a square-root transform applied to the colour scale, we choose to have a linear scale by setting **trans = FALSE**. Finally, to show a ‘disk’, the **offset** is set at 80. Increasing

```
# weighted mean SO2 concentrations
polarFreq(mydata, pollutant = "so2", type = "year",
           statistic = "weighted.mean", min.bin = 2)
```

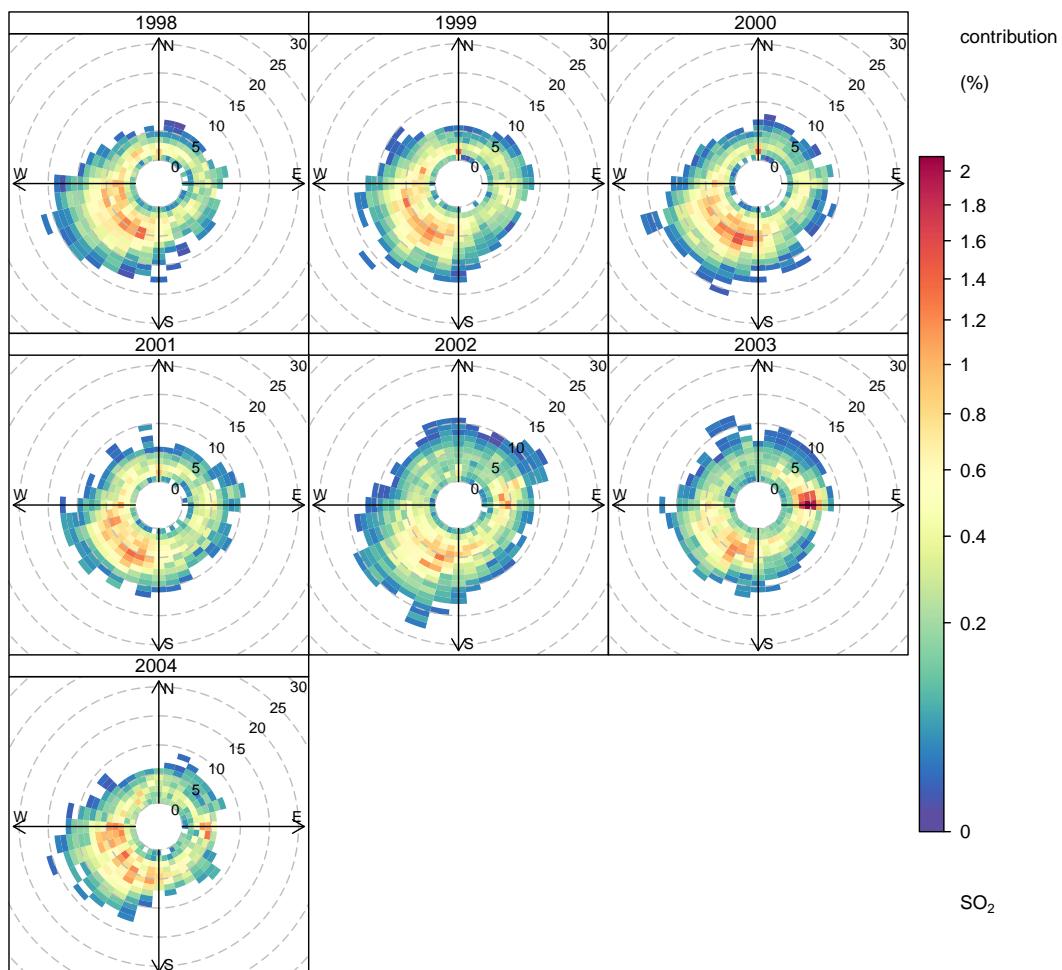


FIGURE 14.4 Use of **polarFreq** function to plot weighted mean SO₂ concentrations (ppb) by wind speed/directions and year.

the value of the offset will narrow the disk.

While Figure 14.5 is useful — e.g. it clearly shows that concentrations of NO_x at this site are totally dominated by south-westerly winds, the use of **pollutionRose** for this type of plot is more effective, as shown in Section 12.

```
polarFreq(mydata, pollutant = "nox", ws.int = 30, statistic = "weighted.mean",
offset = 80, trans = FALSE, col = "heat")
```

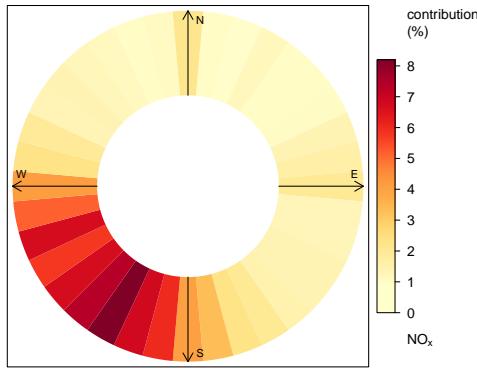


FIGURE 14.5 The percentage contribution to overall mean concentrations of NO_x at Marylebone Road.

15 The **polarPlot** and **polarCluster** functions

15.1 Purpose

see also
polarFreq
polarAnnulus
percentileRose
pollutionRose

The **polarPlot** function plots a *bivariate polar plot* of concentrations. Concentrations are shown to vary by wind speed and wind direction. In many respects they are similar to the plots shown in (§14) but include some additional enhancements. These enhancements include: plots are shown as a continuous surface and surfaces are calculated through modelling using smoothing techniques. These plots are not entirely new as others have considered the joint wind speed-direction dependence of concentrations (see for example Yu et al. (2004)). However, plotting the data in polar coordinates and for the purposes of source identification is new. Furthermore, the basic polar plot is since been enhanced in many ways as described below. Recent publications that describe or use the technique are Carslaw et al. (2006) and Westmoreland et al. (2007). These plots have proved to be useful for quickly gaining a graphical impression of potential sources influences at a location.

The **polarPlot** function is described in more detail in Carslaw et al. (2006) where it is used to triangulate sources in an airport setting, Carslaw and Beevers (2013) where it is used with a clustering technique to identify features in a polar plot with similar characteristics and Uriarte-Tellaetxe and Carslaw (2014) where it is extended to include a conditional probability function to extract more information from the plots. The latter paper is Open Access and can be downloaded from the [openair](#) website.

For many, maybe most situations, increasing wind speed generally results in lower concentrations due to increased dilution through advection and increased mechanical turbulence. There are, however, many processes that can lead to interesting concentration-wind speed dependencies and we will provide a more theoretical treatment of this in due course. However, below are a few reasons why concentrations can change with increasing wind speeds.

- Buoyant plumes from tall stacks can be brought down to ground-level resulting in high concentrations under high wind speed conditions.
- Particle suspension increases with increasing wind speeds e.g. PM_{10} from spoil heaps and the like.

- ‘Particle’ suspension can be important close to coastal areas where higher wind speeds generate more sea spray.
- The wind speed dependence of concentrations in a street canyon can be very complex: higher wind speeds do not always result in lower concentrations due to re-circulation. Bivariate polar plots are very good at revealing these complexities.
- As Carslaw et al. (2006) showed, aircraft emissions have an unusual wind speed dependence and this can help distinguish them from other sources. If several measurement sites are available, polar plots can be used to triangulate different sources.
- Concentrations of NO₂ can increase with increasing wind speed — or at least not decline steeply due to increased mixing. This mixing can result in O₃-rich air converting NO to NO₂.

The function has been developed to allow variables other than wind speed to be plotted with wind direction in polar coordinates. The key issue is that the other variable plotted against wind direction should be discriminating in some way. For example, temperature can help reveal high-level sources brought down to ground level in unstable atmospheric conditions, or show the effect a source emission dependent on temperature e.g. biogenic isoprene. For research applications where many more variables could be available, discriminating sources by these other variables could be very insightful.

Bivariate polar plots are constructed in the following way. First, wind speed, wind direction and concentration data are partitioned into wind speed-direction bins and the mean concentration calculated for each bin. Testing on a wide range of data suggests that wind direction intervals at 10 degrees and 30 wind speed intervals capture sufficient detail of the concentration distribution. The wind direction data typically available are generally rounded to 10 degrees and for typical surface measurements of wind speed in the range 0 to 20 to 30 m s⁻¹, intervals greater than 30 would be difficult to justify based on a consideration of the accuracy of the instruments. Binning the data in this way is not strictly necessary but acts as an effective data reduction technique without affecting the fidelity of the plot itself. Furthermore, because of the inherent wind direction variability in the atmosphere, data from several weeks, months or years typically used to construct a bivariate polar plot tends to be diffuse and does not vary abruptly with either wind direction or speed and more finely resolved bin sizes or working with the raw data directly does not yield more information.

The wind components, u and v are calculated i.e.

$$u = \bar{u} \cdot \sin\left(\frac{2\pi}{\theta}\right), v = \bar{u} \cdot \cos\left(\frac{2\pi}{\theta}\right) \quad (1)$$

with \bar{u} is the mean hourly wind speed and θ is the mean wind direction in degrees with 90 degrees as being from the east.

The calculations above provides a u, v, concentration (C) surface. While it would be possible to work with this surface data directly a better approach is to apply a model to the surface to describe the concentration as a function of the wind components u and v to extract real source features rather than noise. A flexible framework for fitting a surface is to use a Generalized Additive Model (GAM) e.g. (Hastie and Tibshirani 1990; Wood 2006). GAMs are a useful modelling framework with respect to air pollution prediction because typically the relationships between variables are non-linear and variable interactions are important, both of which issues can be addressed in a GAM framework. GAMs can be expressed as shown in Equation 2:

$$\sqrt{C_i} = \beta_0 + \sum_{j=1}^n s_j(x_{ij}) + e_i \quad (2)$$

where C_i is the i th pollutant concentration, β_0 is the overall mean of the response, $s_j(x_{ij})$ is the smooth function of i th value of covariate j , n is the total number of covariates, and e_i is the i th residual. Note that C_i is square-root transformed as the transformation generally produces better model diagnostics e.g. normally distributed residuals.

The model chosen for the estimate of the concentration surface is given by [Equation 3](#). In this model the square root-transformed concentration is a smooth function of the bivariate wind components u and v . Note that the smooth function used is *isotropic* because u and v are on the same scales. The isotropic smooth avoids the potential difficulty of smoothing two variables on different scales e.g. wind speed and direction, which introduces further complexities.

$$\sqrt{C_i} = s(u, v) + e_i \quad (3)$$

15.2 Options available

The `polarPlot` function has the following options:

mydata A data frame minimally containing `wd`, another variable to plot in polar coordinates (the default is a column “ws” — wind speed) and a pollutant. Should also contain `date` if plots by time period are required.

pollutant Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. `pollutant = "nox"`. There can also be more than one pollutant specified e.g. `pollutant = c("nox", "no2")`. The main use of using two or more pollutants is for model evaluation where two species would be expected to have similar concentrations. This saves the user stacking the data and it is possible to work with columns of data directly. A typical use would be `pollutant = c("obs", "mod")` to compare two columns “obs” (the observations) and “mod” (modelled values). When pair-wise statistics such as Pearson correlation and regression techniques are to be plotted, `pollutant` takes two elements too. For example, `pollutant = c("bc", "pm25")` where “bc” is a function of “pm25”.

x Name of variable to plot against wind direction in polar coordinates, the default is wind speed, “ws”.

wd Name of wind direction field.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

statistic The statistic that should be applied to each wind speed/direction bin. Because of the smoothing involved, the colour scale for some of these statistics is only to provide an indication of overall pattern and should not be interpreted in concentration units e.g. for `statistic = "weighted.mean"` where the bin mean is multiplied by the bin frequency and divided by the total frequency. In many cases using `polarFreq` will be better. Setting `statistic = "weighted.mean"` can be useful because it provides an indication of the concentration * frequency of occurrence and will highlight the wind speed/direction conditions that dominate the overall mean. Can be:

- “mean” (default), “median”, “max” (maximum), “frequency”, “stdev” (standard deviation), “weighted.mean”
- `statistic = "cpf"` the conditional probability function (CPF) is plotted and a single (usually high) percentile level is supplied. The CPF is defined as $CPF = my/ny$, where my is the number of samples in the y bin (by default a wind direction, wind speed interval) with mixing ratios greater than the *overall* percentile concentration, and ny is the total number of samples in the same wind sector (see Ashbaugh et al., 1985). Note that percentile intervals can also be considered; see `percentile` for details.
- When `statistic = "r"`, the Pearson correlation coefficient is calculated for *two* pollutants. The calculation involves a weighted Pearson correlation coefficient, which is weighted by Gaussian kernels for wind direction and the radial variable (by default wind speed). More weight is assigned to values close to a wind speed-direction interval. Kernel weighting is used to ensure that all data are used rather than relying on the potentially small number of values in a wind speed-direction interval.
- `"robust.slope"` is another option for pair-wise statistics and `"quantile.slope"`, which uses quantile regression to estimate the slope for a particular quantile level (see also `tau` for setting the quantile level).

resolution Two plot resolutions can be set: “normal” and “fine” (the default), for a smoother plot. It should be noted that plots with a “fine” resolution can take longer to render.

limits The function does its best to choose sensible limits automatically. However, there are circumstances when the user will wish to set different ones. An example would be a series of plots showing each year of data separately. The limits are set in the form `c(lower, upper)`, so `limits = c(0, 100)` would force the plot limits to span 0-100.

exclude.missing Setting this option to `TRUE` (the default) removes points from the plot that are too far from the original data. The smoothing routines will produce predictions at points where no data exist i.e. they predict. By

removing the points too far from the original data produces a plot where it is clear where the original data lie. If set to `FALSE` missing data will be interpolated.

uncertainty Should the uncertainty in the calculated surface be shown? If `TRUE` three plots are produced on the same scale showing the predicted surface together with the estimated lower and upper uncertainties at the 95 Calculating the uncertainties is useful to understand whether features are real or not. For example, at high wind speeds where there are few data there is greater uncertainty over the predicted values. The uncertainties are calculated using the GAM and weighting is done by the frequency of measurements in each wind speed-direction bin. Note that if uncertainties are calculated then the type is set to "default".

percentile If `statistic = "percentile"` then `percentile` is used, expressed from 0 to 100. Note that the percentile value is calculated in the wind speed, wind direction 'bins'. For this reason it can also be useful to set `min.bin` to ensure there are a sufficient number of points available to estimate a percentile. See `quantile` for more details of how percentiles are calculated.

`percentile` is also used for the Conditional Probability Function (CPF) plots. `percentile` can be of length two, in which case the percentile *interval* is considered for use with CPF. For example, `percentile = c(90, 100)` will plot the CPF for concentrations between the 90 and 100th percentiles. Percentile intervals can be useful for identifying specific sources. In addition, `percentile` can also be of length 3. The third value is the 'trim' value to be applied. When calculating percentile intervals many can cover very low values where there is no useful information. The trim value ensures that values greater than or equal to the trim * mean value are considered *before* the percentile intervals are calculated. The effect is to extract more detail from many source signatures. See the manual for examples. Finally, if the trim value is less than zero the percentile range is interpreted as absolute concentration values and subsetting is carried out directly.

cols Colours to be used for plotting. Options include "default", "increment", "heat", "jet" and `RColorBrewer` colours—see the `openair openColours` function for more details. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`. `cols` can also take the values "viridis", "magma", "inferno", or "plasma" which are the viridis colour maps ported from Python's Matplotlib library.

weights At the edges of the plot there may only be a few data points in each wind speed-direction interval, which could in some situations distort the plot if the concentrations are high. `weights` applies a weighting to reduce their influence. For example and by default if only a single data point exists then the weighting factor is 0.25 and for two points 0.5. To not apply any weighting and use the data as is, use `weights = c(1, 1, 1)`.

An alternative to down-weighting these points they can be removed altogether using `min.bin`.

<code>min.bin</code>	The minimum number of points allowed in a wind speed/wind direction bin. The default is 1. A value of two requires at least 2 valid records in each bin and so on; bins with less than 2 valid records are set to NA. Care should be taken when using a value > 1 because of the risk of removing real data points. It is recommended to consider your data with care. Also, the <code>polarFreq</code> function can be of use in such circumstances.
<code>mis.col</code>	When <code>min.bin</code> is > 1 it can be useful to show where data are removed on the plots. This is done by shading the missing data in <code>mis.col</code> . To not highlight missing data when <code>min.bin</code> > 1 choose <code>mis.col = "transparent"</code> .
<code>alpha</code>	The alpha transparency to use for the plotting surface (a value between 0 and 1 with zero being fully transparent and 1 fully opaque). Setting a value below 1 can be useful when plotting surfaces on a map using the package <code>openairmapss</code> .
<code>upper</code>	This sets the upper limit wind speed to be used. Often there are only a relatively few data points at very high wind speeds and plotting all of them can reduce the useful information in the plot.
<code>angle.scale</code>	The wind speed scale is by default shown at a 315 degree angle. Sometimes the placement of the scale may interfere with an interesting feature. The user can therefore set <code>angle.scale</code> to another value (between 0 and 360 degrees) to mitigate such problems. For example <code>angle.scale = 45</code> will draw the scale heading in a NE direction.
<code>units</code>	The units shown on the polar axis scale.
<code>force.positive</code>	The default is <code>TRUE</code> . Sometimes if smoothing data with steep gradients it is possible for predicted values to be negative. <code>force.positive = TRUE</code> ensures that predictions remain positive. This is useful for several reasons. First, with lots of missing data more interpolation is needed and this can result in artifacts because the predictions are too far from the original data. Second, if it is known beforehand that the data are all positive, then this option carries that assumption through to the prediction. The only likely time where setting <code>force.positive = FALSE</code> would be if background concentrations were first subtracted resulting in data that is legitimately negative. For the vast majority of situations it is expected that the user will not need to alter the default option.
<code>k</code>	This is the smoothing parameter used by the <code>gam</code> function in package <code>mgcv</code> . Typically, value of around 100 (the default) seems to be suitable and will resolve important features in the plot. The most appropriate choice of <code>k</code> is problem-dependent; but extensive testing of polar plots for many different problems suggests a value of <code>k</code> of about 100 is suitable. Setting <code>k</code> to higher values will not tend to affect the surface predictions by much but will add to the computation time. Lower values of <code>k</code> will increase smoothing. Sometimes with few data to plot <code>polarPlot</code> will fail. Under these circumstances it can be worth lowering the value of <code>k</code> .
<code>normalise</code>	If <code>TRUE</code> concentrations are normalised by dividing by their mean value. This is done <i>after</i> fitting the smooth surface. This option is particularly useful if one is interested in the patterns of concentrations for several pollutants on different scales e.g. NOx and CO. Often useful if more than one <code>pollutant</code> is chosen.

key.header Adds additional text/labels to the scale key. For example, passing the options `key.header = "header", key.footer = "footer1"` adds addition text above and below the scale key. These arguments are passed to `drawOpenKey` via `quickText`, applying the `auto.text` argument, to handle formatting.

key.footer see `key.footer`.

key.position Location where the scale key is to plotted. Allowed arguments currently include `"top"`, `"right"`, `"bottom"` and `"left"`.

key Fine control of the scale key via `drawOpenKey`. See `drawOpenKey` for further details.

auto.text Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO₂.

ws_spread An integer used for the weighting kernel spread for wind speed when correlation or regression techniques are used. Default is `15`.

wd_spread An integer used for the weighting kernel spread for wind direction when correlation or regression techniques are used. Default is `4`.

kernel Type of kernel used for the weighting procedure for when correlation or regression techniques are used. Only `"gaussian"` is supported but this may be enhanced in the future.

tau The quantile to be estimated when `statistic` is set to `"quantile.slope"`. Default is `0.5` which is equal to the median and will be ignored if `"quantile.slope"` is not used.

... Other graphical parameters passed onto `lattice:levelplot` and `cutData`. For example, `polarPlot` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `levelplot` via `quickText` to handle routine formatting.

15.3 Example of use

We first use the function in its simplest form to make a polar plot of NO_x. The code is very simple as shown in Figure 15.1.

This produces Figure 15.1. The scale is automatically set using whatever units the original data are in. This plot clearly shows highest NO_x concentrations when the wind is from the south-west. Given that the monitor is on the *south* side of the street and the highest concentrations are recorded when the wind is blowing *away* from the monitor is strong evidence of street canyon recirculation.

Figure 15.2 and Figure 15.3 shows polar plots using different defaults and for other pollutants. In the first (Figure 15.2, a different colour scheme is used and some adjustments are made to the key. In Figure 15.3, SO₂ concentrations are shown. What is interesting about this plot compared with either Figure 15.2 or Figure 15.1 is that the concentration pattern is very different i.e. high concentrations with high wind speeds from the east. The most likely source of this SO₂ are industrial sources to the east of

```
polarPlot(mydata, pollutant = "nox")
```

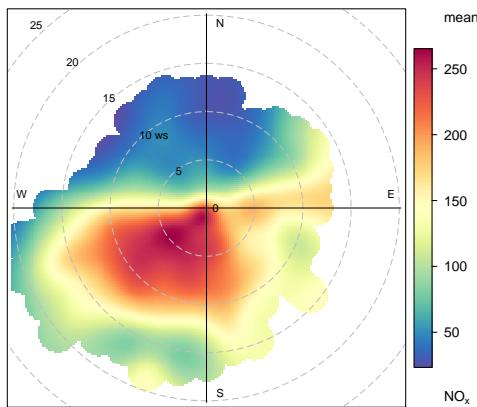


FIGURE 15.1 Default use of the `polarPlot` function applied to Marylebone Road NO_x concentrations.

London. The plot does still however show evidence of a source to the south-west, similar to the plot for NO_x , which implies that road traffic sources of SO_2 can also be detected.

These plots often show interesting features at higher wind speeds. For these conditions there can be very few measurements and therefore greater uncertainty in the calculation of the surface. There are several ways in which this issue can be tackled. First, it is possible to avoid smoothing altogether and use `polarFreq`. The problem with this approach is that it is difficult to know how best to bin wind speed and direction: the choice of interval tends to be arbitrary. Second, the effect of setting a minimum number of measurements in each wind speed-direction bin can be examined through `min.bin`. It is possible that a single point at high wind speed conditions can strongly affect the surface prediction. Therefore, setting `min.bin = 3`, for example, will remove all wind speed-direction bins with fewer than 3 measurements *before* fitting the surface. This is a useful strategy for testing how sensitive the plotted surface is to the number of measurements available. While this is a useful strategy to get a feel for how the surface changes with different `min.bin` settings, it is still difficult to know how many points should be used as a minimum. Third, consider setting `uncertainty = TRUE`. This option will show the predicted surface together with upper and lower 95% confidence intervals, which take account of the frequency of measurements. The uncertainty approach ought to be the most robust and removes any arbitrary setting of other options. There is a close relationship between the amount of smoothing and the uncertainty: more smoothing will tend to reveal less detail and lower uncertainties in the fitted surface and vice-versa.

The default however is to down-weight the bins with few data points when fitting a surface. Weights of 0.25, 0.5 and 0.75 are used for bins containing 1, 2 and 3 data points respectively. The advantage of this approach is that no data are actually removed (which is what happens when using `min.bin`). This approach should be robust in a very wide range of situations and is also similar to the approaches used when trying to locate sources when using back trajectories as described in Section 26. Users can ignore the automatic weighting by supplying the option `weights = c(1, 1, 1)`.

A very useful approach for understanding air pollution is to consider ratios of pollutants. One reason is that pollutant ratios can be largely independent of meteorological variation. In many circumstances it is possible to gain a lot of insight into sources if pollutant ratios are considered. First, it is necessary to calculate a ratio, which is easy in

`polarFreq`
provides an
un-smoothed
surface

```
## NOx plot
polarPlot(mydata, pollutant = "nox", col = "jet", key.position = "bottom",
key.header = "mean nox ( $\mu\text{g}/\text{m}^3$ )", key.footer = NULL)
```

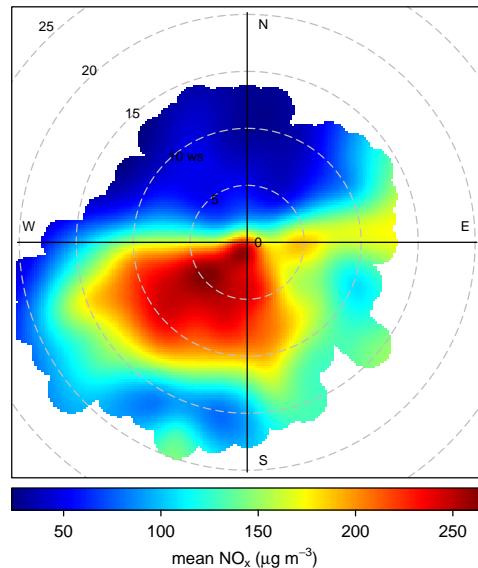


FIGURE 15.2 Example plots using the **polarPlot** function with different options for the mean concentration of NO_x .

```
polarPlot(mydata, pollutant = "so2")
```

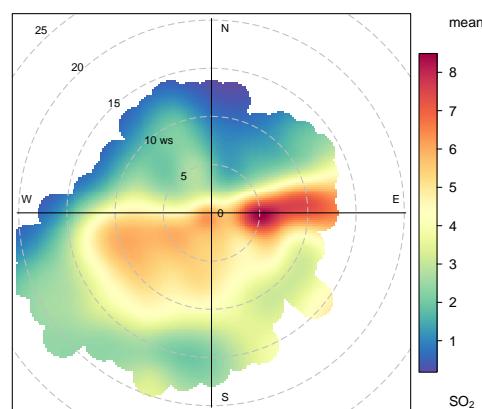


FIGURE 15.3 Example plots using the **polarPlot** function for the mean concentration of SO_2 .

```
polarPlot(mydata, pollutant = "ratio", main = "so2/nox ratio")
```

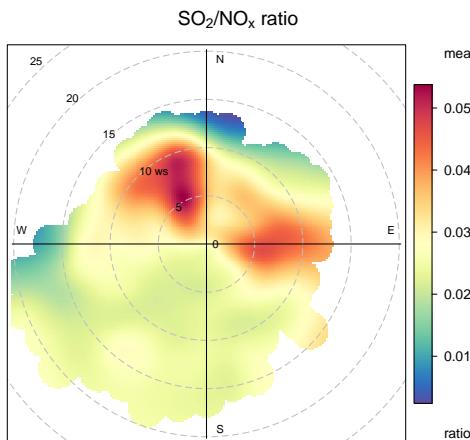


FIGURE 15.4 Bivariate polar plot of the ratio of SO_2/NO_x .

R. In this example we consider the ratio of SO_2/NO_x :

```
mydata <- transform(mydata, ratio = so2 / nox)
```

Working with
ratios of
pollutants

This makes a new variable called `ratio`. Sometimes it can be problematic if there are values equal to zero on the denominator, as is the case here. The mean and maximum value of the ratio is infinite, as shown by the `Inf` in the statistics below. Luckily, R can deal with infinity and the `openair` functions will remove these values before performing calculations. It is very simple therefore to calculate ratios.

```
summary(mydata[, "ratio"])
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
	0.000	0.018	0.024	Inf	0.034	Inf	11782

A polar plot of the SO_2/NO_x ratio is shown in Figure 15.4. The plot highlights some new features not seen before. First, to the north there seems to be evidence that the air tends to have a higher SO_2/NO_x ratio. Also, the source to the east has a higher SO_2/NO_x ratio compared with that when the wind is from the south-west i.e. dominated by road sources. It seems therefore that the easterly source(s), which are believed to be industrial sources have a different SO_2/NO_x ratio compared with road sources. This is a very simple analysis, but ratios can be used effectively in many functions and are particularly useful in the presence of high source complexity.

Sometimes when considering ratios it might be necessary to limit the values in some way; perhaps due to some unusually low value denominator data resulting in a few very high values for the ratio. This is easy to do with the `subset` command. The code below selects ratios less than 0.1.

```
polarPlot(subset(mydata, ratio < 0.1), pollutant = "ratio")
```

The uncertainties in the surface can be calculated by setting the option `uncertainty = TRUE`. The details are described above and here we show the example of SO_2 concentrations (Figure 15.5). In general the uncertainties are higher at high wind speeds i.e. at the ‘fringes’ of a plot where there are fewer data. However, the magnitude depends on

```
polarPlot(mydata, pollutant = "so2", uncertainty = TRUE)
```

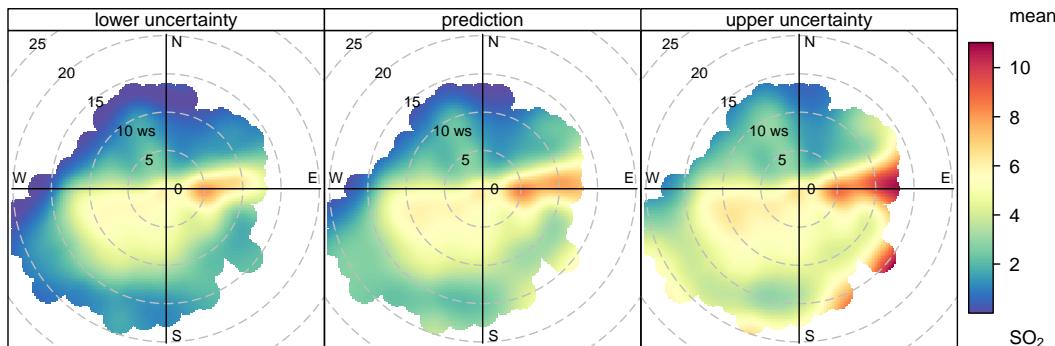


FIGURE 15.5 Bivariate polar plot of SO₂ concentrations at Marylebone Road. Three surfaces are shown: the central prediction (middle) and the lower and upper 95% estimated uncertainties. These plots help to show that in this particular case, some of the concentrations for strong easterly and south-easterly winds are rather uncertain. However, the central feature to the east remains, suggesting this feature is ‘real’ and not an artifact of there being too few data.

both the frequency and magnitude of the concentration close to the points of interest. The pattern of uncertainty is not always obvious and it can differ markedly for different pollutants.

The `polarPlot` function can also produce plots dependent on another variable (see the `type` option). For example, the variation of SO₂ concentrations at Marylebone Road by hour of the day in the code below. The function was called as shown in in this case the minimum number of points in each wind speed/direction was set to 2.

```
polarPlot(mydata, pollutant = "so2", type = "hour", min.bin = 2)
```

This plot shows that concentrations of SO₂ tend to be highest from the east (as also shown in Figure 15.3) and for hours in the morning. Together these plots can help better understand different source types. For example, does a source only seem to be present during weekdays, or winter months etc. In the case of `type = "hour"`, the more obvious presence during the morning hours could be due to meteorological factors and this possibility should be investigated also. In other settings where there are many sources that vary in their source emission and temporal characteristics, the `polarPlot` function should prove to be very useful.

One issue to be aware of is the amount of data required to generate some of these plots; particularly the hourly plots. If only a relatively short time series is available there may not be sufficient information to produce useful plots. Whether this is important or not will depend on the specific circumstances e.g. the prevalence of wind speeds and directions from the direction of interest. When used to produce many plots (e.g. `type = "hour"`), the run time can be quite long.

15.3.1 Conditional Probability Function (CPF) plot

The conditional probability functions (CPF) was described on page 100 in the context of the `percentileRose` function. The CPF was originally used to show the wind directions that dominate a (specified) high concentration of a pollutant; showing the probability of such concentrations occurring by wind direction (Ashbaugh et al. 1985). However, these ideas can very usefully be applied to bivariate polar plots. In this case

the CPF is defined as $CPF = m_{\theta,j}/n_{\theta,j}$, where $m_{\theta,j}$ is the number of samples in the wind sector θ and wind speed interval j with mixing ratios greater than some ‘high’ concentration, and $n_{\theta,j}$ is the total number of samples in the same wind direction-speed interval. Note that j does not have to be wind speed but could be any numeric variable e.g. ambient temperature. CPF analysis is very useful for showing which wind direction, wind speed intervals are dominated by high concentrations and give the probability of doing so. A full explanation of the development and use of the bivariate case of the CPF is described in Uria-Tellaetxe and Carslaw (2014) where it is applied to monitoring data close to a steelworks.

```
polarPlot(mydata, pollutant = "so2", statistic = "cpf", percentile = 90)
```

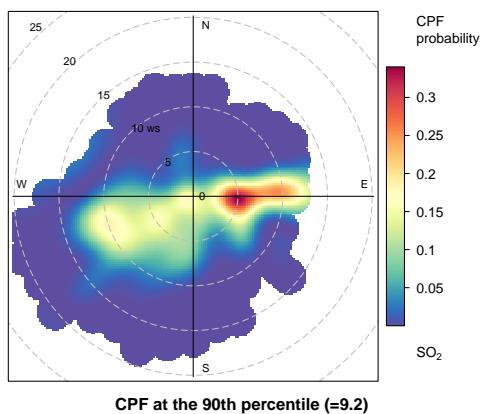


FIGURE 15.6 `polarPlot` of SO₂ concentrations at Marylebone Road based on the CPF function.

An example of a CPF polar plot is shown in Figure 15.6 for the 90th percentile concentration of SO₂. This plot shows that for most wind speed-directions the probability of SO₂ concentrations being greater than the 90th percentile is zero. The clearest areas where the probability is higher is to the east. Indeed, the plot now clearly reveals two potential sources of SO₂, which are not as apparent in the ‘standard’ plot shown in Figure 15.3. Note that Figure 15.6 also gives the calculated percentile at the bottom of the plot (9.2 ppb in this case). Figure 15.6 can also be compared with the CPF plot based only on wind direction shown in Figure 13.4. While Figure 13.4 very clearly shows that easterly wind dominate high concentrations of SO₂, Figure 15.6 provides additional valuable information by also considering wind speed, which in this case is able to discriminate between two sources (or groups of sources) to the east.

The polar CPF plot is therefore potentially very useful for source identification and characterisation. It is, for example, worth also considering other percentile levels and other pollutants. For example, considering the 95th percentile for SO₂ ‘removes’ one of the sources (the one at highest wind speed). This helps to show some maybe important differences between the sources that could easily have been missed. Similarly, considering other pollutants can help build up a good understanding of these sources. A CPF plot for NO₂ at the 90th percentile shows the single dominance of the road source. However, a CPF plot at the 75th percentile level indicates source contributions from the east (likely tall stacks), which again are not as clear in the standard bivariate polar plot. Considering a range of percentile values can therefore help to build up a more complete understanding of source contributions.

However, even more useful information can be gained by considering *intervals* of percentiles e.g. 50–60, 60–70 etc. By considering intervals of percentiles it becomes

```
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(0, 10))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(10, 20))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(20, 30))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(30, 40))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(40, 50))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(50, 60))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(60, 70))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(70, 80))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(80, 90))
polarPlot(mydata, poll = "so2", stati = "cpf", percentile = c(90, 100))
```

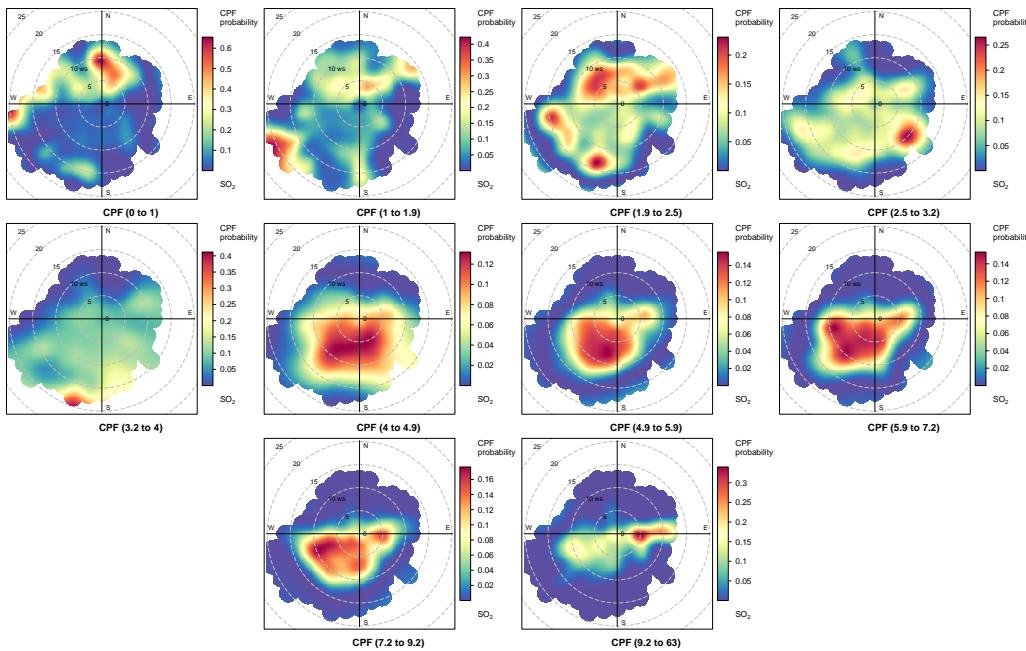


FIGURE 15.7 `polarPlot` of SO₂ concentrations at Marylebone Road based on the CPF function for a range of percentile intervals from 0–10, 10–20, ..., 90–100.

clear that some sources only affect a limited percentile range. `polarPlot` can accept a `percentile` argument of length two e.g. `percentile = c(80, 90)`. In this case concentrations in the range from the lower to upper percentiles will be considered. In Figure 15.7 for example, it is apparent that the road source to the south west is only important between the 60 to 90th percentiles. As mentioned previously, the chimney stacks to the east are important for the higher percentiles (90 to 100). What is interesting though is the emergence of what appears to be other sources at the lower percentile intervals. These potential sources are not apparent in Figure 15.3. The other interesting aspect is that it does seem that specific sources tend to be prominent for specific percentile ranges. If this characteristic is shown to be the case more generally, then CPF intervals could be a powerful way in which to identify many sources. Whether these particular sources are important or not is questionable and depends on the aims of the analysis. However, there is no reason to believe that the potential sources shown in the percentile ranges 0 to 50 are artifacts. They could for example be signals from more distant point sources whose plumes have diluted more over longer distances. Such sources would be ‘washed out’ in an ordinary polar plot. For a fuller example of this approach see Uria-Tellaetxe and Carslaw (2014).

Note that it is easy to work out what the concentration intervals are for the percentiles shown in Figure 15.7:

```
quantile(mydata$so2, probs = seq(0, 1, by = 0.1), na.rm = TRUE)

##   0% 10% 20% 30% 40% 50% 60% 70% 80% 90% 100%
## 0.0 1.0 1.9 2.5 3.2 4.0 4.9 5.9 7.2 9.2 63.2
```

To plot the Figures on one page it is necessary to make the plot objects first and then decide how to plot them. To plot the Figures in a particular layout see [page 63](#).

15.3.2 The **polarCluster** function for feature identification and extraction

The **polarPlot** function will often identify interesting features that would be useful to analyse further. It is possible to select areas of interest based only on a consideration of a plot. Such a selection could be based on wind direction and wind speed intervals for example e.g.

```
subdata <- subset(mydata, ws >3 & wd >= 180 & wd <=270)
```

which would select wind speeds $>3 \text{ m s}^{-1}$ and wind directions from 180 to 270 degrees from **mydata**. That subset of data, **subdata**, could then be analysed using other functions. While this approach may be useful in many circumstances it is rather arbitrary. In fact, the choice of ‘interesting feature’ in the first place can even depend on the colour scale used, which is not very robust. Furthermore, many interesting patterns can be difficult to select and won’t always fall into convenient intervals of other variables such as wind speed and direction.

A better approach is to use a method that can select group similar features together. One such approach is to use *cluster analysis*. **openair** uses k-means clustering as a way in which bivariate polar plot features can be identified and grouped. The main purpose of grouping data in this way is to identify records in the original time series data by cluster to enable post-processing to better understand potential source characteristics. The process of grouping data in k-means clustering proceeds as follows. First, k points are randomly chosen from the space represented by the objects that are being clustered into k groups. These points represent initial group centroids. Each object is assigned to the group that has the closest centroid. When all objects have been assigned, the positions of the k centroids are re-calculated. The previous two steps are repeated until the centroids no longer move. This produces a separation of the objects into groups from which the metric to be minimised can be calculated.

Central to the idea of clustering data is the concept of distance i.e. some measure of similarity or dissimilarity between points. Clusters should be comprised of points separated by small distances relative to the distance between the clusters. Careful consideration is required to define the distance measure used because the effectiveness of clustering itself fundamentally depends on its choice. The similarity of concentrations shown in [Figure 15.1](#) for example is determined by three variables: the u and v wind components and the concentration. All three variables are equally important in characterising the concentration-location information, but they exist on different scales i.e. a wind speed-direction measure and a concentration. Let $X = \{x_i\}, i = 1, \dots, n$ be a set of n points to be clustered into K clusters, $C = \{c_k, k = 1, \dots, K\}$. The basic k-means algorithm for K clusters is obtained by minimising:

$$\sum_{k=1}^K \sum_{x_i \in c_k} \|x_i - \mu_k\|^2 \quad (4)$$

where $\|x_i - \mu_k\|^2$ is a chosen distance measure, μ_k is the mean of cluster c_k .

The distance measure is defined as the Euclidean distance:

$$d_{x,y} = \left(\sum_{j=1}^J (x_j - y_j)^2 \right)^{1/2} \quad (5)$$

Where x and y are two J -dimensional vectors, which have been standardized by subtracting the mean and dividing by the standard deviation. In the current case J is of length three i.e. the wind components u and v and the concentration C , each of which is standardized e.g.:

$$x_j = \left(\frac{x_j - \bar{x}}{\sigma_x} \right) \quad (6)$$

Standardization is necessary because the wind components u and v are on different scales to the concentration. In principle, more weight could be given to the concentration rather than the u and v components, although this would tend to identify clusters with similar concentrations but different source origins.

`polarCluster` can be thought of as the ‘local’ version of clustering of back trajectories. Rather than using air mass origins, wind speed, wind direction and concentration are used to group similar conditions together. Section 26.3 provides the details of clustering back trajectories in `openair`. A fuller description of the clustering approach is described in Carslaw and Beevers (2013).

The `polarCluster` function has the following options.

- mydata** A data frame minimally containing `wd`, another variable to plot in polar coordinates (the default is a column “ws” — wind speed) and a pollutant. Should also contain `date` if plots by time period are required.
- pollutant** Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. `pollutant = "nox"`. Only one pollutant can be chosen.
- x** Name of variable to plot against wind direction in polar coordinates, the default is wind speed, “ws”.
- wd** Name of wind direction field.
- n.clusters** Number of clusters to use. If `n.clusters` is more than length 1, then a `lattice` panel plot will be output showing the clusters identified for each one of `n.clusters`.
- cols** Colours to be used for plotting. Useful options for categorical data are available from `RColorBrewer` colours — see the `openair openColours` function for more details. Useful schemes include “Accent”, “Dark2”, “Paired”, “Pastel1”, “Pastel2”, “Set1”, “Set2”, “Set3” — but see `?brewer.pal` for the maximum useful colours in each. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`.
- angle.scale** The wind speed scale is by default shown at a 315 degree angle. Sometimes the placement of the scale may interfere with an interesting feature. The user can therefore set `angle.scale` to another value (between 0 and 360 degrees) to mitigate such problems. For example `angle.scale = 45` will draw the scale heading in a NE direction.

- units** The units shown on the polar axis scale.
- auto.text** Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO₂.
- ...** Other graphical parameters passed onto `polarPlot`, `lattice:levelplot` and `cutData`. Common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed via `quickText` to handle routine formatting.

The use of the `polarCluster` is very similar to the use of all `openair` functions. While there are many techniques available to try and find the optimum number of clusters, it is difficult for these approaches to work in a consistent way for identifying features in bivariate polar plots. For this reason it is best to consider a range of solutions that covers a number of clusters.

Cluster analysis is computationally intensive and the `polarCluster` function can take a comparatively long time to run. The basic idea is to calculate the solution to several cluster levels and then choose one that offers the most appropriate solution for post-processing.

The example given below is for concentrations of SO₂, shown in Figure 15.3 and the aim is to identify features in that plot. A range of numbers of clusters will be calculated — in this case from two to ten.

The real benefit of `polarCluster` is being able to identify clusters in the original data frame. To do this, the results from the analysis must be read into a new variable, as in Figure 15.9, where the results are read into a data frame `results`. Now it is possible to use this new information. In the 8-cluster solution to Figure 15.9, cluster 6 seems to capture the elevated SO₂ concentrations to the east well (see Figure 15.3 for comparison), while cluster 5 will strongly represent the road contribution.

The results are here:

```
head(results[["data"]])
##      .id          date   ws   wd nox no2 o3 pm10 so2   co pm25 ws2 wd2
## 1    1 1998-01-01 00:00:00 0.6 280 285 39   1    29 4.7 3.37   NA 3.7 327
## 2    10 1998-01-01 09:00:00 4.0 170 113 39   2    12 2.9 1.21   NA 4.8 154
## 3   100 1998-01-05 03:00:00 8.8 240  53 24  16    12 1.6 0.54   NA 9.6 318
## 4  1000 1998-02-11 15:00:00 7.2 230 372 92   2    49 8.6 3.72   NA 8.5 270
## 5 10000 1999-02-21 15:00:00 6.4 290  71 32  20    10 2.6 1.05     8 8.2 281
## 6 10001 1999-02-21 16:00:00 7.7 290   85 39  16    12 2.9 1.48     9 10.8 327
##      ratio cluster
## 1  0.017       4
## 2  0.026       5
## 3  0.030       4
## 4  0.023       4
## 5  0.037       7
## 6  0.034       7
```

Note that there is an additional column `cluster` that gives the cluster a particular row belongs to and that this is a *character* variable. It might be easier to read these results into a new data frame:

```
results <- results[["data"]]
```

It is easy to find out how many points are in each cluster:

```
polarCluster(mydata, pollutant="so2", n.clusters=2:10, cols= "Set2")
```

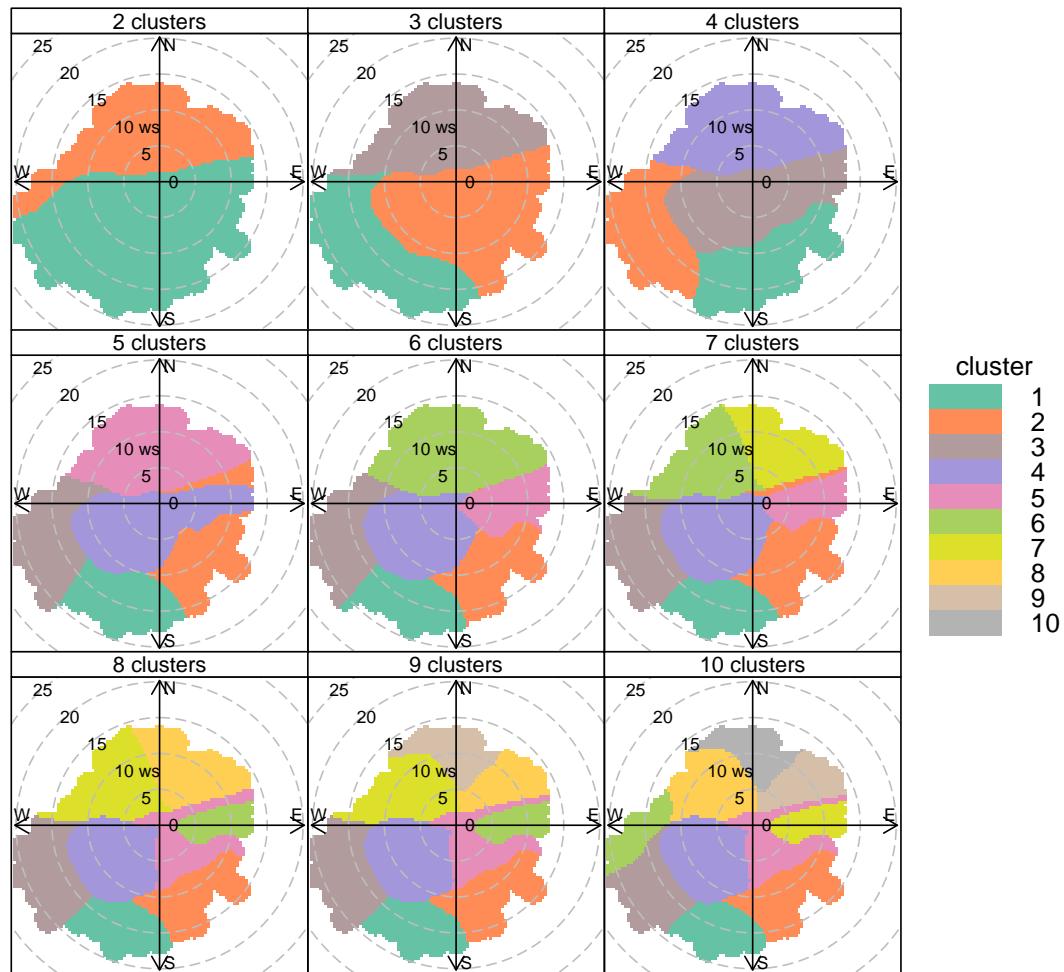


FIGURE 15.8 Use of the **polarCluster** function applied to SO₂ concentrations at Marylebone Road. In this case 2 to 10 clusters have been chosen.

```
results <- polarCluster(mydata, pollutant="so2", n.clusters=8, cols= "Set2")
```

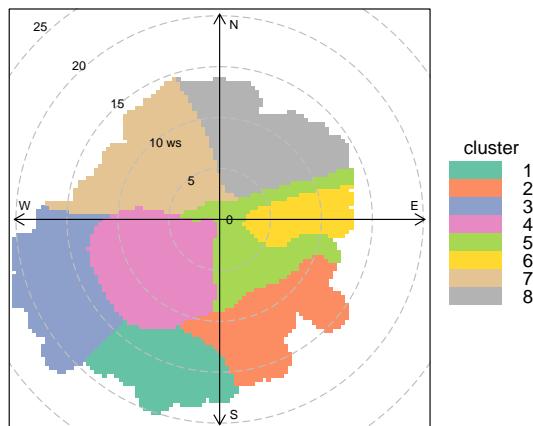


FIGURE 15.9 Use of the **polarCluster** function applied to SO₂ concentrations at Marylebone Road. In this case 8 clusters have been chosen.

```
timeProp(selectByDate(results, year = 2003), pollutant = "so2", avg.time = "day",
proportion = "cluster", col = "Set2", key.position = "top",
key.columns = 8, date.breaks = 10, ylab = "so2 (ug/m3)")
```

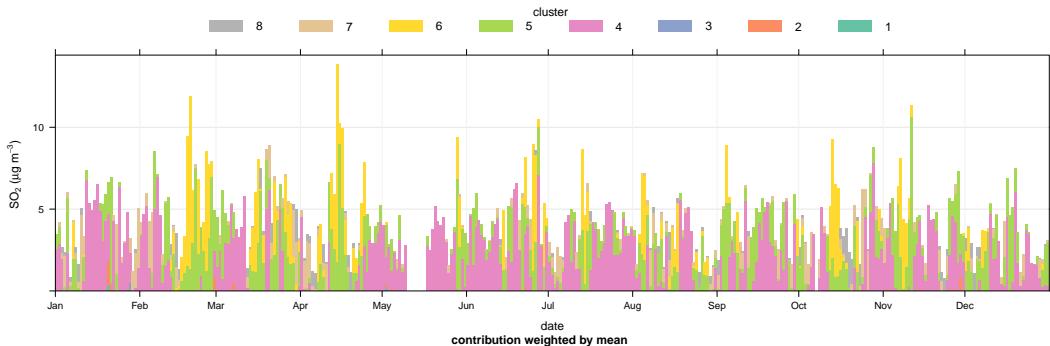


FIGURE 15.10 Temporal variation in daily SO₂ concentration at the Marylebone Road site show by contribution of each cluster for 2003.

