# 01- Reading Files

R Workshop
- Data Formatting and Reshaping -
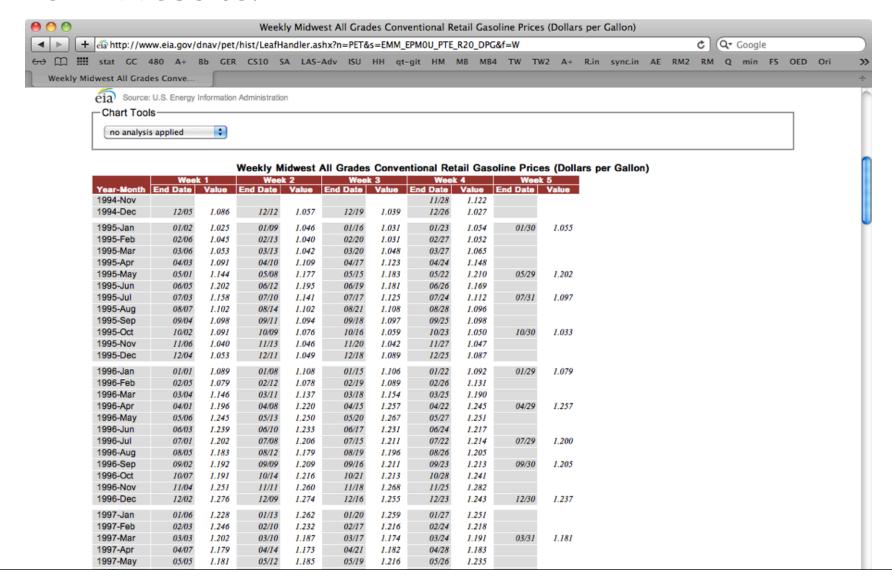
# Outline

- Reading files: Excel and R

- packages gdata and foreign

- reading SAS xport files

# Data in Excel

- Formats xls and csv - what's the difference?

- File extensions xls or xlsx are proprietary Excel formats, they are binary files

- csv is the extension for comma separated value files. They are text files - and directly readable.

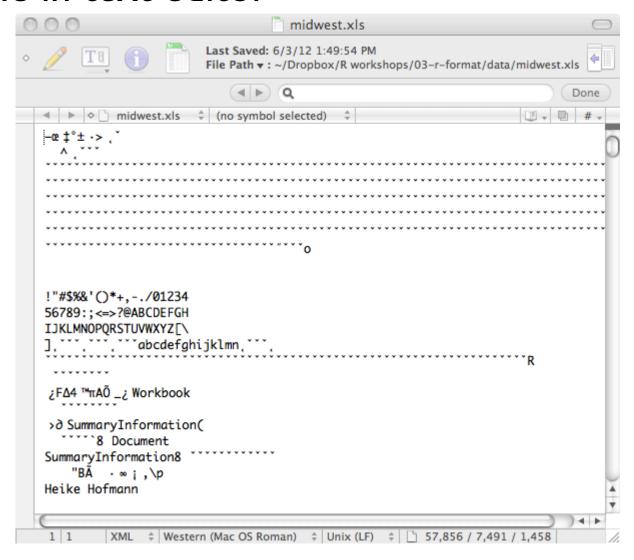- Example: gas prices in the Midwest since 1994 (from data.gov and EIA)

# Gas Prices

xls file in Website:

# Gas Prices

xls file in Excel:



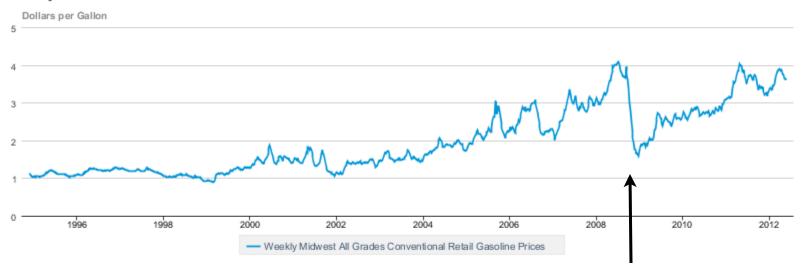| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| | | Week 1 | | Week 2 | | Week 3 |
| | Year-Month | End Date | Value | End Date | Value | End Date |
| 3 | 1994-Nov | | | | | |
| 4 | 1994-Dec | 5-Dec | 1.086 | 12-Dec | 1.057 | 19- |
| 5 | 1995-Jan | 2-Jan | 1.025 | 9-Jan | 1.046 | 16 |
| 6 | 1995-Feb | 6-Feb | 1.045 | 13-Feb | 1.04 | 20- |
| 7 | 1995-Mar | 6-Mar | 1.053 | 13-Mar | 1.042 | 20- |
| 8 | 1995-Apr | 3-Apr | 1.091 | 10-Apr | 1.109 | 17 |
| 9 | 1995-May | 1-May | 1.144 | 8-May | 1.177 | 15- |
| 10 | 1995-Jun | 5-Jun | 1.202 | 12-Jun | 1.195 | 19 |
| 11 | 1995-Jul | 3-Jul | 1.158 | 10-Jul | 1.141 | 17 |
| 12 | 1995-Aug | 7-Aug | 1.102 | 14-Aug | 1.102 | 21- |
| 13 | 1995-Sep | 4-Sep | 1.098 | 11-Sep | 1.094 | 18- |
| 14 | 1995-Oct | 2-Oct | 1.091 | 9-Oct | 1.076 | 16 |
| 15 | 1995-Nov | 6-Nov | 1.04 | 13-Nov | 1.046 | 20- |
| 16 | 1995-Dec | 4-Dec | 1.053 | 11-Dec | 1.049 | 18- |
| 17 | 1996-Jan | 1-Jan | 1.089 | 8-Jan | 1.108 | 15 |
| 18 | 1996-Feb | 5-Feb | 1.079 | 12-Feb | 1.078 | 19- |
| 19 | 1996-Mar | 4-Mar | 1.146 | 11-Mar | 1.137 | 18- |
| 20 | 1996-Apr | 1-Apr | 1.196 | 8-Apr | 1.22 | 15 |
| 21 | 1996-May | 6-May | 1.245 | 13-May | 1.25 | 20- |
| 22 | 1996-Jun | 3-Jun | 1.239 | 10-Jun | 1.233 | 17 |

# Gas Prices

- xls file in text editor

# Gas Prices

- csv file in text editor

# … what we'd like to do with the data …



Weekly Midwest All Grades Conventional Retail Gasoline Prices

Oct 2008

# Reading Files in R

- Textfiles: usually comma separated (or tabular separated)

```
read.csv(file, header = TRUE, sep = ",", quote = "\\"",
         dec = ".",  fill = TRUE, comment.char = "", ...)

read.table (file, header = FALSE, sep = "", quote = "\\"'", dec = ".",
    row.names, col.names, as.is = !stringsAsFactors, na.strings = "NA",
    colClasses = NA, nrows = -1, skip = 0, check.names = TRUE,
    fill = !blank.lines.skip, strip.white = FALSE, blank.lines.skip = TRUE,
    comment.char = "#", allowEscapes = FALSE, flush = FALSE,
    stringsAsFactors = default.stringsAsFactors(), fileEncoding = "",
    encoding = "