```
table(results[, "cluster"])

##
##      1     2     3     4     5     6     7     8
##    206   412   160 24133 16049  2590  7839 2918
```

Now other **openair** analysis functions can be used to analyse the results. For example, to consider the temporal variations by cluster:

```
timeVariation(results, pollutant = "so2", group = "cluster",
col = "Set2", ci = FALSE, lwd = 3)
```

Or if we just want to plot a couple of clusters (5 and 6) *using the same colours* as in Figure 15.9:

```
timeVariation(subset(results, cluster %in% c("C5", "C6")), pollutant = "so2",
group = "cluster", col = openColours("Set2", 8)[5:6], lwd = 3)
```

polarCluster will work on any surface that can be plotted by **polarPlot** e.g. the radial variable does not have to be wind speed but could be another variable such as temperature. While it is not always possible for **polarCluster** to identify all features in a surface it certainly makes it easier to post-process **polarPlots** using other **openair** functions or indeed other analyses altogether.

Another useful way of understanding the clusters is to use the **timeProp** function, which can display a time series as a bar chart split by a categorical variable (in this case the cluster). In this case it is useful to plot the time series of SO₂ and show how much of the concentration is contributed to by each cluster. Such a plot is shown in Figure 15.10. It is now easy to see for example that many of the peaks in SO₂ are associated with cluster 6 (power station sources from the east), seen in Figure 15.9. Cluster 6 is particularly prominent during springtime, but those sources also make important contributions through the whole year.

16 The `polarAnnulus` function

16.1 Purpose

see also
[polarFreq](#)
[polarPlot](#)
[percentileRose](#)
[pollutionRose](#)

The `polarAnnulus` function provides a way in which to consider the temporal aspects of a pollutant concentration by wind direction. This is another means of visualising diurnal, day of week, seasonal and trend variations. Plotting as an annulus, rather than a circle avoids to some extent the difficulty in interpreting values close to the origin. These plots have the capacity to display potentially important information regarding sources; particularly if more than one pollutant is available.

16.2 Options available

The `polarAnnulus` function has the following options:

mydata A data frame minimally containing `date`, `wd` and a pollutant.

pollutant Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. `pollutant = "nox"`. There can also be more than one pollutant specified e.g. `pollutant = c("nox", "no2")`. The main use of using two or more pollutants is for model evaluation where two species would be expected to have similar concentrations. This saves the user stacking the data and it is possible to work with columns of data directly. A typical use would be `pollutant = c("obs", "mod")` to compare two columns “obs” (the observations) and “mod” (modelled values).

resolution Two plot resolutions can be set: “normal” and “fine” (the default).

local.tz Should the results be calculated in local time that includes a treatment of daylight savings time (DST)? The default is not to consider DST issues, provided the data were imported without a DST offset. Emissions activity tends to occur at local time e.g. rush hour is at 8 am every day. When the clocks go forward in spring, the emissions are effectively released into the atmosphere typically 1 hour earlier during the summertime i.e. when DST applies. When plotting diurnal profiles, this has the effect of “smearing-out” the concentrations. Sometimes, a useful approach is to express time as local time. This correction tends to produce better-defined diurnal profiles of concentration (or other variables) and allows a better comparison to be made with emissions/activity data. If set to `FALSE` then GMT is used. Examples of usage include `local.tz = "Europe/London"`, `local.tz = "America/New_York"`. See `cutData` and `import` for more details.

period This determines the temporal period to consider. Options are “hour” (the default, to plot diurnal variations), “season” to plot variation throughout the year, “weekday” to plot day of the week variation and “trend” to plot the trend by wind direction.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If `type` is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

`Type` can be up length two e.g. `type = c("season", "site")` will produce a 2x2 plot split by season and site. The use of two types is mostly meant for situations where there are several sites. Note, when two types are provided the first forms the columns and the second the rows.

Also note that for the `polarAnnulus` function some type/period combinations are forbidden or make little sense. For example, `type = "season"` and `period = "trend"` (which would result in a plot with too many gaps in it for sensible smoothing), or `type = "weekday"` and `period = "weekday"`.

statistic The statistic that should be applied to each wind speed/direction bin. Can be “mean” (default), “median”, “max” (maximum), “frequency”, “stdev” (standard deviation), “weighted.mean” or “cpf” (Conditional Probability Function). Because of the smoothing involved, the colour scale for some of these statistics is only to provide an indication of overall pattern and should not be interpreted in concentration units e.g. for `statistic = "weighted.mean"` where the bin mean is multiplied by the bin frequency and divided by the total frequency. In many cases using `polarFreq` will be better. Setting `statistic = "weighted.mean"` can be useful because it provides an indication of the concentration * frequency of occurrence and will highlight the wind speed/direction conditions that dominate the overall mean.

percentile If `statistic = "percentile"` or `statistic = "cpf"` then `percentile` is used, expressed from 0 to 100. Note that the percentile value is calculated in the wind speed, wind direction ‘bins’. For this reason it can also be useful to set `min.bin` to ensure there are a sufficient number of points available to estimate a percentile. See `quantile` for more details of how percentiles are calculated.

limits Limits for colour scale.

cols Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and user defined. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`

width The width of the annulus; can be “normal” (the default), “thin” or “fat”.

min.bin The minimum number of points allowed in a wind speed/wind direction bin. The default is 1. A value of two requires at least 2 valid records in each bin and so on; bins with less than 2 valid records are set to NA. Care should be taken when using a value > 1 because of the risk of removing real data points. It is recommended to consider your data with care. Also, the `polarFreq` function can be of use in such circumstances.

exclude.missing Setting this option to `TRUE` (the default) removes points from the plot that are too far from the original data. The smoothing routines will

produce predictions at points where no data exist i.e. they predict. By removing the points too far from the original data produces a plot where it is clear where the original data lie. If set to `FALSE` missing data will be interpolated.

date.pad For `type = "trend"` (default), `date.pad = TRUE` will pad-out missing data to the beginning of the first year and the end of the last year. The purpose is to ensure that the trend plot begins and ends at the beginning or end of year.

force.positive The default is `TRUE`. Sometimes if smoothing data with steep gradients it is possible for predicted values to be negative. `force.positive = TRUE` ensures that predictions remain positive. This is useful for several reasons. First, with lots of missing data more interpolation is needed and this can result in artifacts because the predictions are too far from the original data. Second, if it is known beforehand that the data are all positive, then this option carries that assumption through to the prediction. The only likely time where setting `force.positive = FALSE` would be if background concentrations were first subtracted resulting in data that is legitimately negative. For the vast majority of situations it is expected that the user will not need to alter the default option.

k The smoothing value supplied to `gam` for the temporal and wind direction components, respectively. In some cases e.g. a trend plot with less than 1-year of data the smoothing with the default values may become too noisy and affected more by outliers. Choosing a lower value of `k` (say 10) may help produce a better plot.

normalise If `TRUE` concentrations are normalised by dividing by their mean value. This is done *after* fitting the smooth surface. This option is particularly useful if one is interested in the patterns of concentrations for several pollutants on different scales e.g. NOx and CO. Often useful if more than one `pollutant` is chosen.

key.header Adds additional text/labels to the scale key. For example, passing the options `key.header = "header"`, `key.footer = "footer1"` adds addition text above and below the scale key. These arguments are passed to `drawOpenKey` via `quickText`, applying the `auto.text` argument, to handle formatting.

key.footer see `key.header`.

key.position Location where the scale key is to plotted. Allowed arguments currently include “top”, “right”, “bottom” and “left”.

key Fine control of the scale key via `drawOpenKey`. See `drawOpenKey` for further details.

auto.text Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO₂.

... Other graphical parameters passed onto `lattice:levelplot` and `cutData`. For example, `polarAnnulus` passes the option `hemisphere = "southern"`

on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `levelplot` via `quickText` to handle routine formatting.

16.3 Example of use

We apply the four variations of the `polarAnnulus` plot to PM₁₀ concentrations at Marylebone Road. Figure 16.1 shows the different temporal components. Similar to other analyses for PM₁₀, the trend plot show that concentrations are dominated by southerly winds and there is little overall change in concentrations from 1998 to 2005, as shown by the red colouring over the period. The seasonal plot shows that February/March is important for easterly winds, while the summer/late summer period is more important for southerly and south-westerly winds. The day of the week plot clearly shows concentrations to be elevated for during weekdays but not weekends — for all wind directions. Finally, the diurnal plot highlights that higher concentrations are observed from 6 am to 6 pm.

Interestingly, the plot for NO_x and CO (not shown, but easily produced) did not show such a strong contribution for south-easterly winds. This raises the question as to whether the higher particle concentrations seen for these wind directions are dominated by different sources (i.e. not the road itself). One explanation is that during easterly flow, concentrations are strongly affected by long-range transport. However, as shown in the diurnal plot in Figure 16.1, the contribution from the south-east also has a sharply defined profile — showing very low concentrations at night, similar to the likely contribution from the road. This type of profile might not be expected from a long-range source where emissions are well-mixed and secondary particle formation has had time to occur. The same is also true for the day of the week plot, where there is little evidence of ‘smeared-out’ long-range transport sources. These findings may suggest a different, local source of PM₁₀ that is not the road itself. Clearly, a more detailed analysis would be required to confirm the patterns shown; but it does highlight the benefit of being able to analyse data in different ways.

Where there is interest in considering the wind direction dependence of concentrations, it can be worth filtering for wind speeds. At low wind speed with wind direction becomes highly variable (and is often associated with high pollutant concentrations). Therefore, for some situations it might be worth considering removing the very low wind speeds. The code below provides two ways of doing this using the `subset` function. The first selects data where the wind speed is > 2 m s⁻¹. The second part shows how to select wind speeds greater than the 10th percentile, using the `quantile` function. The latter way of selecting is quite useful, because it is known how much data are selected i.e. in this case 90 %. It is worth experimenting with different values because it is also important not to lose information by ignoring wind speeds that provide useful information.

```
## wind speed >2
polarAnnulus(subset(mydata, ws > 2), poll="pm10", type = "hour")
## wind speed > 10th percentile
polarAnnulus(subset(mydata, ws > quantile(ws, probs = 0.1, na.rm = TRUE)),
              poll="pm10", type = "hour")
```

```
data(mydata)
polarAnnulus(mydata, poll = "pm10", period = "trend", main = "Trend")
polarAnnulus(mydata, poll = "pm10", period = "season", main = "Season")
polarAnnulus(mydata, poll = "pm10", period = "weekday", main = "Weekday")
polarAnnulus(mydata, poll = "pm10", period = "hour", main = "Hour")
```

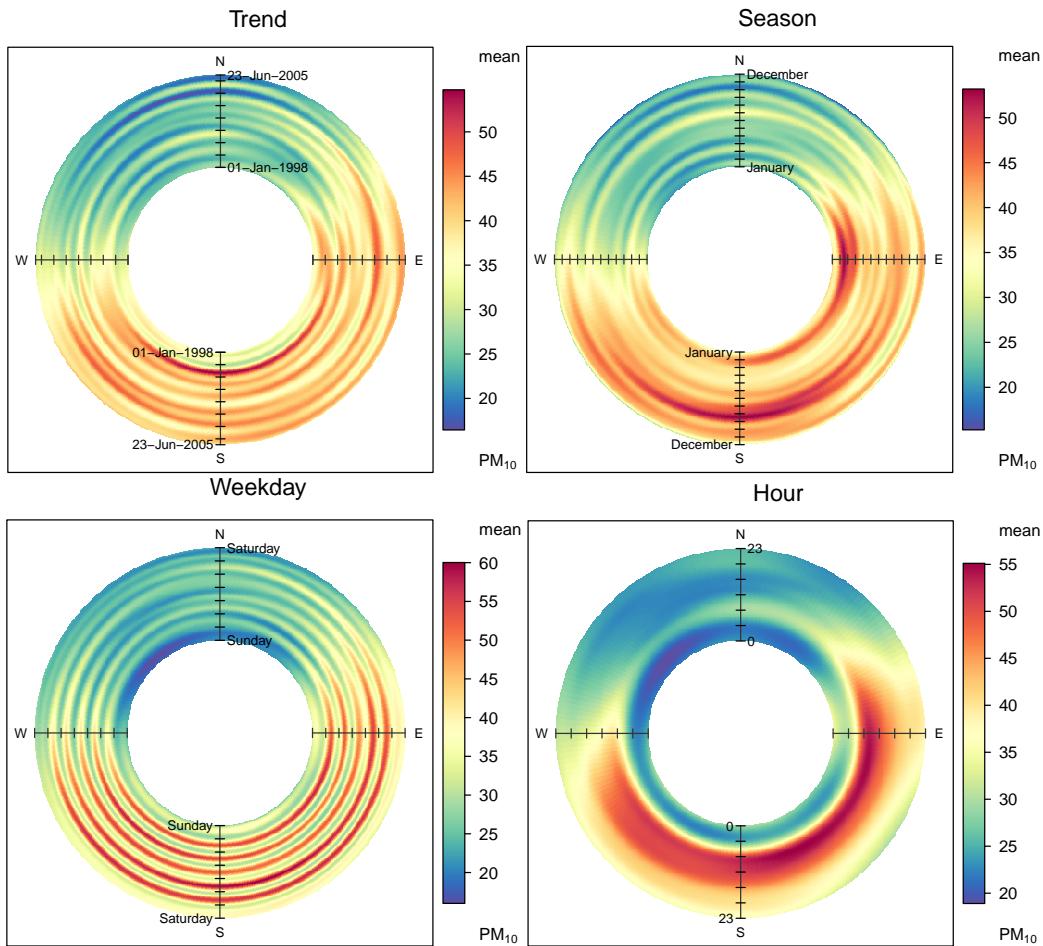


FIGURE 16.1 Examples of the **polarAnnulus** function applied to Marylebone Road

17 The **timePlot** function

17.1 Purpose

see also
[smoothTrend](#)
[TheilSen](#)
[timeVariation](#)
[scatterPlot](#)

The **timePlot** function is designed to quickly plot time series of data, perhaps for several pollutants or variables. This is, or should be, a very common task in the analysis of air pollution. In doing so, it is helpful to be able to plot several pollutants at the same time (and maybe other variables) and quickly choose the time periods of interest. It will plot time series of type **Date** and hourly and high time resolution data.

The function offers fine control over many of the plot settings such as line type, colour and width. If more than one pollutant is selected, then the time series are shown in a compact way in different panels with different scales. Sometimes it is useful to get an idea of whether different variables 'go up and down' together. Such comparisons in **timePlot** are made easy by setting **group = TRUE**, and maybe also **normalise = "mean"**. The latter setting divides each variable by its mean value, thus enabling several variables to be plotted together using *the same scale*. The **normalise** option will

also take a date as a string (in British format dd/mm/YYYY), in which case all data are normalise to equal 100 at that time. Normalising data like this makes it easy to compare time series on different scales e.g. emissions and ambient measurements.

`timePlot` works very well in conjunction with `selectByDate`, which makes it easy to select specific time series intervals. See (§31.1) for examples of how to select parts of a data frame based on the date.

Another useful feature of `timePlot` is the ability to average the data in several ways. This makes it easy, for example, to plot daily or monthly means from hourly data, or hourly means from 15-minute data. See the option `avg.time` for more details and (§31.4) where a full description of time averaging of data frames is given.

17.2 Options available

mydata A data frame of time series. Must include a `date` field and at least one variable to plot.

pollutant Name of variable to plot. Two or more pollutants can be plotted, in which case a form like `pollutant = c("nox", "co")` should be used.

group If more than one pollutant is chosen, should they all be plotted on the same graph together? The default is `FALSE`, which means they are plotted in separate panels with their own scaled. If `TRUE` then they are plotted on the same plot with the same scale.

stack If `TRUE` the time series will be stacked by year. This option can be useful if there are several years worth of data making it difficult to see much detail when plotted on a single plot.

normalise Should variables be normalised? The default is not to normalise the data. `normalise` can take two values, either “mean” or a string representing a date in UK format e.g. “1/1/1998” (in the format dd/mm/YYYY). If `normalise = "mean"` then each time series is divided by its mean value. If a date is chosen, then values at that date are set to 100 and the rest of the data scaled accordingly. Choosing a date (say at the beginning of a time series) is very useful for showing how trends diverge over time. Setting `group = TRUE` is often useful too to show all time series together in one panel.

avg.time This defines the time period to average to. Can be “sec”, “min”, “hour”, “day”, “DSTday”, “week”, “month”, “quarter” or “year”. For much increased flexibility a number can precede these options followed by a space. For example, a timeAverage of 2 months would be `period = "2 month"`. See function `timeAverage` for further details on this.

data.thresh The data capture threshold to use (%) when aggregating the data using `avg.time`. A value of zero means that all available data will be used in a particular period regardless of the number of values available. Conversely, a value of 100 will mean that all data will need to be present for the average to be calculated, else it is recorded as `NA`. Not used if `avg.time = "default"`.

statistic The statistic to apply when aggregating the data; default is the mean. Can be one of “mean”, “max”, “min”, “median”, “frequency”, “sd”, “percentile”. Note that “sd” is the standard deviation and “frequency” is the number

(frequency) of valid records in the period. “percentile” is the percentile level (%) between 0-100, which can be set using the “percentile” option - see below. Not used if `avg.time = "default"`.

percentile The percentile level in % used when `statistic = "percentile"` and when aggregating the data with `avg.time`. More than one percentile level is allowed for `type = "default"` e.g. `percentile = c(50, 95)`. Not used if `avg.time = "default"`.

date.pad Should missing data be padded-out? This is useful where a data frame consists of two or more “chunks” of data with time gaps between them. By setting `date.pad = TRUE` the time gaps between the chunks are shown properly, rather than with a line connecting each chunk. For irregular data, set to `FALSE`. Note, this should not be set for `type` other than `default`.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Only one `type` is currently allowed in `timePlot`.

cols Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and `RColorBrewer` colours— see the `openair openColours` function for more details. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`

plot.type The `lattice` plot type, which is a line (`plot.type = "l"`) by default. Another useful option is `plot.type = "h"`, which draws vertical lines.

key Should a key be drawn? The default is `TRUE`.

log Should the y-axis appear on a log scale? The default is `FALSE`. If `TRUE` a well-formatted log10 scale is used. This can be useful for plotting data for several different pollutants that exist on very different scales. It is therefore useful to use `log = TRUE` together with `group = TRUE`.

windflow This option allows a scatter plot to show the wind speed/direction as an arrow. The option is a list e.g. `windflow = list(col = "grey", lwd = 2, scale = 0.1)`. This option requires wind speed (`ws`) and wind direction (`wd`) to be available.

The maximum length of the arrow plotted is a fraction of the plot dimension with the longest arrow being `scale` of the plot x-y dimension. Note, if the plot size is adjusted manually by the user it should be re-plotted to ensure the correct wind angle. The list may contain other options to `panel.arrows` in the `lattice` package. Other useful options include

`length`, which controls the length of the arrow head and `angle`, which controls the angle of the arrow head.

This option works best where there are not too many data to ensure overplotting does not become a problem.

`smooth` Should a smooth line be applied to the data? The default is `FALSE`.

`ci` If a smooth fit line is applied, then `ci` determines whether the 95% confidence intervals are shown.

`y.relation` This determines how the y-axis scale is plotted. "same" ensures all panels use the same scale and "free" will use panel-specific scales. The latter is a useful setting when plotting data with very different values.

`ref.x` See `ref.y` for details. In this case the correct date format should be used for a vertical line e.g. `ref.x = list(v = as.POSIXct("2000-06-15"), lty = 5)`.

`ref.y` A list with details of the horizontal lines to be added representing reference line(s). For example, `ref.y = list(h = 50, lty = 5)` will add a dashed horizontal line at 50. Several lines can be plotted e.g. `ref.y = list(h = c(50, 100), lty = c(1, 5), col = c("green", "blue"))`. See `panel.abline` in the `lattice` package for more details on adding/controlling lines.

`key.columns` Number of columns to be used in the key. With many pollutants a single column can make the key too wide. The user can thus choose to use several columns by setting `columns` to be less than the number of pollutants.

`key.position` Location where the scale key is to be plotted. Can include "top", "bottom", "right" and "left".

`name.pol` This option can be used to give alternative names for the variables plotted. Instead of taking the column headings as names, the user can supply replacements. For example, if a column had the name "nox" and the user wanted a different description, then setting `name.pol = "nox before change"` can be used. If more than one pollutant is plotted then use `c` e.g. `name.pol = c("nox here", "o3 there")`.

`date.breaks` Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of `date.breaks` up or down.

`date.format` This option controls the date format on the x-axis. While `timePlot` generally sets the date format sensibly there can be some situations where the user wishes to have more control. For format types see `strptime`. For example, to format the date like "Jan-2012" set `date.format = "%b-%Y"`.

`auto.text` Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO₂.

... Other graphical parameters are passed onto `cutData` and `lattice:xyplot`. For example, `timePlot` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, most common plotting parameters, such as `layout` for panel arrangement and `pch` and `cex` for plot symbol type and size and `lty` and `lwd` for line type and width, as passed to `xyplot`, although some maybe locally managed by `openair` on route, e.g. axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed via `quickText` to handle routine formatting. See examples below.

17.3 Example of use

A full set of examples is shown in the help pages — see `?timePlot` for details. At the basic level, concentrations are shown using a simple call e.g. to plot time series of NO_x and O_3 in separate panels with their own scales.

```
timePlot(mydata, pollutant = c("nox", "o3"), y.relation = "free")
```

Often it is necessary to only consider part of a time series and using the `openair` function `selectByDate` makes it easy to do this. Some examples are shown below.

To plot data only for 2003:

```
timePlot(selectByDate(mydata, year = 2003), pollutant = c("nox", "o3"), y.relation = "free")
```

Plots for several pollutants for August 2003, are shown in [Figure 17.1](#).

Some other examples (not plotted) are:

```
## plot monthly means of ozone and no2
timePlot(mydata, pollutant = c("o3", "no2"), avg.time = "month",
          y.relation = "free")

## plot 95th percentile monthly concentrations
timePlot(mydata, pollutant = c("o3", "no2"), avg.time = "month",
          statistic = "percentile", percentile = 95,
          y.relation = "free")

## plot the number of valid records in each 2-week period
timePlot(mydata, pollutant = c("o3", "no2"), avg.time = "2 week",
          statistic = "frequency", y.relation = "free")
```

An example of normalising data is shown in [Figure 17.2](#). In this plot we have:

- Averaged the data to annual means;
- Chosen to normalise to the beginning of 2008;
- Set the line width to 4 and the line type to 1 (continuous line);
- Chosen to group the data in one panel.

[Figure 17.2](#) shows that concentrations of NO_2 and O_3 have increased over the period 1998–2005; SO_2 and CO have shown the greatest reductions (by about 60%), whereas NO_x concentrations have decreased by about 20%.

Another example is grouping pollutants from *several sites* on one plot. It is easy to import data from several sites and to plot the data in separate panels e.g.

```
data(mydata)
timePlot(selectByDate(mydata, year = 2003, month = "aug"),
          pollutant = c("nox", "o3", "pm25", "pm10", "ws"),
          y.relation = "free")
```

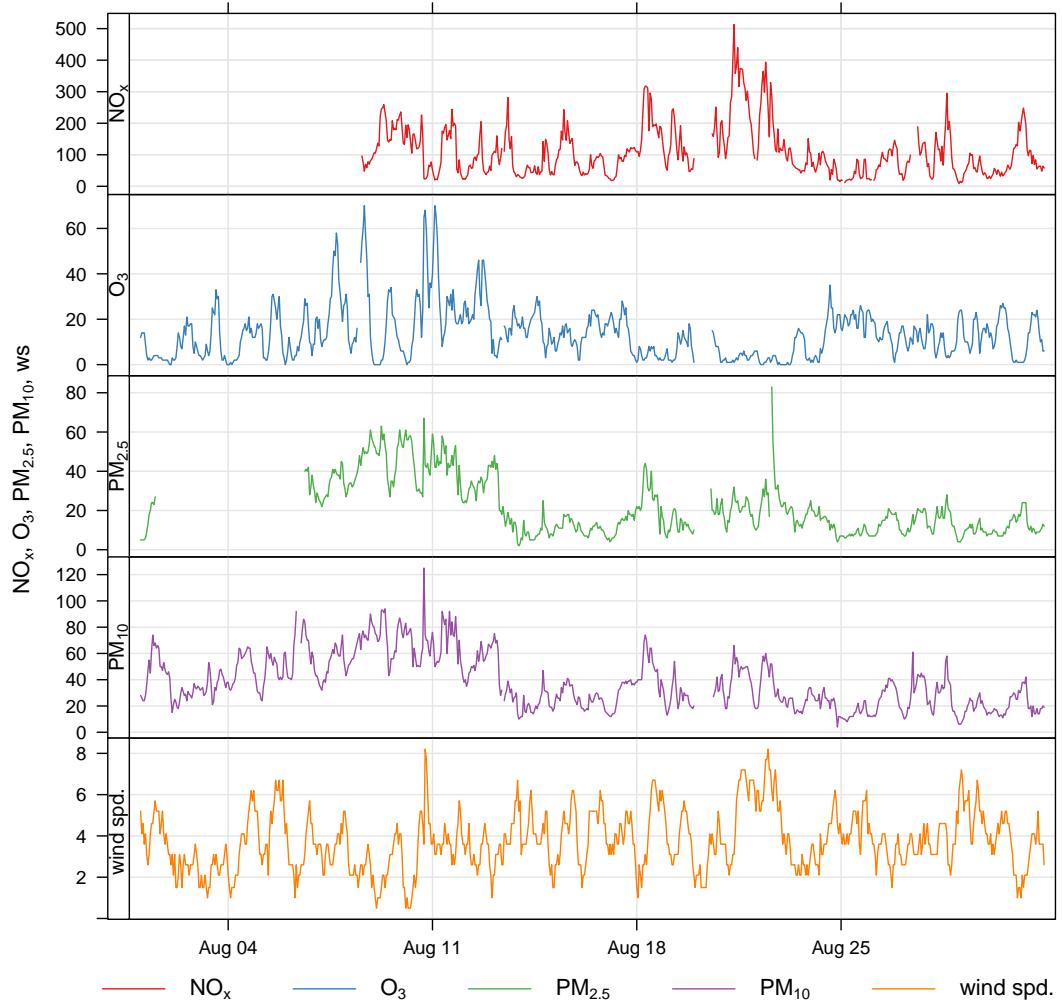


FIGURE 17.1 Time series for several variables using the `timePlot` and the `selectByDate` functions. The data shown are for August 2003.

```
## import data from 3 sites
thedata <- importAURN(site = c("kc1", "my1", "nott"), year = 2005:2010)

## plot it
timePlot(thedata, pollutant = "nox", type = "site", avg.time = "month")
```

Using the code above it is also possible to include several species. But what if we wanted to plot NO_x concentrations across all sites in one panel? To do this we need to re-organise the data, as described in Section 5.4. An example of how to do this is shown below. Note, in order to make referring to the columns easier, we will drop the full (long) site name and use the site code instead.

```
timePlot(mydata, pollutant = c("nox", "no2", "co", "so2", "pm10"),
         avg.time = "year", normalise = "1/1/1998", lwd = 4, lty = 1,
         group = TRUE, ylim = c(0, 120))
```

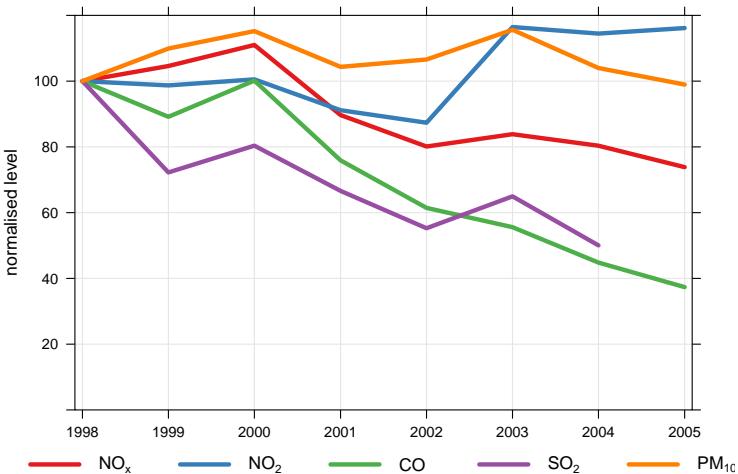


FIGURE 17.2 An example of normalising time series data to fix values to equal 100 at the beginning of 1998.

```
## first drop site name
thedata <- subset(thedata, select = -site)
## now reshape the data using the reshape package
thedata <- melt(thedata, id.vars = c("date", "code"))
thedata <- dcast(thedata, ... ~ code + variable)
```

The final step will make columns of each site/pollutant combination e.g. ‘KC1_nox’, ‘KC1_pm10’ and so on. It is then easy to use any of these names to make the plot:

```
timePlot(thedata, pollutant = c("KC1_nox", "MY1_nox", "NOTT_nox"),
         avg.time = "month", group = TRUE)
```

An alternative way of selecting all columns containing the character ‘nox’ is to use the `grep` command (see page 31). For example:

```
timePlot(thedata, pollutant = names(thedata)[grep(pattern = "nox",
                                                 names(thedata))],
         avg.time = "month", group = TRUE)
```

If wind speed (`ws`) and wind direction (`wd`) are available they can be used in plots and shown as ‘wind vectors’. Plotting data in this way conveys more information in an easy to understand way, which works best for relatively short time periods e.g. a pollution episode lasting a few days. As an example Figure 17.3 shows the first 48 hours of NO_x and NO₂ data with wind arrows shown. The arrows are controlled by a list of option that control the length, shape and colour of the arrows. The maximum length of the arrow plotted is a fraction of the plot dimension with the longest arrow being `scale` of the plot x-y dimension. Note, if the plot size is adjusted manually by the user it should be re-plotted to ensure the correct wind angle. The list may contain other options to `panel.arrows` in the `lattice` package. Other useful options include `length`, which controls the length of the arrow head and `angle`, which controls the angle of the arrow head. Wind vector arrows can also be used with the `scatterPlot` function.

```
timePlot(head(mydata, 48), pollutant = c("nox", "no2"),
         windflow = list(scale = 0.1, lwd = 2, col = "pink"),
         lwd = 3, group = FALSE,
         ylab = "concentration (ug/m3)")
```

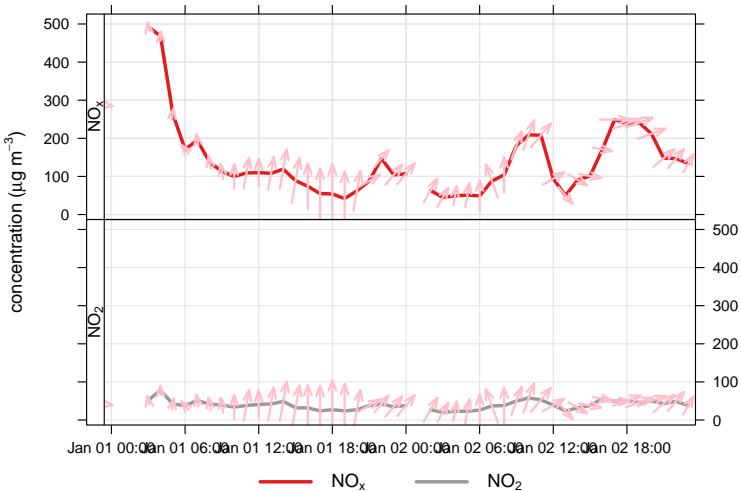


FIGURE 17.3 An example of using the `windflow` option in `timePlot`.

18 The `timeProp` function

The `timeProp` ('time proportion') function shows time series plots as stacked bar charts. For a particular time, proportions of a chosen variable are shown as a stacked bar chart. The different categories in the bar chart are made up from a character or factor variable in a data frame. The function is primarily developed to support the plotting of cluster analysis output from `polarCluster` (see Section 15) and `trajCluster` (see Section 26.3) that consider local and regional (back trajectory) cluster analysis respectively. However, the function has more general use for understanding time series data. In order to plot time series in this way, some sort of time aggregation is needed, which is controlled by the option `avg.time`.

The plot shows the value of `pollutant` on the y-axis (averaged according to `avg.time`). The time intervals are made up of bars split according to `proportion`. The bars therefore show how the total value of `pollutant` is made up for any time interval.

The `timeProp` function has the following options:

mydata A data frame containing the fields `date`, `pollutant` and a splitting variable `proportion`

pollutant Name of the pollutant to plot contained in `mydata`.

proportion The splitting variable that makes up the bars in the bar chart e.g. `proportion = "cluster"` if the output from `polarCluster` is being analysed. If `proportion` is a numeric variable it is split into 4 quantiles (by default) by `cutData`. If `proportion` is a factor or character variable then the categories are used directly.

avg.time This defines the time period to average to. Can be "sec", "min", "hour", "day", "DSTday", "week", "month", "quarter" or "year". For much increased flexibility a number can precede these options followed by a

space. For example, a timeAverage of 2 months would be `period = "2 month"`. In addition, `avg.time` can equal “season”, in which case 3-month seasonal values are calculated with spring defined as March, April, May and so on.

Note that `avg.time` when used in `timeProp` should be greater than the time gap in the original data. For example, `avg.time = "day"` for hourly data is OK, but `avg.time = "hour"` for daily data is not.

<code>type</code>	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. “season”, “year”, “weekday” and so on. For example, <code>type = "season"</code> will produce four plots — one for each season. It is also possible to choose <code>type</code> as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another. <code>type</code> must be of length one.
<code>statistic</code>	Determines how the bars are calculated. The default (“mean”) will provide the contribution to the overall mean for a time interval. <code>statistic = "frequency"</code> will give the proportion in terms of counts.
<code>normalise</code>	If <code>normalise = TRUE</code> then each time interval is scaled to 100. This is helpful to show the relative (percentage) contribution of the proportions.
<code>cols</code>	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours — see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
<code>date.breaks</code>	Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of <code>date.breaks</code> up or down.
<code>date.format</code>	This option controls the date format on the x-axis. While <code>timePlot</code> generally sets the date format sensibly there can be some situations where the user wishes to have more control. For format types see <code>strptime</code> . For example, to format the date like “Jan-2012” set <code>date.format = "%b-%Y"</code> .
<code>key.columns</code>	Number of columns to be used in the key. With many pollutants a single column can make to key too wide. The user can thus choose to use several columns by setting <code>columns</code> to be less than the number of pollutants.
<code>key.position</code>	Location where the scale key is to plotted. Allowed arguments currently include “top”, “right”, “bottom” and “left”.
<code>key.title</code>	The title of the key.

```
timeProp(selectByDate(mydata, year = 2003), pollutant = "so2", avg.time = "3 day",
proportion = "wd", date.breaks = 10, key.position = "top",
key.columns = 8, ylab = "so2 (ug/m3)")
```

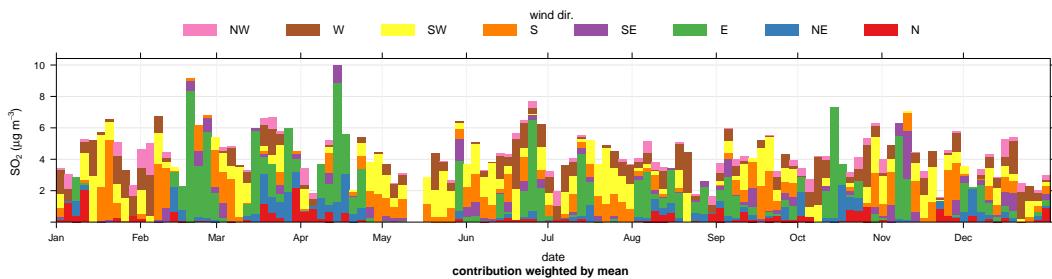


FIGURE 18.1 `timeProp` plot for SO₂ concentrations in 2003. The data are categorised into 8 wind sectors for 3-day averages.

`auto.text` Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels etc. will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO₂.

`...` Other graphical parameters passed onto `timeProp` and `cutData`. For example, `timeProp` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

An example of the `timeProp` function is shown in Figure 18.1. In this example SO₂ concentrations are considered for 2003 (using the `selectByDate` function). The averaging period is set to 3 days and the mean concentration is plotted and the proportion contribution by wind sector is given. Other options are chosen to place the key at the top and choose the number of columns used in the key. It is apparent from Figure 18.1 that the highest SO₂ concentrations are dominated by winds from an easterly sector, but actually occur throughout the year.

Note that `proportion` can be an existing categorical (i.e. factor or character) variable in a data frame. If a numeric variable is supplied, then it is typically cut into four quantile levels. So, for example, the plot below would show four intervals of wind speed, which would help show the wind speed conditions that control high SO₂ concentration — and importantly, when they occur.

One of the key uses of `timeProp` is to post-process cluster analysis data. Users should consider the uses of `timeProp` for cluster analysis shown in Section 15 and Section 26.3. In both these cases the cluster analysis yields a categorical output directly i.e. cluster, which lends itself to analysis using `timeProp`.

```
timeProp(selectByDate(mydata, year = 2003), pollutant = "so2",
avg.time = "3 day", proportion = "ws", date.breaks = 10,
key.position = "top", key.columns = 4)
```

19 The `calendarPlot` function

19.1 Purpose

Sometimes it is useful to visualise data in a familiar way. Calendars are the obvious way to represent data for data on the time scale of days or months. The `calendarPlot` function provides an effective way to visualise data in this way by showing daily concentrations laid out in a calendar format. The concentration of a species is shown by its colour. The data can be shown in different ways. By default `calendarPlot` overlays the day of the month. However, if wind speed and wind direction are available then an arrow can be shown for each day giving the vector-averaged wind direction. In addition, the arrow can be scaled according to the wind speed to highlight both the direction and strength of the wind on a particular day, which can help show the influence of meteorology on pollutant concentrations.

`calendarPlot` can also show the daily mean concentration as a number on each day and can be extended to highlight those conditions where daily mean (or maximum etc.) concentrations are above a particular threshold. This approach is useful for highlighting daily air quality limits e.g. when the daily mean concentration is greater than $50 \mu\text{g m}^{-3}$.

The `calendarPlot` function can also be used to plot categorical scales. This is useful for plotting concentrations expressed as an air quality index i.e. intervals of concentrations that are expressed in ways like ‘very good’, ‘good’, ‘poor’ and so on.

19.2 Options available

<code>mydata</code>	A data frame minimally containing <code>date</code> and at least one other numeric variable. The date should be in either <code>Date</code> format or class <code>POSIXct</code> .
<code>pollutant</code>	Mandatory. A pollutant name corresponding to a variable in a data frame should be supplied e.g. <code>pollutant = "nox"</code> .
<code>year</code>	Year to plot e.g. <code>year = 2003</code> . If not supplied all data potentially spanning several years will be plotted.
<code>month</code>	If only certain month are required. By default the function will plot an entire year even if months are missing. To only plot certain months use the <code>month</code> option where month is a numeric 1:12 e.g. <code>month = c(1, 12)</code> to only plot January and December.
<code>type</code>	Not yet implemented.
<code>annotate</code>	This option controls what appears on each day of the calendar. Can be: “date” — shows day of the month; “wd” — shows vector-averaged wind direction, or “ws” — shows vector-averaged wind direction scaled by wind speed. Finally it can be “value” which shows the daily mean value.
<code>statistic</code>	Statistic passed to <code>timeAverage</code> .
<code>cols</code>	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours — see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
<code>limits</code>	Use this option to manually set the colour scale limits. This is useful in the case when there is a need for two or more plots and a consistent scale

is needed on each. Set the limits to cover the maximum range of the data for all plots of interest. For example, if one plot had data covering 0–60 and another 0–100, then set `limits = c(0, 100)`. Note that data will be ignored if outside the limits range.

<code>lim</code>	A threshold value to help differentiate values above and below <code>lim</code> . It is used when <code>annotate = "value"</code> . See next few options for control over the labels used.
<code>col.lim</code>	For the annotation of concentration labels on each day. The first sets the colour of the text below <code>lim</code> and the second sets the colour of the text above <code>lim</code> .
<code>col.arrow</code>	The colour of the annotated wind direction / wind speed arrows.
<code>font.lim</code>	For the annotation of concentration labels on each day. The first sets the font of the text below <code>lim</code> and the second sets the font of the text above <code>lim</code> . Note that font = 1 is normal text and font = 2 is bold text.
<code>cex.lim</code>	For the annotation of concentration labels on each day. The first sets the size of the text below <code>lim</code> and the second sets the size of the text above <code>lim</code> .
<code>digits</code>	The number of digits used to display concentration values when <code>annotate = "value"</code> .
<code>data.thresh</code>	Data capture threshold passed to <code>timeAverage</code> . For example, <code>data.thresh = 75</code> means that at least 75% of the data must be available in a day for the value to be calculate, else the data is removed.
<code>labels</code>	If a categorical scale is required then these labels will be used. Note there is one less label than break. For example, <code>labels = c("good", "bad", "very bad")</code> . <code>breaks</code> must also be supplied if labels are given.
<code>breaks</code>	If a categorical scale is required then these breaks will be used. For example, <code>breaks = c(0, 50, 100, 1000)</code> . In this case “good” corresponds to values between 0 and 50 and so on. Users should set the maximum value of <code>breaks</code> to exceed the maximum data value to ensure it is within the maximum final range e.g. 100–1000 in this case.
<code>w.shift</code>	Controls the order of the days of the week. By default the plot shows Saturday first (<code>w.shift = 0</code>). To change this so that it starts on a Monday for example, set <code>w.shift = 2</code> , and so on.
<code>remove.empty</code>	Should months with no data present be removed? Default is <code>TRUE</code> .
<code>main</code>	The plot title; default is pollutant and year.
<code>key.header</code>	Adds additional text/labels to the scale key. For example, passing <code>calendarPlot(mydata, key.header = "header", key.footer = "footer")</code> adds addition text above and below the scale key. These arguments are passed to <code>drawOpenKey</code> via <code>quickText</code> , applying the <code>auto.text</code> argument, to handle formatting.
<code>key.footer</code>	see <code>key.header</code> .

```
calendarPlot(mydata, pollutant = "o3", year = 2003)
```

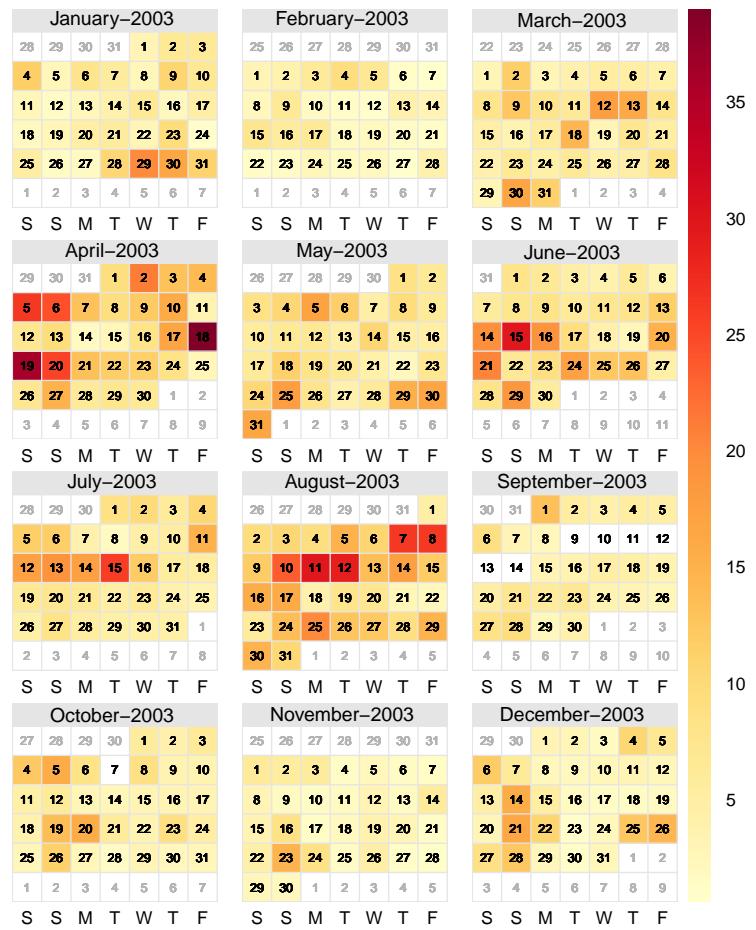


FIGURE 19.1 `calendarPlot` for O₃ concentrations in 2003.

key.position Location where the scale key is to plotted. Allowed arguments currently include "top", "right", "bottom" and "left".

key Fine control of the scale key via `drawOpenKey`. See `drawOpenKey` for further details.

auto.text Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO₂.

... Other graphical parameters are passed onto the `lattice` function `lattice:levelplot`, with common axis and title labelling options (such as `xlab`, `ylab`, `main`) being passed to via `quickText` to handle routine formatting.

19.3 Example of use

The function is called in the usual way. As a minimum, a data frame, pollutant and year is required. So to show O₃ concentrations for each day in 2003 (Figure 19.1). Note that if `year` is not supplied the full data set will be used.

It is sometimes useful to annotate the plots with other information. It is possible to show the daily mean wind angle, which can also be scaled to wind speed. The idea here

```
data(mydata) ## make sure openair 'mydata' loaded fresh
calendarPlot(mydata, pollutant = "pm10", year = 2003, annotate = "value",
             lim = 50, cols = "Purples", col.lim = c("black", "orange"),
             layout = c(4, 3))
```

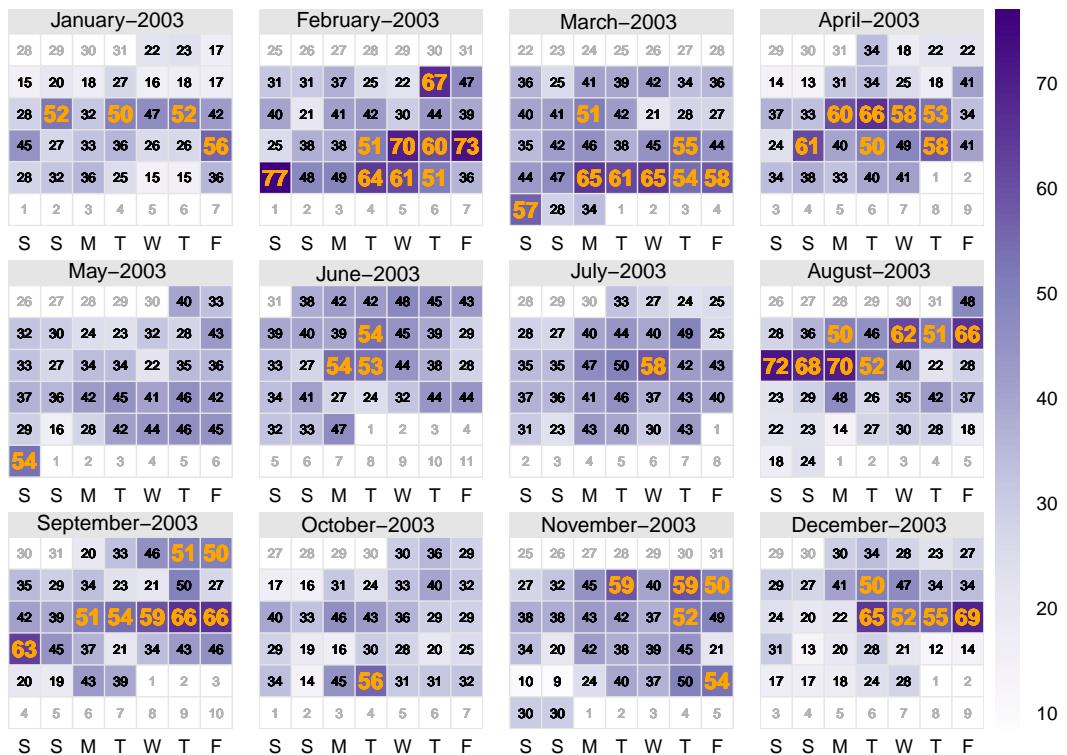


FIGURE 19.2 `calendarPlot` for PM₁₀ concentrations in 2003 with annotations highlighting those days where the concentration of PM₁₀ > 50 $\mu\text{g m}^{-3}$. The numbers show the PM₁₀ concentration in $\mu\text{g m}^{-3}$.

being to provide some information on meteorological conditions on each day. Another useful option is to set `annotate = "value"` in which case the daily concentration will be shown on each day. Furthermore, it is sometimes useful to highlight particular values more clearly. For example, to highlight daily mean PM₁₀ concentrations above 50 $\mu\text{g m}^{-3}$. This is where setting `lim` (a concentration limit) is useful. In setting `lim` the user can then differentiate the values below and above `lim` by colour of text, size of text and type of text e.g. plain and bold.

Figure 19.2 highlights those days where PM₁₀ concentrations exceed 50 $\mu\text{g m}^{-3}$ by making the annotation for those days bigger, bold and orange. Plotting the data in this way clearly shows the days where PM₁₀ > 50 $\mu\text{g m}^{-3}$.

Other `openair` functions can be used to plot other statistics. For example, `rollingMean` could be used to calculate rolling 8-hour mean O₃ concentrations. Then, `calendarPlot` could be used with `statistic = "max"` to show days where the maximum daily rolling 8-hour mean O₃ concentration is greater than a certain threshold e.g. 100 or 120 $\mu\text{g m}^{-3}$.

To show wind angle, scaled to wind speed (Figure 19.3).

Note again that `selectByDate` can be useful. For example, to plot select months:

```
calendarPlot(selectByDate(mydata, year = 2003, month = c("jun", "jul", "aug")),
             pollutant = "o3", year = 2003)
```

```
calendarPlot(mydata, pollutant = "o3", year = 2003, annotate = "ws")
```

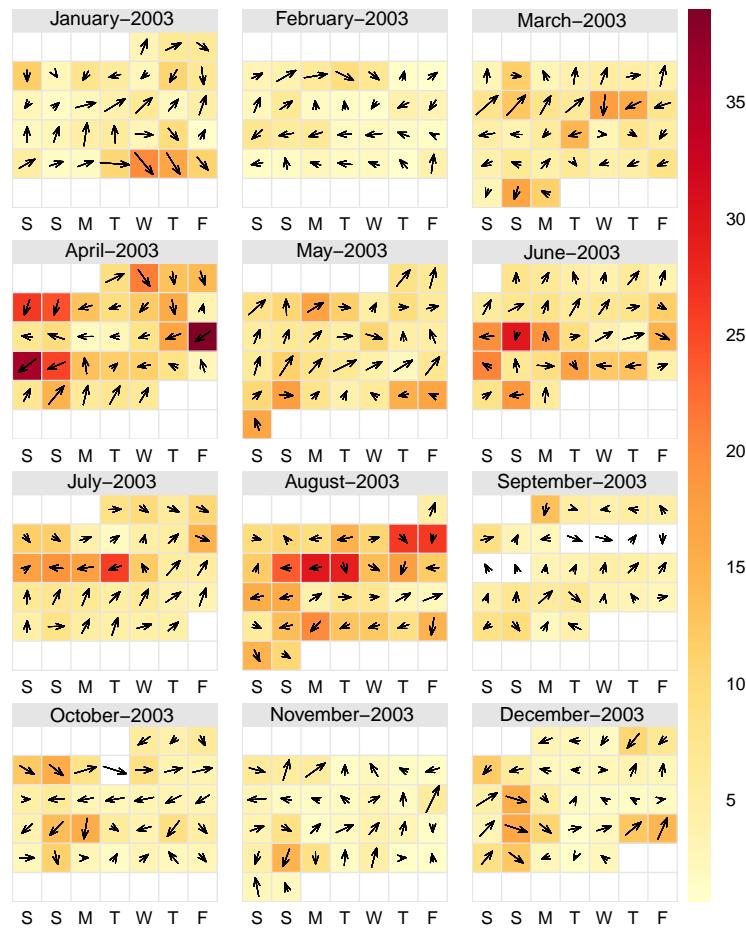


FIGURE 19.3 `calendarPlot` for O₃ concentrations in 2003 with annotations showing wind angle scaled to wind speed i.e. the longer the arrow, the higher the wind speed. It shows for example high O₃ concentrations on the 17 and 18th of April were associated with strong north-easterly winds.

Figure 19.4 shows an example of plotting data with a categorical scale. In this case the options `labels` and `breaks` have been used to define concentration intervals and their descriptions. Note that `breaks` needs to be one longer than `labels`. In the example in Figure 19.4 the first interval ('Very low') is defined as concentrations from 0 to 50 (ppb), 'Low' is 50 to 100 and so on. Note that the upper value of `breaks` should be a number greater than the maximum value contained in the data to ensure that it is encompassed. In the example given in Figure 19.4 the maximum daily concentration is plotted. These types of plots are very useful for considering national or international air quality indexes.

The user can explicitly set each colour interval:

```
calendarPlot(mydata, pollutant = "no2", year = 2003,
  breaks = c(0, 50, 100, 150, 1000),
  labels = c("Very_low", "Low", "High", "Very_High"),
  cols = c("lightblue", "forestgreen", "yellow", "red"),
  statistic = "max")
```

```
calendarPlot(mydata, pollutant = "no2", year = 2003,
  breaks = c(0, 50, 100, 150, 1000),
  labels = c("Very low", "Low", "High", "Very High"),
  cols = "increment", statistic = "max")
```

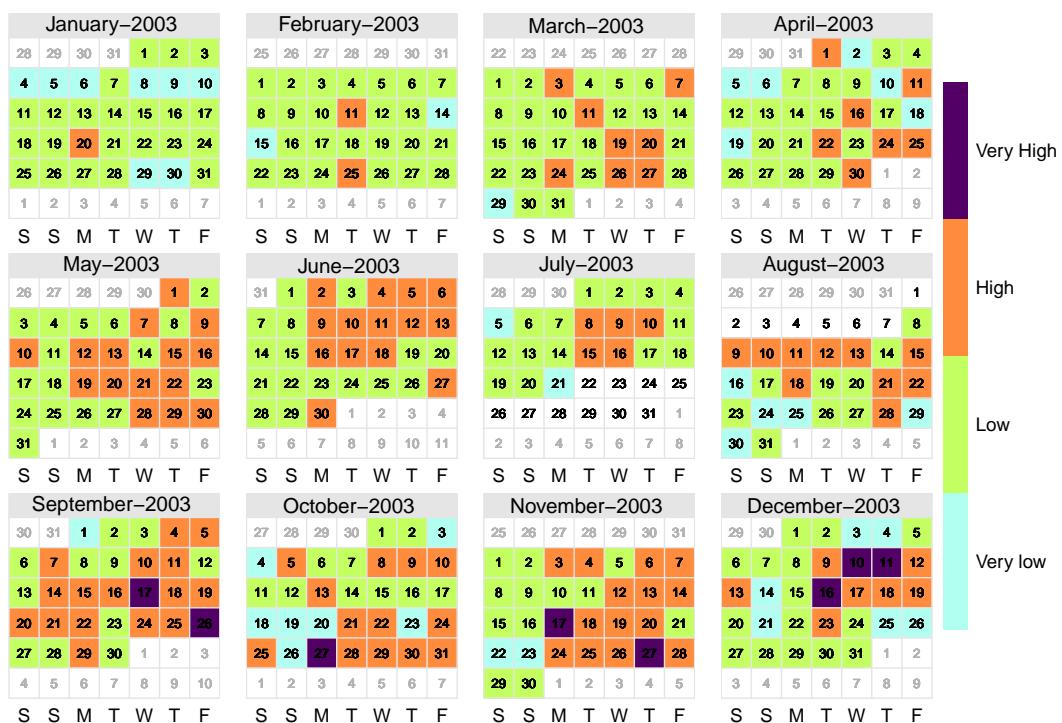


FIGURE 19.4 `calendarPlot` for NO₂ concentrations in 2003 with a user-defined categorical scale.

Note that in the case of categorical scales it is possible to define the breaks and labels first and then make the plot. For example:

```
breaks <- c(0, 34, 66, 100, 121, 141, 160, 188, 214, 240, 500)
labels <- c("Low.1", "Low.2", "Low.3", "Moderate.4", "Moderate.5", "Moderate.6",
          "High.7", "High.8", "High.9", "Very High.10")

calendarPlot(mydata, pollutant = "no2", year = 2003,
             breaks = breaks, labels = labels,
             cols = "jet", statistic = "max")
```

It is also possible to first use `rollingMean` to calculate statistics. For example, if one was interested in plotting the maximum daily rolling 8-hour mean concentration, the data could be prepared and plotted as follows.

```
## makes a new field 'rolling8o3'
dat <- rollingMean(dat, pollutant = "o3", hours = 8)
breaks <- c(0, 34, 66, 100, 121, 141, 160, 188, 214, 240, 500)
labels <- c("Low.1", "Low.2", "Low.3", "Moderate.4", "Moderate.5", "Moderate.6",
          "High.7", "High.8", "High.9", "Very High.10")

calendarPlot(mydata, pollutant = "rolling8o3", year = 2003,
             breaks = breaks, labels = labels,
             cols = "jet", statistic = "max")
```

The UK has an air quality index for O₃, NO₂, PM₁₀ and PM_{2.5} described in detail at <http://uk-air.defra.gov.uk/air-pollution/daqi> and COMEAP (2011). The air quality index is shown in Table 19.1. The index is most relevant to air quality forecasting, but is used widely for public information. Most other countries have similar indexes. Note that the indexes are calculated for different averaging times dependent on the pollutant: rolling 8-hour mean for O₃, hourly means for NO₂ and a fixed 24-hour mean for PM₁₀ and PM_{2.5}.

TABLE 19.1 The UK daily air quality index (values in $\mu\text{g m}^{-3}$).

Band	Description	NO ₂	O ₃	PM ₁₀	PM _{2.5}
1	Low	0–66	0–33	0–16	0–11
2	Low	67–133	34–65	17–33	12–23
3	Low	134–199	66–99	34–49	24–34
4	Moderate	200–267	100–120	50–58	35–41
5	Moderate	268–334	121–140	59–66	42–46
6	Moderate	335–399	141–159	67–74	47–52
7	High	400–467	160–187	75–83	53–58
8	High	468–534	188–213	84–91	59–64
9	High	535–599	214–239	92–99	65–69
10	Very High	600 or more	240 or more	100 or more	70 or more

In the code below the labels and breaks are defined for each pollutant in Table 19.1 to make it easier to use the index in the `calendarPlot` function.

```

## the labels - same for all species
labels <- c("1 - Low", "2 - Low", "3 - Low", "4 - Moderate", "5 - Moderate",
          "6 - Moderate", "7 - High", "8 - High", "9 - High", "10 - Very High")
o3.breaks <- c(0, 34, 66, 100, 121, 141, 160, 188, 214, 240, 500)
no2.breaks <- c(0, 67, 134, 200, 268, 335, 400, 468, 535, 600, 1000)
pm10.breaks <- c(0, 17, 34, 50, 59, 67, 75, 84, 92, 100, 1000)
pm25.breaks <- c(0, 12, 24, 35, 42, 47, 53, 59, 65, 70, 1000)

```

Remember it is necessary to use the correct averaging time. Assuming data are imported using `importAURN` or `importKCL` then the units will be in $\mu\text{g m}^{-3}$ — if not the user should ensure this is the case. Note that rather than showing the day of the month (the default), `annotate = "value"` can be used to show the actual numeric value on each day. In this way, the colours represent the categorical interval the concentration on a day corresponds to and the actual value itself is shown.

```

## import test data
dat <- importAURN(site = "kc1", year = 2010)

## no2 index example
calendarPlot(dat, year = 2010, pollutant = "no2", labels = labels,
              breaks = no2.breaks, statistic = "max", cols = "jet")

## for PM10 or PM2.5 we need the daily mean concentration
calendarPlot(dat, year = 2010, pollutant = "pm10", labels = labels,
              breaks = pm10.breaks, statistic = "mean", cols = "jet")

## for ozone, need the rolling 8-hour mean
dat <- rollingMean(dat, pollutant = "o3", hours = 8)
calendarPlot(dat, year = 2010, pollutant = "rolling8o3", labels = labels,
              breaks = o3.breaks, statistic = "max", cols = "jet")

```

20 The TheilSen function

20.1 Purpose

see also
`smoothTrend`
`timePlot`

Calculating trends for air pollutants is one of the most important and common tasks that can be undertaken. Trends are calculated for all sorts of reasons. Sometimes it is useful to have a general idea about how concentrations might have changed. On other occasions a more definitive analysis is required; for example, to establish statistically whether a trend is significant or not. The whole area of trend calculation is a complex one and frequently trends are calculated with little consideration as to their validity. Perhaps the most common approach is to apply linear regression and not think twice about it. However, there can be many pitfalls when using ordinary linear regression, such as the assumption of normality, autocorrelation etc.

One commonly used approach for trend calculation in studies of air pollution is the non-parametric *Mann-Kendall* approach (Hirsch et al. 1982). Wilcox (2010) provides an excellent case for using ‘modern methods’ for regression including the benefits of non-parametric approaches and bootstrap simulations. Note also that the all the regression parameters are estimated through bootstrap resampling.

The Theil-Sen method dates back to 1950, but the basic idea pre-dates 1950 (Theil 1950; Sen 1968). It is one of those methods that required the invention of fast computers to be practical. The basic idea is as follows. Given a set of n x, y pairs, the slopes between all pairs of points are calculated. Note, the number of slopes can increase by $\approx n^2$ so that the number of slopes can increase rapidly as the length of the data set increases. The Theil-Sen estimate of the slope is the median of all these slopes. The advantage of

the using the Theil-Sen estimator is that it tends to yield accurate confidence intervals even with non-normal data and heteroscedasticity (non-constant error variance). It is also resistant to outliers — both characteristics can be important in air pollution. As previously mentioned, the estimates of these parameters can be made more robust through bootstrap-resampling, which further adds to the computational burden, but is not an issue for most time series which are expressed either as monthly or annual means. Bootstrap resampling also provides the estimate of p for the slope.

An issue that can be very important for time series is dependence or *autocorrelation* in the data. Normal (in the statistical sense) statistics assume that data are independent, but in time series this is rarely the case. The issue is that neighbouring data points are similar to one another (correlated) and therefore not independent. Ignoring this dependence would tend to give an overly optimistic impression of uncertainties. However, taking account of it is far from simple. A discussion of these issues is beyond the aims of this report and readers are referred to standard statistical texts on the issue. In **openair** we follow the suggestion of Kunsch (1989) of setting the block length to $n^{1/3}$ where n is the length of the time series.

There is a temptation when considering trends to use all the available data. Why? Often it is useful to consider specific periods. For example, is there any evidence that concentrations of NO_x have decreased since 2000? Clearly, the time period used depends on both the data and the questions, but it is good to be aware that considering subsets of data can be very insightful.

Another aspect is that almost all trends are shown as mean concentration versus time; typically by year. Such analyses are very useful for understanding how concentrations have changed through time and for comparison with air quality limits and regulations. However, if one is interested in *understanding* why trends are as they are, it can be helpful to consider how concentrations vary in other ways. The trend functions in **openair** do just this. Trends can be plotted by day of the week, month, hour of the day, by wind direction sector and by different wind speed ranges. All these capabilities are easy to use and their effectiveness will depend on the situation in question. One of the reasons that trends are not considered in these many different ways is that there can be a considerable overhead in carrying out the analysis, which is avoided by using these functions. Few, for example, would consider a detailed trend analysis by hour of the day, ensuring that robust statistical methods were used and uncertainties calculated. However, it can be useful to consider how concentrations vary in this way. It may be, for example, that the hours around midday are dominated by heavy vehicle emissions rather than by cars — so is the trend for a pollutant different for those hours compared with say, hours dominated by other vehicle types? Similarly, a much more focussed trend analysis can be done by considering different wind direction, as this can help isolate different source influences.

The **TheilSen** function is typically used to determine trends in pollutant concentrations over several years. However, it can be used to calculate the trend in any numeric variable. It calculates monthly mean values from daily, hourly or higher time resolution data, as well as working directly with monthly means. Whether it is meaningful to calculate trends over shorter periods of time (e.g. 2 years) depends very much on the data. It may well be that statistically significant trends can be detected over relatively short periods but it is another matter whether it matters. Because seasonal effects can be important for monthly data, there is the option to deseasonalise the data first. The **timeVariation** function are both useful to determine whether there is a seasonal cycle that should be removed.

Note also that the symbols shown next to each trend estimate relate to how statistically significant the trend estimate is: $p < 0.001 = ***$, $p < 0.01 = **$, $p < 0.05 = *$ and $p <$

$0.1 = +.$

20.2 Options available

The `TheilSen` function has the following options:

- mydata** A data frame containing the field `date` and at least one other parameter for which a trend test is required; typically (but not necessarily) a pollutant.
- pollutant** The parameter for which a trend test is required. Mandatory.
- deseason** Should the data be de-seasonalized first? If `TRUE` the function `stl` is used (seasonal trend decomposition using loess). Note that if `TRUE` missing data are first imputed using a Kalman filter and Kalman smooth.
- type** `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.
It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.
Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.
- avg.time** Can be “month” (the default), “season” or “year”. Determines the time over which data should be averaged. Note that for “year”, six or more years are required. For “season” the data are split up into spring: March, April, May etc. Note that December is considered as belonging to winter of the following year.
- statistic** Statistic used for calculating monthly values. Default is “mean”, but can also be “percentile”. See `timeAverage` for more details.
- percentile** Single percentile value to use if `statistic = "percentile"` is chosen.
- data.thresh** The data capture threshold to use (the data using `avg.time`. A value of zero means that all available data will be used in a particular period regardless if of the number of values available. Conversely, a value of 100 will mean that all data will need to be present for the average to be calculated, else it is recorded as `NA`.
- alpha** For the confidence interval calculations of the slope. The default is 0.05. To show 99% confidence intervals for the value of the trend, choose alpha = 0.01 etc.
- dec.place** The number of decimal places to display the trend estimate at. The default is 2.

xlab	x-axis label, by default "year".
lab.frac	Fraction along the y-axis that the trend information should be printed at, default 0.99.
lab.cex	Size of text for trend information.
x.relation	This determines how the x-axis scale is plotted. "same" ensures all panels use the same scale and "free" will use panel-specific scales. The latter is a useful setting when plotting data with very different values.
y.relation	This determines how the y-axis scale is plotted. "same" ensures all panels use the same scale and "free" will use panel-specific scales. The latter is a useful setting when plotting data with very different values.
data.col	Colour name for the data
trend	list containing information on the line width, line type and line colour for the main trend line and confidence intervals respectively.
text.col	Colour name for the slope/uncertainty numeric estimates
slope.text	The text shown for the slope (default is 'units/year').
cols	Predefined colour scheme, currently only enabled for "greyscale".
shade	The colour used for marking alternate years. Use "white" or "transparent" to remove shading.
auto.text	Either TRUE (default) or FALSE . If TRUE titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO2.
autocor	Should autocorrelation be considered in the trend uncertainty estimates? The default is FALSE . Generally, accounting for autocorrelation increases the uncertainty of the trend estimate — sometimes by a large amount.
slope.percent	Should the slope and the slope uncertainties be expressed as a percentage change per year? The default is FALSE and the slope is expressed as an average units/year change e.g. ppb. Percentage changes can often be confusing and should be clearly defined. Here the percentage change is expressed as $100 * (C.end/C.start - 1) / (\text{end.year} - \text{start.year})$. Where C.start is the concentration at the start date and C.end is the concentration at the end date. For avg.time = "year" ($\text{end.year} - \text{start.year}$) will be the total number of years - 1. For example, given a concentration in year 1 of 100 units and a percentage reduction of 5 units but the actual time span will be 6 years i.e. year 1 is used as a reference year. Things are slightly different for monthly values e.g. avg.time = "month" , which will use the total number of months as a basis of the time span and is therefore able to deal with partial years. There can be slight differences in the depending on whether monthly or annual values are considered.
date.breaks	Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work

```
TheilSen(mydata, pollutant = "o3", ylab = "ozone (ppb)", deseason = TRUE)
## [1] "Taking bootstrap samples. Please wait."
```

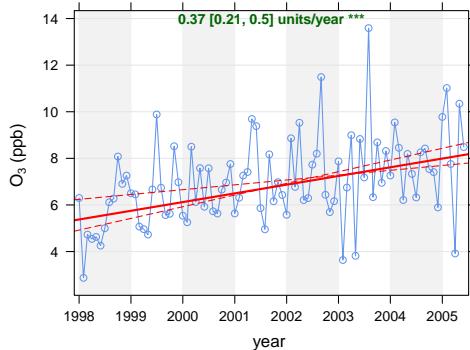


FIGURE 20.1 Trends in ozone at Marylebone Road. The plot shows the deseasonalised monthly mean concentrations of O₃. The solid red line shows the trend estimate and the dashed red lines show the 95 % confidence intervals for the trend based on resampling methods. The overall trend is shown at the top-left as 0.38 (ppb) per year and the 95 % confidence intervals in the slope from 0.21–0.51 ppb/year. The * * * show that the trend is significant to the 0.001 level.

as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of `date.breaks` up or down.

- `plot` Should a plot be produced. `FALSE` can be useful when analysing data to extract trend components and plotting them in other ways.
- `silent` When `FALSE` the function will give updates on trend-fitting progress.
- `...` Other graphical parameters passed onto `cutData` and `lattice:xyplot`. For example, `TheilSen` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

20.3 Example of use

We first show the use of the `TheilSen` function by applying it to concentrations of O₃. The function is called as shown in Figure 20.1.

Because the function runs simulations to estimate the uncertainty in the slope, it can take a little time for all the calculations to finish. These printed results show that in this case the trend in O₃ was +0.38 units (i.e. ppb) per year as an average over the entire period. It also shows the 95 % confidence intervals in the trend ranged between 0.21 to 0.51 ppb/year. Finally, the significance level in this case is very high; providing very strong evidence that concentrations of O₃ increased over the period. The plot together with the summary results is shown in Figure 20.1. Note that if one wanted to display the confidence intervals in the slope at the 99 % confidence intervals, the code would be Figure 20.2.

```
TheilSen(mydata, pollutant = "o3", ylab = "ozone (ppb)", alpha = 0.01)
```

Sometimes it is useful to consider a subset of data, perhaps by excluding some years. This is easy with the `subset` function. The following code calculates trends for years greater than 1999 i.e. from 2000 onwards.

```
TheilSen(subset(mydata, format(date, "%Y") > 1999), pollutant = "o3",
        ylab = "ozone (ppb)")
```

It is also possible to calculate trends in many other ways e.g. by wind direction. Considering how trends vary by wind direction can be extremely useful because the influence of different sources invariably depends on the direction of the wind. The `TheilSen` function splits the wind direction into 8 sectors i.e. N, NE, E etc. The Theil-Sen slopes are then calculated for each direction in turn. This function takes rather longer to run because the simulations need to be run eight times in total. Considering concentrations of O₃ again, the output is shown in Figure 20.2. Note that this plot is specifically laid out to assist interpretation, with each panel located at the correct point on the compass. This makes it easy to see immediately that there is essentially no trend in O₃ for southerly winds i.e. where the road itself has the strongest influence. On the other hand the strongest evidence of increasing O₃ are for northerly winds, where the influence of the road is much less. The reason that there is no trend in O₃ for southerly winds is that there is always a great excess of NO, which reacts with O₃ to form NO₂. At this particular location it will probably take many more years before O₃ concentrations start to increase when the wind direction is southerly. Nevertheless, there will always be *some* hours that do not have such high concentrations of NO.

The option `slope.percent` can be set to express slope estimates as a percentage change per year. This is useful for comparing slopes for sites with very different concentration levels and for comparison with emission inventories. The percentage change uses the concentration at the beginning and end months to express the mean slope.

The trend, T is defined as:

$$T[\%.\text{yr}^{-1}] = 100 \cdot \left(\frac{C_{\text{End}}}{C_{\text{Start}}} - 1 \right) / N_{\text{years}} \quad (7)$$

where C_{End} and C_{Start} are the mean concentrations for the end and start date, respectively. N_{years} is the number of years (or fractions of) the time series spans.

```
TheilSen(mydata, pollutant = "o3", deseason = TRUE,
        slope.percent = TRUE)
```

The `TheilSen` function was written to work with hourly data, which is then averaged into monthly or annual data. However, it is realised that users may already have data that is monthly or annual. The function can therefore accept as input monthly or annual data directly. *However, it is necessary to ensure the date field is in the correct format.* Assuming data in an Excel file in the format dd/mm/YYYY (e.g. 23/11/2008), it is necessary to convert this to a date format understood by R, as shown below. Similarly, if annual data were available, get the dates in formats like '2005-01-01', '2006-01-01'... and make sure the date is again formatted using `as.Date`. Note that if dates are pre-formatted as YYYY-mm-dd, then it is sufficient to use `as.Date` without providing any format information because it is already in the correct format.

```
mydata$date = as.Date(mydata$date, format = "%d/%m/%Y")
```

Finally, the `TheilSen` function can consider trends at different sites, provided the input data are correctly formatted. For input, a data frame with three columns is required: date, pollutant and *site*. The call would then be, for example:

```
TheilSen(mydata, pollutant = "o3", type = "wd", deseason = TRUE,
ylab = "ozone (ppb)")
```

```
## [1] "Taking bootstrap samples. Please wait."
```

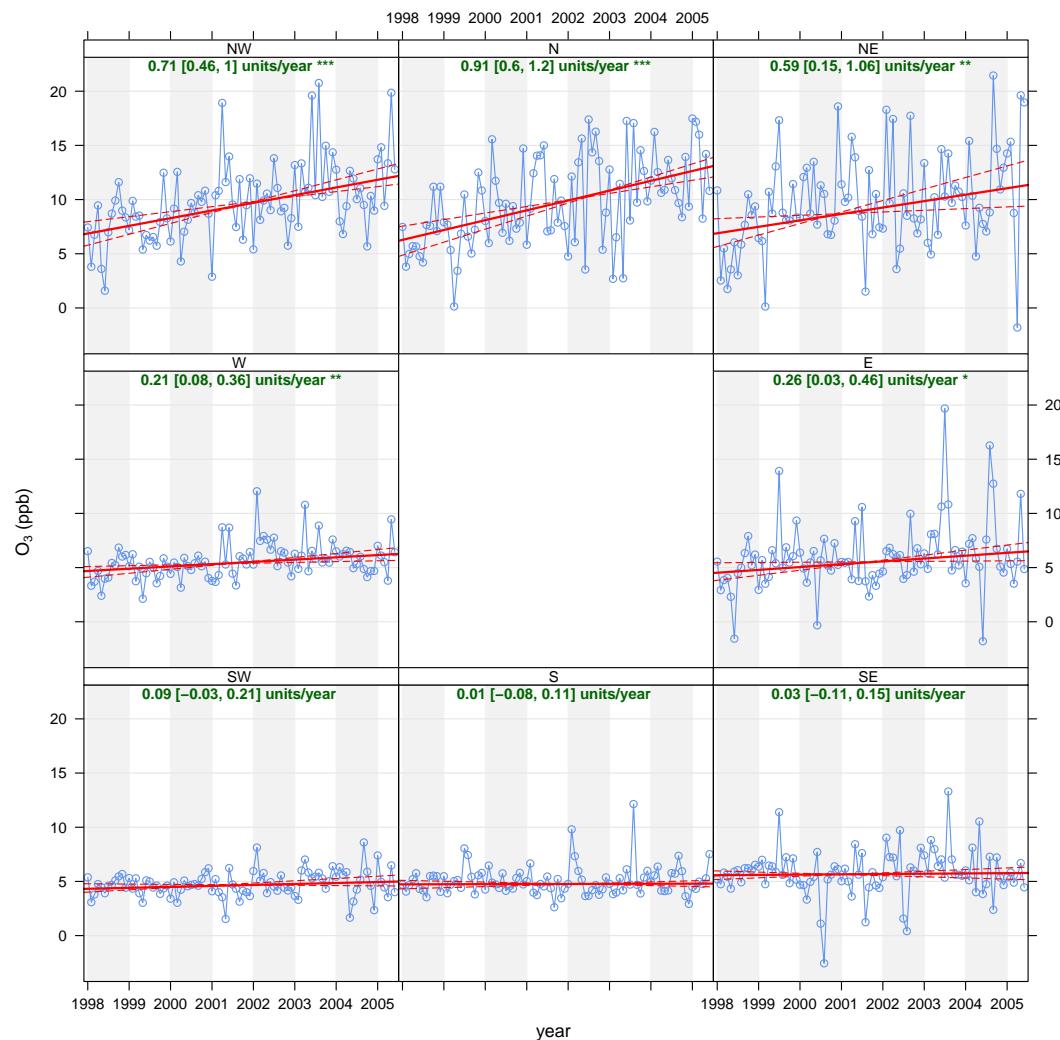


FIGURE 20.2 Trends in ozone at Marylebone Road split by eight wind sectors. The **TheilSen** function will automatically organise the separate panels by the different compass directions.

```
TheilSen(mydata, pollutant = "no2", type = "site")
```

20.4 Output

The **TheilSen** function provides lots of output data for further analysis or adding to a report. To obtain it, it is necessary to read it into a variable:

```
MKresults <- TheilSen(mydata, pollutant = "o3", deseason = TRUE, type = "wd")

## [1] "Taking bootstrap samples. Please wait."
```

This returns a list of two data frames containing all the monthly mean values and trend statistics and an aggregated summary. The first 6 lines are shown next:

```
head(MKresults$data[[1]])

##   wd      date    conc       a       b upper.a     upper.b lower.a
## 1 E 1998-01-01 5.552314 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
## 2 E 1998-02-01 2.919583 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
## 3 E 1998-03-01 3.849336 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
## 4 E 1998-04-01 4.051648 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
## 5 E 1998-05-01 2.304723 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
## 6 E 1998-06-01 -1.560546 -2.824218 0.0007184427 -
## 18363 0.001271986 4.738197
##   lower.b      p.p.stars      slope intercept intercept.lower
## 1 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
## 2 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
## 3 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
## 4 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
## 5 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
## 6 6.918724e-05 0.02337229 * 0.2622316 -2.824218 4.738197
##   intercept.upper      lower      upper slope.percent lower.percent upper.percent
## 1      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
## 2      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
## 3      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
## 4      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
## 5      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
## 6      -9.18363 0.02525334 0.464275 5.797357 0.4637236 12.13799
```

Often only the trend statistics are required and not all the monthly values. These can be obtained by:

```
MKresults$data[[2]]

##   wd p.stars      date     conc       a       b    upper.a
## 1 E * 2001-09-15 5.989971 -2.824218 7.184427e-04 -9.1836297
## 2 N *** 2001-09-15 9.786267 -19.142101 2.485907e-03 -28.8217389
## 3 NE ** 2001-09-15 9.728994 -9.725419 1.624683e-03 -24.0956107
## 4 NW *** 2001-09-15 9.785850 -12.965114 1.939117e-03 -22.3399248
## 5 S 2001-09-15 5.052729 4.472587 2.509363e-05 1.3410740
## 6 SE 2001-09-15 5.780651 4.822725 7.357450e-05 0.9208026
## 7 SW 2001-09-15 4.761100 1.818339 2.444784e-04 -1.8442920
## 8 W ** 2001-09-15 5.618726 -1.178749 5.736018e-04 -6.1035872
##   upper.b lower.a    lower.b      p    slope intercept
## 1 0.0012719864 4.738197 6.918724e-05 0.023372287 0.262231587 -2.824218
## 2 0.0032938473 -9.376480 1.654787e-03 0.0000000000 0.907355931 -19.142101
## 3 0.0029095076 3.896026 4.243103e-04 0.006677796 0.593009472 -9.725419
## 4 0.0027505150 -4.983341 1.265576e-03 0.0000000000 0.707777816 -12.965114
## 5 0.0002908397 7.340138 -2.193123e-04 0.864774624 0.009159175 4.472587
## 6 0.0004164597 9.005264 -2.969094e-04 0.767946578 0.026854691 4.822725
## 7 0.0005728771 5.556577 -7.560627e-05 0.150250417 0.089234619 1.818339
## 8 0.0009987953 2.988595 2.063957e-04 0.003338898 0.209364674 -1.178749
##   intercept.lower intercept.upper      lower      upper slope.percent
## 1        4.738197      -9.1836297 0.02525334 0.4642750 5.7973571
## 2       -9.376480     -28.8217389 0.60399712 1.2022542 14.4454295
## 3        3.896026     -24.0956107 0.15487326 1.0619703 8.6065403
## 4       -4.983341     -22.3399248 0.46193522 1.0039380 10.3080868
## 5        7.340138      1.3410740 -0.08004900 0.1061565 0.1936720
## 6        9.005264      0.9208026 -0.10837193 0.1520078 0.4816837
## 7       5.556577      -1.8442920 -0.02759629 0.2091002 2.0662764
## 8        2.988595     -6.1035872 0.07533442 0.3645603 4.4664681
##   lower.percent upper.percent
## 1        0.4637236 12.137988
## 2        8.0031175 24.715177
## 3        1.8805688 18.762980
## 4        5.8034222 17.340392
## 5       -1.5704407 2.459893
## 6       -1.8156490 2.934550
## 7       -0.5769237 5.208593
## 8        1.4773182 8.867723
```

In the results above the ‘lower’ and ‘upper’ fields provide the 95% (or chosen confidence interval using the **alpha** option) of the trend and ‘slope’ is the trend estimate expressed in units/year.

21 The **smoothTrend** function

21.1 Purpose

see also
TheilSen
timePlot

The **smoothTrend** function calculates smooth trends in the monthly mean concentrations of pollutants. In its basic use it will generate a plot of monthly concentrations and fit a smooth line to the data and show the 95 % confidence intervals of the fit. The smooth line is essentially determined using Generalized Additive Modelling using the **mgcv** package. This package provides a comprehensive and powerful set of methods for modelling data. In this case, however, the *model* is a relationship between time and pollutant concentration i.e. a trend. One of the principal advantages of this approach is that the amount of smoothness in the trend is optimised in the sense that it is neither too smooth (therefore missing important features) nor too variable (perhaps fitting ‘noise’ rather than real effects). Some background information on the use of this approach in an air quality setting can be found in Carslaw et al. (2007).

Appendix C considers smooth trends in more detail and considers how different models can be developed that can be quite sophisticated. Readers should consider this section if they are considering trend analysis in more depth.

The user can select to deseasonalise the data first to provide a clearer indication of the overall trend on a monthly basis. The data are deseasonalised using the `The stl` function. The user may also select to use bootstrap simulations to provide an alternative method of estimating the uncertainties in the trend. In addition, the simulated estimates of uncertainty can account for autocorrelation in the residuals using a block bootstrap approach.

21.2 Options available

The `smoothTrend` function has the following options:

- mydata** A data frame containing the field `date` and at least one other parameter for which a trend test is required; typically (but not necessarily) a pollutant.
- pollutant** The parameter for which a trend test is required. Mandatory.
- deseason** Should the data be de-seasonalized first? If `TRUE` the function `stl` is used (seasonal trend decomposition using loess). Note that if `TRUE` missing data are first imputed using a Kalman filter and Kalman smooth.
- type** `type` determines how the data are split i.e. conditioned, and then plotted. The default will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.
It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.
Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.
- statistic** Statistic used for calculating monthly values. Default is “mean”, but can also be “percentile”. See `timeAverage` for more details.
- avg.time** Can be “month” (the default), “season” or “year”. Determines the time over which data should be averaged. Note that for “year”, six or more years are required. For “season” the data are split up into spring: March, April, May etc. Note that December is considered as belonging to winter of the following year.
- percentile** Percentile value(s) to use if `statistic = "percentile"` is chosen. Can be a vector of numbers e.g. `percentile = c(5, 50, 95)` will plot the 5th, 50th and 95th percentile values together on the same plot.

data.thresh	The data capture threshold to use (the data using <code>avg.time</code> . A value of zero means that all available data will be used in a particular period regardless if of the number of values available. Conversely, a value of 100 will mean that all data will need to be present for the average to be calculated, else it is recorded as <code>NA</code> . Not used if <code>avg.time = "default"</code> .
simulate	Should simulations be carried out to determine the Mann-Kendall tau and p-value. The default is <code>FALSE</code> . If <code>TRUE</code> , bootstrap simulations are undertaken, which also account for autocorrelation.
n	Number of bootstrap simulations if <code>simulate = TRUE</code> .
autocor	Should autocorrelation be considered in the trend uncertainty estimates? The default is <code>FALSE</code> . Generally, accounting for autocorrelation increases the uncertainty of the trend estimate sometimes by a large amount.
cols	Colours to use. Can be a vector of colours e.g. <code>cols = c("black", "green")</code> or pre-defined openair colours — see <code>openColours</code> for more details.
shade	The colour used for marking alternate years. Use “white” or “transparent” to remove shading.
xlab	x-axis label, by default “year”.
y.relation	This determines how the y-axis scale is plotted. “same” ensures all panels use the same scale and “free” will use panel-specific scales. The latter is a useful setting when plotting data with very different values. See <code>ref.y</code> for details. In this case the correct date format should be used for a vertical line e.g. <code>ref.x = list(v = as.POSIXct("2000-06-15"), lty = 5)</code> .
ref.x	See <code>ref.y</code> .
ref.y	A list with details of the horizontal lines to be added representing reference line(s). For example, <code>ref.y = list(h = 50, lty = 5)</code> will add a dashed horizontal line at 50. Several lines can be plotted e.g. <code>ref.y = list(h = c(50, 100), lty = c(1, 5), col = c("green", "blue"))</code> . See <code>panel.abline</code> in the <code>lattice</code> package for more details on adding/controlling lines.
key.columns	Number of columns used if a key is drawn when using the option <code>statistic = "percentile"</code> .
name.pol	Names to be given to the pollutant(s). This is useful if you want to give a fuller description of the variables, maybe also including subscripts etc.
ci	Should confidence intervals be plotted? The default is <code>FALSE</code> .
alpha	The alpha transparency of shaded confidence intervals - if plotted. A value of 0 is fully transparent and 1 is fully opaque.
date.breaks	Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of <code>date.breaks</code> up or down.

`auto.text` Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in `NO2`.

`k` This is the smoothing parameter used by the `gam` function in package `mgcv`. By default it is not used and the amount of smoothing is optimised automatically. However, sometimes it is useful to set the smoothing amount manually using `k`.

... Other graphical parameters are passed onto `cutData` and `lattice:xyplot`. For example, `smoothTrend` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common graphical arguments, such as `xlim` and `ylim` for plotting ranges and `pch` and `cex` for plot symbol type and size, are passed on `xyplot`, although some local modifications may be applied by `openair`. For example, axis and title labelling options (such as `xlab`, `ylab` and `main`) are passed to `xyplot` via `quickText` to handle routine formatting. One special case here is that many graphical parameters can be vectors when used with `statistic = "percentile"` and a vector of `percentile` values, see examples below.

21.3 Example of use

We apply the function to concentrations of O_3 and NO_2 using the code below. The first plot shows the smooth trend in raw O_3 concentrations, which shows a very clear seasonal cycle. By removing the seasonal cycle of O_3 , a better indication of the trend is given, shown in the second plot. Removing the seasonal cycle is more effective for pollutants (or locations) where the seasonal cycle is stronger e.g. for ozone and background sites. Figure 21.1 shows the results of the simulations for NO_2 without the seasonal cycle removed. It is clear from this plot that there is little evidence of a seasonal cycle. The principal advantage of the smoothing approach compared with the Theil-Sen method is also clearly shown in this plot. Concentrations of NO_2 first decrease, then increase strongly. The trend is therefore not monotonic, violating the Theil-Sen assumptions. Finally, the last plot shows the effects of first deseasonalising the data: in this case with little effect.

The `smoothTrend` function share many of the functionalities of the `TheilSen` function. Figure 21.2 shows the result of applying this function to O_3 concentrations. The code that produced Figure 21.2 was:

The `smoothTrend` function can easily be used to gain a large amount of information on trends easily. For example, how do trends in NO_2 , O_3 and PM_{10} vary by season and wind sector. There are 8 wind sectors and four seasons i.e. 32 plots. In Figure 21.3 all three pollutants are chosen and two types (season and wind direction). We also reduce the number of axis labels and the line to improve clarity. There are numerous combinations of analyses that could be produced here and it is very easy to explore the data in a wide number of ways.

```

smoothTrend(mydata, pollutant = "o3", ylab = "concentration (ppb)",
               main = "monthly mean o3")

smoothTrend(mydata, pollutant = "o3", deseason = TRUE, ylab = "concentration (ppb)",
               main = "monthly mean deseasonalised o3")

smoothTrend(mydata, pollutant = "no2", simulate = TRUE, ylab = "concentration (ppb)",
               main = "monthly mean no2 (bootstrap uncertainties)")

## [1] "Taking bootstrap samples. Please wait..."

smoothTrend(mydata, pollutant = "no2", deseason = TRUE, simulate =TRUE,
               ylab = "concentration (ppb)",
               main = "monthly mean deseasonalised no2 (bootstrap uncertainties)")

## [1] "Taking bootstrap samples. Please wait..."

```

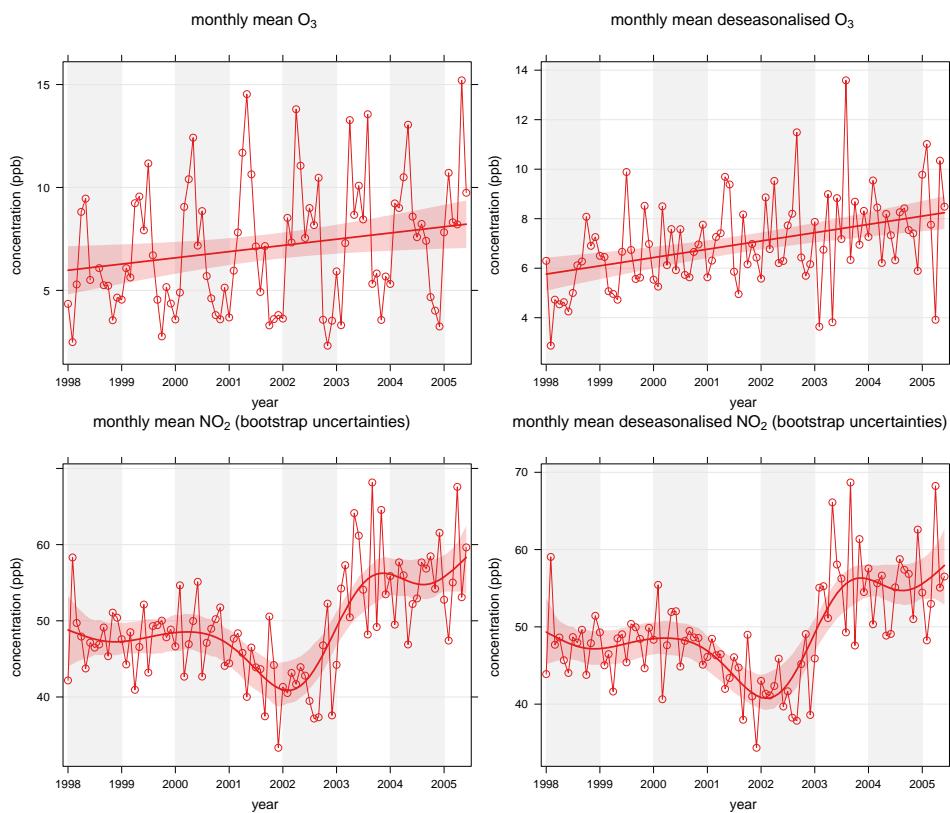


FIGURE 21.1 Examples of the `smoothTrend` function applied to Marylebone Road

```
smoothTrend(mydata, pollutant = "o3", deseason = TRUE,
              type = "wd")
```

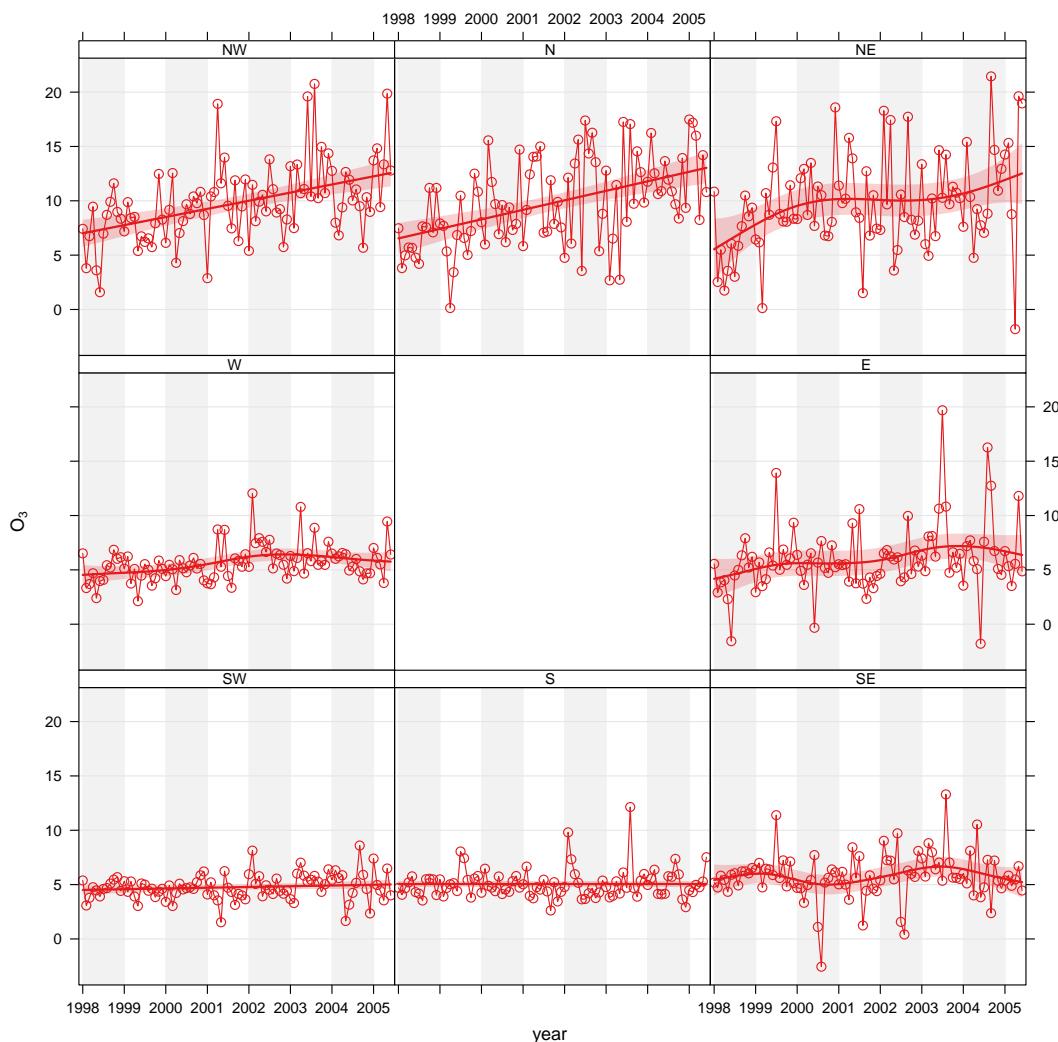


FIGURE 21.2 Trends in O_3 using the **smoothTrend** function applied to Marylebone Road. The shading shows the estimated 95 % confidence intervals. This plot can usefully be compared with Figure 20.2.

```
smoothTrend(mydata, pollutant = c("no2", "pm10", "o3"), type = c("wd", "season"),
date.breaks = 3, lty = 0)
```

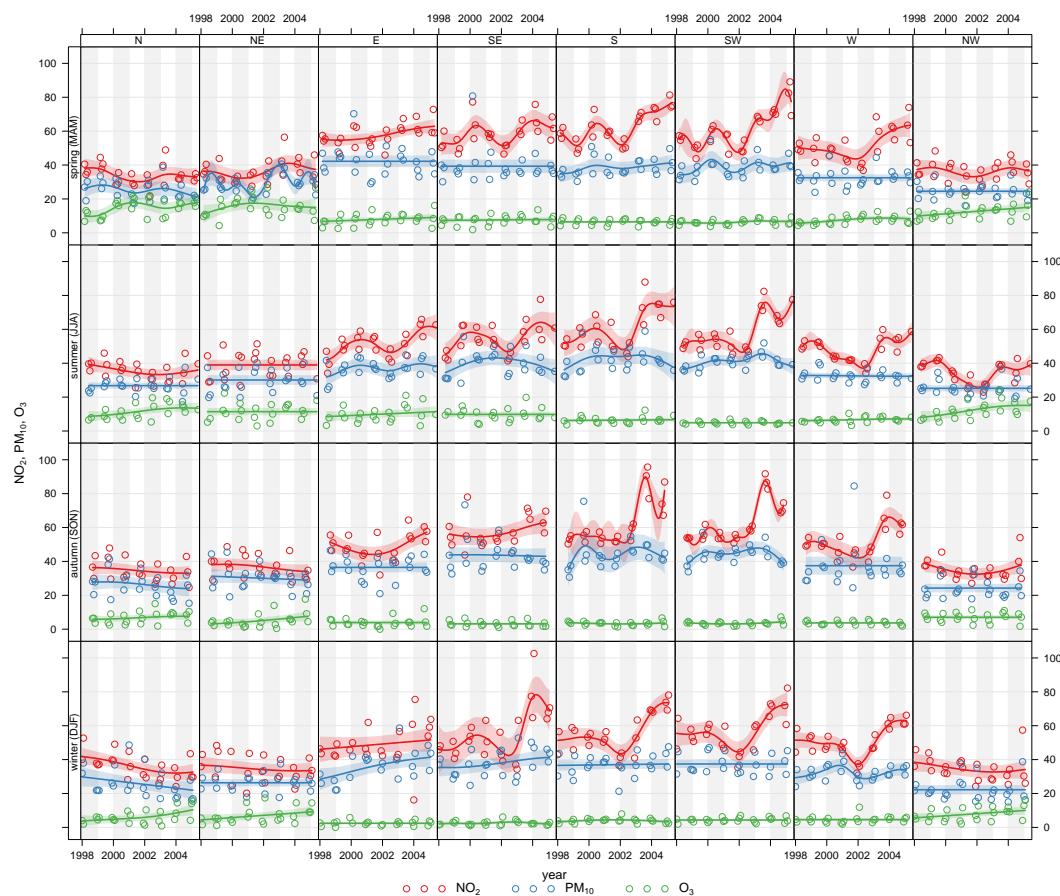


FIGURE 21.3 The **smoothTrend** function applied to three pollutants, split by wind sector and season.

22 The **timeVariation** function

22.1 Purpose

see also
linearRelation

In air pollution, the variation of a pollutant by time of day and day of week can reveal useful information concerning the likely sources. For example, road vehicle emissions tend to follow very regular patterns both on a daily and weekly basis. By contrast some industrial emissions or pollutants from natural sources (e.g. sea salt aerosol) may well have very different patterns.

The **timeVariation** function produces four plots: day of the week variation, mean hour of day variation and a combined hour of day – day of week plot and a monthly plot. Also shown on the plots is the 95 % confidence interval in the mean. These uncertainty limits can be helpful when trying to determine whether one candidate source is different from another. The uncertainty intervals are calculated through bootstrap re-sampling, which will provide better estimates than the application of assumptions based on normality, particularly when there are few data available. The function can consider one or two input variables. In addition, there is the option of ‘normalising’ concentrations (or other quantities). Normalising is very useful for comparing the patterns of two different pollutants, which often cover very different ranges in concentration. Normalising is achieved by dividing the concentration of a pollutant by its mean value. Note also that any other variables besides pollutant concentrations can be considered e.g. meteorological or traffic data.

There is also an option **difference** which is very useful for considering the difference in two time series and how they vary over different temporal resolutions. Again, bootstrap re-sampling methods are used to estimate the uncertainty of the difference in two means.

Averaging wind direction

Care has been taken to ensure that wind direction (**wd**) is vector-averaged. Less obvious though is the uncertainty in wind direction. A pragmatic approach has been adopted here that considers how wind direction *changes*. For example, consider the following wind directions: 10, 10, 10, 180, 180, 180°. The standard deviation of these numbers is 93°. However, what actually occurs is the wind direction is constant at 10° then switches to 180°. In terms of changes there is a sequence of numbers: 0, 0, 170, 0, 0 with a standard deviation of 76°. We use the latter method as a basis of calculating the 95% confidence intervals in the mean.

There are also problems with simple averaging—for example, what is the average of 20 and 200°. It can't be known. In some situations where the wind direction is bimodal with differences around 180°, the mean can be ‘unstable’. For example, wind that is funnelled along a valley forcing it to be either easterly or westerly. Consider for example the mean of 0° and 179° (89.5°), but a small change in wind direction to 181° gives a mean of 270.5°. Some care should be exercised therefore when averaging wind direction. It is always a good idea to use the **windRose** function with type set to ‘month’ or ‘hour’.

The **timeVariation** function is probably one of the most useful functions that can be used for the analysis of air pollution. Here are a few uses/advantages:

- Variations in time are one of the most useful ways of characterising air pollution for

a very wide range of pollutants including local urban pollutants and tropospheric background concentrations of ozone and the like.

- The function works well in conjunction with other functions such as `polarPlot`, where the latter may identify conditions of interest (say a wind speed/direction range). By sub-setting for those conditions in `timeVariation` the temporal characteristics of a particular source could be characterised and perhaps contrasted with another subset of conditions.
- The function can be used to compare a wide range of variables, if available. Suggestions include meteorological e.g. boundary layer height and traffic flows.
- The function can be used for comparing pollutants over different sites. See §(31.7) for examples of how to do this.
- The function can be used to compare one part of a time series with another. This is often a very powerful thing to do, particularly if concentrations are normalised. For example, there is often interest in knowing how diurnal/weekday/seasonal patterns vary with time. If a pollutant showed signs of an increase in recent years, then splitting the data set and comparing each part together can provide information on what is driving the change. Is there, for example, evidence that morning rush hour concentrations have become more important, or Sundays have become relatively more important? An example is given below using the `splitByDate` function.
- `timeVariation` can be used to consider the differences between two time series, which will have multiple benefits. For example, for model evaluation it can be very revealing to consider the difference between observations and modelled values over different time scales. Considering such differences can help reveal the character and some of the reasons for why a model departs from reality.

22.2 Options available

The `timeVariation` function has the following options:

<code>mydata</code>	A data frame of hourly (or higher temporal resolution data). Must include a <code>date</code> field and at least one variable to plot.
<code>pollutant</code>	Name of variable to plot. Two or more pollutants can be plotted, in which case a form like <code>pollutant = c("nox", "co")</code> should be used.
<code>local.tz</code>	Should the results be calculated in local time that includes a treatment of daylight savings time (DST)? The default is not to consider DST issues, provided the data were imported without a DST offset. Emissions activity tends to occur at local time e.g. rush hour is at 8 am every day. When the clocks go forward in spring, the emissions are effectively released into the atmosphere typically 1 hour earlier during the summertime i.e. when DST applies. When plotting diurnal profiles, this has the effect of “smearing-out” the concentrations. Sometimes, a useful approach is to express time as local time. This correction tends to produce better-defined diurnal profiles of concentration (or other variables) and allows a better comparison to be made with emissions/activity data. If set to <code>FALSE</code> then GMT is used. Examples of usage include <code>local.tz = "Europe/London"</code> , <code>local.tz = "America/New_York"</code> . See <code>cutData</code> and <code>import</code> for more details.

normalise Should variables be normalised? The default is `FALSE`. If `TRUE` then the variable(s) are divided by their mean values. This helps to compare the shape of the diurnal trends for variables on very different scales.

xlab x-axis label; one for each sub-plot.

name.pol Names to be given to the pollutant(s). This is useful if you want to give a fuller description of the variables, maybe also including subscripts etc.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Only one `type` is allowed in `timeVariation`.

group This sets the grouping variable to be used. For example, if a data frame had a column `site` setting `group = "site"` will plot all sites together in each panel. See examples below.

difference If two pollutants are chosen then setting `difference = TRUE` will also plot the difference in means between the two variables as `pollutant[2] - pollutant[1]`. Bootstrap 95% confidence intervals of the difference in means are also calculated. A horizontal dashed line is shown at $y = 0$.

statistic Can be “mean” (default) or “median”. If the statistic is ‘mean’ then the mean line and the 95% confidence interval in the mean are plotted by default. If the statistic is ‘median’ then the median line is plotted together with the 5/95 and 25/75th quantiles are plotted. Users can control the confidence intervals with `conf.int`.

conf.int The confidence intervals to be plotted. If `statistic = "mean"` then the confidence intervals in the mean are plotted. If `statistic = "median"` then the `conf.int` and `1 - conf.int` quantiles are plotted. `conf.int` can be of length 2, which is most useful for showing quantiles. For example `conf.int = c(0.75, 0.99)` will yield a plot showing the median, 25/75 and 5/95th quantiles.

B Number of bootstrap replicates to use. Can be useful to reduce this value when there are a large number of observations available to increase the speed of the calculations without affecting the 95% confidence interval calculations by much.

ci Should confidence intervals be shown? The default is `TRUE`. Setting this to `FALSE` can be useful if multiple pollutants are chosen where over-lapping confidence intervals can over complicate plots.

<code>cols</code>	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours — see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
<code>ref.y</code>	A list with details of the horizontal lines to be added representing reference line(s). For example, <code>ref.y = list(h = 50, lty = 5)</code> will add a dashed horizontal line at 50. Several lines can be plotted e.g. <code>ref.y = list(h = c(50, 100), lty = c(1, 5), col = c("green", "blue"))</code> . See <code>panel.abline</code> in the <code>lattice</code> package for more details on adding/controlling lines.
<code>key</code>	By default <code>timeVariation</code> produces four plots on one page. While it is useful to see these plots together, it is sometimes necessary just to use one for a report. If <code>key</code> is <code>TRUE</code> , a key is added to all plots allowing the extraction of a single plot with <code>key</code> . See below for an example.
<code>key.columns</code>	Number of columns to be used in the key. With many pollutants a single column can make the key too wide. The user can thus choose to use several columns by setting <code>columns</code> to be less than the number of pollutants.
<code>start.day</code>	What day of the week should the plots start on? The user can change the start day by supplying an integer between 0 and 6. Sunday = 0, Monday = 1, ...For example to start the weekday plots on a Saturday, choose <code>start.day = 6</code> .
<code>auto.text</code>	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO ₂ .
<code>alpha</code>	The alpha transparency used for plotting confidence intervals. 0 is fully transparent and 1 is opaque. The default is 0.4
<code>...</code>	Other graphical parameters passed onto <code>lattice:xyplot</code> and <code>cutData</code> . For example, in the case of <code>cutData</code> the option <code>hemisphere = "southern"</code> .

22.3 Example of use

We apply the `timeVariation` function to PM₁₀ concentrations and take the opportunity to filter the data to maximise the signal from the road. The `polarPlot` function described in (§15) is very useful in this respect in highlighting the conditions under which different sources have their greatest impact. A subset of data is used filtering for wind speeds > 3 m s⁻¹ and wind directions from 100–270 degrees. The code used is:

The results are shown in Figure 22.1. The plot shown at the top-left shows the diurnal variation of concentrations for all days. It shows for example that PM₁₀ concentrations tend to peak around 9 am. The shading shows the 95 % confidence intervals of the mean. The plot at the top-right shows how PM₁₀ concentrations vary by day of the week. Here there is strong evidence that PM₁₀ is much lower at the weekends and that there is a significant difference compared with weekdays. It also shows that concentrations tend to increase during the weekdays. Finally, the plot at the bottom shows both sets of information together to provide an overview of how concentrations vary.

Note that the plot need not just consider pollutant concentrations. Other useful variables (if available) are meteorological and traffic flow or speed data. Often, the combination of several sets of data can be very revealing.

`timeVariation`
is also very
useful for other
variables such
as traffic and
meteorological
data

```
timeVariation(subset(mydata, ws > 3 & wd > 100 & wd < 270),
pollutant = "pm10", ylab = "pm10 (ug/m3)")
```

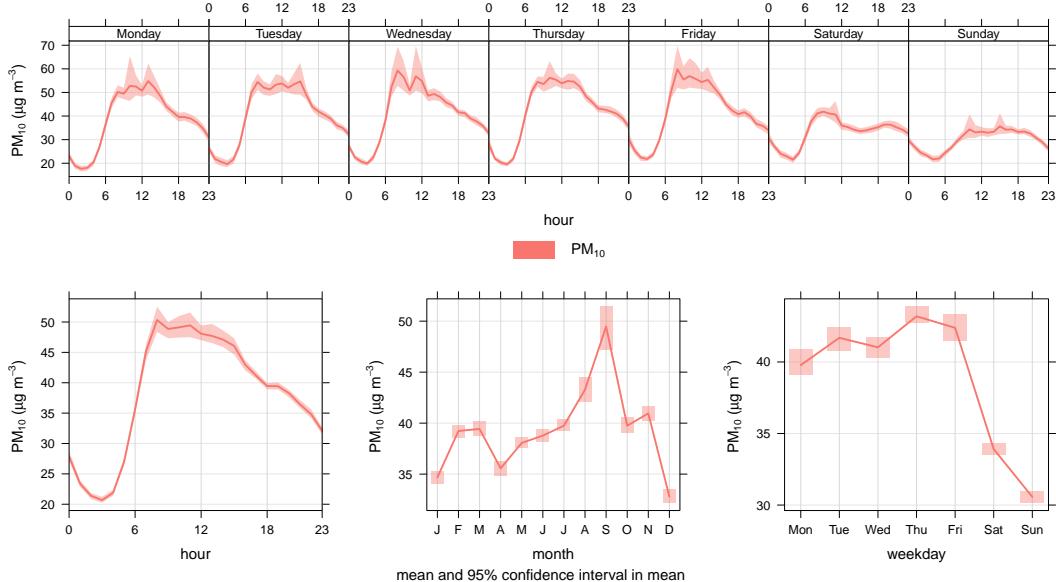


FIGURE 22.1 Example plot using the `timeVariation` function to plot PM₁₀ concentrations at Marylebone Road.

The `subset` function is extremely useful in this respect. For example, if it were believed that a source had an effect under specific conditions; they can be isolated with the `subset` function. It is also useful if it is suspected that two or more sources are important that they can be isolated to some degree and compared. This is where the uncertainty intervals help — they provide an indication whether the behaviour of one source differs significantly from another.

Figure 22.2 shows the function applied to concentrations of NO_x, CO, NO₂ and O₃ concentrations. In this case the concentrations have been normalised. The plot clearly shows the markedly different temporal trends in concentration. For CO, there is a very pronounced increase in concentrations during the peak pm rush hour. The other important difference is on Sundays when CO concentrations are relatively much higher than NO_x. This is because flows of cars (mostly petrol) do not change that much by day of the week, but flows of vans and HGVs (diesel vehicles) are much less on Sundays. Note, however, that the monthly trend is very similar in each case — which indicates very similar source origins. Taken together, the plots highlight that traffic emissions dominate this site for CO and NO_x, but there are important differences in how these emissions vary by hour of day and day of week.

Also shown in the very different behaviour of O₃. Because O₃ reacts with NO, concentrations of NO_x and O₃ tend to be anti-correlated. Note also the clear peak in O₃ in April/May, which is due to higher northern hemispheric background concentrations in the spring. Even at a busy roadside site in central London this influence is clear to see.

Another example is splitting the data set by time. We use the `splitByDate` function to divide up the data into dates before January 2003 and after January 2003. This time the option `difference` is used to highlight how NO₂ concentrations have changed over these two periods. The results are shown in Figure 22.3. There is some indication in this plot that data after 2003 seem to show more of a double peak in the diurnal plots; particularly in the morning rush hour. Also, the difference line does more clearly

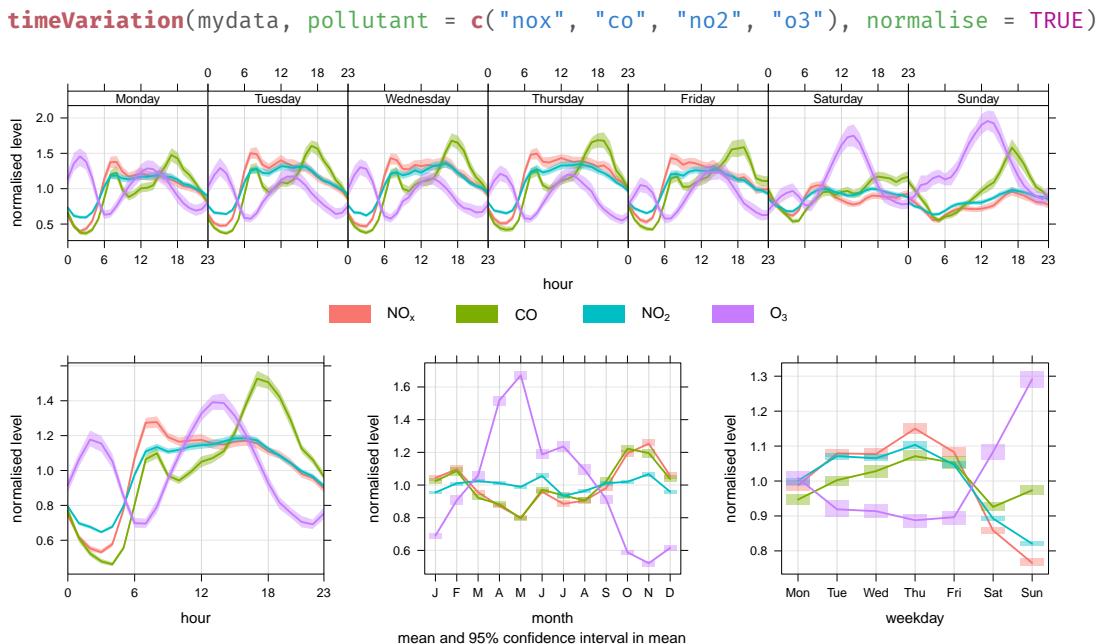


FIGURE 22.2 Example plot using the `timeVariation` function to plot NO_x , CO, NO_2 and O_3 concentrations at Marylebone Road. In this plot, the concentrations are normalised.

highlight a more substantial change over weekdays and weekends. Given that cars are approximately constant at this site each day, the change may indicate a change in vehicle emissions from other vehicle types. Given that it is known that primary NO_2 emissions are known to have increased sharply from the beginning of 2003 onwards, this perhaps provides clues as to the principal cause.

In the next example it is shown how to compare one subset of data of interest with another. Again, there can be many reasons for wanting to do this and perhaps the data set at Marylebone Road is not the most interesting to consider. Nevertheless, the code below shows how to approach such a problem. The scenario would be that one is interested in a specific set of conditions and it would be useful to compare that set, with another set. A good example would be from an analysis using the `polarPlot` function where a ‘feature’ of interest has been identified—maybe an indication of a different source. But does this potentially different source behave in a different way in terms of temporal variation? If it does, then maybe that provides evidence to support that it is a different source. In a wider context, this approach could be used in many different ways depending on available data. A good example is the analysis of model output where many diagnostic meteorological data are available. This is an area that will be developed.

The approach here is to first make a new variable called ‘feature’ and fill it with the value ‘other’. A subset of data is defined and the associated locations in the data frame identified. The subset of data is then used to update the ‘feature’ field with a new description. This approach could be extended to some quite complex situations.

There are a couple of things to note in Figure 22.2. There seems to be evidence that for easterly winds $> 4 \text{ m s}^{-1}$ that concentrations of SO_2 are lower at night. Also, there is some evidence that concentrations for these conditions are also lower at weekends. This *might* reflect that SO_2 concentrations for these conditions tend to be dominated by tall stack emissions that have different activities to road transport sources. This technique will be returned to with different data sets in future.

```
## split data into two periods (see Utilities section for more details)
mydata <- splitByDate(mydata, dates = "1/1/2003",
                      labels = c("before Jan. 2003", "After Jan. 2003"))

timeVariation(mydata, pollutant = "no2", group = "split.by", difference = TRUE)
```

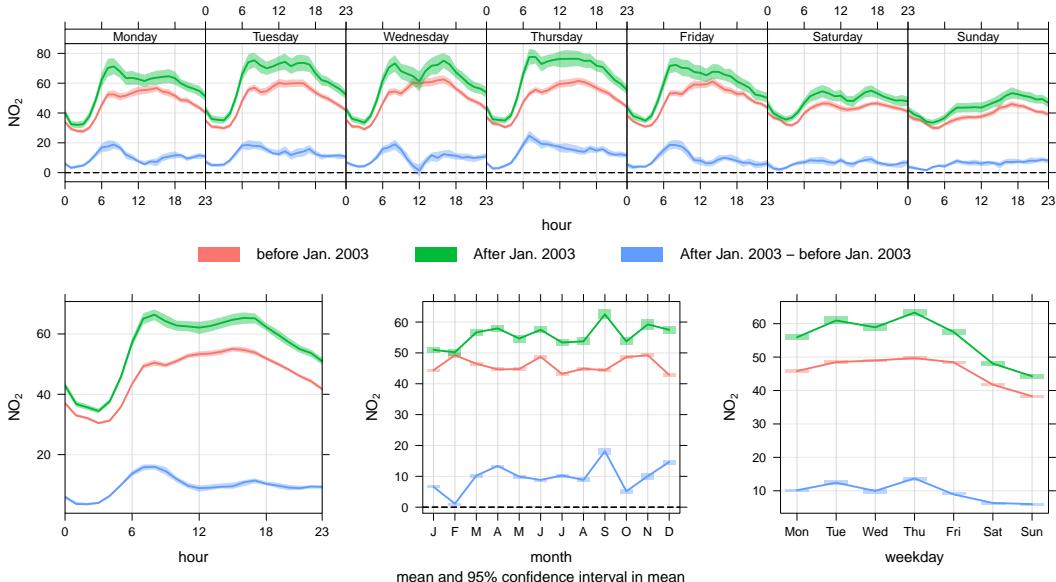


FIGURE 22.3 Example plot using the `timeVariation` function to plot NO₂ concentrations at Marylebone Road. In this plot, the concentrations are shown before and after January 2003.

By default `timeVariation` shows the mean variation in different temporal components and the 95% confidence interval in the mean. However, it is also possible to show how the data are distributed by using a different option for `statistic`. When `statistic = "median"` the median line is shown together with the 25/75th and 5/95th quantile values. Users can control the quantile values shown by setting the `conf.int`. For example, `conf.int = c(0.25, 0.99)` will show the median, 25/75th and 1/99th quantile values. The `statistic = "median"` option is therefore very useful for showing how the data are distributed — somewhat similar to a box and whisker plot. Note that it is expected that only one pollutant should be shown when `statistic = "median"` is used due to potential over-plotting; although the function will display several species if required. An example is shown in Figure 22.5 for PM₁₀ concentrations.

```
## make a field called "site" and fill: make all values = "other"
mydata$feature <- "other"

## now find which indexes correspond to easterly conditions > 4m/s ws
id <- which(with(mydata, ws > 4 & wd > 0 & wd <= 180))

## use the ids to update the site column
## there are now two values in site: "other" and "easterly"
mydata$feature[id] <- "easterly"
timeVariation(mydata, pollutant = "so2", group = "feature", ylab = "so2 (ppb)",
               difference = TRUE)
```

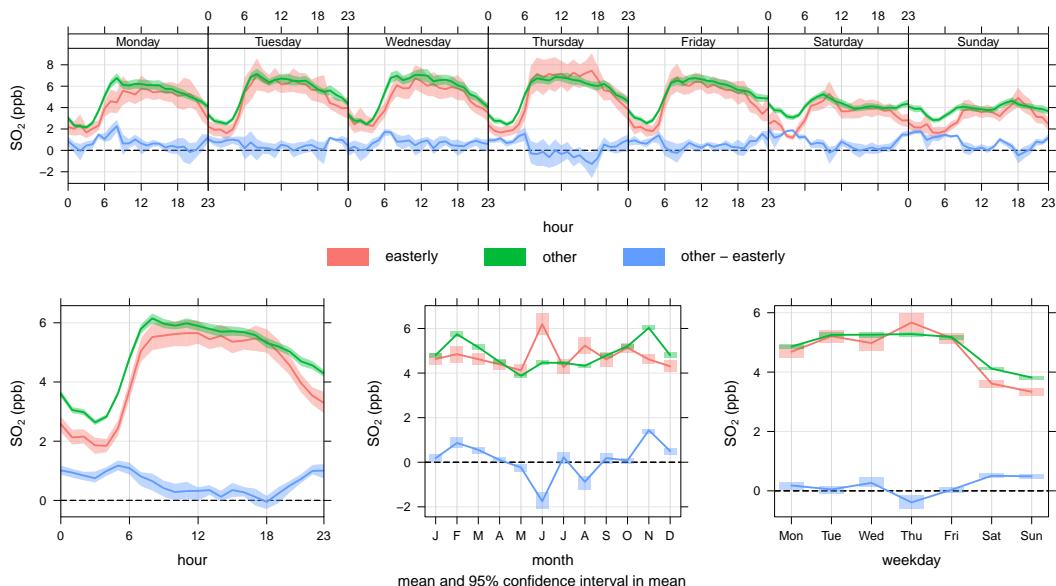


FIGURE 22.4 Example plot using the **timeVariation** function to plot SO₂ concentrations at Marylebone Road. In this plot, the concentrations are shown for a subset of easterly conditions and everything else. Note that the uncertainty in the mean values for easterly winds is greater than 'other'. This is mostly because the sample size is much lower for 'easterly' compared with 'other'.

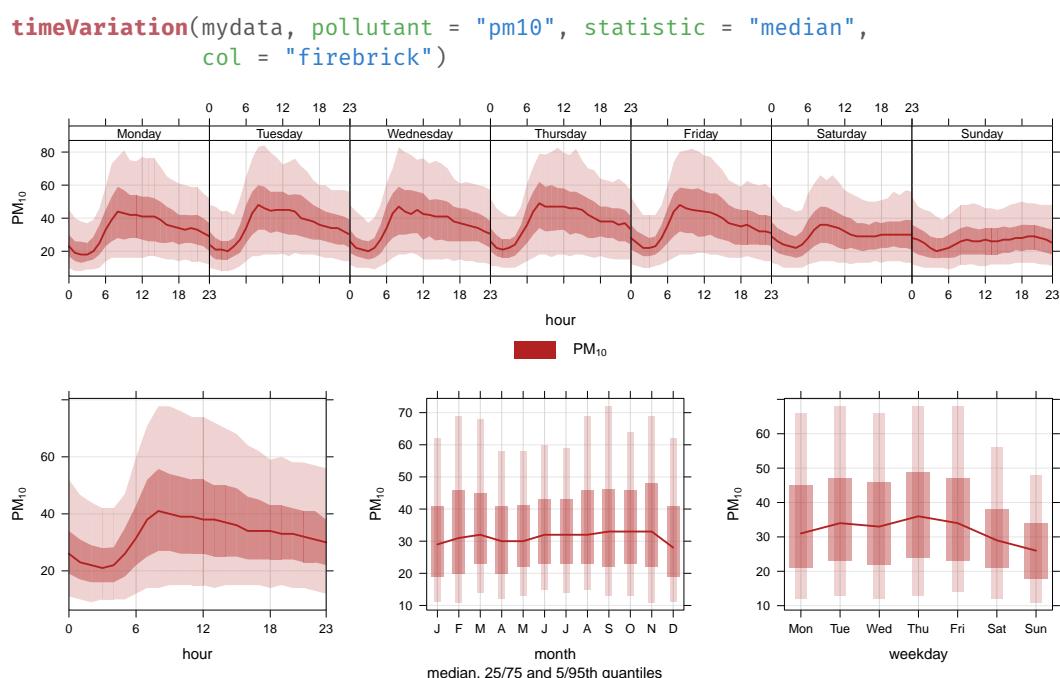


FIGURE 22.5 Example plot using the **timeVariation** function to show the variation in the median, 25/75th and 5/95th quantile values for PM_{10} . The shading shows the extent to the 25/75th and 5/95th quantiles.

22.4 Output

The **timeVariation** function produces several outputs that can be used for further analysis or plotting. It is necessary to read the output into a variable for further processing. The code below shows the different objects that are returned and the code shows how to access them.

```
myOutput <- timeVariation(mydata, pollutant = "so2")
## show the first part of the day/hour variation
## note that value = mean, and Upper/Lower the 95% confid. intervals
head(myOutput$data$day.hour)

## # A tibble: 6 x 8
## # Groups: ci [1]
##   variable wkday     hour default          Mean Lower Upper   ci
##   <fct>    <ord>    <int> <fct>        <dbl> <dbl> <dbl> <dbl>
## 1 so2      Monday     0 01 January 1998 to 23 June 2~  2.93  2.62  3.23  0.95
## 2 so2      Tuesday    0 01 January 1998 to 23 June 2~  3.21  2.92  3.52  0.95
## 3 so2      Wednesday  0 01 January 1998 to 23 June 2~  3.35  3.13  3.57  0.95
## 4 so2      Thursday   0 01 January 1998 to 23 June 2~  3.22  2.99  3.48  0.95
## 5 so2      Friday     0 01 January 1998 to 23 June 2~  3.64  3.32  3.87  0.95
## 6 so2      Saturday   0 01 January 1998 to 23 June 2~  4.25  3.97  4.57  0.95

## can make a new data frame of this data e.g.
day.hour <- myOutput$data$day.hour
head(day.hour)

## # A tibble: 6 x 8
## # Groups: ci [1]
##   variable wkday     hour default          Mean Lower Upper   ci
##   <fct>    <ord>    <int> <fct>        <dbl> <dbl> <dbl> <dbl>
## 1 so2      Monday     0 01 January 1998 to 23 June 2~  2.93  2.62  3.23  0.95
## 2 so2      Tuesday    0 01 January 1998 to 23 June 2~  3.21  2.92  3.52  0.95
## 3 so2      Wednesday  0 01 January 1998 to 23 June 2~  3.35  3.13  3.57  0.95
## 4 so2      Thursday   0 01 January 1998 to 23 June 2~  3.22  2.99  3.48  0.95
## 5 so2      Friday     0 01 January 1998 to 23 June 2~  3.64  3.32  3.87  0.95
## 6 so2      Saturday   0 01 January 1998 to 23 June 2~  4.25  3.97  4.57  0.95
```

All the numerical results are given by:

```
myOutput$data$day.hour ## are the weekday and hour results
myOutput$data$hour ## are the diurnal results
myOutput$data$day ## are the weekday results
myOutput$data$month ## are the monthly results
```

It is also possible to plot the individual plots that make up the (four) plots produced by **timeVariation**:

```
## just the diurnal variation
plot(myOutput, subset = "hour")
## day and hour
plot(myOutput, subset = "day.hour")
## weekday variation
plot(myOutput, subset = "day")
## monthly variation
plot(myOutput, subset = "month")
```

23 The `scatterPlot` function

23.1 Purpose

Scatter plots are extremely useful and a very commonly used analysis technique for considering how variables relate to one another. R does of course have many capabilities for plotting data in this way. However, it can be tricky to add linear relationships, or split scatter plots by levels of other variables etc. The purpose of the `scatterPlot` function is to make it straightforward to consider how variables are related to one another in a way consistent with other `openair` functions. We have added several capabilities that can be used just by setting different options, some of which are shown below.

- A smooth fit is automatically added to help reveal the underlying relationship between two variables together with the estimated 95% confidence intervals of the fit. This is in general an extremely useful thing to do because it helps to show the (possibly) non-linear relationship between variables in a very robust way — or indeed whether the relationship is linear.
- It is easy to add a linear regression line. The resulting equation is shown on the plot together with the R^2 value.
- For large data sets there is the possibility to ‘bin’ the data using hexagonal binning or kernel density estimates. This approach is very useful when there is considerable over-plotting.
- It is easy to show how two variables are related to one another dependent on levels of a third variable. This capability is very useful for exploring how different variables depend on one another and can help reveal the underlying important relationships.
- A plot of two variables can be colour-coded by a continuous colour scale of a third variable.
- It can handle date/time x-axis formats to provide an alternative way of showing time series, which again can be colour-coded by a third variable.

The `scatterPlot` function isn’t really specific to atmospheric sciences, in the same way as other plots. It is more a function for convenience, written in a style that is consistent with other `openair` functions. Nevertheless, along with the `timePlot` function they do form an important part of `openair` because of the usefulness of understanding how variables relate to one another. Furthermore, there are many options to make it easy to explore data in an interactive way without worrying about processing data or formatting plots.

23.2 Options available

`mydata` A data frame containing at least two numeric variables to plot.

`x` Name of the x-variable to plot. Note that `x` can be a date field or a factor. For example, `x` can be one of the `openair` built in types such as "`year`" or "`season`".

`y` Name of the numeric y-variable to plot.

<code>z</code>	Name of the numeric z-variable to plot for <code>method = "scatter"</code> or <code>method = "level"</code> . Note that for <code>method = "scatter"</code> points will be coloured according to a continuous colour scale, whereas for <code>method = "level"</code> the surface is coloured.
<code>method</code>	Methods include “scatter” (conventional scatter plot), “hexbin” (hexagonal binning using the <code>hexbin</code> package). “level” for a binned or smooth surface plot and “density” (2D kernel density estimates).
<code>group</code>	The grouping variable to use, if any. Setting this to a variable in the data frame has the effect of plotting several series in the same panel using different symbols/colours etc. If set to a variable that is a character or factor, those categories or factor levels will be used directly. If set to a numeric variable, it will split that variable in to quantiles.
<code>avg.time</code>	This defines the time period to average to. Can be “sec”, “min”, “hour”, “day”, “DSTday”, “week”, “month”, “quarter” or “year”. For much increased flexibility a number can precede these options followed by a space. For example, a timeAverage of 2 months would be <code>period = "2 month"</code> . See function <code>timeAverage</code> for further details on this. This option is useful as one method by which the number of points plotted is reduced i.e. by choosing a longer averaging time.
<code>data.thresh</code>	The data capture threshold to use (the data using <code>avg.time</code> . A value of zero means that all available data will be used in a particular period regardless of the number of values available. Conversely, a value of 100 will mean that all data will need to be present for the average to be calculated, else it is recorded as <code>NA</code> . Not used if <code>avg.time = "default"</code> .
<code>statistic</code>	The statistic to apply when aggregating the data; default is the mean. Can be one of “mean”, “max”, “min”, “median”, “frequency”, “sd”, “percentile”. Note that “sd” is the standard deviation and “frequency” is the number (frequency) of valid records in the period. “percentile” is the percentile level (using the “percentile” option - see below. Not used if <code>avg.time = "default"</code> .
<code>percentile</code>	The percentile level in % used when <code>statistic = "percentile"</code> and when aggregating the data with <code>avg.time</code> . The default is 95. Not used if <code>avg.time = "default"</code> .
<code>type</code>	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. “season”, “year”, “weekday” and so on. For example, <code>type = "season"</code> will produce four plots — one for each season. It is also possible to choose <code>type</code> as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another. Type can be up length two e.g. <code>type = c("season", "weekday")</code> will produce a 2x2 plot split by season and day of the week. Note, when

two types are provided the first forms the columns and the second the rows.

<code>smooth</code>	A smooth line is fitted to the data if <code>TRUE</code> ; optionally with 95% confidence intervals shown. For <code>method = "level"</code> a smooth surface will be fitted to binned data.
<code>spline</code>	A smooth spline is fitted to the data if <code>TRUE</code> . This is particularly useful when there are fewer data points or when a connection line between a sequence of points is required.
<code>linear</code>	A linear model is fitted to the data if <code>TRUE</code> ; optionally with 95% confidence intervals shown. The equation of the line and R2 value is also shown.
<code>ci</code>	Should the confidence intervals for the smooth/linear fit be shown?
<code>mod.line</code>	If <code>TRUE</code> three lines are added to the scatter plot to help inform model evaluation. The 1:1 line is solid and the 1:0.5 and 1:2 lines are dashed. Together these lines help show how close a group of points are to a 1:1 relationship and also show the points that are within a factor of two (FAC2). <code>mod.line</code> is appropriately transformed when x or y axes are on a log scale.
<code>cols</code>	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours—see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
<code>plot.type</code>	<code>lattice</code> plot type. Can be “p” (points — default), “l” (lines) or “b” (lines and points).
<code>key</code>	Should a key be drawn? The default is <code>TRUE</code> .
<code>key.title</code>	The title of the key (if used).
<code>key.columns</code>	Number of columns to be used in the key. With many pollutants a single column can make to key too wide. The user can thus choose to use several columns by setting <code>columns</code> to be less than the number of pollutants.
<code>key.position</code>	Location where the scale key is to plotted. Allowed arguments currently include “top”, “right”, “bottom” and “left”.
<code>strip</code>	Should a strip be drawn? The default is <code>TRUE</code> .
<code>log.x</code>	Should the x-axis appear on a log scale? The default is <code>FALSE</code> . If <code>TRUE</code> a well-formatted log10 scale is used. This can be useful for checking linearity once logged.
<code>log.y</code>	Should the y-axis appear on a log scale? The default is <code>FALSE</code> . If <code>TRUE</code> a well-formatted log10 scale is used. This can be useful for checking linearity once logged.
<code>x.inc</code>	The x-interval to be used for binning data when <code>method = "level"</code> .
<code>y.inc</code>	The y-interval to be used for binning data when <code>method = "level"</code> .

<code>limits</code>	For <code>method = "level"</code> the function does its best to choose sensible limits automatically. However, there are circumstances when the user will wish to set different ones. The limits are set in the form <code>c(lower, upper)</code> , so <code>limits = c(0, 100)</code> would force the plot limits to span 0-100.
<code>windflow</code>	This option allows a scatter plot to show the wind speed/direction shows as an arrow. The option is a list e.g. <code>windflow = list(col = "grey", lwd = 2, scale = 0.1)</code> . This option requires wind speed (<code>ws</code>) and wind direction (<code>wd</code>) to be available. The maximum length of the arrow plotted is a fraction of the plot dimension with the longest arrow being <code>scale</code> of the plot x-y dimension. Note, if the plot size is adjusted manually by the user it should be re-plotted to ensure the correct wind angle. The list may contain other options to <code>panel.arrows</code> in the <code>lattice</code> package. Other useful options include <code>length</code> , which controls the length of the arrow head and <code>angle</code> , which controls the angle of the arrow head. This option works best where there are not too many data to ensure overplotting does not become a problem.
<code>y.relation</code>	This determines how the y-axis scale is plotted. “same” ensures all panels use the same scale and “free” will use panel-specific scales. The latter is a useful setting when plotting data with very different values.
<code>x.relation</code>	This determines how the x-axis scale is plotted. “same” ensures all panels use the same scale and “free” will use panel-specific scales. The latter is a useful setting when plotting data with very different values.
<code>ref.x</code>	See <code>ref.y</code> for details.
<code>ref.y</code>	A list with details of the horizontal lines to be added representing reference line(s). For example, <code>ref.y = list(h = 50, lty = 5)</code> will add a dashed horizontal line at 50. Several lines can be plotted e.g. <code>ref.y = list(h = c(50, 100), lty = c(1, 5), col = c("green", "blue"))</code> . See <code>panel.abline</code> in the <code>lattice</code> package for more details on adding/controlling lines.
<code>k</code>	Smoothing parameter supplied to <code>gam</code> for fitting a smooth surface when <code>method = "level"</code> .
<code>dist</code>	When plotting smooth surfaces (<code>method = "level"</code> and <code>smooth = TRUE</code> , <code>dist</code> controls how far from the original data the predictions should be made. See <code>exclude.too.far</code> from the <code>mgcv</code> package. Data are first transformed to a unit square. Values should be between 0 and 1.
<code>map</code>	Should a base map be drawn? This option is under development.
<code>auto.text</code>	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO ₂ .
<code>...</code>	Other graphical parameters are passed onto <code>cutData</code> and an appropriate <code>lattice</code> plot function (<code>xyplot</code> , <code>levelplot</code> or <code>hexbinplot</code> depending on <code>method</code>). For example, <code>scatterPlot</code> passes the option <code>hemisphere</code>

```
data2003 <- selectByDate(mydata, year = 2003)
scatterPlot(data2003, x = "nox", y = "no2")
```

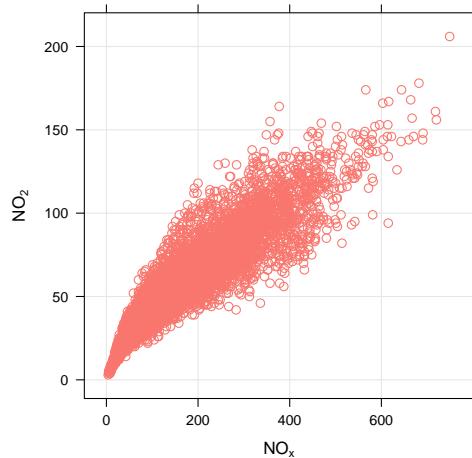


FIGURE 23.1 Scatter plot of hourly NO_x vs. NO_2 at Marylebone Road for 2003.

`= "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, for the default case `method = "scatter"` common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting. Other common graphical parameters, e.g. `layout` for panel arrangement, `pch` for plot symbol and `lwd` and `lty` for line width and type, as also available (see examples below).

For `method = "hexbin"` it can be useful to transform the scale if it is dominated by a few very high values. This is possible by supplying two functions: one that applies the transformation and the other that inverses it. For log scaling (the default) for example, `trans = function(x) log(x)` and `inv = function(x) exp(x)`. For a square root transform use `trans = sqrt` and `inv = function(x) x^2`. To not carry out any transformation the options `trans = NULL` and `inv = NULL` should be used.

23.3 Example of use

We provide a few examples of use and as usual, users are directed towards the help pages (type `?scatterPlot`) for more extensive examples.

First we select a subset of data (2003) using the `openair selectByDate` function and plot NO_x vs. NO_2 (Figure 23.1).

Often with several years of data, points are over-plotted and it can be very difficult to see what the underlying relationship looks like. One very effective method to use in these situations is to ‘bin’ the data and to colour the intervals by the number of counts of occurrences in each bin. There are various ways of doing this, but ‘hexagonal binning’ is particularly effective because of the way hexagons can be placed next to one another.¹⁴ To use hexagonal binning it will be necessary to install the `hexbin` package:

¹⁴In fact it is not possible to have a shape with more than 6 sides that can be used to form a lattice without gaps.

```
scatterPlot(data2003, x = "nox", y = "no2", method = "hexbin", col = "jet")
```

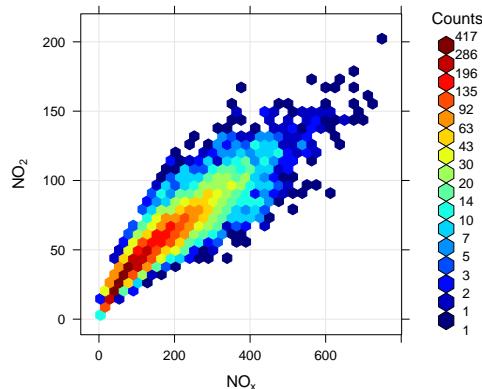


FIGURE 23.2 Scatter plot of hourly NO_x vs. NO_2 at Marylebone Road using hexagonal binning. The number of occurrences in each bin is colour-coded (not on a linear scale). It is now possible to see where most of the data lie and a better indication of the relationship between NO_x and NO_2 is revealed.

```
install.packages("hexbin")
```

Now it should be possible to make the plot by setting the method option to `method = "hexbin"`, as shown in Figure 23.2. The benefit of hexagonal binning is that it works equally well with enormous data sets e.g. several million records. In this case Figure 23.2 provides a clearer indication of the relationship between NO_x and NO_2 than Figure 23.1 because it reveals where most of the points lie, which is not apparent from Figure 23.1. Note that For `method = "hexbin"` it can be useful to transform the scale if it is dominated by a few very high values. This is possible by supplying two functions: one that applies the transformation and the other that inverses it. For log scaling for example (the default), `trans = function(x) log(x)` and `inv = function(x) exp(x)`. For a square root transform use `trans = sqrt` and `inv = function(x) x^2`. To not apply any transformation `trans = NULL` and `inv = NULL` should be used.

Note that when `method = "hexbin"` there are various options that are useful e.g. a border around each bin and the number of bins. For example, to place a grey border around each bin and set the bin size try:

```
scatterPlot(mydata, x = "nox", y = "no2", method = "hexbin", col = "jet",
           border = "grey", xbin = 15)
```

The hexagonal binning and other binning methods are useful but often the choice of bin size is somewhat arbitrary. Another useful approach is to use a kernel density estimate to show where most points lie. This is possible in `scatterPlot` with the `method = "density"` option. Such a plot is shown in Figure 23.3.

Sometimes it is useful to consider how the relationship between two variables varies by levels of a third. In `openair` this approach is possible by setting the option `type`. When `type` is another numeric variables, four plots are produced for different quantiles of that variable. We illustrate this point by considering how the relationship between NO_x and NO_2 varies with different levels of O_3 . We also take the opportunity to not plot the smooth line, but plot a linear fit instead and force the layout to be a 2 by 2 grid.

```
scatterPlot(selectByDate(mydata, year = 2003), x = "nox", y = "no2",
           method = "density", col = "jet")

## (loaded the KernSmooth namespace)
```

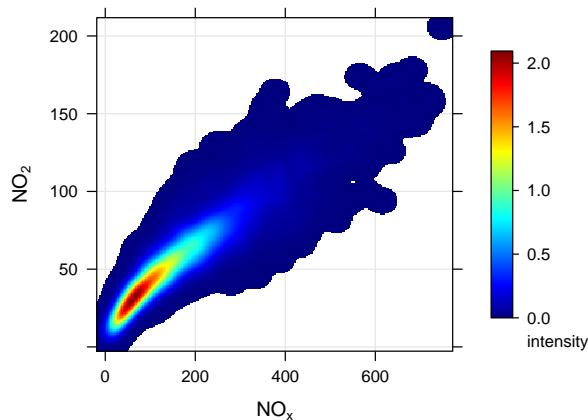


FIGURE 23.3 Scatter plot of hourly NO_x vs. NO_2 at Marylebone Road using a kernel density estimate to show where most of the points lie. The ‘intensity’ is a measure of how many points lie in a unit area of NO_x and NO_2 concentration.

```
scatterPlot(data2003, x = "nox", y = "no2", type = "o3", smooth = FALSE,
           linear = TRUE, layout = c(2, 2))
```

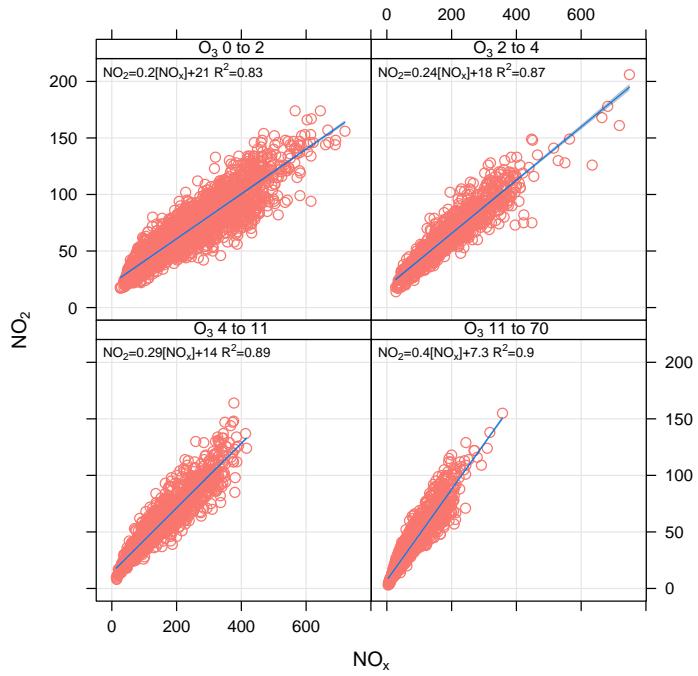


FIGURE 23.4 Scatter plot of hourly NO_x vs. NO_2 at Marylebone Road by different levels of O_3 .

Below is an extended example that brings together data manipulation, refined plot options and linear fitting of two variables with NO_x . The aim is to plot the weekly concentration of NO_x against PM_{10} and $\text{PM}_{2.5}$ and fit linear equations to both relationships. To do this we need the x variable as NO_x and the y variable as PM_{10} or $\text{PM}_{2.5}$, which

means we also need a column that will act as a grouping column i.e. identifies whether the y is PM₁₀ or PM_{2.5}.

```
# load the packages we need
library(dplyr)
library(reshape2)

# select the variables of interest
subdat <- select(mydata, date, nox, pm10, pm25)

# calculate weekly averages
subdat <- timeAverage(subdat, avg.time = "week")

# reshape so we have two variable columns
subdat <- melt(subdat, id.vars = c("date", "nox"))
head(subdat)

##           date nox variable value
## 1 1997-12-29 128     pm10    22
## 2 1998-01-05 189     pm10    34
## 3 1998-01-12 203     pm10    29
## 4 1998-01-19 142     pm10    29
## 5 1998-01-26 158     pm10    29
## 6 1998-02-02 308     pm10    41
```

Now we will plot weekly NO_x versus PM₁₀ and PM_{2.5} and fit a linear equation to both — and adjust some of the symbols (shown in Figure 23.5).

```
scatterPlot(subdat, x = "nox", y = "value", group = "variable",
            pch = 21:22, cex = 1.6, fill = c("dodgerblue", "tomato"),
            col = "white",
            linear = TRUE,
            xlab = "nox (ppb)",
            ylab = "PM concentration ( $\mu\text{g m}^{-3}$ )")
```

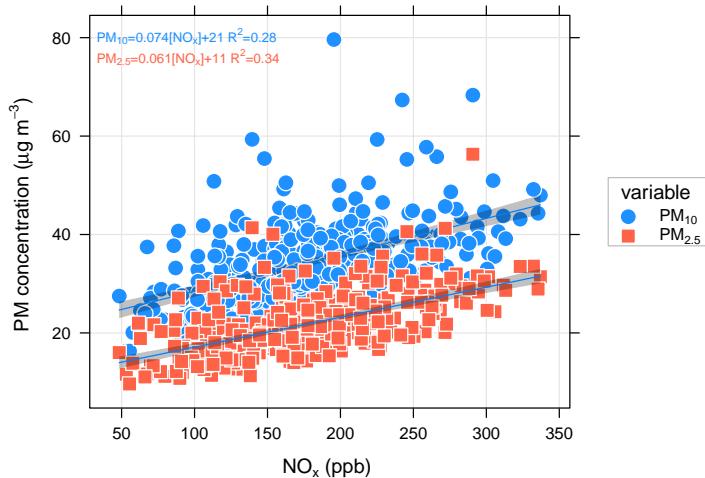


FIGURE 23.5 Scatter plot of weekly NO_x vs. PM₁₀ and PM_{2.5} at Marylebone Road with linear equations shown and plot symbols modified.

To gain a better idea of where the data lie and the linear fits, adding some transparency helps:

```
scatterPlot(subdat, x = "nox", y = "value", group = "variable",
            pch = 21:22, cex = 1.6, fill = c("dodgerblue", "tomato"),
            col = "white",
            linear = TRUE,
            xlab = "nox (ppb)",
            ylab = "PM concentration (ug/m3)",
            alpha = 0.2)
```

The above example will also work with `type`. For example, to consider how NO_x againsts PM_{10} and $\text{PM}_{2.5}$ varies by season:

```
scatterPlot(subdat, x = "nox", y = "value", group = "variable",
            pch = 21:22, cex = 2, fill = c("dodgerblue", "tomato"),
            col = "white", linear = TRUE,
            xlab = "nox (ppb)",
            ylab = "PM concentration (ug/m3)",
            type = "season")
```

Finally, we show how to plot a continuous colour scale for a third *numeric* variable setting the value of `z` to the third variable. Figure 23.6 shows again the relationship between NO_x and NO_2 but this time colour-coded by the concentration of O_3 . We also take the opportunity to split the data into seasons and weekday/weekend by setting `type = c("season", "weekend")`. There is an enormous amount of information that can be gained from plots such as this. Differences between weekdays and the weekend can highlight changes in emission sources, splitting by seasons can show seasonal influences in meteorology and background O_3 and colouring the data by the concentration of O_3 helps to show how O_3 concentrations affect NO_2 concentrations. For example, consider the summertime-weekday panel where it clearly shows that the higher NO_2 concentrations are associated with high O_3 concentrations. Indeed there are some hours where NO_2 is >100 ppb at quite low concentrations of NO_x (≈ 200 ppb). It would also be interesting instead of using O_3 concentrations from Marylebone Road to use O_3 from a background site.

Figure 23.6 was very easily produced but contains a huge amount of useful information showing the relationship between NO_x and NO_2 dependent upon the concentration of O_3 , the season and the day of the week. There are of course numerous other plots that are equally easily produced.

```
scatterPlot(selectByDate(data2003, month = 8), x = "date", y = "so2",
           z = "wd")
```

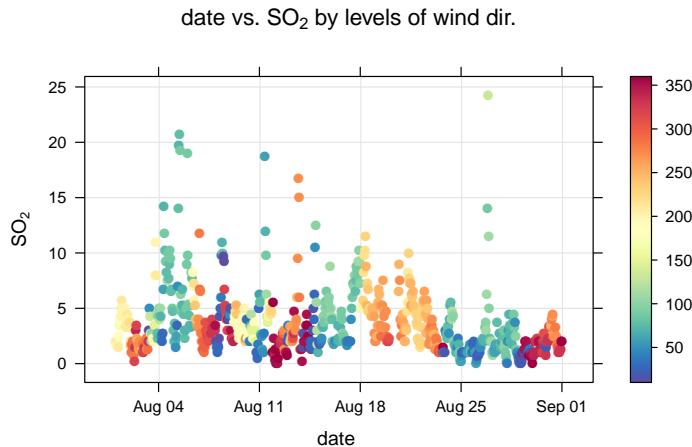


FIGURE 23.7 Scatter plot of date vs. SO₂ at Marylebone Road by different levels of wind direction for August 2003.

```
scatterPlot(data2003, x = "nox", y = "no2", z = "o3", type = c("season", "weekend"),
           limits = c(0, 30))
```

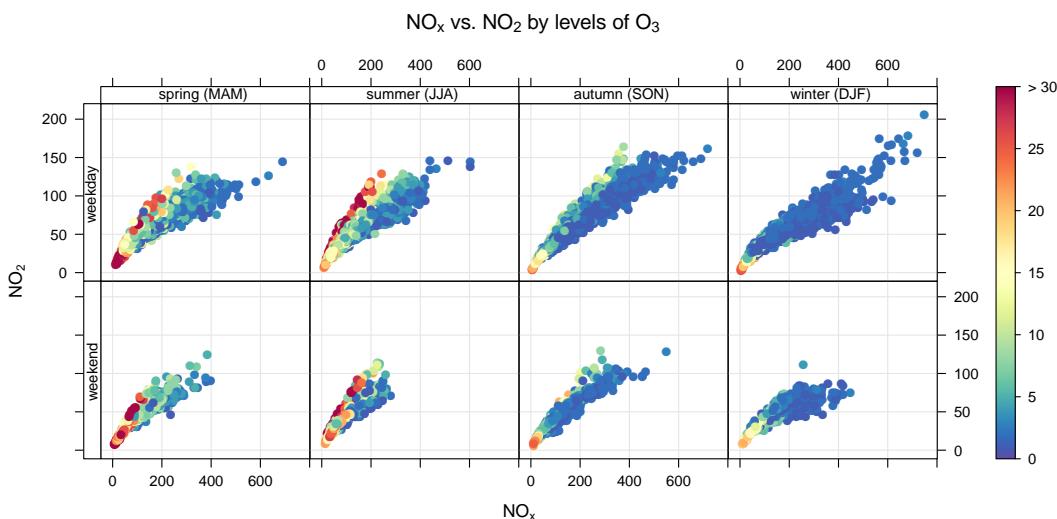


FIGURE 23.6 Scatter plot of hourly NO_x vs. NO₂ at Marylebone Road by different levels of O₃ split by season and weekday-weekend.

Figure 23.7 shows that `scatterPlot` can also handle dates on the x-axis; in this case shown for SO₂ concentrations coloured by wind direction for August 2003.

Similar to `timePlot`, `scatterPlot` can also plot wind vector arrows if wind speed and wind direction are available in the data frame. Figure 23.8 shows an example of using the `windflow` option. The Figure also sets many other options including showing the concentration of O₃ as a colour, setting the colour scale used and selecting a few days of interest using the `selectByDate` function. Figure 23.8 shows that when the wind direction changes to northerly, the concentration of NO₂ decreases and that of O₃ increases.

```
scatterPlot(selectByDate(mydata, start = "1/6/2001", end = "5/6/2001"),
  x = "date", y = "no2", z = "o3",
  col = "increment",
  windflow = list(scale = 0.15),
  key.footer = "o3\n(ppb)", main = NULL, ylab = "no2 (ppb)")
```

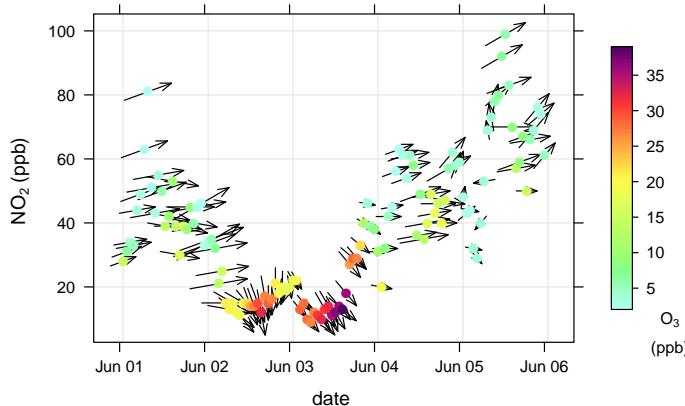


FIGURE 23.8 Scatter plot of date vs. NO₂ with the colour scale representing O₃. The wind flow vectors are also shown.

24 The **linearRelation** function

see also
[timeVariation](#)
[calcFno2](#)

This function considers linear relationships between two pollutants. The relationships are calculated on different times bases using a linear model. The slope and 95 % confidence interval in slope relationships by time unit are plotted in several different ways. The function is particularly useful when considering whether relationships are consistent with emissions inventories.

The relationships between pollutants can yield some very useful information about source emissions and how they change. A scatter plot between two pollutants is the usual way to investigate the relationship. A linear regression is useful to test the strength of the relationship. However, considerably more information can be gleaned by considering different time periods, such as how the relationship between two pollutants vary over time, by day of the week, diurnally and so on. The **linearRelation** function does just that — it fits a linear relationship between two pollutants over a wide range of time periods determined by **period**.

Consider the relationship between NO_x and NO₂. It is best to think of the relationship as:

$$y = m \cdot x + c \quad (8)$$

i.e.

$$NO_2 = m \cdot NO_x + c \quad (9)$$

In which case x corresponds to NO_x and y corresponds to NO₂. The plots show the gradient, m in what ever units the original data were in. For comparison with emission inventories it makes sense to have all the units expressed as mass. By contrast, oxidant slopes are best calculated in volume units e.g. ppb.

linearRelation function is particularly useful if background concentrations are first removed from roadside concentrations, as the increment will relate more directly

with changes in emissions. In this respect, using `linearRelation` can provide valuable information on how emissions may have changed over time, by hour of the day etc. Using the function in this way will require users to do some basic manipulation with their data first.

If a data frame is supplied that contains `nox`, `no2` and `o3`, the `y` can be chosen as `y = "ox"`. In function will therefore consider total oxidant slope (sum of $\text{NO}_2 + \text{O}_3$), which can provide valuable information on likely vehicle primary NO emissions. Note, however, that most roadside sites do not have ozone measurements and `calcFno2` is the alternative.

24.1 Options available

- mydata** A data frame minimally containing `date` and two pollutants.
- x** First pollutant that when plotted would appear on the x-axis of a relationship e.g. `x = "nox"`.
- y** Second pollutant that when plotted would appear on the y-axis of a relationship e.g. `y = "pm10"`.
- period** A range of different time periods can be analysed. The default is `month` but can be `year` and `week`. For increased flexibility an integer can be used e.g. for 3-month values `period = "3 month"`. Other cases include "`hour`" will show the diurnal relationship between `x` and `y` and "`weekday`" the day of the week relationship between `x` and `y`. "`day.hour`" will plot the relationship by weekday and hour of the day.
- condition** For `period = "hour"`, `period = "day"` and `period = "day.hour"`, setting `condition = TRUE` will plot the relationships split by year. This is useful for seeing how the relationships may be changing over time.
- n** The minimum number of points to be sent to the linear model. Because there may only be a few points e.g. hours where two pollutants are available over one week, `n` can be set to ensure that at least `n` points are sent to the linear model. If a period has hours < `n` that period will be ignored.
- rsq.thresh** The minimum correlation coefficient (R^2) allowed. If the relationship between `x` and `y` is not very good for a particular period, setting `rsq.thresh` can help to remove those periods where the relationship is not strong. Any R^2 values below `rsq.thresh` will not be plotted.
- ylab** y-axis title, specified by the user.
- auto.text** Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the '2' in NO_2 .
- cols** Colour for the points and uncertainty intervals.
- date.breaks** Number of major x-axis intervals to use. The function will try and choose a sensible number of dates/times as well as formatting the date/time appropriately to the range being considered. This does not always work as desired automatically. The user can therefore increase or decrease the number of intervals by adjusting the value of `date.breaks` up or down.
- ...** Other graphical parameters. A useful one to remove the strip with the date range on at the top of the plot is to set `strip = FALSE`.

```
linearRelation(mydata, x = "nox", y = "so2")
```

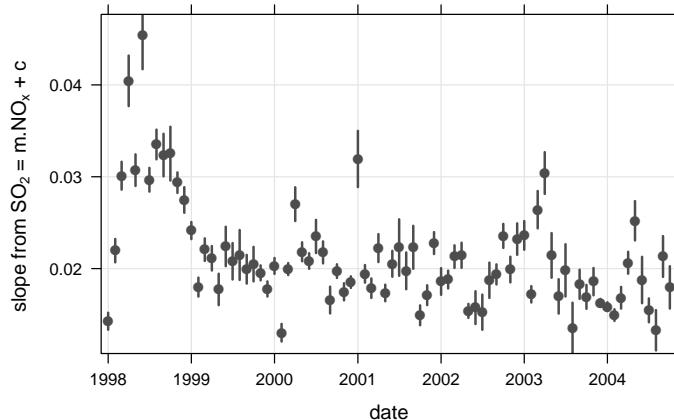


FIGURE 24.1 Relationship between NO_x and SO_2 using the `linearRelation` function. Note that the units of both pollutants are in ppb. The uncertainty in the slope of the hourly relationship between SO_2 and NO_x on a monthly basis is shown at 95 % confidence intervals. The smooth line and shaded area show the general trend using a loess smooth.

```
linearRelation(mydata, x = "nox", y = "ox", period = "day.hour")
```

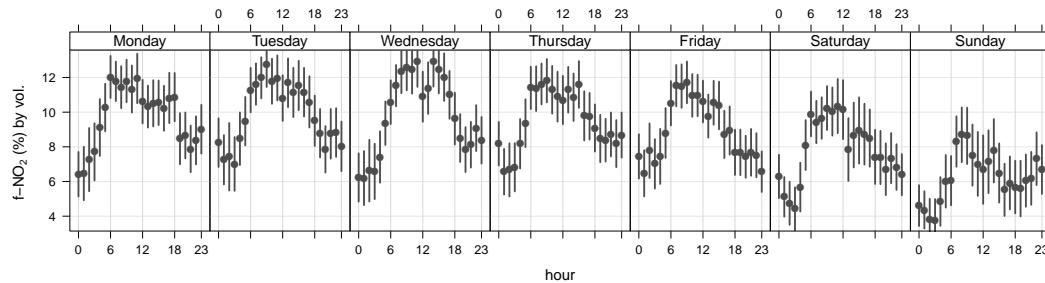


FIGURE 24.2 Oxidant slope by day of the week and hour of the day.

24.2 Example of use

Some examples of the `linearRelation` function are given in this section. The first example considers the ratio of SO_2/NO_x , which is plotted in Figure 24.1.

Figure 24.1 shows the relationship between NO_x and SO_2 . Early in the series (pre-1999) the ratio of SO_2/NO_x was relatively high (about 3.5 in volume units). However, from 1999 onwards the relationship has been relatively constant. One (probable) explanation for the higher ratio at the beginning of the series is due to a higher fuel sulphur content of petrol and diesel. There are many other examples shown in the package itself, type `?linearRelation` to see them.

One of the useful applications of this function is to consider the ‘oxidant’ (sum of NO_2 and O_3) slope where there are measurements of NO_x , NO_2 and O_3 at a site. At roadside sites the oxidant slope provides a good indication of the likely ratio of NO_2/NO_x in vehicle exhausts. Because there are few sites that measure O_3 at the roadside, the `calcFno2` function provides an alternative method of estimation. Figure 24.2 shows how the oxidant slope (an estimate of $f\text{-NO}_2$) varies by day of the week and hour of the day.

25 The `trendLevel` function

25.1 Purpose

The `trendLevel` function provides a way of rapidly showing a large amount of data in a condensed way. It is particularly useful for plotting the level of a value against two categorical variables. These categorical variables can pre-exist in a data set or be made on the fly using `openair`. By default it will show the mean value of a variable against two categorical variables but can also consider a wider range of statistics e.g. the maximum, frequency, or indeed a user-defined function. The function is much more flexible than this by showing temporal data and can plot ‘heat maps’ in many flexible ways. Both continuous colour scales and user-defined categorical scales can be used.

The `trendLevel` function shows how the value of a variable varies according to intervals of two other variables. The x and y variables can be categorical (factor or character) or numeric. The third variable (z) must be numeric and is coloured according to its value. Despite being called `trendLevel` the function is flexible enough to consider a wide range of plotting variables.

If the x and y variables are not categorical they are made so by splitting the data into quantiles (using `cutData`).

25.2 Options available

<code>mydata</code>	The openair data frame to use to generate the <code>trendLevel</code> plot.
<code>pollutant</code>	The name of the data series in <code>mydata</code> to sample to produce the <code>trendLevel</code> plot.
<code>x</code>	The name of the data series to use as the <code>trendLevel</code> x-axis. This is used with the <code>y</code> and <code>type</code> options to bin the data before applying <code>statistic</code> (see below). Other data series in <code>mydata</code> can also be used. (Note: <code>trendLevel</code> does not allow duplication in <code>x</code> , <code>y</code> and <code>type</code> options within a call.)
<code>y</code>	The names of the data series to use as the <code>trendLevel</code> y-axis and for additional conditioning, respectively. As <code>x</code> above.
<code>type</code>	See <code>y</code> .
<code>rotate.axis</code>	The rotation to be applied to <code>trendLevel</code> x and y axes. The default, <code>c(90, 0)</code> , rotates the x axis by 90 degrees but does not rotate the y axis. (Note: If only one value is supplied, this is applied to both axes; if more than two values are supplied, only the first two are used.)
<code>n.levels</code>	The number of levels to split <code>x</code> , <code>y</code> and <code>type</code> data into if numeric. The default, <code>c(10, 10, 4)</code> , cuts numeric <code>x</code> and <code>y</code> data into ten levels and numeric <code>type</code> data into four levels. (Notes: This option is ignored for date conditioning and factors. If less than three values are supplied, three values are determined by recursion; if more than three values are supplied, only the first three are used.)
<code>limits</code>	The colour scale range to use when generating the <code>trendLevel</code> plot.
<code>cols</code>	The colour set to use to colour the <code>trendLevel</code> surface. <code>cols</code> is passed to <code>openColours</code> for evaluation. See <code>?openColours</code> for more details.

`auto.text` Automatic routine text formatting. `auto.text = TRUE` passes common `lattice` labelling terms (e.g. `xlab` for the x-axis, `ylab` for the y-axis and `main` for the title) to the plot via `quickText` to provide common text formatting. The alternative `auto.text = FALSE` turns this option off and passes any supplied labels to the plot without modification.

`key.header, key.footer` Adds additional text labels above and/or below the scale key, respectively. For example, passing the options `key.header = "", key.footer = c("mean", "nox")` adds the addition text as a scale footer. If enabled (`auto.text = TRUE`), these arguments are passed to the scale key (`drawOpenKey`) via `quickText` to handle formatting. The term "`get.stat.name`", used as the default `key.header` setting, is reserved and automatically adds statistic function names or defaults to "`level`" when unnamed functions are requested via `statistic`.

`key.position` Location where the scale key should be plotted. Allowed arguments currently include "top", "right", "bottom" and "left".

`key` Fine control of the scale key via `drawOpenKey`. See `?drawOpenKey` for further details.

`labels` If a categorical colour scale is required then these labels will be used. Note there is one less label than break. For example, `labels = c("good", "bad", "very bad")`. `breaks` must also be supplied if labels are given.

`breaks` If a categorical colour scale is required then these breaks will be used. For example, `breaks = c(0, 50, 100, 1000)`. In this case "good" corresponds to values between 0 and 50 and so on. Users should set the maximum value of `breaks` to exceed the maximum data value to ensure it is within the maximum final range e.g. 100–1000 in this case. `labels` must also be supplied.

`statistic` The statistic method to be use to summarise locally binned `pollutant` measurements with. Three options are currently encoded: "mean" (default), "max" and "frequency". (Note: Functions can also be sent directly via `statistic`. However, this option is still in development and should be used with caution. See Details below.)

`stat.args` Additional options to be used with `statistic` if this is a function. The extra options should be supplied as a list of named parameters. (see Details below.)

`stat.safe.mode` An addition protection applied when using functions direclty with `statistic` that most users can ignore. This option returns `NA` instead of running `statistic` on binned subsamples that are empty. Many common functions terminate with an error message when applied to an empty dataset. So, this option provides a mechanism to work with such functions. For a very few cases, e.g. for a function that counted missing entries, it might need to be set to `FALSE` (see Details below.)

`drop.unused.types` Hide unused/empty `type` conditioning cases. Some conditioning options may generate empty cases for some data sets, e.g. a hour of the day when no measurements were taken. Empty `x` and `y` cases generate 'holes' in individual plots. However, empty `type` cases would produce

```
trendLevel(mydata, pollutant = "nox")
```

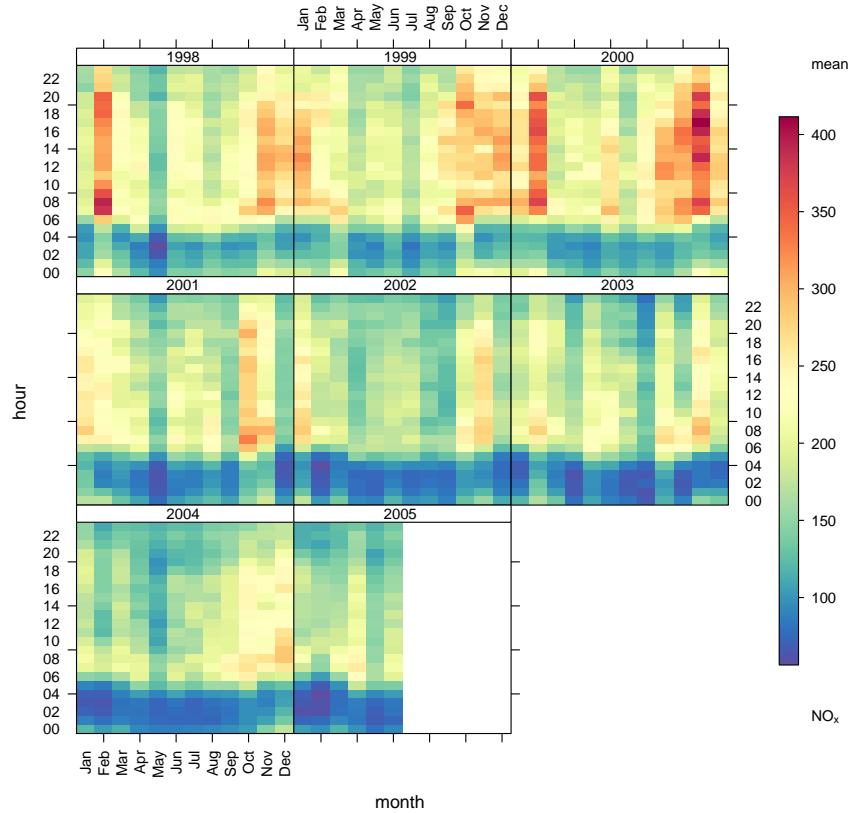


FIGURE 25.1 Standard `trendLevel` output.

blank panels if plotted. Therefore, the default, `TRUE`, excludes these empty panels from the plot. The alternative `FALSE` plots all `type` panels.

- `col.na` Colour to be used to show missing data.
- `...` Addition options are passed on to `cutData` for `type` handling and `levelplot` in `lattice` for finer control of the plot itself.

25.3 Example of use

The standard output from `trendLevel` is shown in Figure 25.1, which shows the variation in NO_x concentrations by year and hour of the day.

Figure 25.3 indicates that the highest NO_x concentrations most strongly associate with wind sectors about 200 degrees, appear to be decreasing over the years, but do not appear to associate with an SO₂ rich NO_x source. Using `type = "so2"` would have conditioned by absolute SO₂ concentration. As both a moderate contribution from an SO₂ rich source and a high contribution from an SO₂ poor source could generate similar SO₂ concentrations, such conditioning can sometimes blur interpretations. The use of this type of ‘over pollutant’ ratio reduces this blurring by focusing conditioning on cases when NO_x concentrations (be they high or low) associate with relatively high or low SO₂ concentrations.

The plot can be used in much more flexible ways. Here are some examples (not plotted):

```
trendLevel(mydata, pollutant = "nox", y = "wd", border = "white",
           cols = "jet")
```

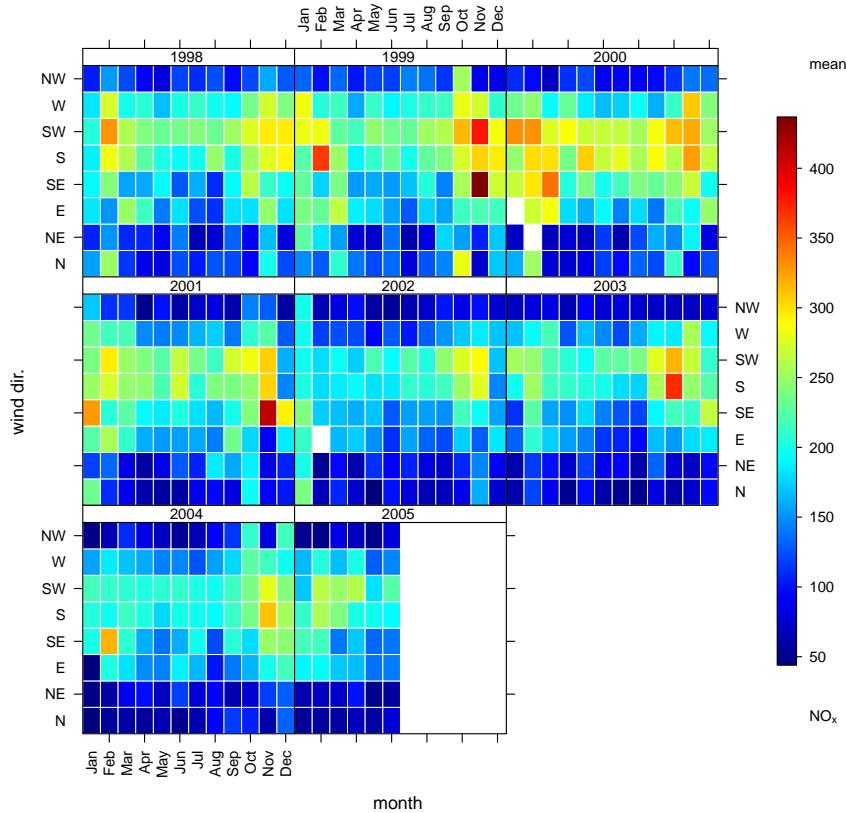


FIGURE 25.2 `trendLevel` output with wind direction as `y` axis. This plot also shows the effect of setting the border colour to white.

A plot of mean O₃ concentration shown by season and by daylight/nighttime hours.

```
trendLevel(mydata, x = "season", y = "daylight", pollutant = "o3")
```

Or by season and hour of the day:

```
trendLevel(mydata, x = "season", y = "hour", pollutant = "o3",
           cols = "increment")
```

How about NO_x versus NO₂ coloured by the concentration of O₃? `scatterPlot` could also be used to produce such a plot. However, one interesting difference with using `trendLevel` is that the data are split into quantiles where equal numbers of data exist in each interval. This approach can make it a bit easier to see the underlying relationship between variables. A scatter plot may have too much data to be clear and also outliers (or regions with relatively few data) that make it harder to see what is going on. The plot generated by the command below makes it a bit easier to see that it is the higher quantiles of NO₂ that are associated with higher O₃ concentration (as well as low NO_x and NO₂ concentrations).

```
trendLevel(mydata, x = "nox", y = "no2", pollutant = "o3", border = "white",
           n.levels = 10, statistic = "max", limits = c(0, 50))
```

```
## new field: so2/nox ratio
mydata$ratio <- mydata$so2/mydata$nox

## condition by mydata$new
trendLevel(mydata, "nox", x = "year", y = "wd", type = "ratio")
```

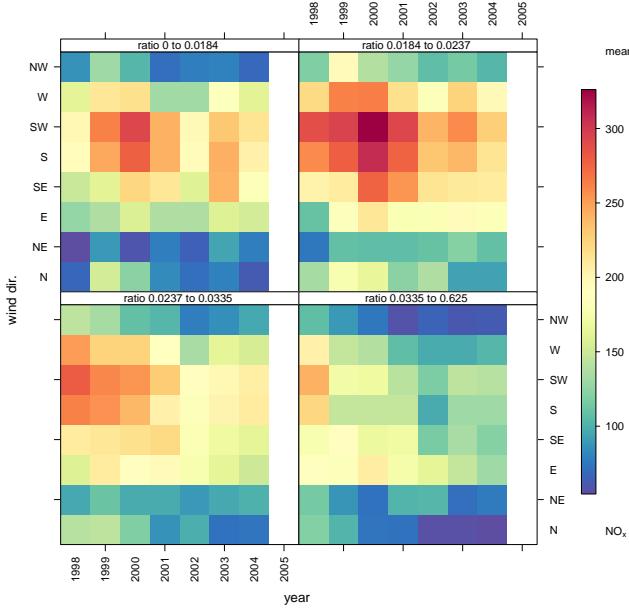


FIGURE 25.3 `trendLevel` output with SO₂: NO_x ratio `type` conditioning.

The plot can also be shown by wind direction sector, this time showing how O₃ varies by weekday, wind direction sector and NO_x quantile.

```
trendLevel(mydata, x = "nox", y = "weekday", pollutant = "o3",
           border = "white", n.levels = 10, statistic = "max",
           limits = c(0, 50), type = "wd")
```

By default `trendLevel` subsamples the plotted `pollutant` data by the supplied `x`, `y` and `type` parameters and in each case calculates the mean. The option `statistic` has always let you apply other statistics. For example, `trendLevel` also calculated the maximum via the option `statistic = "max"`. The user may also use their own statistic function.

As a simple example, consider the above plot which summarises by mean. This tells us about average concentrations. It might also be useful to consider a particular percentile of concentrations. This can be done by defining one's own function as shown in Figure 25.4.

This type of flexibility really opens up the potential of the function as a screening tool for the early stages of data analysis. Increased control of `x`, `y`, `type` and `statistic` allow you to very quickly explore your data and develop an understanding of how different parameters interact. Patterns in `trendLevel` plots can also help to direct your `openair` analysis. For example, possible trends in data conditioned by year would suggest that functions like `smoothTrend` or `TheilSen` could provide further insight. Likewise, `windRose` or `polarPlot` could be useful next steps if wind speed and direct conditioning produces interesting features. However, perhaps most interestingly, novel conditioning or the incorporation of novel parameters in this type of highly flexible

```
## function to estimate 95th percentile
percentile <- function(x) quantile(x, probs = 0.95, na.rm = TRUE)

## apply to present plot
trendLevel(mydata, "nox", x = "year", y = "wd", type = "ratio",
            statistic = percentile)
```

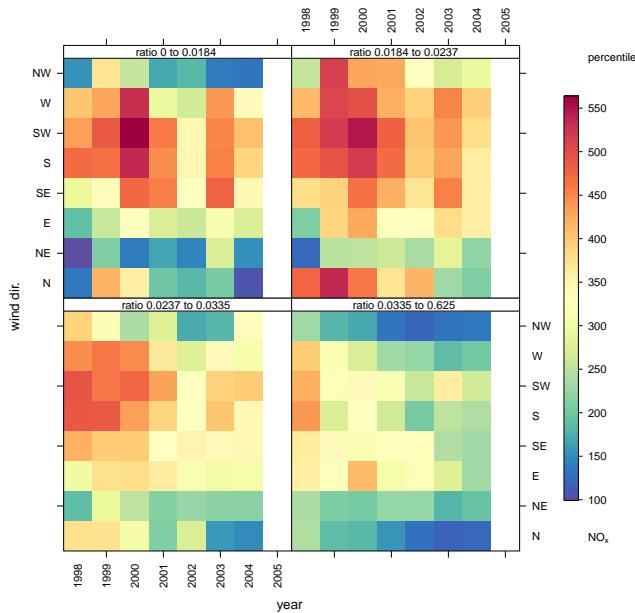


FIGURE 25.4 `trendLevel` using locally defined `statistic`.

function provides a means of developing new data visualisation and analysis methods.

`trendLevel` can also be used with user defined discrete colour scales as shown in Figure 25.5. In this case the default x and y variables are chosen (month and hour) split by `type` (year).

```
trendLevel(mydata, pollutant = "no2",
            border = "white", statistic = "max",
            breaks = c(0, 50, 100, 500),
            labels = c("low", "medium", "high"),
            cols = c("forestgreen", "yellow", "red"))
```

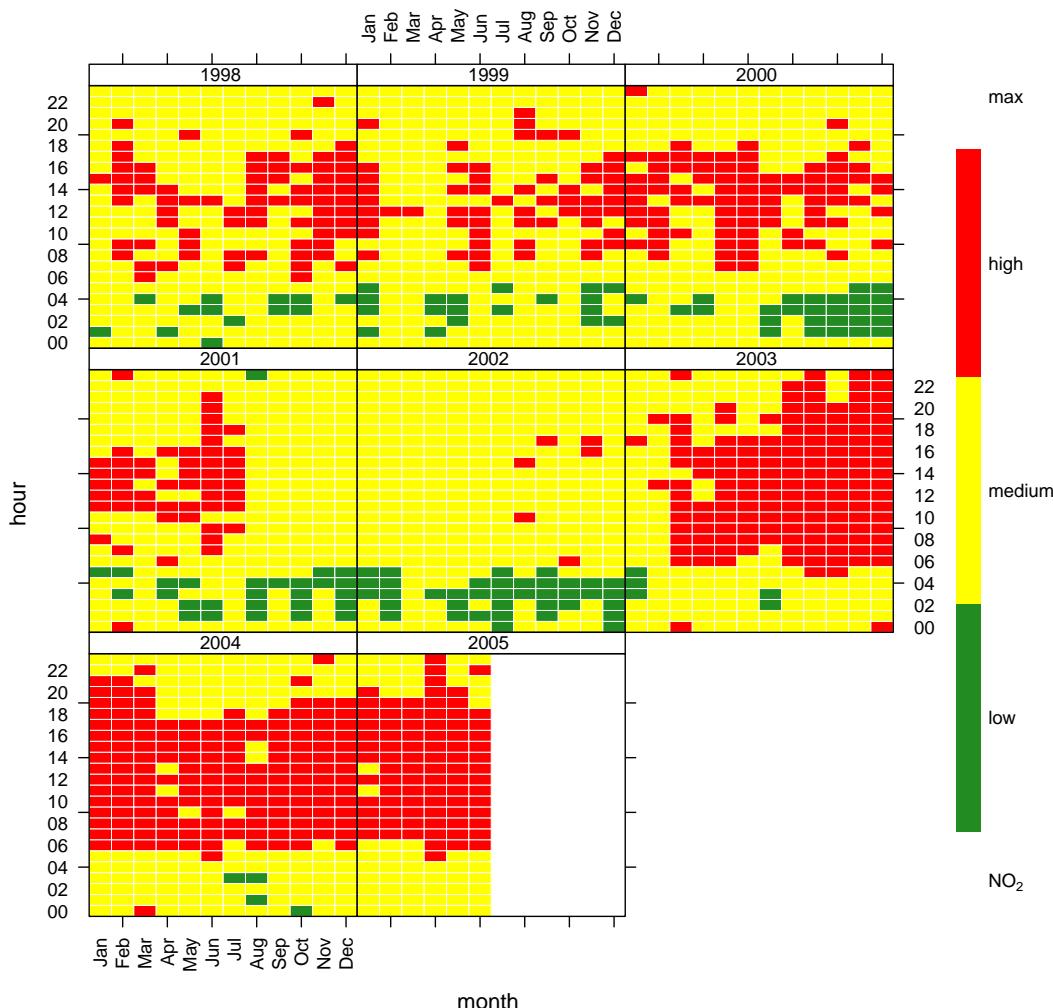


FIGURE 25.5 **trendLevel** plot for maximum NO₂ concentrations using a user-defined discrete colour scale.

26 openair back trajectory functions

Back trajectories are extremely useful in air pollution and can provide important information on air mass origins. Despite the clear usefulness of back trajectories, their use tends to be restricted to the research community. Back trajectories are used for many purposes from understanding the origins of air masses over a few days to undertaking longer term analyses. They are often used to filter air mass origins to allow for more refined analyses of air pollution — for example trends in concentration by air mass origin. They are often also combined with more sophisticated analyses such as cluster analysis to help group similar type of air mass by origin.

Perhaps one of the reasons why back trajectory analysis is not carried out more often is that it can be time consuming to do. This is particularly so if one wants to consider several years at several sites. It can also be difficult to access back trajectory data. In an

attempt to overcome some of these issues and expand the possibilities for data analysis, **openair** makes several functions available to access and analyse pre-calculated back trajectories.

Currently these functions allow for the import of pre-calculated back trajectories are several pre-define locations and some trajectory plotting functions. In time all of these functions will be developed to allow more sophisticated analyses to be undertaken. Also it should be recognised that these functions are in their early stages of development and will may continue to change and be refined.

This **importTraj** function imports pre-calculated back trajectories using the HYSPLIT trajectory model (Hybrid Single Particle Lagrangian Integrated Trajectory Model <http://ready.arl.noaa.gov/HYSPLIT.php>). Trajectories are run at 3-hour intervals and stored in yearly files (see below). The trajectories are started at ground-level (10m) and propagated backwards in time. The data are stored on web-servers at King's College London in a similar way to **importKCL**, which makes it very easy to import pre-processed trajectory data for a range of locations and years. **Note — the back trajectories have been pre-calculated for specific locations and stored as .RData objects. Users should contact David Carslaw to request the addition of other locations.** So far only a few receptors are available to users but in time the number will increase. It should be feasible for example to run back trajectories for the past 20 years at all the EMEP sites in Europe.¹⁵

Users may for various reasons wish to run HYSPLIT themselves e.g. for different starting heights, longer periods or more locations. Code and instructions have been provided in **Appendix D** for users wishing to do this. Users can also use different means of calculating back trajectories e.g. ECMWF and plot them in **openair** provided a few basic fields are present: **date** (POSIXct), **lat** (decimal latitude), **lon** (decimal longitude) and **hour.inc** the hour offset from the arrival date (i.e. from zero decreasing to the length of the back trajectories). See **?importTraj** for more details.

These trajectories have been calculated using the Global NOAA-NCEP/NCAR reanalysis data archives. The global data are on a latitude-longitude grid (2.5 degree). Note that there are many different meteorological data sets that can be used to run HYSPLIT e.g. including ECMWF data. However, in order to make it practicable to run and store trajectories for many years and sites, the NOAA-NCEP/NCAR reanalysis data is most useful. In addition, these archives are available for use widely, which is not the case for many other data sets e.g. ECMWF. HYSPLIT calculated trajectories based on archive data may be distributed without permission (see http://ready.arl.noaa.gov/HYSPLIT_agreement.php). For those wanting, for example, to consider higher resolution meteorological data sets it may be better to run the trajectories separately.

openair uses the **mapproj** package to allow users to user different map projections. By default the projection used is Lambert conformal, which is a conic projection best used for mid-latitude areas. The HYSPLIT model itself will use any one of three different projections depending on the latitude of the origin. If the latitude greater than 55.0 (or less than -55.0) then a polar stereographic projection is used, if the latitude greater than -25.0 and less than 25.0 the mercator projection is used and elsewhere (the mid-latitudes) the Lambert projection. All these projections (and many others) are available in the **mapproj** package.

Users should see the help file for **importTraj** to get an up to date list of receptors where back trajectories have been calculated.

As an example, we will import trajectories for London in 2010. Importing them is easy:

¹⁵It takes about 15 hours to run 20 years of 96-hour back trajectories at 3-hour intervals.

```
traj <- importTraj(site = "london", year = 2010)
```

The file itself contains lots of information that is of use for plotting back trajectories:

```
head(traj)

##   receptor year month day hour hour.inc lat      lon height pressure
## 1          1 2010     1    1    9        0 52 -0.100    10    995
## 2          1 2010     1    1    8       -1 52  0.057    10    995
## 3          1 2010     1    1    7       -2 52  0.250    10    995
## 4          1 2010     1    1    6       -3 52  0.488    11    995
## 5          1 2010     1    1    5       -4 53  0.767    11    995
## 6          1 2010     1    1    4       -5 53  1.065    11    996
##                     date2           date
## 1 2010-01-01 09:00:00 2010-01-01 09:00:00
## 2 2010-01-01 08:00:00 2010-01-01 09:00:00
## 3 2010-01-01 07:00:00 2010-01-01 09:00:00
## 4 2010-01-01 06:00:00 2010-01-01 09:00:00
## 5 2010-01-01 05:00:00 2010-01-01 09:00:00
## 6 2010-01-01 04:00:00 2010-01-01 09:00:00
```

The **traj** data frame contains among other things the latitude and longitude of the back trajectory, the height (m) and pressure (Pa) of the trajectory. The **date** field is the *arrival* time of the air-mass and is useful for linking with ambient measurement data.

The **trajPlot** function is used for plotting back trajectory lines and density plots and has the following options:

mydata Data frame, the result of importing a trajectory file using **importTraj**.

lon Column containing the longitude, as a decimal.

lat Column containing the latitude, as a decimal.

pollutant Pollutant to be plotted. By default the trajectory height is used.

type **type** determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in **cutData** e.g. "season", "year", "weekday" and so on. For example, **type = "season"** will produce four plots — one for each season.

It is also possible to choose **type** as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

type can be up length two e.g. **type = c("season", "weekday")** will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

map Should a base map be drawn? If **TRUE** the world base map from the **maps** package is used.

group It is sometimes useful to group and colour trajectories according to a grouping variable. See example below.

- map.fill** Should the base map be a filled polygon? Default is to fill countries.
- map.res** The resolution of the base map. By default the function uses the ‘world’ map from the **maps** package. If **map.res = "hires"** then the (much) more detailed base map ‘worldHires’ from the **mapdata** package is used. Use **library(mapdata)**. Also available is a map showing the US states. In this case **map.res = "state"** should be used.
- map.cols** If **map.fill = TRUE** **map.cols** controls the fill colour. Examples include **map.fill = "grey40"** and **map.fill = openColours("default", 10)**. The latter colours the countries and can help differentiate them.
- map.alpha** The transparency level of the filled map which takes values from 0 (full transparency) to 1 (full opacity). Setting it below 1 can help view trajectories, trajectory surfaces etc. and a filled base map.
- projection** The map projection to be used. Different map projections are possible through the **mapproj** package. See **?mapproject** for extensive details and information on setting other parameters and orientation (see below).
- parameters** From the **mapproj** package. Optional numeric vector of parameters for use with the projection argument. This argument is optional only in the sense that certain projections do not require additional parameters. If a projection does not require additional parameters then set to null i.e. **parameters = NULL**.
- orientation** From the **mapproj** package. An optional vector c(latitude, longitude, rotation) which describes where the “North Pole” should be when computing the projection. Normally this is c(90, 0), which is appropriate for cylindrical and conic projections. For a planar projection, you should set it to the desired point of tangency. The third value is a clockwise rotation (in degrees), which defaults to the midrange of the longitude coordinates in the map.
- grid.col** The colour of the map grid to be used. To remove the grid set **grid.col = "transparent"**.
- npoints** A dot is placed every **npoints** along each full trajectory. For hourly back trajectories points are plotted every **npoint** hours. This helps to understand where the air masses were at particular times and get a feel for the speed of the air (points closer together correspond to slower moving air masses). If **npoints = NA** then no points are added.
- origin** If true a filled circle dot is shown to mark the receptor point.
- ...** other arguments are passed to **cutData** and **scatterPlot**. This provides access to arguments used in both these functions and functions that they in turn pass arguments on to. For example, **plotTraj** passes the argument **cex** on to **scatterPlot** which in turn passes it on to the **lattice** function **xyplot** where it is applied to set the plot symbol size.

Next, we consider how to plot back trajectories with a few simple examples. The first example will consider a potentially interesting period when the Icelandic volcano, Eyjafjallajökull erupted in April 2010. The eruption of Eyjafjallajökull resulted in a flight-ban that lasted six days across many European airports. In Figure 26.1 **selectByDate** is

```
trajPlot(selectByDate(traj, start = "15/4/2010", end = "21/4/2010"),
         map.cols = openColours("hue", 10),
         col = "grey30")
```

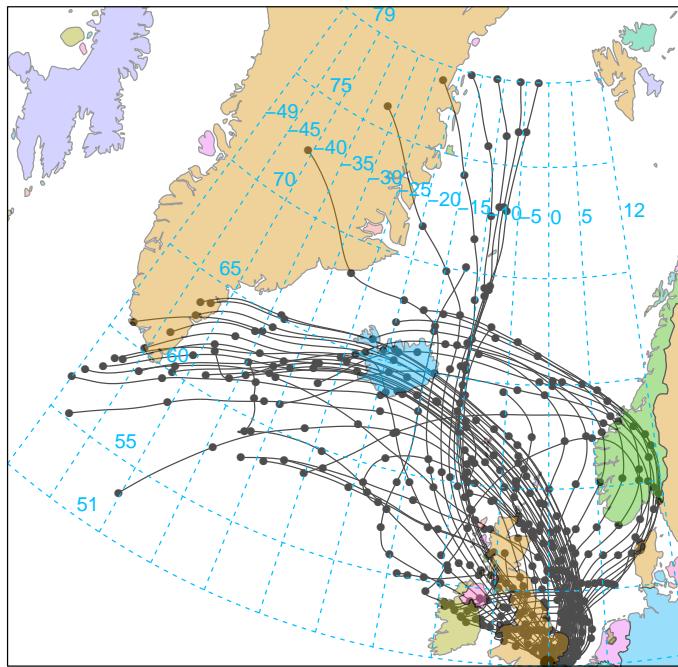


FIGURE 26.1 96-hour HYSPLIT back trajectories centred on London for 7 days in April 2010. Note the additional option to vary the country colours using `map.cols`. By default the map colours for all countries are grey.

used to consider the 7 days of interest and we choose to plot the back trajectories as lines rather than points (the default). Figure 26.1 does indeed show that many of the back trajectories originated from Iceland over this period. Note also the plot automatically includes a world base map. The base map itself is not at very high resolution by default but is useful for the sorts of spatial scales that back trajectories exist over. The base map is also global, so provided that there are pre-calculated back trajectories, these maps can be generated anywhere in the world. By default the function uses the ‘world’ map from the `maps` package. If `map.res = "hires"` then the (much) more detailed base map ‘worldHires’ from the `mapdata` package is used.¹⁶

Note that `trajPlot` will only plot *full length* trajectories. This can be important when plotting something like a single month e.g. by using `selectByDate` when on partial sections of some trajectories may be selected.

There are a few other ways of representing the data shown in Figure 26.1. For example, it might be useful to plot the trajectories for each day. To do this we need to make a new column ‘day’ which can be used in the plotting. The first example considers plotting the back trajectories in separate panels (Figure 26.2).

Another way of plotting the data is to group the trajectories by day and colour them. This time we also set a few other options to get the layout we want — shown in Figure 26.3.

So far the plots have provided information on where the back trajectories come from, grouped or split by day. It is also possible, in common with most other `openair` functions to split the trajectories by many other variables e.g. month, season and so on. However,

¹⁶You will need to load the `mapdata` package i.e. `library(mapdata)`.

```
## make a day column
traj$day <- as.Date(traj$date)

## plot it choosing a specific layout
trajPlot(selectByDate(traj, start = "15/4/2010", end = "21/4/2010"),
         type = "day", layout = c(7, 1))
```

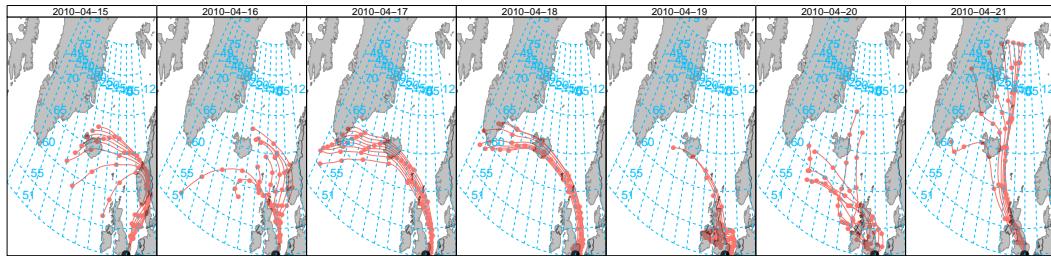


FIGURE 26.2 96-hour HYSPLIT back trajectories centred on London for 7 days in April 2010, shown separately for each day.

```
trajPlot(selectByDate(traj, start = "15/4/2010", end = "21/4/2010"),
         group = "day", col = "jet", lwd = 2, key.pos = "top",
         key.col = 4)
```

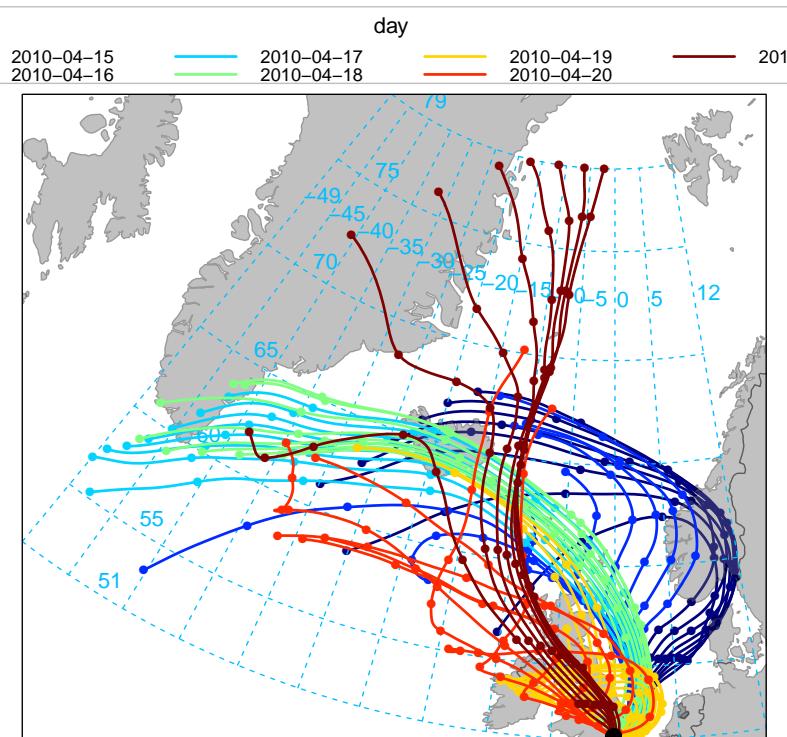


FIGURE 26.3 96-hour HYSPLIT back trajectories centred on London for 7 days in April 2010, shown grouped for each day and coloured accordingly.

perhaps one of the most useful approaches is to link the back trajectories with the concentrations of a pollutant. As mentioned previously, the back trajectory data has a column ‘date’ representing the arrival time of the air mass that can be used to link with concentration measurements. A couple of steps are required to do this using the **merge** function.

```
## import data for North Kensington
kc1 <- importAURN("kc1", year = 2010)
# now merge with trajectory data by 'date'
traj <- merge(traj, kc1, by = "date")
## look at first few lines
head(traj)

##                               date receptor year month      day hour hour.inc lat    lon
## 1 2010-01-01 09:00:00          1 2010     1 2010-01-01    9      0 52 -
0.100
## 2 2010-01-01 09:00:00          1 2010     1 2010-01-01    8      - 
1 52  0.057
## 3 2010-01-01 09:00:00          1 2010     1 2010-01-01    7      - 
2 52  0.250
## 4 2010-01-01 09:00:00          1 2010     1 2010-01-01    6      - 
3 52  0.488
## 5 2010-01-01 09:00:00          1 2010     1 2010-01-01    5      - 
4 53  0.767
## 6 2010-01-01 09:00:00          1 2010     1 2010-01-01    4      - 
5 53  1.065
##   height pressure                  date2 code           site o3 no2 co so2
## 1     10        995 2010-01-01 09:00:00 KC1 London N. Kensington 46 29 0.3 0
## 2     10        995 2010-01-01 08:00:00 KC1 London N. Kensington 46 29 0.3 0
## 3     10        995 2010-01-01 07:00:00 KC1 London N. Kensington 46 29 0.3 0
## 4     11        995 2010-01-01 06:00:00 KC1 London N. Kensington 46 29 0.3 0
## 5     11        995 2010-01-01 05:00:00 KC1 London N. Kensington 46 29 0.3 0
## 6     11        996 2010-01-01 04:00:00 KC1 London N. Kensington 46 29 0.3 0
##   pm10 nox no pm2.5 nv2.5 v2.5 nv10 v10 ws wd
## 1     8 38  6   NA   NA   NA   8   0 NA NA
## 2     8 38  6   NA   NA   NA   8   0 NA NA
## 3     8 38  6   NA   NA   NA   8   0 NA NA
## 4     8 38  6   NA   NA   NA   8   0 NA NA
## 5     8 38  6   NA   NA   NA   8   0 NA NA
## 6     8 38  6   NA   NA   NA   8   0 NA NA
```

This time we can use the option **pollutant** in the function **trajPlot**, which will plot the back trajectories coloured by the concentration of a pollutant. Figure 26.4 does seem to show elevated PM₁₀ concentrations originating from Iceland over the period of interest. In fact, these elevated concentrations occur on two days as shown in Figure 26.2. However, care is needed when interpreting such data because other analysis would need to rule out other reasons why PM₁₀ could be elevated; in particular due to local sources of PM₁₀. There are lots of **openair** functions that can help here e.g. **timeVariation** or **timePlot** to see if NO_x concentrations were also elevated (which they seem to be). It would also be worth considering other sites for back trajectories that could be less influenced by local emissions.

However, it is possible to account for the PM that is local to some extent by considering the relationship between NO_x and PM₁₀ (or PM_{2.5}). For example, using **scatterPlot** (not shown):

```
scatterPlot(kc1, x = "nox", y = "pm2.5", avg = "day", linear = TRUE)
```

which suggests a gradient of 0.084. Therefore we can remove the PM₁₀ that is associated NO_x in **kc1** data, making a new column **pm.new**:

```
trajPlot(selectByDate(traj, start = "15/4/2010", end = "21/4/2010"),
         pollutant = "pm10", col = "jet", lwd = 2)
```

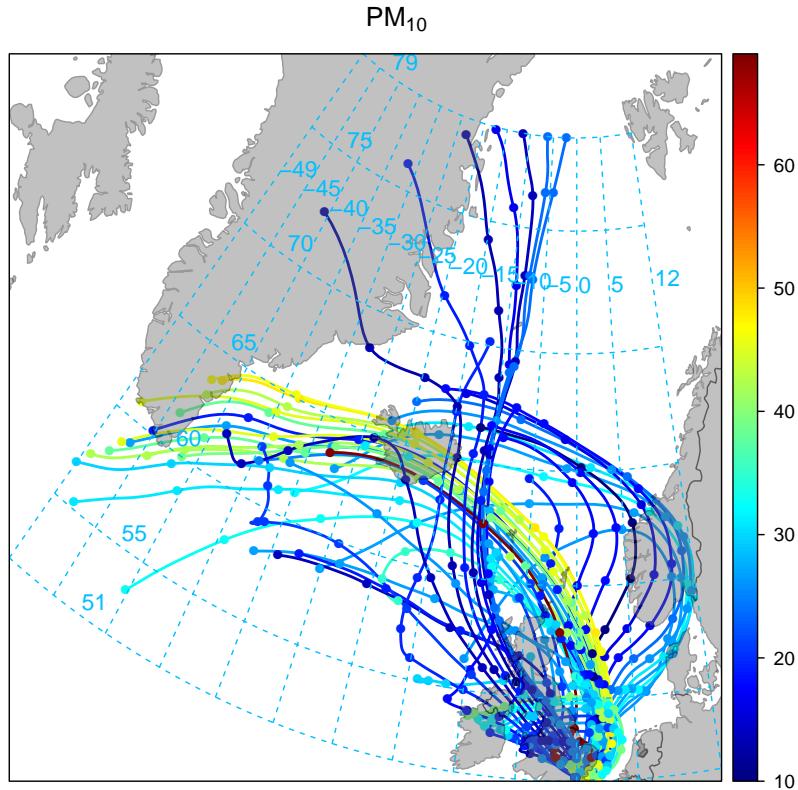


FIGURE 26.4 96-hour HYSPLIT back trajectories centred on London for 7 days in April 2010, coloured by the concentration of PM_{10} ($\mu\text{g m}^{-3}$).

```
kc1 <- transform(kc1, pm.new = pm10 - 0.084 * nox)
```

We have already merged `kc1` with `traj`, so to keep things simple we import `traj` again and merge it with `kc1`. Note that if we had thought of this initially, `pm.new` would have been calculated first before merging with `traj`.

```
traj <- importTraj(site = "london", year = 2010)
traj <- merge(traj, kc1, by = "date")
```

Now it is possible to plot the trajectories:

```
trajPlot(selectByDate(traj, start = "15/4/2010", end = "21/4/2010"),
         pollutant = "pm.new", col = "jet", lwd = 2)
```

Which, interestingly still clearly shows elevated PM_{10} concentrations for those two days that cross Iceland. The same is also true for $\text{PM}_{2.5}$. However, as mentioned previously, checking other sites in more rural areas would be a good idea.

26.1 Trajectory gridded frequencies

The HYSPLIT model itself contains various analysis options for gridding trajectory data. Similar capabilities are also available in **openair** where the analyses can be extended using other **openair** capabilities. It is useful to gain an idea of where trajectories come

from. Over the course of a year representing trajectories as lines or points results in a lot of over-plotting. Therefore it is useful to grid the trajectory data and calculate various statistics by considering latitude-longitude intervals.

The first analysis considers the number of unique trajectories in a particular grid square. This is achieved by using the `trajLevel` function and setting the `statistic` option to “frequency”. Figure 26.5 shows the frequency of back trajectory crossings for the North Kensington data. In this case it highlights that most trajectory origins are from the west and north for 2010 at this site. Note that in this case, `pollutant` can just be the trajectory height (or another numeric field) rather than an actual pollutant because only the frequencies are considered.

The `trajLevel` function has the following options:

<code>mydata</code>	Data frame, the result of importing a trajectory file using <code>importTraj</code>
<code>lon</code>	Column containing the longitude, as a decimal.
<code>lat</code>	Column containing the latitude, as a decimal.
<code>pollutant</code>	Pollutant to be plotted. By default the trajectory height is used.
<code>type</code>	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. “season”, “year”, “weekday” and so on. For example, <code>type = "season"</code> will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

`type` can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

`smooth` Should the trajectory surface be smoothed?

`statistic` For `trajLevel`. By default the function will plot the trajectory frequencies.

For `trajLevel`, the argument `method = "hexbin"` can be used. In this case hexagonal binning of the trajectory *points* (i.e. a point every three hours along each back trajectory). The plot then shows the trajectory frequencies uses hexagonal binning. This is an alternative way of viewing trajectory frequencies compared with `statistic = "frequency"`.

There are also various ways of plotting concentrations.

It is also possible to set `statistic = "difference"`. In this case trajectories where the associated concentration is greater than `percentile` are compared with the the full set of trajectories to understand the differences in freqeuncies of the origin of air masses. The comparsion is made by comparing the percentage change in gridded frequencies. For example,

such a plot could show that the top 10% of concentrations of PM10 tend to originate from air-mass origins to the east.

If `statistic = "pscf"` then a Potential Source Contribution Function map is produced. If `statistic = "cwt"` then concentration weighted trajectories are plotted.

If `statistic = "cwt"` then the Concentration Weighted Trajectory approach is used. See details.

percentile For `trajLevel`. The percentile concentration of `pollutant` against which the all trajectories are compared.

map Should a base map be drawn? If `TRUE` the world base map from the `maps` package is used.

lon.inc The longitude-interval to be used for binning data for `trajLevel`.

lat.inc The latitude-interval to be used for binning data when `trajLevel`.

min.bin For `trajLevel` the minimum number of unique points in a grid cell. Counts below `min.bin` are set as missing. For `trajLevel` gridded outputs.

map.fill Should the base map be a filled polygon? Default is to fill countries.

map.res The resolution of the base map. By default the function uses the ‘world’ map from the `maps` package. If `map.res = "hires"` then the (much) more detailed base map ‘worldHires’ from the `mapdata` package is used. Use `library(mapdata)`. Also available is a map showing the US states. In this case `map.res = "state"` should be used.

map.cols If `map.fill = TRUE` `map.cols` controls the fill colour. Examples include `map.fill = "grey40"` and `map.fill = openColours("default", 10)`. The latter colours the countries and can help differentiate them.

map.alpha The transparency level of the filled map which takes values from 0 (full transparency) to 1 (full opacity). Setting it below 1 can help view trajectories, trajectory surfaces etc. and a filled base map.

projection The map projection to be used. Different map projections are possible through the `mapproj` package. See `?mapproject` for extensive details and information on setting other parameters and orientation (see below).

parameters From the `mapproj` package. Optional numeric vector of parameters for use with the projection argument. This argument is optional only in the sense that certain projections do not require additional parameters. If a projection does not require additional parameters then set to null i.e. `parameters = NULL`.

orientation From the `mapproj` package. An optional vector `c(latitude, longitude, rotation)` which describes where the “North Pole” should be when computing the projection. Normally this is `c(90, 0)`, which is appropriate for cylindrical and conic projections. For a planar projection, you should set it to the desired point of tangency. The third value is a clockwise rotation (in degrees), which defaults to the midrange of the longitude coordinates in the map.

```
trajLevel(traj, statistic = "frequency")
```

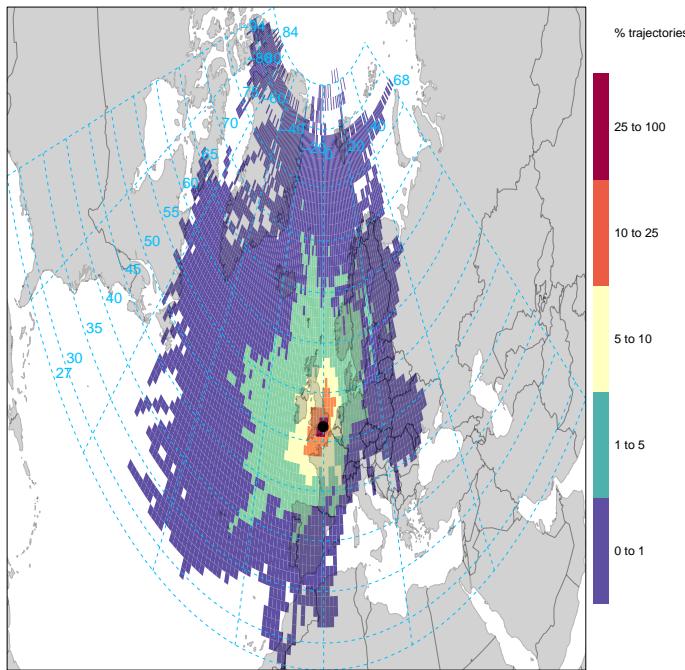


FIGURE 26.5 Gridded back trajectory frequencies. The `border = NA` option removes the border around each grid cell.

grid.col The colour of the map grid to be used. To remove the grid set `grid.col = "transparent"`.

origin should the receptor origin be shown by a black dot?

... other arguments are passed to `cutData` and `scatterPlot`. This provides access to arguments used in both these functions and functions that they in turn pass arguments on to. For example, `plotTraj` passes the argument `cex` on to `scatterPlot` which in turn passes it on to the `lattice` function `xyplot` where it is applied to set the plot symbol size.

It is also possible to use hexagonal binning to gain an idea about trajectory frequencies. In this case each 3-hour point along each trajectory is used in the counting. The code below focuses more on Europe and uses the hexagonal binning method. Note that the effect of the very few high number of points at the origin has been diminished by plotting the data on a log scale — see [page 177](#) for details.

26.2 Trajectory source contribution functions

Back trajectories offer the possibility to undertake receptor modelling to identify the location of major emission sources. When many back trajectories (over months to years) are analysed in specific ways they begin to show the geographic origin most associated with elevated concentrations. With enough (dissimilar) trajectories those locations leading to the highest concentrations begin to be revealed. When a whole year of back trajectory data is plotted the individual back trajectories can extend 1000s of km. There are many approaches using back trajectories in this way and Fleming et al.

```
trajLevel(subset(traj, lat > 30 & lat < 70 & lon > -30 & lon < 20),
           method = "hexbin", col = "jet",
           xbin = 40)
```

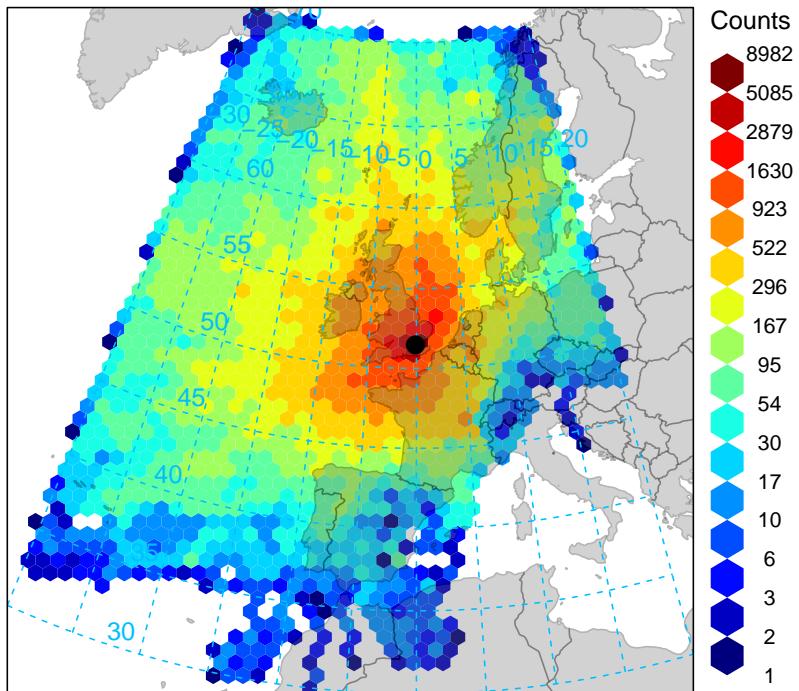


FIGURE 26.6 Gridded back trajectory frequencies with hexagonal binning.

(2012) provide a good overview of the methods available. **openair** has implemented a few of these techniques and over time these will be refined and extended.

26.2.1 Identifying the contribution of high concentration back trajectories

A useful analysis to undertake is to consider the pattern of frequencies for two different conditions. In particular, there is often interest in the origin of high concentrations for different pollutants. For example, compared with data over a whole year, how do the frequencies of occurrence differ? Figure 26.7 shows an example of such an analysis for PM₁₀ concentrations. By default the function will compare concentrations >90th percentile with the full year. The percentile level is controlled by the option `percentile`. Note also there is an option `min.bin` that will exclude grid cells where there are fewer than `min.bin` data points. The analysis compares the percentage of time the air masses are in particular grid squares for all data and a subset of data where the concentrations are greater than the given percentile. The graph shows the absolute percentage difference between the two cases i.e. high – base.

Figure 26.7 shows that compared with the whole year, high PM₁₀ concentrations (>90th percentile) are more prevalent when the trajectories originate from the east, which is seen by the positive values in the plot. Similarly there are relatively fewer occurrences of these high concentration back trajectories when they originate from the west. This analysis is in keeping with the highest PM₁₀ concentrations being largely controlled by secondary aerosol formation from air-masses originating during anticyclonic conditions from mainland Europe.

Note that it is also possible to use conditioning with these plots. For example to split

```
trajLevel(traj, pollutant = "pm10", statistic = "difference",
          col = c("skyblue", "white", "tomato"), min.bin = 50, border = NA,
          xlim = c(-20, 20), ylim = c(40, 70))
```

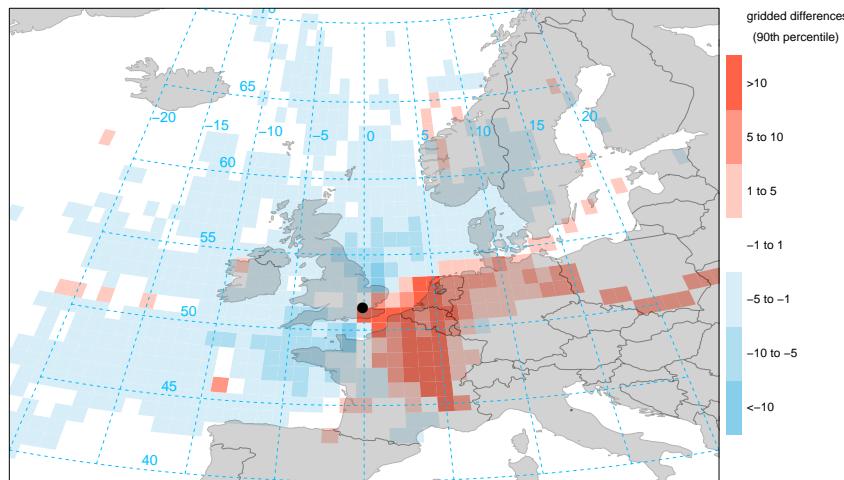


FIGURE 26.7 Gridded back trajectory frequencies showing the percentage difference in occurrence for high PM₁₀ concentrations (90th percentile) compared with conditions over the full year.

the frequency results by season:

```
trajLevel(traj, pollutant = "pm10", statistic = "frequency", col = "heat",
          type = "season")
```

26.2.2 Allocating trajectories to different wind sectors

One of the key aspects of trajectory analysis is knowing something about where air masses have come from. Cluster analysis can be used to group trajectories based on their origins and this is discussed in [Section 26.3](#). A simple approach is to consider different wind sectors e.g. N, NE, E and calculate the proportion of time a particular back trajectory resides in a specific sector. It is then possible to allocate a particular trajectory to a sector based on some assumption about the proportion of time it is in that sector — for example, assume a trajectory is from the west sector if it spends at least 50% of its time in that sector or otherwise record the allocation as ‘unallocated’. The code below can be used as the basis of such an approach.

First we import the trajectories, which in this case are for London in 2010:

```
traj <- importTraj(site = "london", year = 2010)
```

```

## need start/end lat and lon to work out angles
id <- which(traj$hour.inc == 0)
y0 <- traj$lat[id[1]]
x0 <- traj$lon[id[1]]

## calculate angle and then assign sector
traj <- transform(traj, angle = atan2(lon - x0, lat - y0) * 360 / 2 / pi)
ids <- which(traj$angle < 0)
traj$angle[ids] <- traj$angle[ids] + 360

traj$sector <- cut(traj$angle, breaks = seq(22.5, 382.5, 45),
                     labels = c("NE", "E", "SE", "S", "SW", "W",
                               "NW", "N"))
traj[ , "sector"][is.na(traj[ , "sector"])] <- "N" # for wd < 22.5

## count frequencies of sectors for each trajectory and the maximum
alloc <- tapply(traj$sector, traj$date, table)
alloc <- as.data.frame(do.call(rbind, alloc))
alloc$max <- apply(alloc, 1, max)

## identify the most frequent sector
alloc$sec <- sapply(1:nrow(alloc),
                     function(x) colnames(alloc)[which.max(alloc[x, ])])

## assign to most frequent sector, or label unallocated
## below assumes at least 50 out of 96 hours for assignment
alloc$sec <- ifelse(alloc$max > 50, alloc$sec, "unallocated")
alloc$date <- ymd_hms(rownames(alloc))

alloc <- alloc[, c("date", "sec")]

## merge with original data
traj <- merge(traj, alloc, by = "date", all = TRUE)

```

Now it is possible to post-process the data. **traj** now has the angle, sector and allocation (**sec**).

```

head(traj)

##           date receptor year month day hour hour.inc lat   lon height
## 1 2010-01-01 09:00:00      1 2010    1   1    9     0 52 -
0.100    10
## 2 2010-01-01 09:00:00      1 2010    1   1    8     - 
1 52  0.057    10
## 3 2010-01-01 09:00:00      1 2010    1   1    7     - 
2 52  0.250    10
## 4 2010-01-01 09:00:00      1 2010    1   1    6     - 
3 52  0.488    11
## 5 2010-01-01 09:00:00      1 2010    1   1    5     - 
4 53  0.767    11
## 6 2010-01-01 09:00:00      1 2010    1   1    4     - 
5 53  1.065    11
##   pressure          date2 angle sector      sec
## 1      995 2010-01-01 09:00:00    0      N unallocated
## 2      995 2010-01-01 08:00:00   31      NE unallocated
## 3      995 2010-01-01 07:00:00   33      NE unallocated
## 4      995 2010-01-01 06:00:00   36      NE unallocated
## 5      995 2010-01-01 05:00:00   39      NE unallocated
## 6      996 2010-01-01 04:00:00   42      NE unallocated

```

First, merge the air quality data from North Kensington:

```
traj <- merge(traj, kc1, by = "date")
```

We can work out the mean concentration by allocation, which shows the clear importance for the east and south-east sectors.

```
tapply(traj$pm2.5, traj$sec, mean, na.rm = TRUE)
```

	E	N	NE	NW	S	SE
##	22	12	13	14	14	29
##	SW unallocated		W			
##	11	15	12			

Finally, the percentage of the year in each sector can be calculated as follows:

```
100 * prop.table(table(traj$sec))

##
##          E          N          NE         NW          S          SE
##        6.7        5.5       10.3       9.2        2.2       2.1
##    SW unallocated
##        8.8       29.5      25.8
```

26.2.3 Potential Source Contribution Function (PSCF)

If `statistic = "pscfc"` then the Potential Source Contribution Function (PSCF) is plotted. The PSCF calculates the probability that a source is located at latitude i and longitude j (Fleming et al. 2012; Pekney et al. 2006). The PSCF is somewhat analogous to the CPF function described on [page 100](#) that considers local wind direction probabilities. In fact, the two approaches have been shown to work well together (Pekney et al. 2006). The PSCF approach has been widely used in the analysis of air mass back trajectories. Ara Begum et al. (2005) for example assessed the method against the known locations of wildfires and found it performed well for PM_{2.5}, EC (elemental carbon) and OC (organic carbon) and that other (non-fire related) species such as sulphate had different source origins. The basis of PSCF is that if a source is located at (i, j), an air parcel back trajectory passing through that location indicates that material from the source can be collected and transported along the trajectory to the receptor site. PSCF solves

$$\text{PSCF} = \frac{m_{ij}}{n_{ij}} \quad (10)$$

where n_{ij} is the number of times that the trajectories passed through the cell (i, j) and m_{ij} is the number of times that a source concentration was high when the trajectories passed through the cell (i, j). The criterion for determining m_{ij} is controlled by `percentile`, which by default is 90. Note also that cells with few data have a weighting factor applied to reduce their effect.

An example of a PSCF plot is shown in [Figure 26.8](#) for PM_{2.5} for concentrations >90th percentile. This Figure gives a very clear indication that the principal (high) sources are dominated by source origins in mainland Europe — particularly around the Benelux countries.

26.2.4 Concentration Weighted Trajectory (CWT)

A limitation of the PSCF method is that grid cells can have the same PSCF value when sample concentrations are either only slightly higher or much higher than the criterion (Hsu et al. 2003). As a result, it can be difficult to distinguish moderate sources from

```
trajLevel(subset(traj, lon > -20 & lon < 20 & lat > 45 & lat < 60),
          pollutant = "pm2.5", statistic = "pscf", col = "increment",
          border = NA)
```

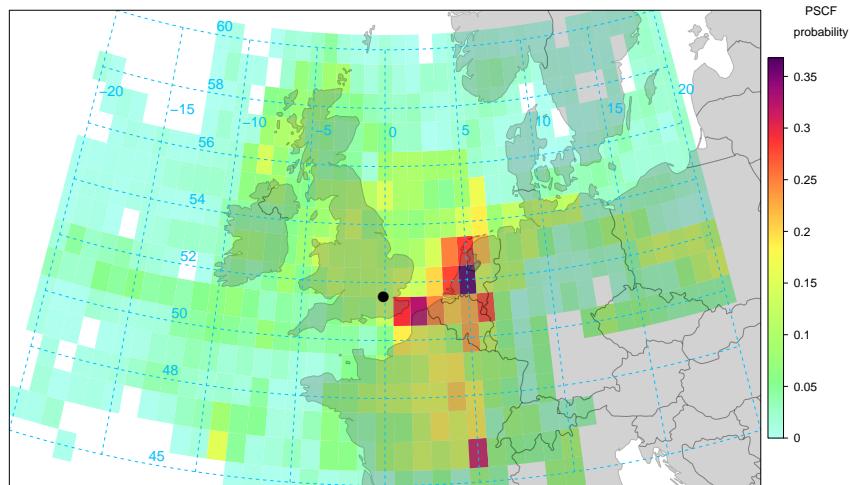


FIGURE 26.8 PSCF probabilities for PM_{2.5} concentrations (90th percentile).

strong ones. Seibert et al. (1994) computed concentration fields to identify source areas of pollutants. This approach is sometimes referred to as the CWT or CF (concentration field). A grid domain was used as in the PSCF method. For each grid cell, the mean (CWT) or logarithmic mean (used in the Residence Time Weighted Concentration (RTWC) method) concentration of a pollutant species was calculated as follows:

$$\ln(\bar{C}_{ij}) = \frac{1}{\sum_{k=1}^N \tau_{ijk}} \sum_{k=1}^N \ln(c_k) \tau_{ijk} \quad (11)$$

where i and j are the indices of grid, k the index of trajectory, N the total number of trajectories used in analysis, c_k the pollutant concentration measured upon arrival of trajectory k , and τ_{ijk} the residence time of trajectory k in grid cell (i, j) . A high value of \bar{C}_{ij} means that, air parcels passing over cell (i, j) would, on average, cause high concentrations at the receptor site.

Figure 26.9 shows the situation for PM_{2.5} concentrations. It was calculated by recording the associated PM_{2.5} concentration for each point on the back trajectory based on the arrival time concentration using 2010 data. The plot shows the geographic areas most strongly associated with high PM_{2.5} concentrations i.e. to the east in continental Europe. Both the CWT and PSCF methods have been shown to give similar results and each have their advantages and disadvantages (Lupu and Maenhaut 2002; Hsu et al. 2003). Figure 26.9 can be compared with Figure 26.8 to compare the overall identification of source regions using the CWT and PSCF techniques. Overall the agreement is good in that similar geographic locations are identified as being important for PM_{2.5}.

Figure 26.9 is useful, but it can be clearer if the trajectory surface is smoothed, which has been done for PM_{2.5} concentrations shown in Figure 26.10.

In common with most other **openair** functions, the flexible ‘type’ option can be used to split the data in different ways. For example, to plot the smoothed back trajectories for PM_{2.5} concentrations by season.

```
trajLevel(subset(traj, lon > -20 & lon < 20 & lat > 45 & lat < 60),
          pollutant = "pm2.5", statistic="cwt", col = "increment",
          border = "white")
```

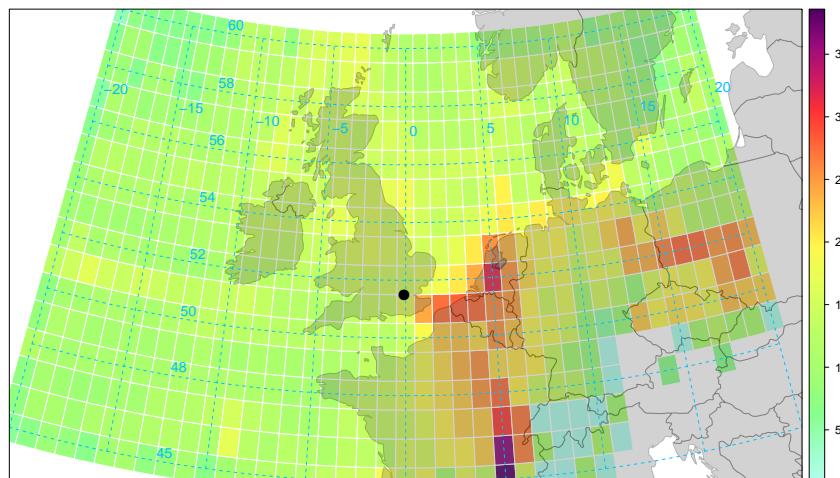


FIGURE 26.9 Gridded back trajectory concentrations showing mean PM_{2.5} concentrations using the CWT approach.

```
trajLevel(subset(traj, lat > 45 & lat < 60 & lon >-20 & lon <20),
          pollutant ="pm2.5", statistic = "cwt", smooth = TRUE,
          col = "increment")
```

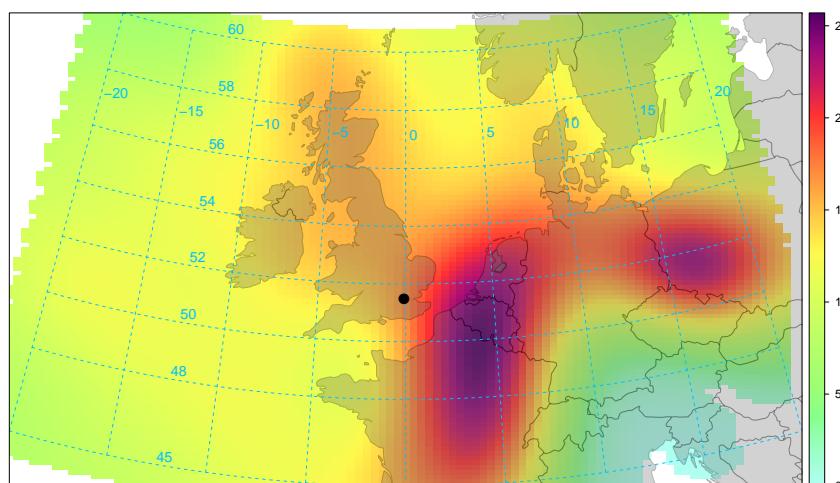


FIGURE 26.10 Gridded and smoothed back trajectory concentrations showing mean PM_{2.5} concentrations using the CWT approach.

```
trajLevel(subset(traj, lat > 40 & lat < 70 & lon >-20 & lon <20),
pollutant = "pm2.5", type = "season", statistic = "pscf",
layout = c(4, 1))
```

It should be noted that it makes sense to analyse back trajectories for pollutants that have a large regional component — such as particles or O₃. It makes little sense to analyse pollutants that are known to have local impacts e.g. NO_x. However, a species such as NO_x can be helpful to exclude ‘fresh’ emissions from the analysis.

26.3 Back trajectory cluster analysis with the **trajCluster** function

Often it is useful to use cluster analysis on back trajectories to group similar air mass origins together. The principal purpose of clustering back trajectories is to post-process data according to cluster origin. By grouping data with similar geographic origins it is possible to gain information on pollutant species with similar chemical histories. There are several ways in which clustering can be carried out and several measures of the similarity of different clusters. A key issue is how the *distance matrix* is calculated, which determines the similarity (or dissimilarity) of different back trajectories. The simplest measure is the Euclidean distance. However, an angle-based measure is also often used. The two distance measures are defined below. In **openair** the distance matrices are calculated using C++ code because their calculation is computationally intensive. Note that these calculations can also be performed directly in the HYSPLIT model itself.

The Euclidean distance between two trajectories is given by [Equation 12](#). Where X₁, Y₁ and X₂, Y₂ are the latitude and longitude coordinates of back trajectories 1 and 2, respectively. n is the number of back trajectory points (96 hours in this case).

$$d_{1,2} = \left(\sum_{i=1}^n ((X_{1i} - X_{2i})^2 + (Y_{1i} - Y_{2i})^2) \right)^{1/2} \quad (12)$$

The *angle* distance matrix is a measure of how similar two back trajectory points are in terms of their angle from the origin i.e. the starting location of the back trajectories. The angle-based measure will often capture some of the important circulatory features in the atmosphere e.g. situations where there is a high pressure located to the east of the UK. However, the most appropriate distance measure will be application dependent and is probably best tested by the extent to which they are able to differentiate different air-mass characteristics, which can be tested through post-processing. The angle-based distance measure is defined as:

$$d_{1,2} = \frac{1}{n} \sum_{i=1}^n \cos^{-1} \left(0.5 \frac{A_i + B_i + C_i}{\sqrt{A_i B_i}} \right) \quad (13)$$

where

$$A_i = (X_1(i) - X_0)^2 + (Y_1(i) - Y_0)^2 \quad (14)$$

$$B_i = (X_2(i) - X_0)^2 + (Y_2(i) - Y_0)^2 \quad (15)$$

$$C_i = (X_2(i) - X_1(i))^2 + (Y_2(i) - Y_1(i))^2 \quad (16)$$

where X₀ and Y₀ are the coordinates of the location being studied i.e. the starting location of the trajectories.

The **trajCluster** function has the following options:

traj	An openair trajectory data frame resulting from the use of <code>importTraj</code> .
method	Method used to calculate the distance matrix for the back trajectories. There are two methods available: “Euclid” and “Angle”.
n.cluster	Number of clusters to calculate.
plot	Should a plot be produced?
type	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. “season”, “year”, “weekday” and so on. For example, <code>type = "season"</code> will produce four plots — one for each season. Note that the cluster calculations are separately made of each level of “type”.
cols	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “jet” and <code>RColorBrewer</code> colours — see the <code>openair openColours</code> function for more details. For user defined the user can supply a list of colour names recognised by R (type <code>colours()</code> to see the full list). An example would be <code>cols = c("yellow", "green", "blue")</code>
split.after	For <code>type</code> other than “default” e.g. “season”, the trajectories can either be calculated for each level of <code>type</code> independently or extracted after the cluster calculations have been applied to the whole data set.
map.fill	Should the base map be a filled polygon? Default is to fill countries.
map.cols	If <code>map.fill = TRUE</code> <code>map.cols</code> controls the fill colour. Examples include <code>map.fill = "grey40"</code> and <code>map.fill = openColours("default", 10)</code> . The latter colours the countries and can help differentiate them.
map.alpha	The transparency level of the filled map which takes values from 0 (full transparency) to 1 (full opacity). Setting it below 1 can help view trajectories, trajectory surfaces etc. and a filled base map.
projection	The map projection to be used. Different map projections are possible through the <code>mapproj</code> package. See <code>?mapproject</code> for extensive details and information on setting other parameters and orientation (see below).
parameters	From the <code>mapproj</code> package. Optional numeric vector of parameters for use with the projection argument. This argument is optional only in the sense that certain projections do not require additional parameters. If a projection does require additional parameters, these must be given in the parameters argument.
orientation	From the <code>mapproj</code> package. An optional vector <code>c(latitude,longitude,rotation)</code> which describes where the “North Pole” should be when computing the projection. Normally this is <code>c(90,0)</code> , which is appropriate for cylindrical and conic projections. For a planar projection, you should set it to the desired point of tangency. The third value is a clockwise rotation (in degrees), which defaults to the midrange of the longitude coordinates in the map.
by.type	The percentage of the total number of trajectories is given for all data by default. Setting <code>by.type = TRUE</code> will make each panel add up to 100.

```
clust <- trajCluster(traj, method = "Angle", n.cluster = 6, col = "Set2",
                      map.cols = openColours("Paired", 10))
```

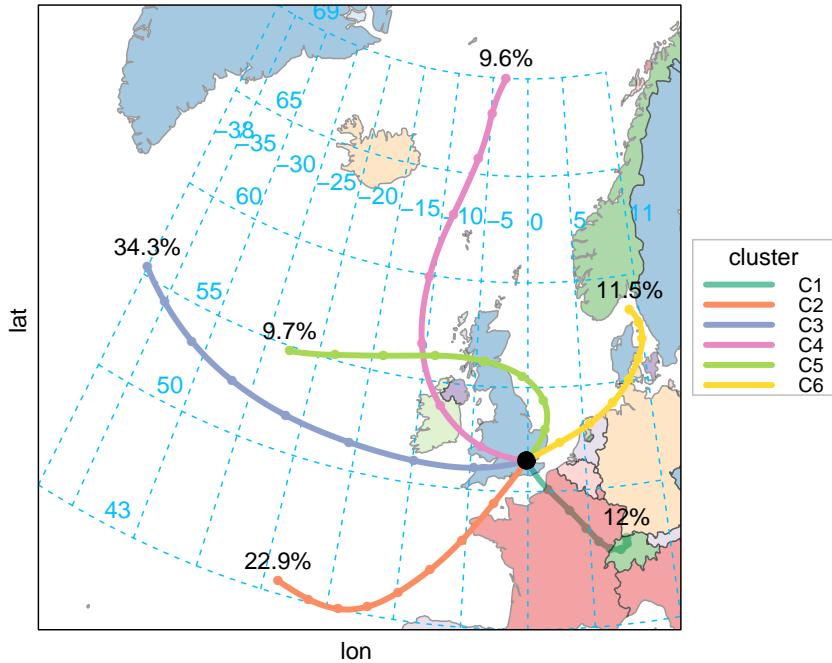


FIGURE 26.11 The 6-cluster solution to back trajectories calculated for the London North Kensington site for 2011 showing the mean trajectory for each cluster.

origin If **TRUE** a filled circle dot is shown to mark the receptor point.

... Other graphical parameters passed onto **lattice::levelplot** and **cutData**. Similarly, common axis and title labelling options (such as **xlab**, **ylab**, **main**) are passed to **levelplot** via **quickText** to handle routine formatting.

As an example we will consider back trajectories for London in 2011.

First, the back trajectory data for London is imported together with the air pollution data for the North Kensington site (KC1).

```
traj <- importTraj(site = "london", year = 2011)
kc1 <- importKCL(site = "kc1", year = 2011)
```

The clusters are straightforward to calculate. In this case the back trajectory data (**traj**) is supplied and the angle-based distance matrix is used. Furthermore, we choose to calculate 6 clusters and choose a specific colour scheme. In this case we read the output from **trajCluster** into a variable **clust** so that the results can be post-processed.

clust returns all the back trajectory information together with the cluster (as a character). This data can now be used together with other data to analyse results further. However, first it is possible to show all trajectories coloured by cluster, although for a year of data there is significant overlap and it is difficult to tell them apart.

```
trajPlot(clust$data, group = "cluster")
```

A useful way in which to see where these air masses come from by trajectory is to produce a frequency plot by cluster. Such a plot (not shown, but code below) provides a good indication of the spread of the different trajectory clusters as well as providing an indication of where air masses spend most of their time. For the London 2011 data it can be seen cluster 1 is dominated by air from the European mainland to the south.

```
trajLevel(clust$data, type = "cluster", col = "increment", border = NA)
```

Perhaps more useful is to merge the cluster data with measurement data. In this case the data at North Kensington site are used. Note that in merging these two data frames it is not necessary to retain all 96 back trajectory hours and for this reason we extract only the first hour.

```
kc1 <- merge(kc1, subset(clust$data, hour.inc == 0), by = "date")
```

Now **kc1** contains air pollution data identified by cluster. The size of this data frame is about a third of the original size because back trajectories are only run every 3 hours.

The numbers of each cluster are given by:

```
table(kc1[["cluster"]])
```

```
##  
##   C1   C2   C3   C4   C5   C6  
## 347 661 989 277 280 333
```

i.e. is dominated by clusters 3 and 2 from west and south-west (Atlantic).

Now it is possible to analyse the concentration data according to the cluster. There are numerous types of analysis that can be carried out with these results, which will depend on what the aims of the analysis are in the first place. However, perhaps one of the first things to consider is how the concentrations vary by cluster. As the summary results below show, there are distinctly different mean concentrations of most pollutants by cluster. For example, clusters 1 and 6 are associated with much higher concentrations of PM₁₀—approximately double that of other clusters. Both of these clusters originate from continental Europe. Cluster 5 is also relatively high, which tends to come from the rest of the UK. Other clues concerning the types of air-mass can be gained from the mean pressure. For example, cluster 5 is associated with the highest pressure (1014 kPa), and as is seen in [Figure 26.11](#) the shape of the line for cluster 5 is consistent with air-masses associated with a high pressure system (a clockwise-type sweep).

```
library(dplyr)  
  
group_by(kc1, cluster) %>%  
  summarise_if(is.numeric, mean, na.rm = TRUE)  
  
## # A tibble: 6 x 23  
##   cluster    nox    no2     o3    so2     co pm10_raw pm10    pm25    v2.5 nv2.5  
##   <chr>    <dbl>  <dbl>  <dbl>  <dbl>  <dbl>    <dbl>  <dbl>  <dbl>  <dbl>  <dbl>  
## 1 C1        89.9  51.4  32.9  3.34  0.388    32.3  35.8  31.3  8.02 23.3  
## 2 C2        40.9  30.9  39.6  1.25  0.223    18.1  17.7  11.6  3.13 8.42  
## 3 C3        47.2  32.3  39.8  1.72  0.189    19.4  17.9  11.3  2.80 8.46  
## 4 C4        44.1  30.6  40.1  1.73  0.188    19.8  17.6  11.0  2.17 8.82  
## 5 C5        57.0  39.2  41.4  2.12  0.226    24.4  25.5  16.5  4.50 12.0  
## 6 C6        64.9  42.4  46.4  2.84  0.290    31.7  35.7  29.7  7.57 22.1  
## # ... with 12 more variables: receptor <dbl>, year <dbl>, month <dbl>,  
## #   day <dbl>, hour <dbl>, hour.inc <dbl>, lat <dbl>, lon <dbl>, height <dbl>,  
## #   pressure <dbl>, traj_len <dbl>, len <dbl>
```

```
trendLevel(kc1, pollutant = "v2.5", type = "cluster", layout = c(6, 1),
           cols = "increment")
```

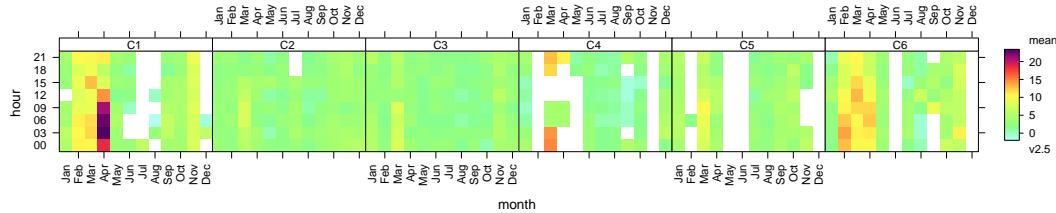


FIGURE 26.12 Some of the temporal characteristics of the volatile PM_{2.5} component plotted by month and hour of the day and by cluster for the London North Kensington site for 2011.

```
timeProp(kc1, pollutant="pm25", avg.time = "day", proportion = "cluster",
         col="Set2", key.position = "top", key.columns = 6)
```

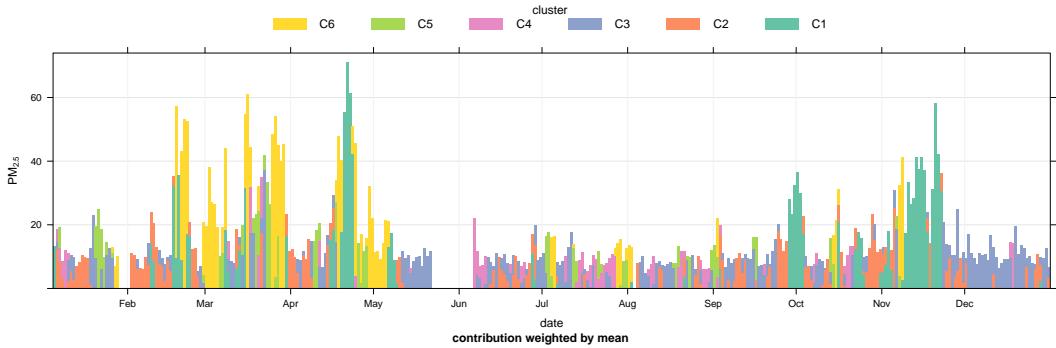


FIGURE 26.13 Temporal variation in daily PM_{2.5} concentrations at the North Kensington site show by contribution of each cluster.

Simple plots can be generated from these results too. For example, it is easy to consider the temporal nature of the volatile component of PM_{2.5} concentrations (`v2.5` in the `kc1` data frame). Figure 26.12 for example shows how the concentration of the volatile component of PM_{2.5} concentrations varies by cluster by plotting the hour of day-month variation. It is clear from Figure 26.12 that the clusters associated with the highest volatile PM_{2.5} concentrations are clusters 1 and 6 (European origin) and that these concentrations peak during spring. There is less data to see clearly what is going on with cluster 5. Nevertheless, the cluster analysis has clearly separated different air mass characteristics which allows for more refined analysis of different air-mass types.

Similarly, as considered in Section 15, the `timeVariation` function can also be used to consider the temporal components.

Another useful plot to consider is `timeProp`, which can show how the concentration of a pollutant is comprised. In this case it is useful to plot the time series of PM_{2.5} and show how much of the concentration is contributed to by each cluster. Such a plot is shown in Figure 26.13. It is now easy to see for example that during the spring months many of the high concentration events were due to clusters 1 and 6, which correspond to European origin air-masses as shown in Figure 26.11.

27 Model evaluation — the **modStats** function

27.1 Purpose

The **modStats** function provides key model evaluation statistics for comparing models against measurements and models against other models.

There are a very wide range of evaluation statistics that can be used to assess model performance. There is, however, no single statistic that encapsulates all aspects of interest. For this reason it is useful to consider several performance statistics and also to understand the sort of information or insight they might provide.

In the following definitions, O_i represents the i th observed value and M_i represents the i th modelled value for a total of n observations.

Fraction of predictions within a factor or two, FAC2

The fraction of modelled values within a factor of two of the observed values are the fraction of model predictions that satisfy:

$$0.5 \leq \frac{M_i}{O_i} \leq 2.0 \quad (17)$$

Mean bias, MB

The mean bias provides a good indication of the mean over or under estimate of predictions. Mean bias in the same units as the quantities being considered.

$$MB = \frac{1}{n} \sum_{i=1}^N M_i - O_i \quad (18)$$

Mean Gross Error, MGE

The mean gross error provides a good indication of the mean error regardless of whether it is an over or under estimate. Mean gross error is in the same units as the quantities being considered.

$$MGE = \frac{1}{n} \sum_{i=1}^N |M_i - O_i| \quad (19)$$

Normalised mean bias, NMB

The normalised mean bias is useful for comparing pollutants that cover different concentration scales and the mean bias is normalised by dividing by the observed concentration.

$$NMB = \frac{\sum_{i=1}^n M_i - O_i}{\sum_{i=1}^n O_i} \quad (20)$$

Normalised mean gross error, NMGE

The normalised mean gross error further ignores whether a prediction is an over or under estimate.

$$NMGE = \frac{\sum_{i=1}^n |M_i - O_i|}{\sum_{i=1}^n O_i} \quad (21)$$

Root mean squared error, RMSE

The RMSE is a commonly used statistic that provides a good overall measure of how close modelled values are to predicted values.

$$RMSE = \left(\frac{\sum_{i=1}^n (M_i - O_i)^2}{n} \right)^{1/2} \quad (22)$$

Correlation coefficient, r

The (Pearson) correlation coefficient is a measure of the strength of the linear relationship between two variables. If there is perfect linear relationship with positive slope between the two variables, $r = 1$. If there is a perfect linear relationship with negative slope between the two variables $r = -1$. A correlation coefficient of 0 means that there is no linear relationship between the variables. Note that `modStats` accepts an option `method`, which can be set to “kendall” and “spearman” for alternative calculations of r .

$$r = \frac{1}{(n-1)} \sum_{i=1}^n \left(\frac{M_i - \bar{M}}{\sigma_M} \right) \left(\frac{O_i - \bar{O}}{\sigma_O} \right) \quad (23)$$

Coefficient of Efficiency, COE

The *Coefficient of Efficiency* based on Legates and McCabe (2012) and Legates and McCabe Jr (1999). There have been many suggestions for measuring model performance over the years, but the COE is a simple formulation which is easy to interpret.

A perfect model has a COE = 1. As noted by Legates and McCabe although the COE has no lower bound, a value of COE = 0.0 has a fundamental meaning. It implies that the model is no more able to predict the observed values than does the observed mean. Therefore, since the model can explain no more of the variation in the observed values than can the observed mean, such a model can have no predictive advantage.

For negative values of COE, the model is less effective than the observed mean in predicting the variation in the observations.

$$COE = 1.0 - \frac{\sum_{i=1}^n |M_i - O_i|}{\sum_{i=1}^n |O_i - \bar{O}|} \quad (24)$$

Index of Agreement, IOA

The *Index of Agreement*, IOA is commonly used in model evaluation (Willmott et al. 2011). It spans between -1 and $+1$ with values approaching $+1$ representing better model performance. An IOA of 0.5, for example, indicates that the sum of the error-magnitudes is one half of the sum of the observed-deviation magnitudes. When IOA =

0.0, it signifies that the sum of the magnitudes of the errors and the sum of the observed-deviation magnitudes are equivalent. When $\text{IOA} = -0.5$, it indicates that the sum of the error-magnitudes is twice the sum of the perfect model-deviation and observed-deviation magnitudes. Values of IOA near -1.0 can mean that the model-estimated deviations about O are poor estimates of the observed deviations; but, they also can mean that there simply is little observed variability — so some caution is needed when the IOA approaches -1 . It is defined as (with $c = 2$):

$$\text{IOA} = \begin{cases} 1.0 - \frac{\sum_{i=1}^n |M_i - O_i|}{c \sum_{i=1}^n |O_i - \bar{O}|}, & \text{when } \sum_{i=1}^n |M_i - O_i| \leq c \sum_{i=1}^n |O_i - \bar{O}| \\ \frac{\sum_{i=1}^n |O_i - \bar{O}|}{\sum_{i=1}^n |M_i - O_i|} - 1.0, & \text{when } \sum_{i=1}^n |M_i - O_i| > c \sum_{i=1}^n |O_i - \bar{O}| \end{cases} \quad (25)$$

27.2 Options available

mydata A data frame.

mod Name of a variable in `mydata` that represents modelled values.

obs Name of a variable in `mydata` that represents measured values.

statistic The statistic to be calculated. See details below for a description of each.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce statistics using the entire data. `type` can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four sets of statistics — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

More than one type can be considered e.g. `type = c("season", "weekday")` will produce statistics split by season and day of the week.

rank.name Simple model ranking can be carried out if `rank.name` is supplied. `rank.name` will generally refer to a column representing a model name, which is to be ranked. The ranking is based the COE performance, as that indicator is arguably the best single model performance indicator available.

... Other arguments to be passed to `cutData` e.g. `hemisphere = "southern"`

27.3 Example of use

The function can be called very simply and only requires two numeric fields to compare. To show how the function works, some synthetic data will be generated for 5 models.

```
## observations; 100 random numbers
set.seed(10)
obs <- 100 * runif(100)
mod1 <- data.frame(obs, mod = obs + 10, model = "model 1")
mod2 <- data.frame(obs, mod = obs + 20 * rnorm(100), model = "model 2")
mod3 <- data.frame(obs, mod = obs - 10 * rnorm(100), model = "model 3")
mod4 <- data.frame(obs, mod = obs / 2 + 10 * rnorm(100), model = "model 4")
mod5 <- data.frame(obs, mod = obs * 1.5 + 3 * rnorm(100), model = "model 5")
modData <- rbind(mod1, mod2, mod3, mod4, mod5)
head(modData)

##      obs mod   model
## 1 50.7 61 model 1
## 2 30.7 41 model 1
## 3 42.7 53 model 1
## 4 69.3 79 model 1
## 5  8.5 19 model 1
## 6 22.5 33 model 1
```

We now have a data frame with observations and predictions for 5 models. The evaluation of the statistics is given by:

```
modStats(modData, obs = "obs", mod = "mod", type = "model")

##      model n FAC2     MB    MGE     NMB NMGE RMSE     r     COE    IOA
## 1 model 1 100 0.89 10.00 10.0  0.225 0.22 10.0 1.00  0.5383 0.77
## 2 model 2 100 0.79  0.92 16.6  0.021 0.37 19.3 0.83  0.2339 0.62
## 3 model 3 100 0.88  1.01  7.9  0.023 0.18  9.5 0.94  0.6358 0.82
## 4 model 4 100 0.56 -20.60 21.9 -0.463 0.49 25.8 0.81 -0.0093 0.50
## 5 model 5 100 0.96 22.52 22.6  0.506 0.51 26.1 1.00 -0.0420 0.48
```

It is possible to rank the statistics based on the *Coefficient of Efficiency*, which is a good general indicator of model performance.

```
modStats(modData, obs = "obs", mod = "mod", type = "model", rank.name = "model")

##      model n FAC2     MB    MGE     NMB NMGE RMSE     r     COE    IOA
## 3 model 3 100 0.88  1.01  7.9  0.023 0.18  9.5 0.94  0.6358 0.82
## 1 model 1 100 0.89 10.00 10.0  0.225 0.22 10.0 1.00  0.5383 0.77
## 2 model 2 100 0.79  0.92 16.6  0.021 0.37 19.3 0.83  0.2339 0.62
## 4 model 4 100 0.56 -20.60 21.9 -0.463 0.49 25.8 0.81 -0.0093 0.50
## 5 model 5 100 0.96 22.52 22.6  0.506 0.51 26.1 1.00 -0.0420 0.48
```

The `modStats` function is however much more flexible than indicated above. While it is useful to calculate model evaluation statistics in a straightforward way it can be much more informative to consider the statistics split by different periods.

Data have been assembled from a Defra model evaluation exercise which consists of hourly O₃ predictions at 15 receptor points around the UK for 2006. The aim here is not to identify a particular model that is ‘best’ and for this reason the models are simply referred to as ‘model 1’, ‘model 2’ and so on. We will aim to make the data more widely available. However, data set has this form:

```
load("~/My Drive/openair/Data/modelData.RData")
head(modTest)

##           site          date o3    mod   group
## 1 Aston.Hill 2006-01-01 00:00:00 NA     NA model 1
## 2 Aston.Hill 2006-01-01 01:00:00 74 65.28 model 1
## 3 Aston.Hill 2006-01-01 02:00:00 72 64.64 model 1
## 4 Aston.Hill 2006-01-01 03:00:00 72 64.46 model 1
## 5 Aston.Hill 2006-01-01 04:00:00 70 64.88 model 1
## 6 Aston.Hill 2006-01-01 05:00:00 66 65.80 model 1
```

There are columns representing the receptor location (`site`), the date, measured values (`o3`), model predictions (`mod`) and the model itself (`group`). There are numerous ways in which the statistics can be calculated. However, of interest here is how the models perform at a single receptor by season. The seasonal nature of O₃ is a very important characteristic and it is worth considering in more detail. The statistics are easy enough to calculate as shown below. In this example a subset of the data is selected to consider only the Harwell site. Second, the `type` option is used to split the calculations by season and model. Finally the statistics are grouped by the IOA for each season. It is now very easy how model performance changes by season and which models perform best in each season.

```
options(digits = 2) ## don't display too many decimal places
modStats(subset(modTest, site == "Harwell"), obs = "o3", mod = "mod",
         type = c("season", "group"), rank = "group")

## # A tibble: 36 x 12
## # Groups:   season [4]
##   season group     n  FAC2     MB     MGE     NMB     NMGE    RMSE      r      COE
##   <ord> <fct> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 sprin~ mode~  1905  0.875    8.58  16.5  0.137  0.263  21.9  0.576  0.0747
## 2 sprin~ mode~  1905  0.876   -2.04  16.8 -0.0325  0.268  21.8  0.452  0.0584
## 3 sprin~ mode~  1905  0.872   11.9   17.7  0.190  0.282  23.7  0.519  0.00818
## 4 sprin~ mode~  1905  0.885    4.98  18.1  0.0793  0.289  23.9  0.361 -0.0151
## 5 sprin~ mode~  1905  0.870   10.8   19.5  0.172  0.311  24.0  0.588 -0.0940
## 6 sprin~ mode~  1825  0.821    7.79  20.0  0.123  0.317  25.8  0.480 -0.118
## 7 sprin~ mode~  1905  0.786   -13.7  21.4 -0.218  0.340  26.1  0.531 -0.195
## 8 sprin~ mode~  1905  0.732   -13.4  24.3 -0.213  0.388  29.8  0.312 -0.362
## 9 sprin~ mode~  1905  0.743     1.01  24.9  0.0161  0.396  32.8  0.264 -0.392
## 10 summe~ mode~ 2002  0.918    3.83  15.5  0.0632  0.256  20.9  0.750  0.317
## # ... with 26 more rows, and 1 more variable: IOA <dbl>
```

Note that it is possible to read the results of the `modStats` function into a data frame, which then allows the results to be plotted. This is generally a good idea when there is a lot of numeric data to consider and plots will convey the information better.

The `modStats` function is much more flexible than indicated above and can be used in lots of interesting ways. The `type` option in particular makes it possible to split the statistics in numerous ways. For example, to summarise the performance of models by site, model and day of the week:

```
modStats(modStats, obs = "o3", mod = "mod",
         type = c("site", "weekday", "group"), rank = "group")
```

Similarly, if other data are available e.g. meteorological data or other pollutant species then these variables can also be used to test models against ranges in their values. This capability is potentially very useful because it allows for a much more probing analysis into model evaluation. For example, with wind speed and direction it is easy to consider how model performance varies by wind speed intervals or wind sectors, both of which

could reveal important performance characteristics.

28 Model evaluation — the **TaylorDiagram** function

28.1 Purpose

The *Taylor Diagram* is one of the more useful methods for evaluating model performance. Details of the diagram can be found at http://www-pcmdi.llnl.gov/about/staff/Taylor/CV/Taylor_diagram_primer.pdf and in Taylor (2001). The diagram provides a way of showing how three complementary model performance statistics vary simultaneously. These statistics are the correlation coefficient R, the standard deviation (sigma) and the (centred) root-mean-square error. These three statistics can be plotted on one (2D) graph because of the way they are related to one another which can be represented through the Law of Cosines.

The **openair** version of the Taylor Diagram has several enhancements that increase its flexibility. In particular, the straightforward way of producing conditioning plots should prove valuable under many circumstances (using the **type** option). Many examples of Taylor Diagrams focus on model-observation comparisons for several models using all the available data. However, more insight can be gained into model performance by partitioning the data in various ways e.g. by season, daylight/nighttime, day of the week, by levels of a numeric variable e.g. wind speed or by land-use type etc.

We first show a diagram and then pick apart the different components to understand how to interpret it. The diagram can look overly complex but once it is understood how to interpret the three main characteristics it becomes much easier to understand. A typical diagram is shown in [Figure 28.1](#) for nine anonymised models used for predicting hourly O₃ concentrations at 15 sites around the UK.

The plots shown in [Figure 28.2](#) break the Taylor Diagrams into three components to aid interpretation. The first plot (top left) highlights the comparison of variability in for each model compared with the measurements. The variability is represented by the standard deviation of the observed and modelled values. The plot shows that the observed variability (given by the standard deviation) is about 27 ($\mu\text{g m}^{-3}$) and is marked as ‘observed’ on the x-axis. The magnitude of the variability is measured as the **radial** distance from the origin of the plot (the red line with the arrow shows the standard deviation for model g, which is about 25 $\mu\text{g m}^{-3}$). To aid interpretation the radial dashed line is shown from the ‘observed’ point. Each model is shown in this case by the position of the letters a to i. On this basis it can be seen that models 1, a, b have more variability than the measurements (because they extend beyond the dashed line), whereas the others have less variability than the measurements. Models a and b are also closest to the dashed line and therefore have the closest variability compared with the observations.

The next statistic to consider is the correlation coefficient, R shown by the top-right Figure in [Figure 28.2](#). This is shown on the arc and points that lie closest to the x-axis have the highest correlation. The grey lines help to show this specific correlation coefficients. The red arc shows R=0.7 for model g. The best performing models with the highest R are models b and g with correlation coefficients around 0.7. Two models stand out as having much worse correlations with the observations: models e and i (values of around 0.4).

Finally, the lower plot in [Figure 28.2](#) highlights the *centred* root-mean square error (RMS). It is centred because the mean values of the data (observations and predictions) are subtracted first. The concentric dashed lines emanating from the ‘observed’ point show the value of the RMS error — so points furthest from the ‘observed’ value are

```
TaylorDiagram(modTest, obs = "o3", mod = "mod", group = "group")
```

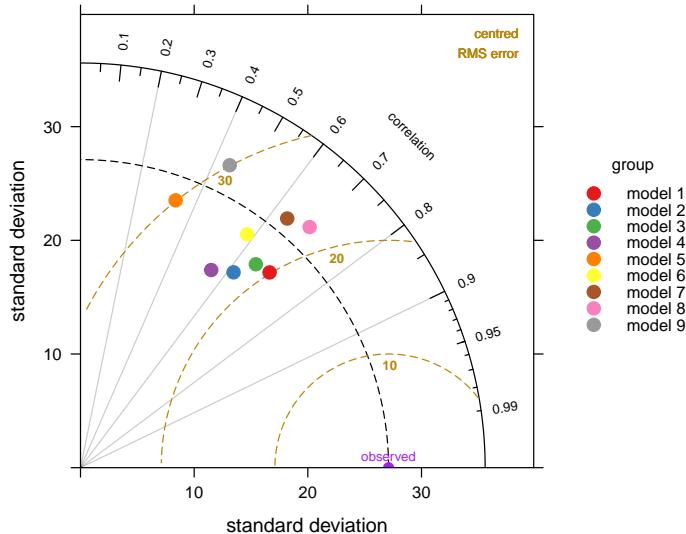


FIGURE 28.1 An example of the use of the `TaylorDiagram` function.

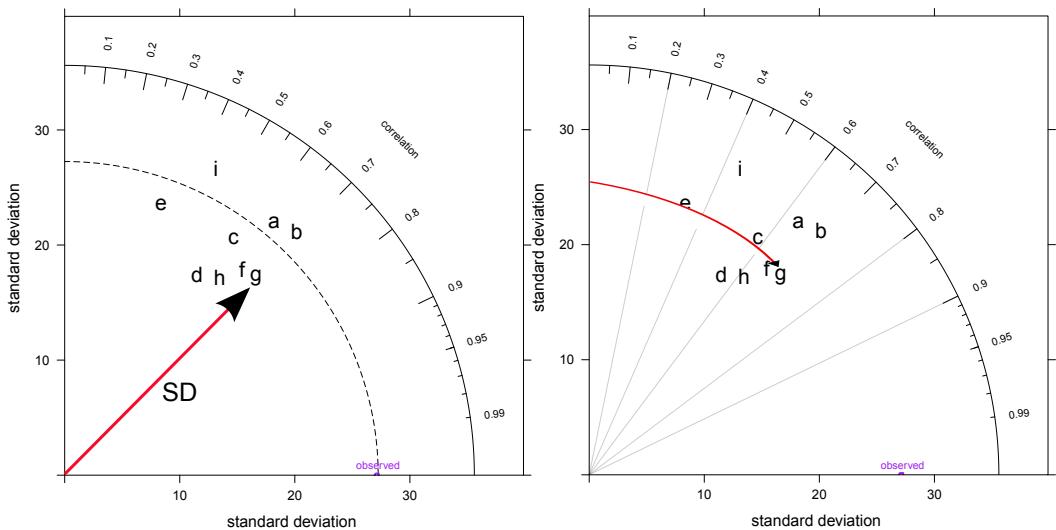
the worst performing models because they have the highest RMS errors. On this basis, model g has the lowest error of about $20 \mu\text{g m}^{-3}$, shown again by the red line. Models e and i are considerably worse because they have RMS errors of around $30 \mu\text{g m}^{-3}$.

So which model is best? Taken as a whole it is probably model g because it has reasonably similar variability compared with the observations, the highest correlation and the least RMS error. However, models f and b also look to be good. Perhaps it is easier to conclude that models e and i are not good

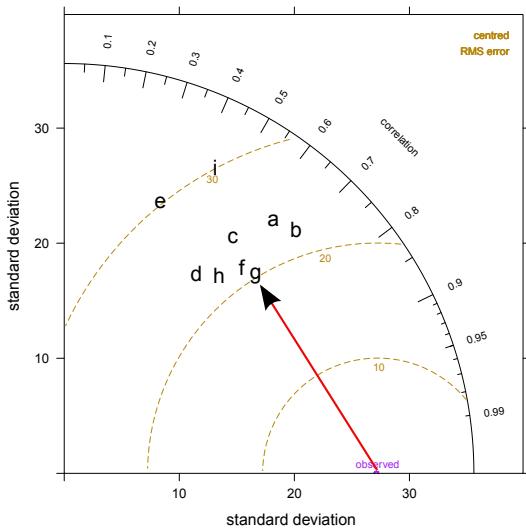
Note that in cases where there is a column ‘site’ it makes sense to use `type = "site"` to ensure that the statistics are calculated on a per site basis and each panel represents a single site.

28.2 Options available

- mydata** A data frame minimally containing a column of observations and a column of predictions.
- obs** A column of observations with which the predictions (`mod`) will be compared.
- mod** A column of model predictions. Note, `mod` can be of length 2 i.e. two lots of model predictions. If two sets of predictions are present e.g. `mod = c("base", "revised")`, then arrows are shown on the Taylor Diagram which show the change in model performance in going from the first to the second. This is useful where, for example, there is interest in comparing how one model run compares with another using different assumptions e.g. input data or model set up. See examples below.
- group** The `group` column is used to differentiate between different models and can be a factor or character. The total number of models compared will be equal to the number of unique values of `group`.



(A) Taylor Diagram highlighting the variation in standard deviation.
(B) Taylor Diagram highlighting the variation in correlation coefficient.



(C) Taylor Diagram highlighting the variation in the centred RMS error.

FIGURE 28.2 Taylor Diagrams broken down to highlight how to interpret the three main statistics. The red line/arrow indicate how to read interpret each of the three statistics.

`group` can also be of length two e.g. `group = c("model", "site")`. In this case all model-site combinations will be shown but they will only be differentiated by colour/symbol by the first grouping variable ("model" in this case). In essence the plot removes the differentiation by the second grouping variable. Because there will be different values of `obs` for each group, `normalise = TRUE` should be used.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. "season", "year", "weekday" and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If

that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

Note that often it will make sense to use `type = "site"` when multiple sites are available. This will ensure that each panel contains data specific to an individual site.

normalise Should the data be normalised by dividing the standard deviation of the observations? The statistics can be normalised (and non-dimensionalised) by dividing both the RMS difference and the standard deviation of the `mod` values by the standard deviation of the observations (`obs`). In this case the “observed” point is plotted on the x-axis at unit distance from the origin. This makes it possible to plot statistics for different species (maybe with different units) on the same plot. The normalisation is done by each `group/type` combination.

cols Colours to be used for plotting. Useful options for categorical data are available from `RColorBrewer` colours — see the `openair openColours` function for more details. Useful schemes include “Accent”, “Dark2”, “Paired”, “Pastel1”, “Pastel2”, “Set1”, “Set2”, “Set3” — but see `?brewer.pal` for the maximum useful colours in each. For user defined the user can supply a list of colour names recognised by R (type `colours()` to see the full list). An example would be `cols = c("yellow", "green", "blue")`.

rms.col Colour for centred-RMS lines and text.

cor.col Colour for correlation coefficient lines and text.

arrow.lwd Width of arrow used when used for comparing two model outputs.

annotate Annotation shown for RMS error.

key Should the key be shown?

key.title Title for the key.

key.columns Number of columns to be used in the key. With many pollutants a single column can make to key too wide. The user can thus choose to use several columns by setting `columns` to be less than the number of pollutants.

key.pos Position of the key e.g. “top”, “bottom”, “left” and “right”. See details in `lattice:xyplot` for more details about finer control.

strip Should a strip be shown?

auto.text Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO₂.

... Other graphical parameters are passed onto `cutData` and `lattice:xyplot`. For example, `TaylorDiagram` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common graphical parameters, such as `layout` for panel arrangement and `pch` and `cex` for plot symbol type and size, are passed on to `xyplot`. Most are passed unmodified, although there are some special cases where `openair` may locally manage this process. For example, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed via `quickText` to handle routine formatting.

28.3 Example of use

The example used here carries on from the previous section using data from a Defra model evaluation exercise. As mentioned previously, the use of the `type` option offers enormous flexibility for comparing models. However, we will only focus on the seasonal evaluation of the models. In the call below, `group` is the column that identified the model and `type` is the conditioning variable that produces in this case four panels — one for each season. Note that in this case we focus on a single site.

Figure 28.3 contains a lot of useful information. Consider the summertime comparison first. All models tend to underestimate the variability of O₃ concentrations because they all lie within the black dashed line. However, models 7 and 9 are close to the observed variability. The general underestimate of the variability for summertime conditions might reflect that the models do not adequately capture regional O₃ episodes when concentrations are high. Models 7 and 8 do best in terms of high correlation with the measurements (around 0.8) and lowest RMS error (around 20–22 µg m⁻³). Models 3, 5 and 6 tend to do worse on all three statistics during the summer.

By contrast, during wintertime conditions models 1 and 3 are clearly best. From an evaluation perspective it would be useful to understand why some models are better for wintertime conditions and others better in summer and this is clearly something that could be investigated further.

There are many other useful comparisons that can be undertaken easily. A few of these are shown below, but not plotted.

```
## by receptor comparison
TaylorDiagram(modTest, obs = "o3", mod = "mod", group = "group", type = "site")

## by month comparison for a SINGLE site
TaylorDiagram(subset(modTest, site == "Harwell"), obs = "o3", mod = "mod",
              group = "group", type = "month")

## By season AND daylight/nighttime
TaylorDiagram(subset(modTest, site == "Harwell"), obs = "o3", mod = "mod",
              group = "group", type = c("season", "daylight"))
```

It is also possible to combine different groups of model results. For example, rather than plot how the models perform at a single site it can be useful to show how they compare at all sites. To do this it is necessary to normalise the data because there will be different values of the observed variable across different sites. In this case we can supply the option `group = c("group", "site")`. This will show the variation by model for *all* sites. The results are shown in Figure 28.4. These results show that in general models tend to predict in similarly good (or bad) ways across all sites as shown by the grouping of points on Figure 28.4.

```
## select a single site
LH <- subset(modTest, site == "Lullington.Heath")
TaylorDiagram(LH, obs = "o3", mod = "mod", group = "group", type = "season")
```

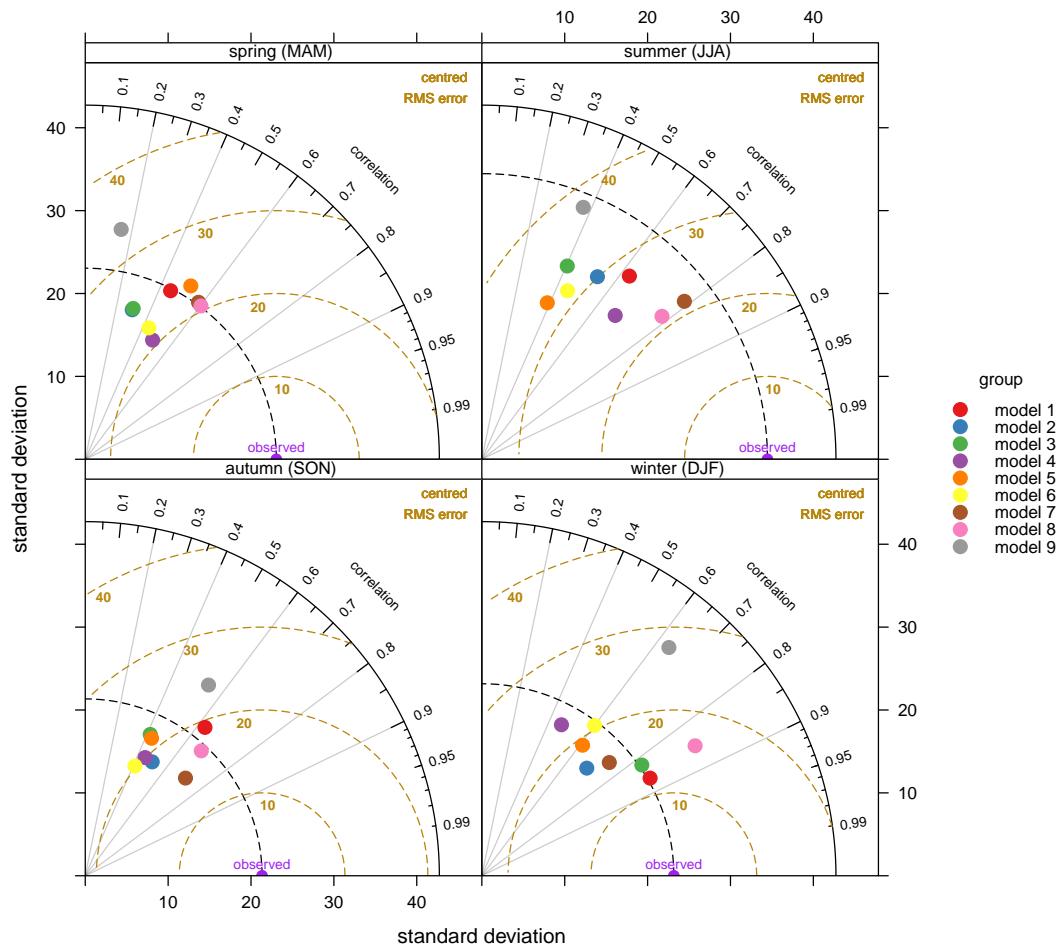


FIGURE 28.3 Use of the `TaylorDiagram` function to show model performance for 9 models used to predict O₃ concentrations at the Lullington Heath site.

```
TaylorDiagram(modTest, obs = "o3", mod = "mod", group = c("group", "site"),
normalise = TRUE, cex = 1, pch = c(15:19, 15:18))
```

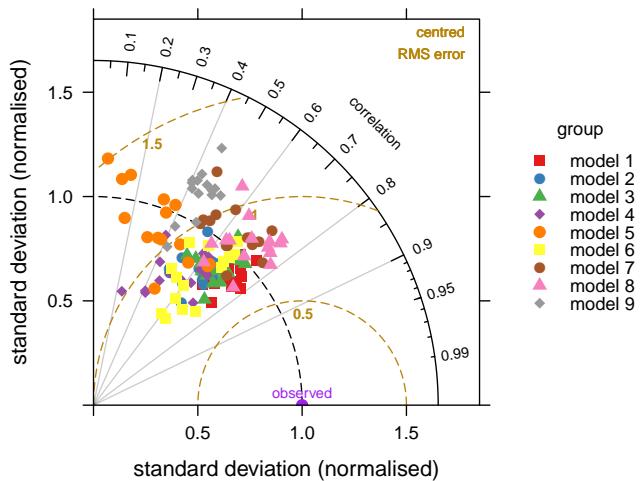


FIGURE 28.4 Use of the `TaylorDiagram` function to show model performance for 9 models used for all sites.

29 Model evaluation — the `conditionalQuantile` and `conditionalEval` functions

29.1 Purpose

Conditional quantiles are a very useful way of considering model performance against observations for continuous measurements Wilks 2005. The conditional quantile plot splits the data into evenly spaced bins. For each predicted value bin e.g. from 0 to $10 \mu\text{g m}^{-3}$ the *corresponding* values of the observations are identified and the median, 25/75th and 10/90 percentile (quantile) calculated for that bin. The data are plotted to show how these values vary across all bins. For a time series of observations and predictions that agree precisely the median value of the predictions will equal that for the observations for each bin.

The conditional quantile plot differs from the quantile-quantile plot (Q-Q plot) that is often used to compare observations and predictions. A Q-Q plot separately considers the distributions of observations and predictions, whereas the conditional quantile uses the corresponding observations for a particular interval in the predictions. Take as an example two time series, the first a series of real observations and the second a lagged time series of the same observations representing the predictions. These two time series will have identical (or very nearly identical) distributions (e.g. same median, minimum and maximum). A Q-Q plot would show a straight line showing perfect agreement, whereas the conditional quantile will not. This is because in any interval of the predictions the corresponding observations now have different values.

Plotting the data in this way shows how well predictions agree with observations and can help reveal many useful characteristics of how well model predictions agree with observations — across the full distribution of values. A single plot can therefore convey a considerable amount of information concerning model performance. The basic function is considerably enhanced by allowing flexible conditioning easily e.g. to evaluate model performance by season, day of the week and so on, as in other `openair` functions.

29.2 Options available

`mydata` A data frame containing the field `obs` and `mod` representing observed and modelled values.

`obs` The name of the observations in `mydata`.

`mod` The name of the predictions (modelled values) in `mydata`.

`type` `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

Type can be up length two e.g. `type = c("season", "weekday")` will produce a 2x2 plot split by season and day of the week. Note, when two types are provided the first forms the columns and the second the rows.

- `bins` Number of bins to be used in calculating the different quantile levels.
- `min.bin` The minimum number of points required for the estimates of the 25/75th and 10/90th percentiles.
- `xlab` label for the x-axis, by default “predicted value”.
- `ylab` label for the y-axis, by default “observed value”.
- `col` Colours to be used for plotting the uncertainty bands and median line. Must be of length 5 or more.
- `key.columns` Number of columns to be used in the key.
- `key.position` Location of the key e.g. “top”, “bottom”, “right”, “left”. See `lattice xyplot` for more details.
- `auto.text` Either `TRUE` (default) or `FALSE`. If `TRUE` titles and axis labels etc. will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO₂.
- `...` Other graphical parameters passed onto `cutData` and `lattice:xyplot`. For example, `conditionalQuantile` passes the option `hemisphere = "southern"` on to `cutData` to provide southern (rather than default northern) hemisphere handling of `type = "season"`. Similarly, common axis and title labelling options (such as `xlab`, `ylab`, `main`) are passed to `xyplot` via `quickText` to handle routine formatting.

29.3 Example of use

To make things more interesting we will use data from a model evaluation exercise organised by Defra in 2010/2011. A large number of models were evaluated as part of the evaluation but we only consider hourly ozone predictions from the CMAQ model being used at King’s College London.

First the data are loaded:

```
load("~/My Drive/openair/Data/CMAQozone.RData")
class(CMAQ.KCL$date) <- c("POSIXct", "POSIXt")
head(CMAQ.KCL)

##           site          date    o3 rollingO3Meas    mod rollingO3Mod group
## 1 Aston.Hill 2006-01-01 00:00:00   NA        NA 92.80      NA CMAQ.KCL
## 2 Aston.Hill 2006-01-01 01:00:00   74        NA 92.18      NA CMAQ.KCL
## 3 Aston.Hill 2006-01-01 02:00:00   72        NA 92.14      NA CMAQ.KCL
## 4 Aston.Hill 2006-01-01 03:00:00   72        NA 91.72      NA CMAQ.KCL
## 5 Aston.Hill 2006-01-01 04:00:00   70        NA 91.50      NA CMAQ.KCL
## 6 Aston.Hill 2006-01-01 05:00:00   66        NA 92.28      NA CMAQ.KCL
```

The data consists of hourly observations of O₃ in $\mu\text{g m}^{-3}$ at 15 rural O₃ sites in the UK together with predicted values.¹⁷ First of all we consider O₃ predictions across all

```
conditionalQuantile(CMAQ.KCL, obs = "o3", mod = "mod")
```

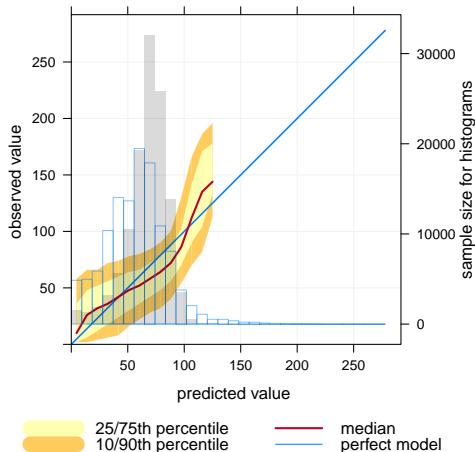


FIGURE 29.1 Example of the use of conditional quantiles applied to the KCL CMAQ model for 15 rural O₃ monitoring sites in 2006, for hourly data. The blue line shows the results for a perfect model. In this case the observations cover a range from 0 to 270 µg m⁻³. The red line shows the median values of the predictions and corresponding observations. The maximum predicted value is 125 µg m⁻³, somewhat less than the maximum observed value. The shading shows the predicted quantile intervals i.e. the 25/75th and the 10/90th. A perfect model would lie on the blue line and have a very narrow spread. There is still some spread because even for a perfect model a specific quantile interval will contain a range of values. However, for the number of bins used in this plot the spread will be very narrow. Finally, the histogram shows the counts of predicted values.

sites to help illustrate the purpose of the function. The results are shown in Figure 29.1. An explanation of the Figure is given in its caption.

A more informative analysis can be undertaken by considering conditional quantiles separately by site, which is easily done using the `type` option. The results are shown in Figure 29.2. It is now easier to see where the model performs best and how it varies by site type. For example, at a remote site in Scotland like Strath Vaich it is clear that the model does not capture either the lowest or highest O₃ concentrations very well.

As with other `openair` functions, the ability to consider conditioning can really help with interpretation. For example, what do the conditional quantiles at Lullington Heath (in south-east England) look like by season? This is easily done by subsetting the data to select that site and setting `type = "season"`, as shown in Figure 29.3. These results show that winter predictions have good coverage i.e. with width of the blue ‘perfect model’ line is the same as the observations. However, the predictions tend to be somewhat lower than observations for most concentrations (the median line is below the blue line) — and the width of the 10/75th and 10/90th percentiles is quite broad. However, the area where the model is less good is in summer and autumn because the predictions have low coverage (the red line only covers less than half of the observation line and the width of the percentiles is wide).

Of course it is also easy to plot by hour of the day, day of the week, by daylight/night-time and so on — easily. All these approaches can help better understand *why* a model does not perform very well rather than just quantifying its performance. Also, these types of analysis are particularly useful when more than one model is involved in a comparison as in the recent Defra model evaluation exercise, which we will come back

¹⁷We thank Dr Sean Beevers and Dr Nutthida Kitwiroon for access to these data.

```
conditionalQuantile(CMAQ.KCL, obs = "o3", mod = "mod", type = "site")
```

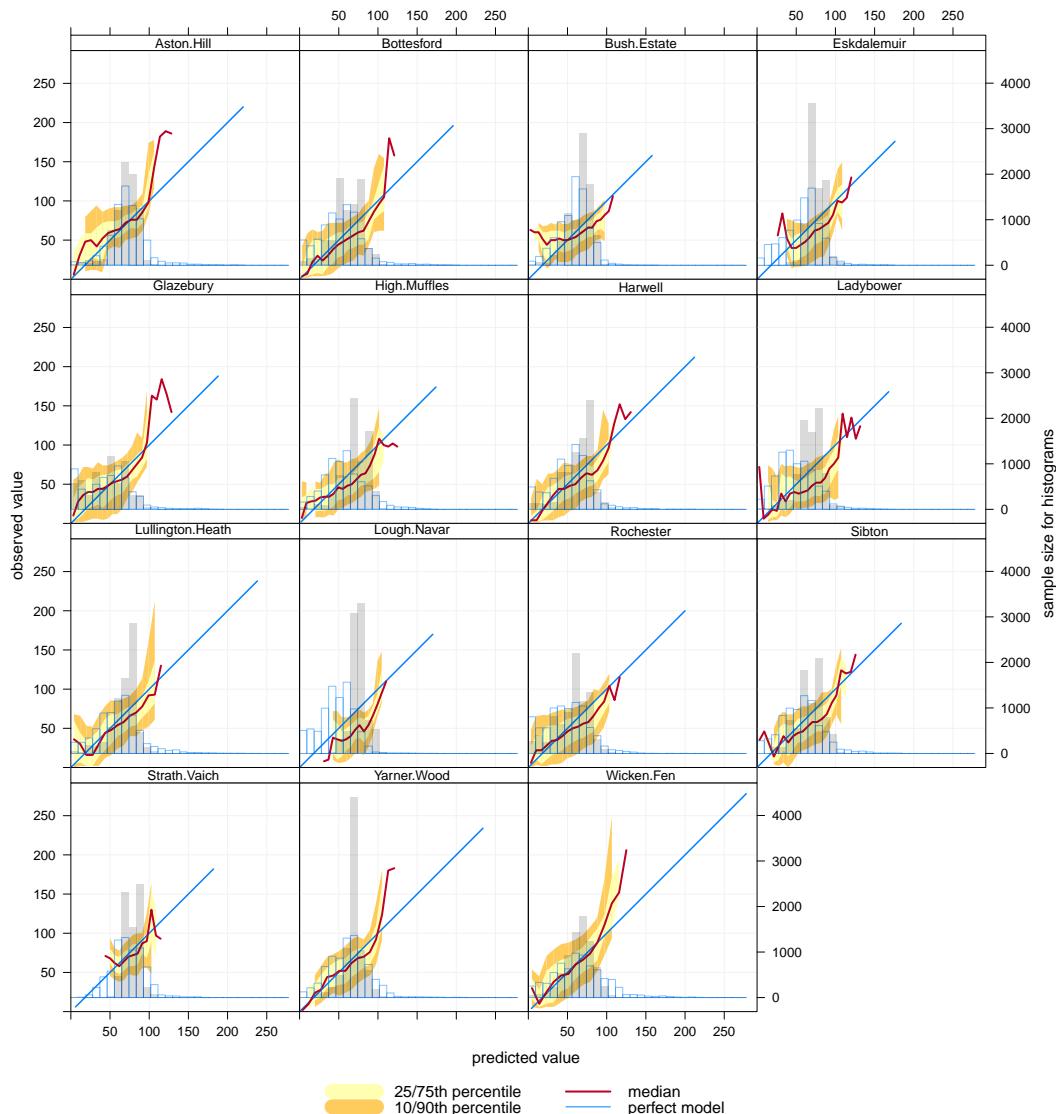


FIGURE 29.2 Conditional quantiles by site for 15 O₃ monitoring sites in the UK.

to later when some of the results are published.

There are numerous ways in which model performance can be assessed, including the use of common statistical measures described in Section 27. These approaches are very useful for comparing models against observations and other models. However, model developers would generally like to know *why* a model may have poor performance under some situations. This is a much more challenging issue to address. However, useful information can be gained by considering how other variables vary simultaneously.

The `conditionalEval` function provides information on how other variables vary across the *same* intervals as shown on the conditional quantile plot. There are two types of variable that can be considered by setting the value of `statistic`. First, `statistic` can be another variable in the data frame. In this case the plot will show the different proportions of `statistic` across the range of predictions. For example `statistic = "season"` will show for each interval of `mod` the proportion of predictions that were spring, summer, autumn or winter. This is useful because if model performance

```
conditionalQuantile(subset(CMAQ.KCL, site == "Lullington.Heath"),
                     obs = "o3",
                     mod = "mod", type = "season")
```

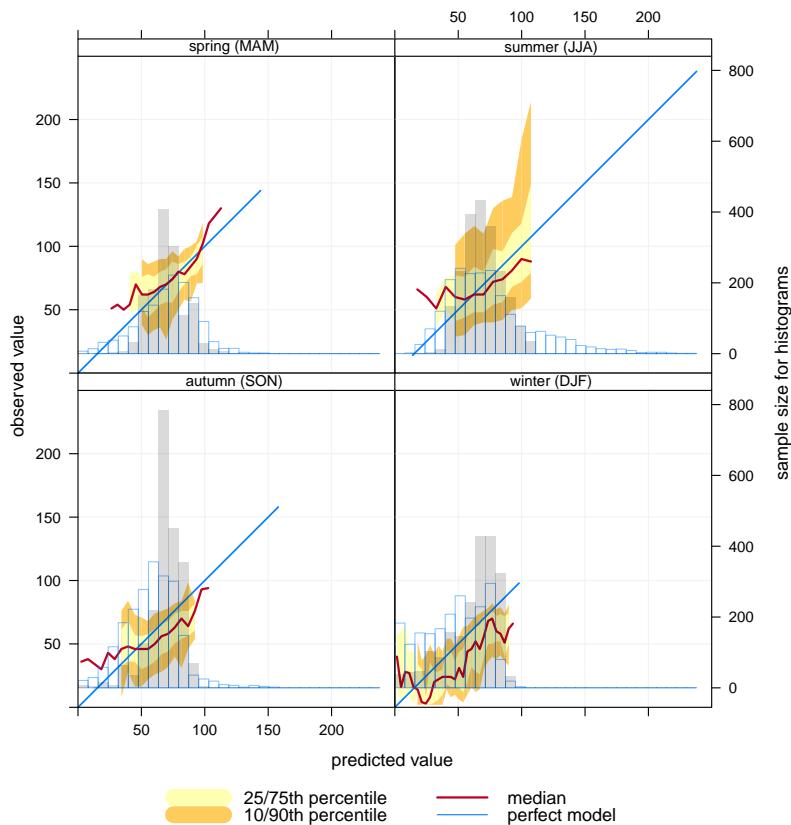


FIGURE 29.3 Conditional quantiles at Lullington Heath conditioned by season.

is worse for example at high concentrations of `mod` then knowing that these tend to occur during a particular season etc. can be very helpful when trying to understand *why* a model fails. See [Section 11](#) for more details on the types of variable that can be `statistic`. Another example would be `statistic = "ws"` (if wind speed were available in the data frame), which would then split wind speed into four quantiles and plot the proportions of each. Again, this would help show whether model performance in predicting concentrations of O₃ for example is related to low to high wind speed conditions.

`conditionalEval` can also simultaneously plot the model performance of other observed/predicted variable *pairs* according to different model evaluation statistics. These statistics derive from the [Section 27](#) function and include *MB*, *NMB*, *r*, *COE*, *MGE*, *NMGE*, *RMSE* and *FAC2*. More than one statistic can be supplied e.g. `statistic = c("NMB", "COE")`. Bootstrap samples are taken from the corresponding values of other variables to be plotted and their statistics with 95% confidence intervals calculated. In this case, the model *performance* of other variables is shown across the same intervals of `mod`, rather than just the values of single variables. In this second case the model would need to provide observed/predicted pairs of other variables.

For example, a model may provide predictions of NO_x and wind speed (for which there are also observations available). The `conditionalEval` function will show how well these other variables are predicted for the *same prediction intervals* of the

main variable assessed in the conditional quantile plot e.g. ozone. In this case, values are supplied to `var.obs` (observed values for other variables) and `var.mod` (modelled values for other variables). For example, to consider how well the model predicts NO_x and wind speed `var.obs = c("nox.obs", "ws.obs")` and `var.mod = c("nox.mod", "ws.mod")` would be supplied (assuming `nox.obs`, `nox.mod`, `ws.obs`, `ws.mod` are present in the data frame). The analysis could show for example, when ozone concentrations are under-predicted, the model may also be shown to over-predict concentrations of NO_x at the same time, or under-predict wind speeds. Such information can thus help identify the underlying causes of poor model performance. For example, an under-prediction in wind speed could result in higher surface NO_x concentrations and lower ozone concentrations. Similarly if wind speed predictions were good and NO_x was over predicted it might suggest an over-estimate of NO_x emissions. One or more additional variables can be plotted.

A special case is `statistic = "cluster"`. In this case a data frame is provided that contains the cluster calculated by `trajCluster` and `importTraj`. Alternatively users could supply their own pre-calculated clusters. These calculations can be very useful in showing whether certain back trajectory clusters are associated with poor (or good) model performance. Note that in the case of `statistic = "cluster"` there will be fewer data points used in the analysis compared with the ordinary statistics above because the trajectories are available for every three hours. Also note that `statistic = "cluster"` cannot be used together with the ordinary model evaluation statistics such as `MB`. The output will be a bar chart showing the proportion of each interval of `mod` by cluster number.

Far more insight can be gained into model performance through conditioning using `type`. For example, `type = "season"` will plot conditional quantiles and the associated model performance statistics of other variables by each season. `type` can also be a factor or character field e.g. representing different models used.

The `conditionalEval` function has the following options:

<code>mydata</code>	A data frame containing the field <code>obs</code> and <code>mod</code> representing observed and modelled values.
<code>obs</code>	The name of the observations in <code>mydata</code> .
<code>mod</code>	The name of the predictions (modelled values) in <code>mydata</code> .
<code>var.obs</code>	Other variable observations for which statistics should be calculated. Can be more than length one e.g. <code>var.obs = c("nox.obs", "ws.obs")</code> . Note that including other variables could reduce the number of data available to plot due to the need of having non-missing data for all variables.
<code>var.mod</code>	Other variable predictions for which statistics should be calculated. Can be more than length one e.g. <code>var.obs = c("nox.obs", "ws.obs")</code> .
<code>type</code>	<code>type</code> determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in <code>cutData</code> e.g. "season", "year", "weekday" and so on. For example, <code>type = "season"</code> will produce four plots — one for each season. It is also possible to choose <code>type</code> as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great

flexibility for understanding the variation of different variables and how they depend on one another.

<code>bins</code>	Number of bins used in <code>conditionalQuantile</code> .
<code>statistic</code>	Statistic(s) to be plotted. Can be “MB”, “NMB”, “r”, “COE”, “MGE”, “NMGE”, “RMSE” and “FAC2”, as described in <code>modStats</code> . When these statistics are chosen, they are calculated from <code>var.mod</code> and <code>var.mod</code> . <code>statistic</code> can also be a value that can be supplied to <code>cutData</code> . For example, <code>statistic = "season"</code> will show how model performance varies by season across the distribution of predictions which might highlight that at high concentrations of NOx the model tends to underestimate concentrations and that these periods mostly occur in winter. <code>statistic</code> can also be another variable in the data frame — see <code>cutData</code> for more information. A special case is <code>statistic = "cluster"</code> if clusters have been calculated using <code>trajCluster</code> .
<code>xlab</code>	label for the x-axis, by default <code>"predicted value"</code> .
<code>ylab</code>	label for the y-axis, by default <code>"observed value"</code> .
<code>col</code>	Colours to be used for plotting the uncertainty bands and median line. Must be of length 5 or more.
<code>col.var</code>	Colours for the additional variables to be compared. See <code>openColours</code> for more details.
<code>var.names</code>	Variable names to be shown on plot for plotting <code>var.obs</code> and <code>var.mod</code> .
<code>auto.text</code>	Either <code>TRUE</code> (default) or <code>FALSE</code> . If <code>TRUE</code> titles and axis labels etc. will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO ₂ .
<code>...</code>	Other graphical parameters passed onto <code>conditionalQuantile</code> and <code>cutData</code> . For example, <code>conditionalQuantile</code> passes the option <code>hemisphere = "southern"</code> on to <code>cutData</code> to provide southern (rather than default northern) hemisphere handling of <code>type = "season"</code> . Similarly, common axis and title labelling options (such as <code>xlab</code> , <code>ylab</code> , <code>main</code>) are passed to <code>xyplot</code> via <code>quickText</code> to handle routine formatting.

As an example, similar data to that described above from CMAQ have been used as an example.

A subset of the data for the North Kensington site can be imported as shown below.

```
condDat <- readRDS(url("https://davidcarslaw.github.io/data/openair/condDat.rds"))
```

The file contains observed and modelled hourly values for O₃, NO_x, wind speed, wind direction, temperature and relative humidity.

```
conditionalEval(condDat, obs = "O3.obs", mod = "O3.mod",
               statistic = "ws.obs",
               col.var = "Set3")
```

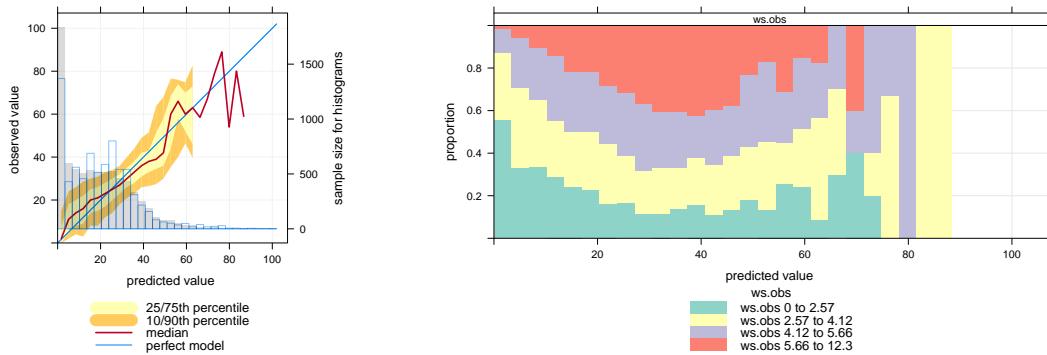


FIGURE 29.4 Conditional quantiles at for O₃ concentrations (left plot). On the right is a plot showing how the wind speed varies across the O₃ prediction intervals.

```
head(condDat)
```

```
##           date O3.obs NOx.obs ws.obs wd.obs temp.obs rh.obs O3.mod
## 5 2006-01-01 00:00:00     10 29.43665 4.6296   190    4.9    89 14.80
## 10 2006-01-01 01:00:00    15 17.55393     NA   210    5.1    90 17.46
## 15 2006-01-01 02:00:00    11 19.64817 2.5720   220    4.9    94 18.31
## 20 2006-01-01 03:00:00    11 19.15393 3.6008   270    5.7    91 18.25
## 25 2006-01-01 04:00:00    11 17.03037 3.0864   270    5.0    94 18.08
## 30 2006-01-01 05:00:00    12 15.98325 3.6008   260    5.8    94 14.87
##           NOx.mod ws.mod wd.mod temp.mod rh.mod
## 5      24.00  2.78   224     3.85  93.16
## 10     19.91  2.63   226     3.85  92.77
## 15     18.25  2.52   236     2.85  99.40
## 20     18.33  2.48   253     2.85  99.19
## 25     18.09  2.24   275     3.85  97.49
## 30     21.38  2.43   285     4.85  94.41
```

The `conditionalEval` function can be used in a straightforward way to provide information on how predictions depend on another variable in general. In this case the option `statistic` can refer to another variable in the data frame to see how the quality of predictions depend on values of that variable. For example, in Figure 29.4 it can be seen how wind speed varies across the O₃ prediction intervals. At low predicted concentrations of O₃ there is a high proportion of low wind speed conditions (0 to 2.57 m s⁻¹). When O₃ is predicted to be around 40 ppb the wind speed tends to be higher — and finally at higher predicted concentrations of O₃ the wind speed tends to decrease again. The important aspect of plotting data in this way is that it can directly relate the prediction performance to values of other variables, which should help develop a much better idea of the conditions that matter most. The user can therefore develop a good feel for the types of conditions where a model performs well or poorly and this might provide clues as to the underlying reasons for specific model behaviour.

In an extension to Figure 29.4 it is possible to derive information on the simultaneous model performance of other variables. Figure 29.5 shows the conditional quantile plot for hourly O₃ predictions. This shows among other things that concentrations of O₃ tend to be under-predicted for concentrations less than about 20 ppb. The Figure on the right shows the simultaneous *model performance* for wind speed and NO_x for the same prediction intervals as shown in the conditional quantile plot. The plot on the

```
conditionalEval(condDat, obs = "O3.obs", mod = "O3.mod",
               var.obs = c("NOx.obs", "ws.obs"),
               var.mod = c("NOx.mod", "ws.mod"),
               statistic = "NMB",
               var.names = c("nox", "wind speed"))
```

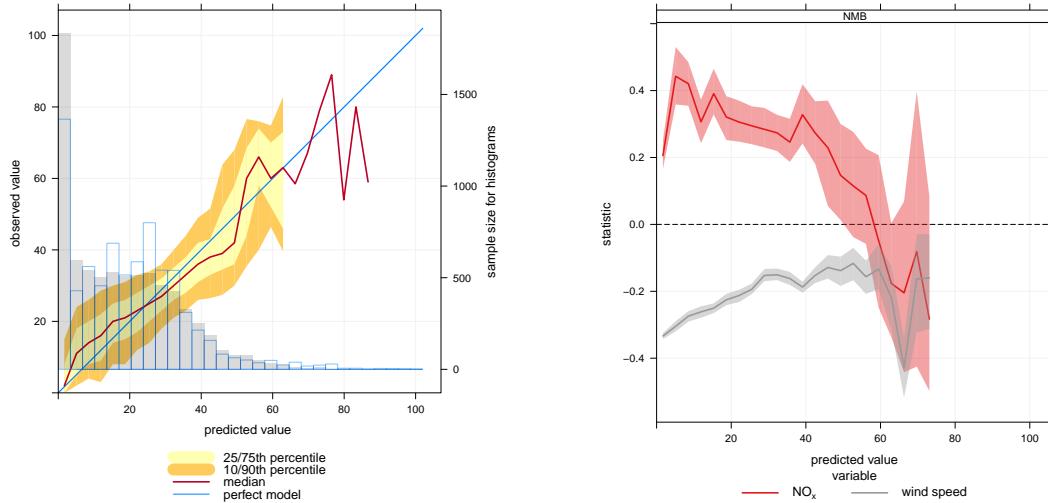


FIGURE 29.5 Conditional quantiles at for O_3 concentrations (left plot). On the right is the model performance for wind speed and NO_x predictions, which in this case is for the Normalised Mean Bias.

right shows that for low concentrations of predicted O_3 there is a tendency for NO_x concentrations to be overestimated ($NMB \approx 0.2$ to 0.4) and wind speeds to be underestimated ($NMB \approx -0.2$ to -0.3). One possible explanation for this behaviour is that the meteorological model tends to produce wind speeds that are too low, which would result in higher concentrations of NO_x , which in turn would result in lower concentrations of O_3 . Note that it is possible to include more than one statistic, which would be plotted in a new panel e.g. `statistic = c("NMB", "r")`.

In essence the `conditionalEval` function provides more information on model performance that can help better diagnose potential problems. Clearly, there are many other ways in which the results can be analysed, which will depend on the data available.

A plot using temperature predictions shows that for most of the range in O_3 predictions there is very little bias in temperature (although there is some negative bias in temperature for very low concentration O_3 predictions):

```
conditionalEval(condDat, obs = "O3.obs", mod = "O3.mod",
               var.obs = c("temp.obs", "ws.obs"),
               var.mod = c("temp.mod", "ws.mod"),
               statistic = "NMB", var.names = c("temperature", "wind speed"))
```

Finally (but not shown) it can be very useful to consider model performance in terms of air mass origin. In the example below, trajectories are imported, a cluster analysis undertaken and then evaluated using `conditionalEval`.

```

## import trajectories for 2006
traj <- importTraj("london", 2006)

## carry out a cluster analysis
cl <- trajCluster(traj, method = "Angle", n.cluster = 5)

## merge with original model eval data
condDat <- merge(condDat, cl, by = "date")

## plot it
conditionalEval(condDat, obs = "O3.obs", mod = "O3.mod",
                 statistic = "cluster",
                 col.var = "Set3")

```

30 The `calcFno2` function—estimating primary NO₂ fractions

30.1 Purpose

see also
`linearRelation`
 for oxidant
 slopes if NO_x,
 NO₂ and O₃
 are available

Recent research has shown that emissions of directly emitted (primary) NO₂ from road vehicles have increased and these increases have had important effects on concentrations of NO₂ (Carslaw 2005; Carslaw and Beevers 2004; Carslaw and Carslaw 2007). Many organisations would like to quantify the level of primary NO₂ from the analysis of ambient monitoring data to help with their air quality management responsibilities. The difficulty is that this is not a straightforward thing to do — it requires some form of modelling. In some situations where NO, NO₂ and O₃ are measured, it is possible to derive an estimate of the primary NO₂ fraction by considering the gradient in ‘total oxidant’ defined as NO₂ + O₃ (Clapp and Jenkin 2001).¹⁸ However, where we most want to estimate primary NO₂ (roadside sites), O₃ is rarely measured and an alternative approach must be used. One approach is using a simple constrained chemistry model described in (Carslaw and Beevers 2005). The `calcFno2` method is based on this work but makes a few simplifying assumptions to make it easier to code.

There are several assumptions that users should be aware of when using this function. First, it is most reliable at estimating f-NO₂ when the roadside concentration is much greater than the background concentration. Second, it is best if the chosen background site is reasonably close to the roadside site and not greatly affected by local sources (which it should not be as a background site). The way the calculations work is to try and distinguish between NO₂ that is directly emitted from vehicles and that derived through the reaction between NO and O₃. During summertime periods when concentrations of NO_x are lower these two influences tend to vary linearly with NO_x, making it difficult for the method to separate them. It can often be useful to select only winter months under these situations (or at least October–March). Therefore, as a simplifying assumption, the time available for chemical reactions to take place, τ , is set to 60 seconds. We have tested the method at roadside sites where O₃ is also measured and $\tau = 60$ s seems to provide a reasonably consistent calculation of f-NO₂.

Note that in some situations it may be worth filtering the data e.g. by wind direction to focus on the road itself. In this respect, the `polarPlot` function described on page 108 can be useful.

¹⁸Note that the volume fraction of NO₂ of total NO_x is termed *f-NO₂*. See the AQEG report for more information (AQEG 2008).

30.2 Options available

- input** A data frame with the following fields. `nox` and `no2` (roadside NOX and NO₂ concentrations), `back_nox`, `back_no2` and `back_o3` (hourly background concentrations of each pollutant). In addition `temp` (temperature in degrees Celsius) and `cl` (cloud cover in Oktas). Note that if `temp` and `cl` are not available, typical means values of 11 deg. C and cloud = 3.5 will be used.
- tau** Mixing time scale. It is unlikely the user will need to adjust this. See details below.
- user.fno2** User-supplied f-NO₂ fraction e.g. 0.1 is a NO₂/NOX ratio of 10 series and is useful for testing "what if" questions.
- main** Title of plot if required.
- xlab** x-axis label.
- ...** Other graphical parameters send to `scatterPlot`.

30.3 Example of use

We apply the technique to roadside data at Marylebone Road, with additional data from a nearby background site at North Kensington and appropriate meteorological variables. The data can be downloaded from the **openair** website (<http://www.openair-project.org>), and remember to change the file path below. The code run is:¹⁹

```
load( "~/My Drive/openair/Data/f-no2Data.RData")
## check first few lines of the file
head(fno2Data)

##           date nox no2 back_nox back_no2 back_o3 temp cl wd
## 1 01/01/1998 285  39      49     34       1  3.3  2 280
## 2 01/01/1998  NA   NA      56     35       0  3.3  5 230
## 3 01/01/1998  NA   NA      41     31       4  4.4  5 190
## 4 01/01/1998 493  52      52     33       1  3.9  5 170
## 5 01/01/1998 468  78      41     30       3  3.9  5 180
## 6 01/01/1998 264  42      36     29       6  3.3  2 190
```

Now apply the function, and collect the results as shown in Figure 30.1, with additional options that are sent to `scatterPlot`.

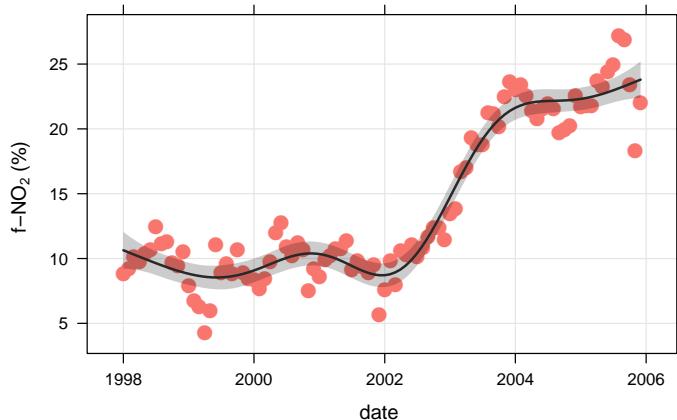
Note that this is a different way to run a function compared with what has been done previously. This time the results are read into the variable `results`, which stores the monthly mean estimated f-NO₂ values.

The results (expressed as a percentage) for f-NO₂ are then available for any other processing or plotting. The function also automatically generates a plot of monthly mean f-NO₂ values as shown in Figure 30.1. It is clear from this Figure that f-NO₂ was relatively stable at around 10 % until the end of 2002, before increasing sharply during 2003 — and remaining at about 20 % until the end of 2007. At this particular site the increases in f-NO₂ are very apparent.

An interesting question is what would NO₂ concentrations have been if f-NO₂ remained at the 10 % level, or indeed any other level. The `calcFno2` function also allows

¹⁹Note that the choice to give the plot a heading is optional.

```
results <- calcFno2(fno2Data, main = "Trends in f-NO2 at Marylebone Road",
                      pch = 16, smooth = TRUE, cex = 1.5)
```



`## the results are a list of two data frame, the first is the f-NO2 results:`

```
head(results$data[[1]])
```

```
##           date fno2
## 1 1.1998 1998-01-01 00:00:00 8.9
## 2 2.1998 1998-02-01 00:00:00 9.2
## 3 3.1998 1998-03-01 00:00:00 10.1
## 4 4.1998 1998-04-01 01:00:00 9.7
## 5 5.1998 1998-05-01 01:00:00 10.3
## 6 6.1998 1998-06-01 01:00:00 10.7
```

`## the second are the hourly nox, no2 and estimated o3:`
`head(results$data[[2]])`

```
##           date nox no2   o3
## 1 1998-01-01 00:00:00 285 56 0.19
## 2 1998-01-01 03:00:00 493 73 0.11
## 3 1998-01-01 04:00:00 468 70 0.35
## 4 1998-01-01 05:00:00 264 54 1.22
## 5 1998-01-01 06:00:00 171 47 0.92
## 6 1998-01-01 07:00:00 195 51 2.18
```

FIGURE 30.1 Plot from the application of the `calcFno2` function applied to Marylebone Road. The plot shows a smooth fit with 95 % confidence intervals.

the user to input their own f-NO₂ level and make new predictions. In this case the code run is slightly different, shown in Figure 30.2.

By providing a value to the option `user.fno2` (expressed as a fraction), the function will automatically calculate NO₂ concentrations with the chosen f-NO₂ value applied to the whole time series. In this case `results` will return a data frame with dates and NO₂ concentrations. In addition a plot is produced as shown in Figure 30.2. The blue line and shading show the measured data and highlight a clear increase in NO₂ concentrations from 2003 onwards. The red line and shading shows the predicted values assuming (in this case) that f-NO₂ was constant at 0.095. Based on these results it is clear that NO₂ concentrations would have been substantially less if it were not for the recent increases in f-NO₂.

```
results <- calcFno2(fno2Data, user.fno2 = 0.095, smooth = TRUE, pch = 16, cex = 1.5)
```

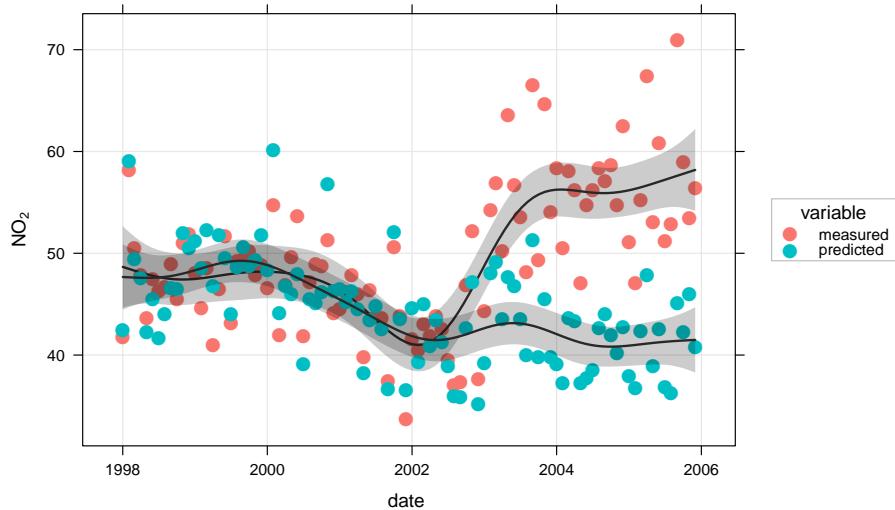


FIGURE 30.2 Plot from the application of the `calcFno2` function applied to Marylebone Road with a forced f-NO₂ values of 0.095 over the whole series. The red line and shading shows the trend in actual measurements and the blue line and shading the predicted trend in NO₂ if the f-NO₂ ratio had remained at 0.095.

31 Utility functions

31.1 Selecting data by date

Selecting by date/time in R can be intimidating for new users—and time consuming for all users. The `selectByDate` function aims to make this easier by allowing users to select data based on the British way of expressing date i.e. d/m/y. This function should be very useful in circumstances where it is necessary to select only part of a data frame.

The function has the following options:

<code>mydata</code>	A data frame containing a <code>date</code> field in hourly or high resolution format.
<code>start</code>	A start date string in the form d/m/yyyy e.g. “1/2/1999” or in ‘R’ format i.e. “YYYY-mm-dd”, “1999-02-01”
<code>end</code>	See <code>start</code> for format.
<code>year</code>	A year or years to select e.g. <code>year = 1998:2004</code> to select 1998-2004 inclusive or <code>year = c(1998, 2004)</code> to select 1998 and 2004.
<code>month</code>	A month or months to select. Can either be numeric e.g. <code>month = 1:6</code> to select months 1-6 (January to June), or by name e.g. <code>month = c("January", "December")</code> . Names can be abbreviated to 3 letters and be in lower or upper case.
<code>day</code>	A day name or days to select. <code>day</code> can be numeric (1 to 31) or character. For example <code>day = c("Monday", "Wednesday")</code> or <code>day = 1:10</code> (to select the 1st to 10th of each month). Names can be abbreviated to 3 letters and be in lower or upper case. Also accepts “weekday” (Monday - Friday) and “weekend” for convenience.

hour An hour or hours to select from 0-23 e.g. `hour = 0:12` to select hours 0 to 12 inclusive.

```
## select all of 1999
data.1999 <- selectByDate(mydata, start = "1/1/1999", end = "31/12/1999")
head(data.1999)

## # A tibble: 6 x 13
##   date           ws    wd   nox   no2    o3   pm10    so2    co   pm25
##   <dttm>        <dbl> <int> <int> <int> <int> <dbl> <dbl> <int>
## 1 1999-01-01 00:00:00  5.04   140    88    35     4    21  3.84  1.02   18
## 2 1999-01-01 01:00:00  4.08   160   132    41     3    17  5.24  2.7   11
## 3 1999-01-01 02:00:00  4.8    160   168    40     4    17  6.51  2.87   8
## 4 1999-01-01 03:00:00  4.92   150    85    36     3    15  4.18  1.62   10
## 5 1999-01-01 04:00:00  4.68   150    93    37     3    16  4.25  1.02   11
## 6 1999-01-01 05:00:00  3.96   160    74    29     5    14  3.88  0.725  NA
## # ... with 3 more variables: split.by <ord>, feature <chr>, ratio <dbl>

tail(data.1999)

## # A tibble: 6 x 13
##   date           ws    wd   nox   no2    o3   pm10    so2    co   pm25
##   <dttm>        <dbl> <int> <int> <int> <int> <dbl> <dbl> <int>
## 1 1999-12-31 18:00:00  4.68   190   226    39    NA    29  5.46  2.38   23
## 2 1999-12-31 19:00:00  3.96   180   202    37    NA    27  4.78  2.15   23
## 3 1999-12-31 20:00:00  3.36   190   246    44    NA    30  5.88  2.45   23
## 4 1999-12-31 21:00:00  3.72   220   231    35    NA    28  5.28  2.22   23
## 5 1999-12-31 22:00:00  4.08   200   217    41    NA    31  4.79  2.17   26
## 6 1999-12-31 23:00:00  3.24   200   181    37    NA    28  3.48  1.78   22
## # ... with 3 more variables: split.by <ord>, feature <chr>, ratio <dbl>

## easier way
data.1999 <- selectByDate(mydata, year = 1999)

## more complex use: select weekdays between the hours of 7 am to 7 pm
sub.data <- selectByDate(mydata, day = "weekday", hour = 7:19)

## select weekends between the hours of 7 am to 7 pm in winter (Dec, Jan, Feb)
sub.data <- selectByDate(mydata, day = "weekend", hour = 7:19,
                           month = c("dec", "jan", "feb"))
```

The function can be used directly in other functions. For example, to make a polar plot using year 2000 data:

```
polarPlot(selectByDate(mydata, year = 2000), pollutant = "so2")
```

31.2 Selecting run lengths of values above a threshold — pollution episodes

A seemingly easy thing to do that has relevance to air pollution episodes is to select run lengths of contiguous values of a pollutant above a certain threshold. For example, one might be interested in selecting O₃ concentrations where there are at least 8 consecutive hours above 90 ppb. In other words, a selection that combines both a threshold and *persistence*. These periods can be very important from a health perspective and it can be useful to study the conditions under which they occur. But how do you select such periods easily? The **selectRunning** utility function has been written to do this. It could be useful for all sorts of situations e.g.

- Selecting hours where primary pollutant concentrations are persistently high — and then applying other **openair** functions to analyse the data in more depth.
- In the study of particle suspension or deposition etc. it might be useful to select hours where wind speeds remain high or rainfall persists for several hours to see how these conditions affect particle concentrations.
- It could be useful in health impact studies to select blocks of data where pollutant concentrations remain above a certain threshold.

The **selectRunning** has the following options:

mydata A data frame with a **date** field and at least one numeric **pollutant** field to analyse.

pollutant Name of variable to process. Mandatory.

run.len Run length for extracting contiguous values of **pollutant** above the **threshold** value.

threshold The threshold value for **pollutant** above which data should be extracted.

As an example we are going to consider O₃ concentrations at a semi-rural site in south-west London (Teddington). The data can be downloaded as follows:

```
ted <- importKCL(site = "td0", year = 2005:2009, met = TRUE)
## see how many rows there are
nrow(ted)
```

We are going to contrast two polar plots of O₃ concentration. The first uses all hours in the data set, and the second uses a subset of hours. The subset of hours is defined by O₃ concentrations above 90 ppb for periods of at least 8-hours i.e. what might be considered as ozone episode conditions.

```
episode <- selectRunning(ted, pollutant = "o3", threshold = 90, run.len = 8)
## see how many rows there are
nrow(episode)

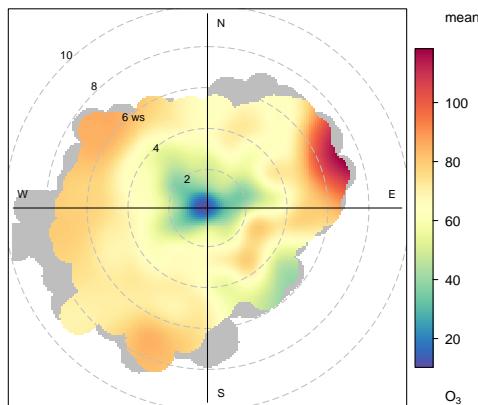
## [1] 1399
```

Now we are going to produce two bivariate polar plots shown in [Figure 31.1](#).

The results are shown in [Figure 31.1](#). The polar plot for all data (left plot of [Figure 31.1](#)) shows that the highest O₃ concentrations tend to occur for high wind speed conditions from almost every direction. Lower concentrations are observed for low wind speeds because concentrations of NO_x are higher, resulting in O₃ destruction. By contrast, a polar plot of the episode conditions (right plot of [Figure 31.1](#)) is very different. In this case there is a clear set of conditions where these criteria are met i.e. lengths of at least 8-hours where the O₃ concentration is at least 90 ppb. It is clear the highest concentrations are dominated by south-easterly conditions i.e. corresponding to easterly flow from continental Europe where there has been time to the O₃ chemistry to take place.

Another interesting test plot is to consider NO_x concentrations at Marylebone Road — see [Figure 15.1](#), which shows that high concentrations are dominated by a swathe of south-westerly wind conditions (even for high wind speeds). However, if a selection is made of episode conditions (defined here as NO_x concentrations >500 ppb for at least 5-hours), then it can be seen that it is actually the low wind speed conditions that

```
polarPlot(ted, pollutant = "o3", min.bin = 2)
```



```
polarPlot(episode, pollutant = "o3", min.bin = 2)
```

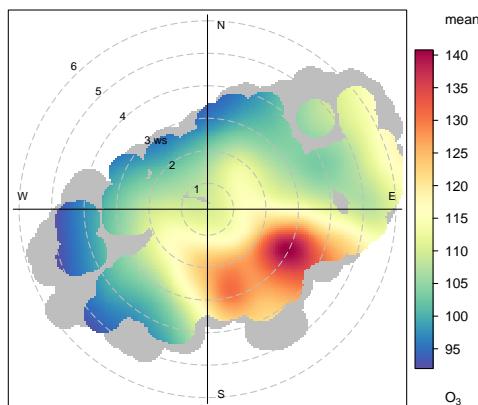


FIGURE 31.1 Example of using the `selectRunning` function to select episode hours to produce bivariate polar plots of O₃ concentration.

dominate. These conditions correspond to low in-canyon wind speeds and low wind speeds across London, which tend to elevate local and background NO_x concentrations. Even though high concentrations of NO_x are observed at high wind speeds, it does not seem that these conditions are as important for episode conditions. Users can run the code below to verify these observations.

```
episode <- selectRunning(mydata, pollutant = "nox", threshold = 800, run.len = 5)


```

31.3 Calculating rolling means

Some air pollution statistics such as for O₃ and particulate matter are expressed as rolling means and it is useful to be able to calculate these. It can also be useful to help smooth-out data for clearer plotting. The `rollingMean` function makes these calculations. One detail that can be important is that for some statistics a mean is only considered valid if there are a sufficient number of valid readings over the averaging period. Often there is a requirement for at least 75 % data capture. For example, with an averaging period of 8 hours and a data capture threshold of 75%, at least 6 hours are required to calculate the mean.

The function is called as follows; in this case to calculate 8-hour rolling mean concentrations of O₃.

```
data(mydata)
mydata <- rollingMean(mydata, pollutant = "o3", hours = 8,
                      new.name = "rollingo3", data.thresh = 75)
tail(mydata)

## # A tibble: 6 x 11
##   date                  ws    wd  nox  no2    o3 pm10    so2    co pm25
##   <dttm>              <dbl> <int> <int> <int> <int> <dbl> <dbl> <dbl>
## 1 2005-06-23 07:00:00    1.5    250    404    156     4    49    NA  1.81   28
## 2 2005-06-23 08:00:00    1.5    260    388    145     6    48    NA  1.64   26
## 3 2005-06-23 09:00:00    1.5    210    404    168     7    58    NA  1.29   34
## 4 2005-06-23 10:00:00    2.6    240    387    175    10    55    NA  1.29   34
## 5 2005-06-23 11:00:00    3.1    220    312    125    15    52    NA  1.29   33
## 6 2005-06-23 12:00:00    3.1    220    287    119    17    55    NA  1.29   35
## # ... with 1 more variable: rollingo3 <dbl>
```

Note that calculating rolling means shortens the length of the data set. In the case of O₃, no calculations are made for the last 7 hours.

Type `help(rollingMean)` into R for more details. Note that the function currently only works with a single site.

31.4 Aggregating data by different time intervals

Aggregating data by different averaging periods is a common and important task. There are many reasons for aggregating data in this way:

1. Data sets may have different averaging periods and there is a need to combine them. For example, the task of combining an hourly air quality data set with a 15-minute average meteorological data set. The need here would be to aggregate the 15-minute data to 1-hour before merging.
2. It is extremely useful to consider data with different averaging times in a straightforward way. Plotting a very long time series of hourly or higher resolution data can hide the main features and it would be useful to apply a specific (but flexible) averaging period to the data for plotting.
3. Those who make measurements during field campaigns (particularly for academic research) may have many instruments with a range of different time resolutions. It can be useful to re-calculate time series with a common averaging period; or maybe help reduce noise.
4. It is useful to calculate statistics other than means when aggregating e.g. percentile values, maximums etc.
5. For statistical analysis there can be short-term autocorrelation present. Being able to choose a longer averaging period is sometimes a useful strategy for minimising autocorrelation.

In aggregating data in this way, there are a couple of other issues that can be useful to deal with at the same time. First, the calculation of proper vector-averaged wind direction is essential. Second, sometimes it is useful to set a minimum number of data points that must be present before the averaging is done. For example, in calculating monthly averages, it may be unwise to not account for data capture if some months only have a few valid points.

When a data capture threshold is set through `data.thresh` it is necessary for `timeAverage` to know what the original time interval of the input time series is. The function will try and calculate this interval based on the most common time gap (and will print the assumed time gap to the screen). This works fine most of the time but there are occasions where it may not e.g. when very few data exist in a data frame. In this case the user can explicitly specify the interval through `interval` in the same format as `avg.time` e.g. `interval = "month"`. It may also be useful to set `start.date` and `end.date` if the time series do not span the entire period of interest. For example, if a time series ended in October and annual means are required, setting `end.date` to the end of the year will ensure that the whole period is covered and that `data.thresh` is correctly calculated. The same also goes for a time series that starts later in the year where `start.date` should be set to the beginning of the year.

All these issues are (hopefully) dealt with by the `timeAverage` function. The options are shown below, but as ever it is best to check the help that comes with the `openair` package.

see also
`timePlot` for
plotting with
different
averaging
times and
statistics

The `timeAverage` function has the following options:

mydata A data frame containing a `date` field . Can be class `POSIXct` or `Date`.

avg.time This defines the time period to average to. Can be “sec”, “min”, “hour”, “day”, “DSTday”, “week”, “month”, “quarter” or “year”. For much increased flexibility a number can precede these options followed by a space. For example, a `timeAverage` of 2 months would be `period = "2 month"`. In addition, `avg.time` can equal “season”, in which case 3-month seasonal values are calculated with spring defined as March, April, May and so on.

Note that `avg.time` can be *less* than the time interval of the original series, in which case the series is expanded to the new time interval. This is useful, for example, for calculating a 15-minute time series from an hourly one where an hourly value is repeated for each new 15-minute period. Note that when expanding data in this way it is necessary to ensure that the time interval of the original series is an exact multiple of `avg.time` e.g. hour to 10 minutes, day to hour. Also, the input time series must have consistent time gaps between successive intervals so that `timeAverage` can work out how much ‘padding’ to apply. To pad-out data in this way choose `fill = TRUE`.

data.thresh The data capture threshold to use (%). A value of zero means that all available data will be used in a particular period regardless if of the number of values available. Conversely, a value of 100 will mean that all data will need to be present for the average to be calculated, else it is recorded as `NA`. See also `interval`, `start.date` and `end.date` to see whether it is advisable to set these other options.

statistic The statistic to apply when aggregating the data; default is the mean. Can be one of “mean”, “max”, “min”, “median”, “frequency”, “sd”, “percentile”. Note that “sd” is the standard deviation, “frequency” is the number (frequency) of valid records in the period and “data.cap” is the percentage data capture. “percentile” is the percentile level (%) between 0-100, which can be set using the “percentile” option — see below. Not used if `avg.time = "default"`.

- type** `type` allows `timeAverage` to be applied to cases where there are groups of data that need to be split and the function applied to each group. The most common example is data with multiple sites identified with a column representing site name e.g. `type = "site"`. More generally, `type` should be used where the date repeats for a particular grouping variable. However, if `type` is not supplied the data will still be averaged but the grouping variables (character or factor) will be dropped.
- percentile** The percentile level in % used when `statistic = "percentile"`. The default is 95.
- start.date** A string giving a start date to use. This is sometimes useful if a time series starts between obvious intervals. For example, for a 1-minute time series that starts “2009-11-29 12:07:00” that needs to be averaged up to 15-minute means, the intervals would be “2009-11-29 12:07:00”, “2009-11-29 12:22:00” etc. Often, however, it is better to round down to a more obvious start point e.g. “2009-11-29 12:00:00” such that the sequence is then “2009-11-29 12:00:00”, “2009-11-29 12:15:00” ...`start.date` is therefore used to force this type of sequence.
- end.date** A string giving an end date to use. This is sometimes useful to make sure a time series extends to a known end point and is useful when `data.thresh > 0` but the input time series does not extend up to the final full interval. For example, if a time series ends sometime in October but annual means are required with a data capture of >75% then it is necessary to extend the time series up until the end of the year. Input in the format yyyy-mm-dd HH:MM.
- interval** The `timeAverage` function tries to determine the interval of the original time series (e.g. hourly) by calculating the most common interval between time steps. The interval is needed for calculations where the `data.thresh > 0`. For the vast majority of regular time series this works fine. However, for data with very poor data capture or irregular time series the automatic detection may not work. Also, for time series such as monthly time series where there is a variable difference in time between months users should specify the time interval explicitly e.g. `interval = "month"`. Users can also supply a time interval to *force* on the time series. See `avg.time` for the format.
- This option can sometimes be useful with `start.date` and `end.date` to ensure full periods are considered e.g. a full year when `avg.time = "year"`.
- vector.ws** Should vector averaging be carried out on wind speed if available? The default is `FALSE` and scalar averages are calculated. Vector averaging of the wind speed is carried out on the u and v wind components. For example, consider the average of two hours where the wind direction and speed of the first hour is 0 degrees and 2m/s and 180 degrees and 2m/s for the second hour. The scalar average of the wind speed is simply the arithmetic average = 2m/s and the vector average is 0m/s. Vector-averaged wind speeds will always be lower than scalar-averaged values.
- fill** When time series are expanded i.e. when a time interval is less than the original time series, data are ‘padded out’ with `NA`. To ‘pad-out’ the

additional data with the first row in each original time interval, choose `fill = TRUE`.

... Additional arguments for other functions calling `timeAverage`.

```
## load in fresh version of mydata
data(mydata)
```

To calculate daily means from hourly (or higher resolution) data:

```
daily <- timeAverage(mydata, avg.time = "day")
head(daily)

## # A tibble: 6 x 10
##   date           ws    wd   nox   no2    o3  pm10    so2    co  pm25
##   <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 1998-01-01 00:00:00  6.84 188. 154. 39.4  6.87 18.2  3.15 2.70  NaN
## 2 1998-01-02 00:00:00  7.07 223. 132. 39.5  6.48 27.8  3.94 1.77  NaN
## 3 1998-01-03 00:00:00 11.0 226. 120. 38.0  8.41 20.2  3.20 1.74  NaN
## 4 1998-01-04 00:00:00 11.5 223. 105. 35.3  9.61 21.0  2.96 1.62  NaN
## 5 1998-01-05 00:00:00  6.61 237. 175. 46.0  4.96 24.2  4.52 2.13  NaN
## 6 1998-01-06 00:00:00  4.38 197. 214. 45.3  1.35 34.6  5.70 2.53  NaN
```

Monthly 95th percentile values:

```
monthly <- timeAverage(mydata, avg.time = "month", statistic = "percentile",
                        percentile = 95)
head(monthly)

## # A tibble: 6 x 10
##   date           ws    wd   nox   no2    o3  pm10    so2    co  pm25
##   <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 1998-01-01 00:00:00 11.2  45.  371.  68.6   14  53   11.1  3.99  NA
## 2 1998-02-01 00:00:00  8.16  16.7  524.   92     7  68.9  17.5  5.63  NA
## 3 1998-03-01 00:00:00 10.6  37.6  417.   85     15  61   18.4  4.85  NA
## 4 1998-04-01 00:00:00  8.16  44.4  384    81.5   20  52   14.6  4.17  NA
## 5 1998-05-01 00:00:00  7.56  40.6  300    80     25  61   12.7  3.55  40
## 6 1998-06-01 00:00:00  8.47  50.7  377   74.2   15  53   12.2  4.28  33.9
```

2-week averages but only calculate if at least 75% of the data are available:

```
twoweek <- timeAverage(mydata, avg.time = "2 week", data.thresh = 75)
head(twoweek)

## # A tibble: 6 x 10
##   date           ws    wd   nox   no2    o3  pm10    so2    co  pm25
##   <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 1997-12-29 00:00:00  6.98 212. 167. 41.4  4.63 29.3  4.47 2.17  NA
## 2 1998-01-12 00:00:00  4.91 221. 173. 42.1  4.70 28.8  5.07 1.86  NA
## 3 1998-01-26 00:00:00  2.78 242. 233. 51.4  2.30 34.9  8.07 2.45  NA
## 4 1998-02-09 00:00:00  4.43 215. 276. 57.1  2.63 43.7  8.98 2.94  NA
## 5 1998-02-23 00:00:00  6.89 237. 248. 56.7  4.99 28.8  9.79 2.57  NA
## 6 1998-03-09 00:00:00  2.97 288. 160. 44.8  5.64 32.7  8.65 1.62  NA
```

Note that `timeAverage` has a `type` option to allow for the splitting of variables by a grouping variable. The most common use for `type` is when data are available for different sites and the averaging needs to be done on a per site basis.

First, retaining by site averages:

```
# import some data for two sites
dat <- importAURN(c("kc1", "my1"), year = 2011:2013)

# annual averages by site
timeAverage(dat, avg.time = "year", type = "site")

## # A tibble: 6 x 16
## # Groups:   site [2]
##   site   date          o3    no2    co    so2   pm10   nox    no pm2.5
##   <fct> <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Lond~ 2011-01-01 00:00:00 39.4  36.1  0.225  2.06  23.7  53.8  11.6  16.3
## 2 Lond~ 2012-01-01 00:00:00 38.5  36.7  0.266  2.03  20.2  57.4  13.3  14.6
## 3 Lond~ 2013-01-01 00:00:00 38.4  36.9  0.250  2.01  23.1  57.9  13.7  14.7
## 4 Lond~ 2011-01-01 00:00:00 18.5  97.2  0.656  6.86  38.4  306.  137.  24.5
## 5 Lond~ 2012-01-01 00:00:00 15.0  94.0  0.589  8.13  30.8  313.  143.  21.5
## 6 Lond~ 2013-01-01 00:00:00 17.7  84.7  0.506  5.98  29.1  281.  128.  20.1
## # ... with 6 more variables: nv2.5 <dbl>, v2.5 <dbl>, nv10 <dbl>, v10 <dbl>,
## #   ws <dbl>, wd <dbl>
```

Retain site name and site code:

```
# can also retain site code
timeAverage(dat, avg.time = "year", type = c("site", "code"))

## # A tibble: 6 x 17
## # Groups:   site, code [2]
##   site   code   date          o3    no2    co    so2   pm10   nox    no pm2.5
##   <fct> <fct> <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Lond~ KC1   2011-01-01 00:00:00 39.4  36.1  0.225  2.06  23.7  53.8  11.6
## 2 Lond~ KC1   2012-01-01 00:00:00 38.5  36.7  0.266  2.03  20.2  57.4  13.3
## 3 Lond~ KC1   2013-01-01 00:00:00 38.4  36.9  0.250  2.01  23.1  57.9  13.7
## 4 Lond~ MY1   2011-01-01 00:00:00 18.5  97.2  0.656  6.86  38.4  306.  137.
## 5 Lond~ MY1   2012-01-01 00:00:00 15.0  94.0  0.589  8.13  30.8  313.  143.
## 6 Lond~ MY1   2013-01-01 00:00:00 17.7  84.7  0.506  5.98  29.1  281.  128.
## # ... with 7 more variables: pm2.5 <dbl>, nv2.5 <dbl>, v2.5 <dbl>, nv10 <dbl>,
## #   v10 <dbl>, ws <dbl>, wd <dbl>
```

Average all data across sites (drops site and code):

```
timeAverage(dat, avg.time = "year")

## # A tibble: 3 x 15
##   date          o3    no2    co    so2   pm10   nox    no pm2.5 nv2.5
##   <dttm>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 2011-01-01 00:00:00 31.5  67.1  0.439  4.31  31.4  181.  74.9  20.5  16.9
## 2 2012-01-01 00:00:00 26.9  64.7  0.424  5.08  25.6  182.  76.7  18.1  14.9
## 3 2013-01-01 00:00:00 28.0  60.7  0.378  3.79  26.8  169.  70.7  17.4  14.1
## # ... with 5 more variables: v2.5 <dbl>, nv10 <dbl>, v10 <dbl>, ws <dbl>,
## #   wd <dbl>
```

`timeAverage` also works the other way in that it can be used to derive higher temporal resolution data e.g. hourly from daily data or 15-minute from hourly data. An example of usage would be the combining of daily mean particle data with hourly meteorological data. There are two ways these two data sets can be combined: either average the meteorological data to daily means or calculate hourly means from the particle data. The `timeAverage` function when used to ‘expand’ data in this way will repeat the original values the number of times required to fill the new time scale. In the example below we calculate 15-minute data from hourly data. As it can be seen, the first line is repeated four times and so on.

```

data15 <- timeAverage(mydata, avg.time = "15 min", fill = TRUE)
head(data15, 20)

## # A tibble: 20 x 10
##   date                  ws    wd  nox  no2   o3 pm10   so2   co pm25
##   <dttm>              <dbl> <int> <int> <int> <dbl> <dbl> <dbl> <dbl>
## 1 1998-01-01 00:00:00  0.6    280   285   39     1    29  4.72  3.37  NA
## 2 1998-01-01 00:15:00  0.6    280   285   39     1    29  4.72  3.37  NA
## 3 1998-01-01 00:30:00  0.6    280   285   39     1    29  4.72  3.37  NA
## 4 1998-01-01 00:45:00  0.6    280   285   39     1    29  4.72  3.37  NA
## 5 1998-01-01 01:00:00  2.16   230    NA    NA     NA   37  NA    NA    NA
## 6 1998-01-01 01:15:00  2.16   230    NA    NA     NA   37  NA    NA    NA
## 7 1998-01-01 01:30:00  2.16   230    NA    NA     NA   37  NA    NA    NA
## 8 1998-01-01 01:45:00  2.16   230    NA    NA     NA   37  NA    NA    NA
## 9 1998-01-01 02:00:00  2.76   190    NA    NA     3    34  6.83  9.60  NA
## 10 1998-01-01 02:15:00 2.76   190    NA    NA     3    34  6.83  9.60  NA
## 11 1998-01-01 02:30:00 2.76   190    NA    NA     3    34  6.83  9.60  NA
## 12 1998-01-01 02:45:00 2.76   190    NA    NA     3    34  6.83  9.60  NA
## 13 1998-01-01 03:00:00 2.16   170   493    52     3    35  7.66  10.2   NA
## 14 1998-01-01 03:15:00 2.16   170   493    52     3    35  7.66  10.2   NA
## 15 1998-01-01 03:30:00 2.16   170   493    52     3    35  7.66  10.2   NA
## 16 1998-01-01 03:45:00 2.16   170   493    52     3    35  7.66  10.2   NA
## 17 1998-01-01 04:00:00 2.4    180   468    78     2    34  8.07  8.91  NA
## 18 1998-01-01 04:15:00 2.4    180   468    78     2    34  8.07  8.91  NA
## 19 1998-01-01 04:30:00 2.4    180   468    78     2    34  8.07  8.91  NA
## 20 1998-01-01 04:45:00 2.4    180   468    78     2    34  8.07  8.91  NA

```

The `timePlot` can apply this function directly to make it very easy to plot data with different averaging times and statistics.

31.5 Calculating percentiles

`calcPercentile` makes it straightforward to calculate percentiles for a single pollutant. It can take account of different averaging periods, data capture thresholds — see Section 31.4 for more details. The function has the following options:

mydata A data frame of data with a `date` field in the format `Date` or `POSIXct`. Must have one variable to apply calculations to.

pollutant Name of variable to process. Mandatory.

avg.time Averaging period to use. See `timeAverage` for details.

percentile A vector of percentile values. For example `percentile = 50` for median values, `percentile = c(5, 50, 95)` for multiple percentile values.

data.thresh Data threshold to apply when aggregating data. See `timeAverage` for details.

start Start date to use - see `timeAverage` for details.

For example, to calculate the 25, 50, 75 and 95th percentiles of O₃ concentration by year:

```
calcPercentile(mydata, pollutant = "o3", percentile = c(25, 50, 75, 95),
               avg.time = "year")

##           date percentile.25 percentile.50 percentile.75 percentile.95
## 1 1998-01-01          2             4             7            16
## 2 1999-01-01          2             4             9            21
## 3 2000-01-01          2             4             9            22
## 4 2001-01-01          2             4            10            24
## 5 2002-01-01          2             4            10            24
## 6 2003-01-01          2             4            11            24
## 7 2004-01-01          2             5            11            23
## 8 2005-01-01          3             7            16            28
```

31.6 The `corPlot` function — correlation matrices

Understanding how different variables are related to one another is always important. However, it can be difficult to easily develop an understanding of the relationships when many different variables are present. One of the useful techniques used is to plot a *correlation matrix*, which provides the correlation between all pairs of data. The basic idea of a correlation matrix has been extended to help visualise relationships between variables by Friendly (2002) and Sarkar (2007).

The `corPlot` shows the correlation coded in three ways: by shape (ellipses), colour and the numeric value. The ellipses can be thought of as visual representations of scatter plot. With a perfect positive correlation a line at 45 degrees positive slope is drawn. For zero correlation the shape becomes a circle — imagine a ‘fuzz’ of points with no relationship between them.

With many different variables it can be difficult to see relationships between variables i.e. which variables tend to behave most like one another. For this reason hierarchical clustering is applied to the correlation matrices to group variables that are most similar to one another (if `cluster = TRUE`.)

It is also possible to use the `openair` type option to condition the data in many flexible ways, although this may become difficult to visualise with too many panels.

The `corPlot` function has the following options:

mydata A data frame which should consist of some numeric columns.

pollutants the names of data-series in `mydata` to be plotted by `corPlot`. The default option `NULL` and the alternative “all” use all available valid (numeric) data.

type `type` determines how the data are split i.e. conditioned, and then plotted. The default is will produce a single plot using the entire data. Type can be one of the built-in types as detailed in `cutData` e.g. “season”, “year”, “weekday” and so on. For example, `type = "season"` will produce four plots — one for each season.

It is also possible to choose `type` as another variable in the data frame. If that variable is numeric, then the data will be split into four quantiles (if possible) and labelled accordingly. If type is an existing character or factor variable, then those categories/levels will be used directly. This offers great flexibility for understanding the variation of different variables and how they depend on one another.

cluster Should the data be ordered according to cluster analysis. If `TRUE` hierarchical clustering is applied to the correlation matrices using `hclust`

to group similar variables together. With many variables clustering can greatly assist interpretation.

dendrogram	Should a dendrogram be plotted? When TRUE a dendrogram is shown on the right of the plot. Note that this will only work for type = "default" .
lower	Should only the lower triangle be plotted?
cols	Colours to be used for plotting. Options include “default”, “increment”, “heat”, “spectral”, “hue”, “greyscale” and user defined (see openColours for more details).
r.thresh	Values of greater than r.thresh will be shown in bold type. This helps to highlight high correlations.
text.col	The colour of the text used to show the correlation values. The first value controls the colour of negative correlations and the second positive.
auto.text	Either TRUE (default) or FALSE . If TRUE titles and axis labels will automatically try and format pollutant names and units properly e.g. by subscripting the ‘2’ in NO ₂ .
...	Other graphical parameters passed onto lattice:levelplot , with common axis and title labelling options (such as xlab , ylab , main) being passed via quickText to handle routine formatting.

An example of the **corPlot** function is shown in [Figure 31.2](#). In this Figure it can be seen the highest correlation coefficient is between PM₁₀ and PM_{2.5} ($r = 0.84$) and that the correlations between SO₂, NO₂ and NO_x are also high. O₃ has a negative correlation with most pollutants, which is expected due to the reaction between NO and O₃. It is not that apparent in [Figure 31.2](#) that the order the variables appear is due to their similarity with one another, through hierarchical cluster analysis. In this case we have chosen to also plot a *dendrogram* that appears on the right of the plot. Dendograms provide additional information to help with visualising how groups of variables are related to one another. Note that dendograms can only be plotted for **type = "default"** i.e. for a single panel plot.

Note also that the **corPlot** accepts a **type** option, so it possible to condition the data in many flexible ways, although this may become difficult to visualise with too many panels. For example:

```
corPlot(mydata, type = "season")
```

When there are a very large number of variables present, the **corPlot** is a very effective way of quickly gaining an idea of how variables are related. As an example (not plotted) it is useful to consider the hydrocarbons measured at Marylebone Road. There is a lot of information in the hydrocarbon plot (about 40 species), but due to the hierarchical clustering it is possible to see that isoprene, ethane and propane behave differently to most of the other hydrocarbons. This is because they have different (non-vehicle exhaust) origins. Ethane and propane results from natural gas leakage whereas isoprene is biogenic in origin (although some is from vehicle exhaust too). It is also worth considering how the relationships change between the species over the years as hydrocarbon emissions are increasingly controlled, or maybe the difference between summer and winter blends of fuels and so on.

```
corPlot(mydata, dendrogram = TRUE)
```

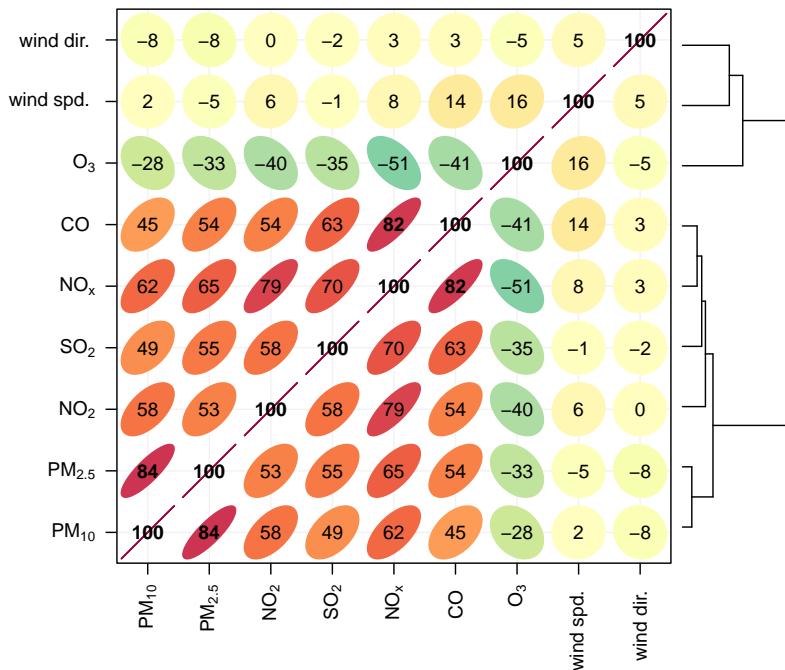


FIGURE 31.2 Example of a correlation matrix showing the relationships between variables.

```
hc <- importAURN(site = "my1", year = 2005, hc = TRUE)
## now it is possible to see the hydrocarbons that behave most
## similarly to one another
corPlot(hc)
```

31.7 Preparing data to compare sites, for model evaluation and intervention analysis

Many of the functions described have the potential to be extremely flexible. Mention has already been made of how to compare different sites in some of the functions. It was stated that the data had to be in a certain format for the functions to work. This section describes a few simple functions to do this — and more.

31.7.1 Intervention analysis

Another common scenario is that there is interest in showing plots by different time intervals *on the same scale*. There could be all sorts of reasons for wanting to do this. A classic example would be to show a before/after plot due to some intervention such as a low emission zone. Again, the function below exploits the flexible ‘site’ option available in many functions.

A small helper function `splitByDate` has been written to simplify chopping up a data set into different defined periods. The function takes three arguments: a data frame to process, a date (or dates) and labels for each period. If there was interest in looking at `mydata` before and after the 1st Jan 2003, it would be split into two periods (before that date and after). In other words, there will always be one more label than there is date. We have made the function easier to use for supplying dates. Dates can be accepted in the form ‘dd/mm/yyyy’ e.g. 13/04/1999 or as ‘yyyy-mm-dd’ e.g. ‘1999-04-13’.

```
polarAnnulus(mydata, pollutant = "no2", type = "split.by", period = "hour",
             layout = c(3, 1))
```

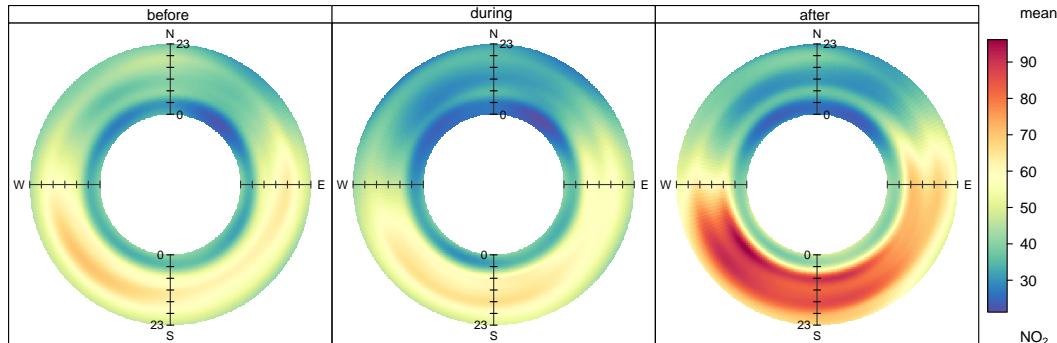


FIGURE 31.3 Example of processing data for use in the `polarAnnulus` function by time period for NO₂ concentrations at Marylebone Road.

The example below chops the data set up into three sections, called ‘before’, ‘during’ and ‘after’. This is done to show how more than one date can be supplied to the function.

```
mydata <- splitByDate(mydata, dates = c("1/1/2000", "1/3/2003"),
                       labels = c("before", "during", "after"))

head(mydata)

## # A tibble: 6 x 11
##   date                  ws    wd    nox    no2    o3    pm10    so2    co    pm25
##   <dttm>              <dbl> <int> <int> <int> <int> <dbl> <dbl> <dbl> <dbl>
## 1 1998-01-01 00:00:00  0.6    280    285    39     1    29  4.72  3.37    NA
## 2 1998-01-01 01:00:00  2.16   230     NA     NA     NA    37  NA    NA    NA
## 3 1998-01-01 02:00:00  2.76   190     NA     NA     3    34  6.83  9.60    NA
## 4 1998-01-01 03:00:00  2.16   170    493    52     3    35  7.66 10.2    NA
## 5 1998-01-01 04:00:00  2.4    180    468    78     2    34  8.07  8.91    NA
## 6 1998-01-01 05:00:00  3     190    264    42     0    16  5.50  3.05    NA
## # ... with 1 more variable: split.by <ord>

tail(mydata)

## # A tibble: 6 x 11
##   date                  ws    wd    nox    no2    o3    pm10    so2    co    pm25
##   <dttm>              <dbl> <int> <int> <int> <int> <dbl> <dbl> <dbl> <dbl>
## 1 2005-06-23 07:00:00  1.5    250    404    156     4    49  NA    1.81   28
## 2 2005-06-23 08:00:00  1.5    260    388    145     6    48  NA    1.64   26
## 3 2005-06-23 09:00:00  1.5    210    404    168     7    58  NA    1.29   34
## 4 2005-06-23 10:00:00  2.6    240    387    175    10    55  NA    1.29   34
## 5 2005-06-23 11:00:00  3.1    220    312    125    15    52  NA    1.29   33
## 6 2005-06-23 12:00:00  3.1    220    287    119    17    55  NA    1.29   35
## # ... with 1 more variable: split.by <ord>
```

As can be seen, there is a new field `split.by` (although the name can be set by the user), where at the beginning of the time series it is labelled ‘before’ and at the end it is labelled ‘after’. Now let us make a polar annulus plot showing the diurnal variation of NO₂ by wind direction:

In some cases it would make sense to have labels that refer to dates. Here is an example:

```
mydata <- splitByDate(mydata, dates = c("1/1/2000", "1/3/2003"),
                      labels = c("before Jan. 2000", "Jan. 2000 - Mar. 2003",
                                   "after Mar. 2003"))
```

31.7.2 Combining lots of sites

A typical example is that imported data have a date field and one or more pollutant fields from one of more sites in a series of columns. The aim would be, for example, to produce a series of plots by site for the same pollutant. If the data contains multiple pollutants and multiple sites, it makes sense to subset the data first.²⁰ For example, if a data frame `mydata` has fields ‘date’, ‘nox.site1’, ‘so2.site1’, ‘nox.site2’, ‘so2.site2’, then just working with the NO_x data can be done by:

```
subdata <- subset(mydata, select = c(date, nox.site1, nox.site2))
```

Rather than import new data, the code below first makes an artificial data set from which to work. In a real situation, the first few lines would not be needed.

²⁰Note that if you are able to use data from the AURN archive using the `importAURN` function, the data will already be in the correct format for direct use by many of the functions—although it may well be necessary to merge some meteorological data first.

```

## load reshape2 package if it is not already loaded...
library(reshape2)
## first make a subset of the data: date and nox, ws, wd
siteData <- subset(mydata, select = c(date, ws, wd, nox))

## look at the first few lines
head(siteData)

## # A tibble: 6 x 4
##   date                  ws     wd    nox
##   <dttm>           <dbl> <int> <int>
## 1 1998-01-01 00:00:00  0.6   280   285
## 2 1998-01-01 01:00:00  2.16  230    NA
## 3 1998-01-01 02:00:00  2.76  190    NA
## 4 1998-01-01 03:00:00  2.16  170   493
## 5 1998-01-01 04:00:00  2.4   180   468
## 6 1998-01-01 05:00:00  3     190   264

## rename the nox field to "site1"
names(siteData)[4] <- "site1"

## now make another field "site2" to be equal to half of site1
siteData$site2 <- siteData$site1 * 0.5
## end of making new data, now let's process it
head(siteData)

## # A tibble: 6 x 5
##   date                  ws     wd site1 site2
##   <dttm>           <dbl> <int> <int> <dbl>
## 1 1998-01-01 00:00:00  0.6   280   285  142.
## 2 1998-01-01 01:00:00  2.16  230    NA   NA
## 3 1998-01-01 02:00:00  2.76  190    NA   NA
## 4 1998-01-01 03:00:00  2.16  170   493  246.
## 5 1998-01-01 04:00:00  2.4   180   468  234
## 6 1998-01-01 05:00:00  3     190   264  132

## now we need to "stack" the data, ready for openair functions
## use the melt function in the reshape2 package (loaded with openair)
siteData <- melt(siteData, measure.vars = c("site1", "site2"))
head(siteData)

##          date  ws  wd variable value
## 1 1998-01-01 0.6 280 site1  285
## 2 1998-01-01 2.2 230 site1    NA
## 3 1998-01-01 2.8 190 site1    NA
## 4 1998-01-01 2.2 170 site1  493
## 5 1998-01-01 2.4 180 site1  468
## 6 1998-01-01 3.0 190 site1  264

## change the variable names (one of them has to be "site")
names(siteData)[4:5] <- c("site", "nox")

```

Now it is possible to run many **openair** functions on this dataset. In this case, let us consider a **polarPlot**. Note that this process would work with many more sites than shown here. Note, however, many functions such as **polarPlot** accept multiple pollutants and the **importAURN** and **importKCL** format multiple site data directly and no additional work is required by the user.

```
polarPlot(siteData, pollutant = "nox", type = "site")
```

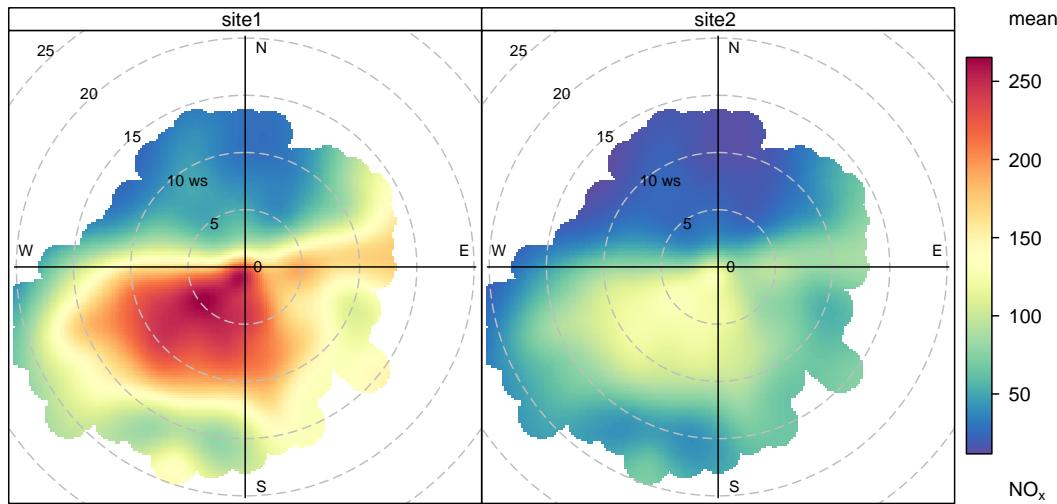


FIGURE 31.4 Example of processing data for use in the `polarPlot` function by site.

Acknowledgements

We are very grateful to a range of organisations that have so far seen this initiative as something worth supporting financially. These include:

- Dr Karl Ropkins of the Institute for Transport Studies at the University of Leeds for his contributions during the NERC project.
- The initial funding supplied by the Faculty of Environment at the University of Leeds to help develop a course on R to analyse and understand air pollution data.
- Sefton Council for their direct funding of data analysis in their borough as part of the Beacon air quality scheme.
- AEA.
- North Lincolnshire Council.
- Defra as part of their work through the AURN.
- The Natural Environment Research Council (NERC) Knowledge Transfer grant NE/G001081/1.

This work would not be possible without the incredible individuals who have given their time freely to develop the R system. This includes those in the R-Core Development Team and all those that contribute to its development (R Core Team 2015).

Further information

For any enquiries related to this document or the `openair` package, please use the contact details below. Please contact us regarding any of the following: bug reports, suggestions for existing functions, suggestions for new functions and offer of code contributions. When reporting potential bugs, it is helpful (sometimes essential) to

submit a reproducible example, which would normally require sending a description of the problem and the data set used. Also, we are interested in developing further funded case studies.

David Carslaw
 Department of Chemistry
 University of York
 Heslington
 York
 YO10 5DD
 UK

e-mail: <mailto:david.carslaw@york.ac.uk>

References

- APPLEQUIST, S. (2012). "Wind Rose Bias Correction". In: *Journal of Applied Meteorology and Climatology* 51.7, pp. 1305–1309 (cit. on pp. 87, 91).
- AQEG (2008). *Trends in primary nitrogen dioxide in the UK*. Air Quality Expert Group. Report prepared by the Air Quality Expert Group for the Department for Environment, Food, Rural Affairs; Scottish Executive; Welsh Assembly Government; and Department of the Environment in Northern Ireland. (cit. on p. 234).
- ARA BEGUM, B., E. KIM, C.-H. JEONG, D.-W. LEE and P. K. HOPKE (2005). "Evaluation of the potential source contribution function using the 2002 Quebec forest fire episode". In: *Atmospheric Environment* 39.20, pp. 3719–3724. DOI: [10.1016/j.atmosenv.2005.03.008](https://doi.org/10.1016/j.atmosenv.2005.03.008) (cit. on p. 205).
- ASHBAUGH, L. L., W. C. MALM and W. Z. SADEH (1985). "A residence time probability analysis of sulfur concentrations at grand Canyon National Park". In: *Atmospheric Environment* (1967) 19.8, pp. 1263–1270. DOI: [10.1016/0004-6981\(85\)90256-2](https://doi.org/10.1016/0004-6981(85)90256-2) (cit. on pp. 100, 118).
- CARSLAW, D. C. (2005). "Evidence of an increasing NO₂/NO_x emissions ratio from road traffic emissions". In: *Atmospheric Environment* 39.26, pp. 4793–4802 (cit. on p. 234).
- CARSLAW, D. C. and S. D. BEEVERS (2004). "Investigating the potential importance of primary NO₂ emissions in a street canyon". In: *Atmospheric Environment* 38.22, pp. 3585–3594 (cit. on p. 234).
- (2005). "Estimations of road vehicle primary NO₂ exhaust emission fractions using monitoring data in London". In: *Atmospheric Environment* 39.1, pp. 167–177 (cit. on p. 234).
 - (2013). "Characterising and understanding emission sources using bivariate polar plots and k-means clustering". In: *Environmental Modelling & Software* 40.0, pp. 325–329. DOI: [10.1016/j.envsoft.2012.09.005](https://doi.org/10.1016/j.envsoft.2012.09.005) (cit. on pp. 108, 122).
- CARSLAW, D. C., S. D. BEEVERS, K. ROPKINS and M. C. BELL (2006). "Detecting and quantifying aircraft and other on-airport contributions to ambient nitrogen oxides in the vicinity of a large international airport". In: *Atmospheric Environment* 40.28, pp. 5424–5434 (cit. on pp. 108, 109).
- CARSLAW, D. C., S. D. BEEVERS and J. E. TATE (2007). "Modelling and assessing trends in traffic-related emissions using a generalised additive modelling approach". In: *Atmospheric Environment* 41.26, pp. 5289–5299 (cit. on p. 155).

- CARSLAW, D. C. and N. CARSLAW (2007). "Detecting and characterising small changes in urban nitrogen dioxide concentrations". In: *Atmospheric Environment* 41.22, pp. 4723–4733 (cit. on p. 234).
- CARSLAW, D. C. and K. ROPKINS (2012). "openair — An R package for air quality data analysis". In: *Environmental Modelling & Software* 27–28.0, pp. 52–61. DOI: [10.1016/j.envsoft.2011.09.008](https://doi.org/10.1016/j.envsoft.2011.09.008) (cit. on p. 7).
- CHAMBERS, J. M. (2007). *Software for Data Analysis: Programming with R*. ISBN 978-0-387-75935-7. New York: Springer (cit. on p. 8).
- CHATFIELD, C. (2004). *The analysis of time series : an introduction / Chris Chatfield*. 6th ed. Boca Raton, FL ; London : Chapman & Hall/CRC (cit. on p. 263).
- CLAPP, L. J. and M. E. JENKIN (2001). "Analysis of the relationship between ambient levels of O₃, NO₂ and NO as a function of NO_x in the UK". In: *Atmospheric Environment* 35.36, pp. 6391–6405 (cit. on p. 234).
- CLEVELAND, W. S. (1985). *The elements of graphing data*. Wadsworth Publ. Co. Belmont, CA, USA (cit. on p. 50).
- (1993). *Visualizing Data*. Hobart Press, Summit, NJ (cit. on p. 50).
- COMEAP (2011). *Review of the UK Air Quality Index: A report by the Committee on the Medical Effects of Air Pollutants* (cit. on p. 146).
- CRAWLEY, M. (2007). *The R book*. Wiley (cit. on p. 12).
- DALGAARD, P. (2008). *Introductory Statistics with R*. 2nd. ISBN 978-0-387-79053-4. Springer, p. 380 (cit. on p. 12).
- DAVISON, A. C. and D. HINKLEY (1997). *Bootstrap methods and their application*. (Anthony Christopher). Cambridge ; New York, NY, USA : Cambridge University Press (cit. on p. 259).
- DROPOPO, J. G. and B. A. NAPIER (2008). "Wind direction bias in generating wind roses and conducting sector-based air dispersion modeling". In: *Journal of the Air & Waste Management Association* 58.7, pp. 913–918 (cit. on p. 87).
- EFRON, B. and R. TIBSHIRANI (1993). *An Introduction to the Bootstrap*. Chapman & Hall (cit. on p. 259).
- FLEMING, Z. L., P. S. MONKS and A. J. MANNING (2012). "Review: Untangling the influence of air-mass history in interpreting observed atmospheric composition". In: *Atmospheric Research* 104-105, pp. 1–39. DOI: [10.1016/j.atmosres.2011.09.009](https://doi.org/10.1016/j.atmosres.2011.09.009) (cit. on pp. 201, 205).
- FRIENDLY, M. (2002). "Corrgrams: Exploratory Displays for Correlation Matrices". In: *The American Statistician* 56.4, pp. 316–325 (cit. on p. 247).
- HASTIE, T. J. and R. TIBSHIRANI (1990). *Generalized additive models*. London: Chapman and Hall (cit. on p. 109).
- HELSEL, D. and R. HIRSCH (2002). *Statistical Methods in Water Resources*. US Geological Survey (cit. on p. 13).
- HENRY, R., G. A. NORRIS, R. VEDANTHAM and J. R. TURNER (2009). "Source Region Identification Using Kernel Smoothing". In: *Environmental Science & Technology* 43.11, 4090–4097. DOI: [{10.1021/es8011723}](https://doi.org/10.1021/es8011723) (cit. on p. 91).
- HIRSCH, R. M., J. R. SLACK and R. A. SMITH (1982). "Techniques Of Trend Analysis For Monthly Water-Quality Data". In: *Water Resources Research* 18.1. ISI Document Delivery No.: NC504, pp. 107–121 (cit. on p. 147).
- HSU, Y.-K., T. M. HOLSEN and P. K. HOPKE (2003). "Comparison of hybrid receptor models to locate PCB sources in Chicago". In: *Atmospheric Environment* 37.4, pp. 545–562. DOI: [10.1016/S1352-2310\(02\)00886-5](https://doi.org/10.1016/S1352-2310(02)00886-5) (cit. on pp. 205, 206).
- KUNSCH, H. R. (1989). "The jackknife and the bootstrap for general stationary observations". In: *Annals of Statistics* 17.3, pp. 1217–1241 (cit. on pp. 148, 260).

- LEGATES, D. R. and G. J. McCABE JR (1999). "Evaluating the use of "goodness-of-fit" measures in hydrologic and hydroclimatic model validation". In: *Water Resources Research* 35.1, pp. 233–241 (cit. on p. 214).
- LEGATES, D. R. and G. J. McCABE (2012). "A refined index of model performance: a rejoinder". In: *International Journal of Climatology* (cit. on p. 214).
- LUPU, A. and W. MAENHAUT (2002). "Application and comparison of two statistical trajectory techniques for identification of source regions of atmospheric aerosol species". In: *Atmospheric Environment* 36, pp. 5607–5618 (cit. on p. 206).
- MAINDONALD, J. and J. BRAUN (2007). *Data Analysis and Graphics Using R: An Example-based Approach*. Cambridge University Press (cit. on p. 12).
- MATLOFF, N. (2011). *The Art of R Programming*. No Starch Press (cit. on p. 13).
- MCHUGH, C. A., D. J. CARRUTHERS and H. A. EDMUND (1997). "ADMS and ADMS-Urban". In: *International Journal of Environment and Pollution* 8.3-6, pp. 438–440 (cit. on p. 59).
- PEKNEY, N. J., C. I. DAVIDSON, L. ZHOU and P. K. HOPKE (2006). "Application of PSCF and CPF to PMF-Modeled Sources of PM 2.5 in Pittsburgh". In: *Aerosol Science and Technology* 40.10, pp. 952–961. doi: [10.1080/02786820500543324](https://doi.org/10.1080/02786820500543324) (cit. on p. 205).
- R CORE TEAM (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria (cit. on p. 253).
- SARKAR, D. (2007). *Lattice Multivariate Data Visualization with R*. ISBN 978-0-387-75968-5. New York: Springer (cit. on pp. 12, 247).
- SEIBERT, P., H. KROMP-KOLB, U. BALTENSPERGER and D. JOST (1994). "Trajectory analysis of high-alpine air pollution data". In: *NATO Challenges of Modern Society* 18, pp. 595–595 (cit. on p. 206).
- SEN, P. K. (1968). "Estimates of regression coefficient based on Kendall's tau". In: *Journal of the American Statistical Association* 63(324), pp. 1379–1389 (cit. on p. 147).
- SPECTOR, P. (2008). *Data Manipulation with R*. ISBN 978-0-387-74730-9. New York: Springer (cit. on p. 12).
- TAYLOR, K. E. (2001). "Summarizing multiple aspects of model performance in a single diagram". In: *Journal of Geophysical Research* 106.D7, pp. 7183–7192 (cit. on p. 218).
- THEIL, H. (1950). "A rank invariant method of linear and polynomial regression analysis, I, II, III". In: *Proceedings of the Koninklijke Nederlandse Akademie Wetenschappen, Series A – Mathematical Sciences* 53, pp. 386–392, 521–525, 1397–1412 (cit. on p. 147).
- TUKEY, J. (1977). *Exploratory data analysis*. Addison-Wesley Series in Behavioral Science: Quantitative Methods, Reading, Mass.: Addison-Wesley, 1977 (cit. on p. 9).
- URIA-TELLAETXE, I. and D. C. CARSLAW (2014). "Conditional bivariate probability function for source identification". In: *Environmental Modelling & Software* 59, pp. 1–9. doi: [10.1016/j.envsoft.2014.05.002](https://doi.org/10.1016/j.envsoft.2014.05.002) (cit. on pp. 108, 119, 120).
- WESTMORELAND, E. J., N. CARSLAW, D. C. CARSLAW, A. GILLAHI and E. BATES (2007). "Analysis of air quality within a street canyon using statistical and dispersion modelling techniques". In: *Atmospheric Environment* 41.39, pp. 9195–9205 (cit. on p. 108).
- WILCOX, R. R. (2010). *Fundamentals of Modern Statistical Methods: Substantially Improving Power and Accuracy*. 2nd. Springer New York (cit. on p. 147).
- WILKS, D. S. (2005). *Statistical Methods in the Atmospheric Sciences, Volume 91, Second Edition (International Geophysics)*. 2nd ed. Academic Press (cit. on p. 225).
- WILLMOTT, C. J., S. M. ROBESON and K. MATSUURA (2011). "A refined index of model performance". In: *International Journal of Climatology* (cit. on p. 214).

- WOOD, S. N. (2006). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC (cit. on pp. 109, 262).
- XIE, Y. (2013a). *Dynamic Documents with R and knitr*. ISBN 978-1482203530. Chapman and Hall/CRC (cit. on p. 9).
- (2013b). *knitr: A general-purpose package for dynamic report generation in R*. R package version 1.1 (cit. on p. 9).
- YU, K., Y. CHEUNG, T. CHEUNG and R. HENRY (2004). “Identifying the impact of large urban airports on local air quality by nonparametric regression”. In: *Atmospheric Environment* 38.27, pp. 4501–4507 (cit. on p. 108).

A Installing and maintaining R

A.1 Downloading and installing R

R can be downloaded from <http://www.r-project.org>, shown in [Figure A.1](#). On the left hand side there is a link to the *Comprehensive R Archive Network* (CRAN), which provides links to local repositories where the software can be downloaded from. The one most relevant to the UK is hosted by the University of Bristol. Follow the link and choose to download ‘Precompiled binary distributions’ from the Download and Install R section. Most likely you will want the Windows version, but versions for Linux and Apple Mac are also available. This will provide an executable file that can be downloaded (something like R-2.12.0-win.exe). A direct link from the UK to the Windows download file is <http://www.stats.bris.ac.uk/R/bin/windows/base/>

Important Information — Internet connections for Windows users

Note that R set up and maintenance works best when users have direct Internet access. Many users that use R through their organisation’s computers may need to install R slightly differently due to the proxy settings used. Rather than accept all the defaults during installation, at the Setup screen choose not to accept the defaults and when offered, choose ‘Internet2’ as the Internet option. This will force R to use the same proxy settings used by Internet Explorer. The defaults for all other options can be accepted.

The other issue on Windows systems is Administrator rights. When it comes to installing or updating packages it may be necessary to ‘Run as Administrator’. To do this, go to R on the Windows menu and right-click, then choose to ‘Run as administrator’.

The installation of R is straightforward. Most users can happily accept all the defaults. However, see box for installation information for users installing R on an organisation’s computer system. Choose where you want to install it (the default is usually appropriate for most systems). It is a good idea to install it where you have write privileges — probably more of an issue for Windows Vista/7 users. Many more details are provided on the R website regarding installation on specific systems.

...then accept all the defaults unless you have a preference for an alternative option.

A.2 Maintenance

One of the advantages of R compared with most commercially available software is that it is regularly updated; fixing any bugs and adding new features. New versions tend to be released every six months and it is wise to keep the version as up to date as possible. You can check the website occasionally, or sign up to be alerted to new releases (recommended). Often, as new versions of the base system are released, certain R packages will also require updating. If there is incompatibility between the base system and an R package you will likely receive a warning message. To keep the packages up to date select Packages | Update packages ..., which will prompt you to choose a CRAN mirror (choose the UK one again). This will check to see if all your packages are up to date, and if not, automatically download and install more recent versions. Note the information about R in the box where it may be necessary to run R as an administrator.

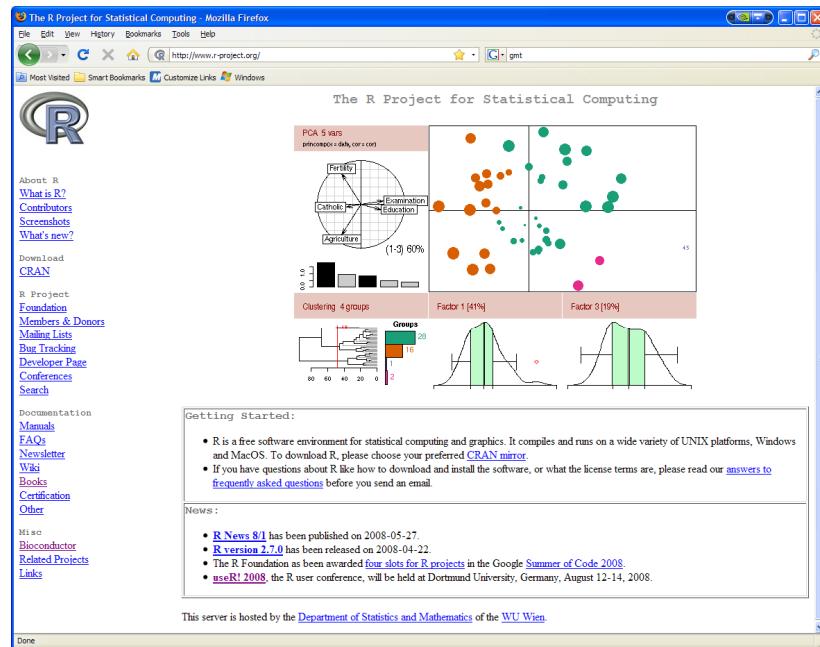


FIGURE A.1 The R-project web pages.

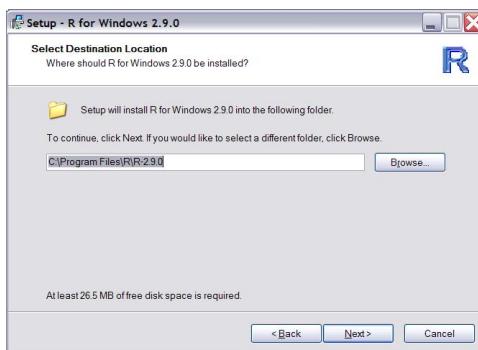


FIGURE A.2 Choose where R should be installed.

B Bootstrap estimates of uncertainty

B.1 The bootstrap

The *bootstrap* is a data-based simulation method for analysing data, including hypothesis testing, standard error and confidence interval estimation. It involves repeatedly drawing random samples from the original data, with replacement (see EFRON and TIBSHIRANI (1993) and DAVISON and HINKLEY (1997) for a detailed history and examples of use of the bootstrap). Each *bootstrap sample* is the same size as the original sample. The ‘with replacement’ bit is important. Sampling with replacement means that after we randomly draw an observation from the original sample we put it back before drawing the next observation i.e. it is possible to draw the same sample more than once. In fact, on average, 37 % of data will not be sampled each time. If one sampled without replacement it would be equivalent to just shuffling the data and no new information is available. Typically, 100s or 1000s of samples are required in order to derive reliable statistics.

The term bootstrap derives from the phrase “to pull oneself up by one’s bootstraps”.

The phrase is based on one of the eighteenth century *Adventures of Baron Munchausen* by Rudolph Erich Raspe. The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps! In a statistical sense it is meant to convey the idea of generating ‘new’ data from the original data set itself, which seems like an implausible thing to do, but has been shown to be valid.

When the bootstrap was discovered in the 1970s it was difficult to apply to many practical problems because computers were not powerful enough to carry out such repetitive and intensive calculations. However, computers are now sufficiently powerful to allow these methods to be used (in most circumstances) easily. This section does not aim to provide an in-depth consideration of statistics and justify the use of these methods, but rather aims to provide some background in their use in **openair**.

When used to estimate confidence intervals, the bootstrap sampling will yield, say, 1000 estimates of the statistic of interest e.g. the slope of a trend. This distribution could be highly skewed and this is one of the principal advantages of the technique: normality is not required. We now have a 1000 bootstrap samples, and 1000 estimates of the statistic of interest, one from each bootstrap sample. If these 1000 bootstrap samples are ordered in increasing value, a bootstrap 95 % confidence interval for the mean would be from the 25th to the 975th largest values. Sometimes, uncertainty estimates are not symmetrical. For example, it may not be possible to report an uncertainty as 100 ± 12 , but $100 [87, 121]$, where 87 and 121 are the lower and upper 95 % confidence intervals, respectfully.

B.2 The block bootstrap

The basic bootstrap assumes that data are independent. However, in time series this is rarely the case due to *autocorrelation* when consecutive points in time are related to one another. For example, for data with a strong seasonal effect, the month of January may tend to have higher values than other months. These effects can, however, be difficult to characterise and model. The motivation for accounting for autocorrelation in this project is mostly to ensure that uncertainty estimates in trends and other statistics are not overly optimistic, which would generally be the case if autocorrelation was not accounted for. These effects can be accounted for by ensuring that the random sampling captures the correlation structure of the data using a *block bootstrap* (KUNSCH 1989). The idea is that if data (or residuals from a model) are sampled in small blocks, the correlation structure is retained, provided there is not significant correlation between the blocks.

The following Figures highlight the importance of accounting for autocorrelation using the **TheilSen** function. [Figure B.1](#) shows a simulated time series comprising of a linear trend + some random noise with no autocorrelation. The Sen-Theil slope and the slope uncertainties are shown (and the slope statistics upper left). By contrast, [Figure B.2](#) shows a similar series but with a high autocorrelation value of 0.8. In this time series it is possible to see that the values seem to fall and rise in ‘chunks’, indicative of autocorrelation. In this plot, no account has been taken of autocorrelation and the uncertainty in the slope is very similar to [Figure B.1](#).

However, if autocorrelation is accounted for using the data shown in [Figure B.2](#), the uncertainty in the slope increases markedly, as shown in [Figure B.3](#). Instead of the 95 % confidence intervals ranging from 0.42–0.54, they now range from 0.36–0.61 — approximately double the uncertainty of the case where no account is taken of autocorrelation.

The block bootstrap has also been applied to models e.g. the Generalized Additive

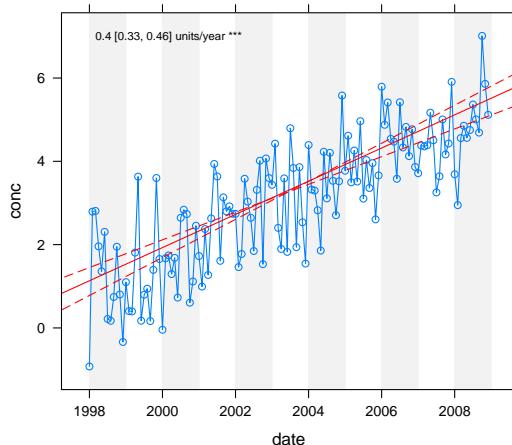


FIGURE B.1 AR1 time random series where the autocorrelation coefficient is zero.

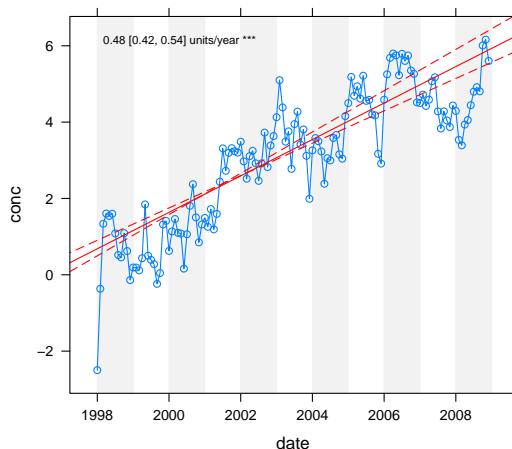


FIGURE B.2 AR1 time random series where the autocorrelation coefficient is 0.8. The uncertainty in the Sen-Theil slope estimate does not account for autocorrelation.

Model (GAM) used in the `smoothTrend` function. There are two options here: the observations can be sampled and the models run many times (called *case resampling*), or the residuals from the model can be sampled and added to the model predictions to make ‘new’ input data and run many times (called *residual resampling*). There are pros and cons with each approach, but often the two methods yield similar results. In the case of a GAM (or specifically the `mgcv` package), which uses cross-validation for model fitting, having duplicate samples through bootstrapping would seem to make little sense. The approach adopted here therefore is to use residual resampling. The effect of taking account of autocorrelation often (but not always) is an increase in the predicted uncertainty intervals, and a smooth function that is less ‘wiggly’ than that derived by not accounting for autocorrelation.

A more ‘robust’ approach is outlined in Appendix C, where models are used to describe the correlation structure of the data.

Clearly, the importance of these issues is data-dependent and there are other issues to consider too. However, if one is interested in drawing important inferences from data, then it would seem wise to account for these effects. It should be noted that these issues are an area of active research and will be revisited from time to time with the aim of improving the robustness of the techniques used.

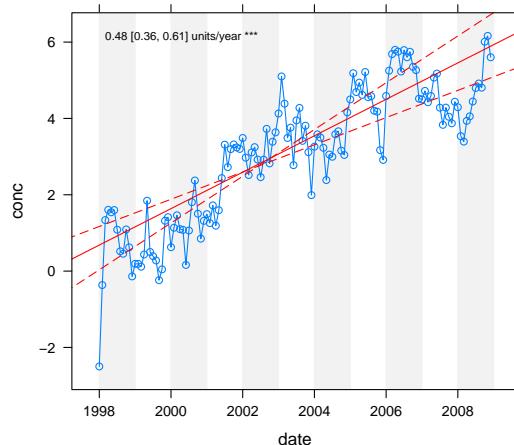


FIGURE B.3 AR1 time random series where the autocorrelation coefficient is 0.8. The uncertainty in the Sen-Theil slope estimate does account for autocorrelation.

C A closer look at trends

Understanding trends is a core component of air quality and the atmospheric sciences in general. `openair` provides two main functions for considering trends (`smoothTrend` and `TheilSen`), the latter useful for linear trend estimates. Understanding trends and quantifying them robustly is not so easy and careful analysis would treat each time series individually and consider a wide range of diagnostics. In this section we take advantage of some of the excellent capabilities that R has to consider fitting trend models. Experience with real atmospheric composition data shows that trends are rarely linear, which is unfortunate given how much of statistics has been built around the linear model.

Generalized Additive Models (GAMs) offer a flexible approach to calculating trends and in particular, the `mgcv` package contains many functions that are very useful for such modelling. Some of the details of this type of model are presented in Wood (2006) and the `mgcv` package itself.

The example considered is 23 years of O₃ measurements at Mace Head on the West coast of Ireland. The example shows the sorts of steps that might be taken to build a model to explain the trend. The data are first imported and then the year, month and ‘trend’ estimated. Note that ‘trend’ here is simply a decimal date that can be used to construct various explanatory models.

First we import the data:

```
library(mgcv)
dat <- importAURN(site = "mh", year = 1988:2010)

## calculate monthly means
monthly <- timeAverage(dat, avg.time = "month")
## now calculate components for the models
monthly$year <- as.numeric(format(monthly$date, "%Y"))
monthly$month <- as.numeric(format(monthly$date, "%m"))
monthly <- transform(monthly, trend = year + (month - 1) / 12)
```

It is *always* a good idea to plot the data first:

Figure C.1 shows that there is a clear seasonal variation in O₃ concentrations, which is certainly expected. Less obvious is whether there is a trend.

Even though it is known there is a seasonal signal in the data, we will first of all

```
timePlot(monthly, pollutant = "o3")
```

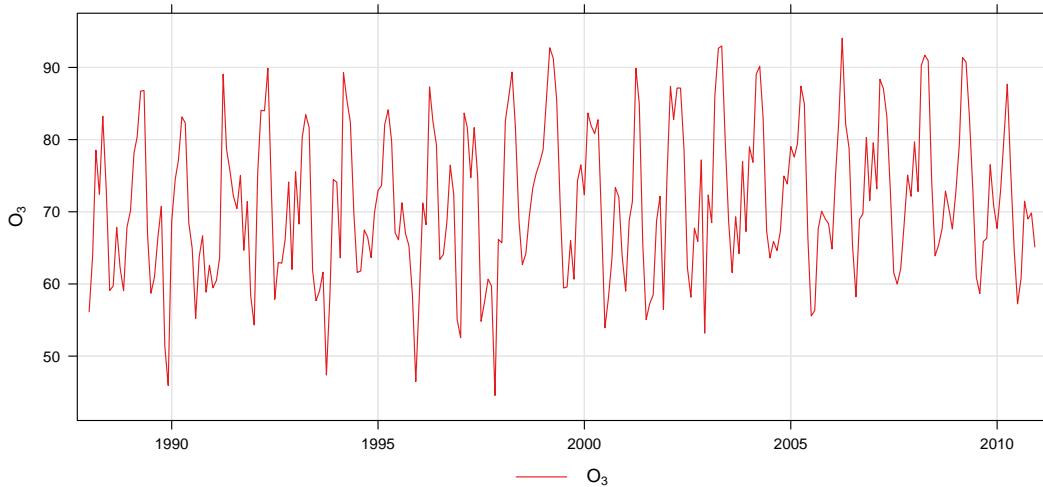


FIGURE C.1 Monthly mean O₃ concentrations at Mace Head, Ireland (1998–2010).

ignore it and build a simple model that only has a trend component (model M0).

```
M0 <- gam(o3 ~ s(trend), data = monthly)
summary(M0)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## o3 ~ s(trend)
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 71.34      0.62     115   <2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##          edf Ref.df F p-value
## s(trend) 1     1 6.96 0.0088 **
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) = 0.0212 Deviance explained = 2.48%
## GCV = 106.82 Scale est. = 106.04 n = 276
```

This model only explains about 2% of the variation as shown by the adjusted r^2 . More of a problem however is that no account has been taken of the seasonal variation. An easy way of seeing the effect of this omission is to plot the autocorrelation function (ACF) of the residuals, shown in Figure C.2. This Figure clearly shows the residuals have a strong seasonal pattern. CHATFIELD (2004) provides lots of useful information on time series modelling.

A refined model should therefore take account of the seasonal variation in O₃ concentrations. Therefore, we add a term taking account of the seasonal variation. Note also that we choose a cyclic spline for the monthly component (**bs** = "cc"), which

```
acf(residuals(M0))
```

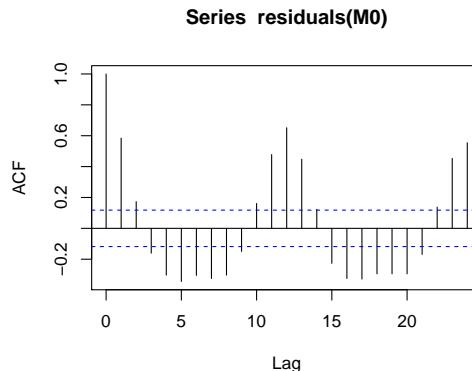


FIGURE C.2 ACF for the residuals of model M0.

joins the first and last points i.e. January and December.

```
M1 <- gam(o3 ~ s(trend) + s(month, bs = "cc"), data = monthly)
summary(M1)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## o3 ~ s(trend) + s(month, bs = "cc")
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 71.343     0.374    191   <2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##          edf Ref.df F p-value
## s(trend) 1.22    1.4 18.0 1.1e-05 ***
## s(month) 6.11    8.0 59.4 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) = 0.644 Deviance explained = 65.4%
## GCV = 39.766 Scale est. = 38.566 n = 276
```

Now we have a model that explains much more of the variation with an r^2 of 0.65. Also, the p-values for the trend and seasonal components are both highly statistically significant. Let's have a look at the separate components for trend and seasonal variation:

The seasonal component shown in Figure C.4 clearly shows the strong seasonal effect on O_3 at this site (peaking in April). The trend component is actually linear in this case and could be modelled as such. This model looks much better, but as is often the case autocorrelation could remain important. The ACF is shown in Figure C.5 and shows there is still some short-term correlation in the residuals.

Note also that there are other diagnostic tests one should consider when comparing these models that are not shown here e.g. such as considering the normality of the residuals. Indeed a consideration of the residuals shows that the model fails to some extent in explaining the very low values of O_3 , which can be seen in Figure C.1. These

```
plot.gam(M1, select = 1, shade = TRUE)
```

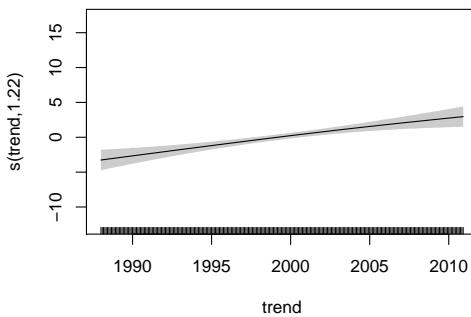


FIGURE C.3 The trend component of model M1.

```
plot.gam(M1, select = 2, shade = TRUE)
```

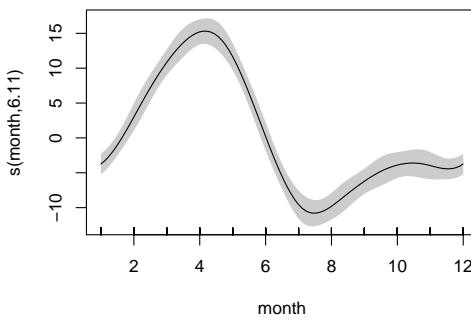


FIGURE C.4 The seasonal component of model M1.

few points (which skew the residuals) may well be associated with air masses from the polluted regions of Europe. Better and more useful models would likely be possible if the data were split my airmass origin, which is something that will be returned to when **openair** includes a consideration of back trajectories.

Further tests, also considering the partial autocorrelation function (PACF) suggest

```
acf(residuals(M1))
```

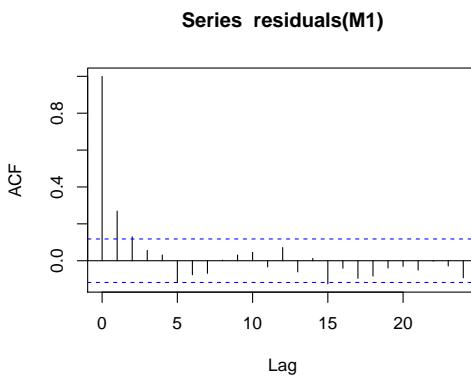


FIGURE C.5 ACF for the residuals of model M1.

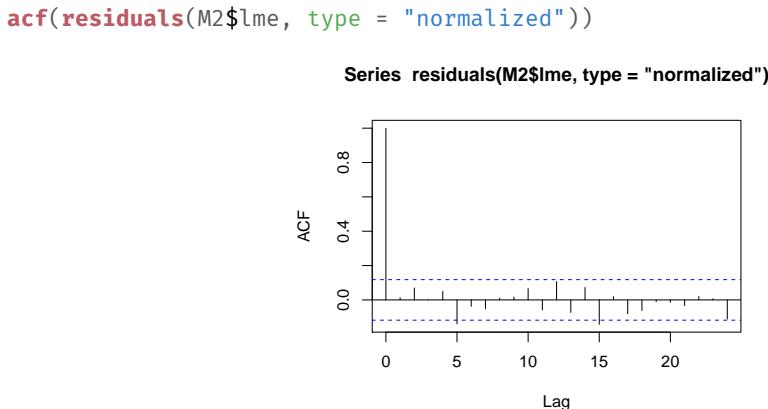


FIGURE C.6 ACF for the residuals of model M2.

that an AR1 model is suitable for modelling this short-term autocorrelation. This is where modelling using a GAMM (Generalized Additive Mixed Model) comes in because it is possible to model the short-term autocorrelation using a linear mixed model. The **gamm** function uses the package **nlme** and the Generalized Linear Mixed Model (GLMM) fitting routine. In the M2 model below the correlation structure is considered explicitly.

```
M2 <- gamm(o3 ~ s(month, bs = "cc") + s(trend), data = monthly,
            correlation = corAR1(form = ~ month | year))
summary(M2$gam)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## o3 ~ s(month, bs = "cc") + s(trend)
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 71.316     0.493    145   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##          edf Ref.df   F p-value
## s(month) 6.91      8 42.2 < 2e-16 ***
## s(trend) 1.00      1 15.1 0.00013 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.643
## Scale est. = 38.863   n = 276
```

The ACF plot is shown in [Figure C.6](#) and shows that the autocorrelation has been dealt with and we can be rather more confident about the trend component (not plotted). Note that in this case we need to use the normalized residuals to get residuals that take account of the fitted correlation structure.

Note that model M2 assumes that the trend and seasonal terms vary independently of one another. However, if the seasonal amplitude and/or phase change over time then a model that accounts for the interaction between the two may be better. Indeed, this

```

plot(M3$gam, select = 2, pers = TRUE, theta = 225, phi = 10, ticktype = "detailed")

## Error in plot(M3$gam, select = 2, pers = TRUE, theta = 225, phi = 10, :
## object 'M3' not found

```

FIGURE C.7 Plot showing the two-way interaction between the trend and seasonal components.

does seem to be the case here, as shown by the improved fit of the model below. This model uses a tensor product smooth (`te`) and the reason for doing this and not using an isotropic smooth (`s`) is that the trend and seasonal components are essentially on different scales. We would not necessarily want to apply the same level of smoothness to both components. An example of covariates on the same scale would be latitude and longitude.

```

M3 <- gamm(o3 ~ s(month, bs = "cc") + te(trend, month), data = monthly,
            correlation = corAR1(form = ~ month | year))

## Error in gamm(o3 ~ s(month, bs = "cc") + te(trend, month), data = monthly,
## : could not find function "gamm"

summary(M3$gam)

## Error in summary(M3$gam): object 'M3' not found

```

It becomes a bit more difficult to plot the two-way interaction between the trend and the month, but it is possible with a surface as shown in Figure C.7. This plot shows for example that during summertime the trends component varies little. However for the autumn and winter months there has been a much greater increase in the trend component for O₃.

While there have been many steps involved in this short analysis, the data at Mace Head are not typical of most air quality data observed, say in urban areas. Much of the data considered in these areas does not appear to have significant autocorrelation in the residuals once the seasonal variation has been accounted for, therefore avoiding the complexities of taking account of the correlation structure of the data. It may be for example that sites like Mace Head and a pollutant such as O₃ are much more prone to larger scale atmospheric processes that are not captured by these models.

D Production of HYSPLIT trajectory files

As discussed in Section 26, `openair` can import pre-calculated trajectory data for specified locations. The data are stored on a Ricardo webserver to make it easy to import 96-hour back trajectory data. Several users have requested how they can run HYSPLIT themselves e.g. for different trajectory start heights or for many locations. This section provides the code necessary to run the HYSPLIT model. The code below assumes that full years are run, but it could be adopted for shorter periods. There are three main parts to producing trajectory files:

1. Download and install the NOAA Hysplit model, somewhere with write access (see below).
2. Download the monthly meteorological (.gbl) files also from the NOAA website.
3. Obtain the code to run Hysplit.

To run back trajectories it is necessary to download the meteorological data files. The easiest way to download the meteorological files is using the function below.

```
getMet <- function (year = 2013, month = 1, path_met = "~/TrajData/") {

  for (i in seq_along(year)) {
    for (j in seq_along(month)) {
      download.file(url = paste0("ftp://arlftp.arlhq.noaa.gov/archives/reanalysis/RP",
                                    year[i], sprintf("%02d", month[j]), ".gbl"),
                     destfile = paste0(path_met, "RP", year[i],
                                       sprintf("%02d", month[j]), ".gbl"), mode = "wb")
    }
  }
}
```

The function will download monthly met files (each about 120 MB) to the chosen directory. Note that the met data files only need be downloaded once. For example, to download files for 2013:

```
getMet(year = 2013, month = 1:12)
```

The original functions have been modified by Stuart Grange (PhD student at the University of York). It is first necessary to ensure that the **stringr** and **devtools** packages are installed. The latter is needed to load some R functions stored as a *Github gist* (some code that can be shared publicly).

```
library(devtools)
source_gist("https://gist.github.com/davidcarslaw/c67e33a04ff6e1be0cd7357796e4bdf5",
            filename = "run_hysplit.R")

## Sourcing https://gist.githubusercontent.com/davidcarslaw/c67e33a04ff6e1be0cd7357796e4bdf5/
## SHA-1 hash of file is e1d37c75528ac18924fd0c6902318702301fa591
```

Now there should be several loaded functions, including **run_hysplit**. To run Hysplit, have a look at the examples here <https://gist.github.com/davidcarslaw/c67e33a04ff6e1be0cd7357796e4bdf5>.

On my Windows machine it is run as follows:

```
data_out <- run_hysplit(
  latitude = 36.134,
  longitude = -5.347,
  runtime = -96,
  start_height = 10,
  model_height = 10000,
  start = 2015,
  end = "2015-01-10",
  hysplit_exec = "~/hysplit4/exec",
  hysplit_input = "~/trajData",
  hysplit_output = "~/temp",
  site = "gibraltar")
```

The **data_out** can then be used directly in **openair** trajectory functions.

Most of the options should be self-explanatory but **hysplit_exec** is the path to the Hysplit executable, **hysplit_input** is the path to the meteorological files (downloaded as described above) and **hysplit_output** is the directory where Hysplit will write its temporary files.

Once run it is then advisable to store the data somewhere. Save it like:

```
saveRDS(data_out, file = "~/trajProc/myTrajData.rds")
```

Then it is easy to read in later and use e.g.

```
traj <- readRDS("~/trajProc/myTrajData.rds")
```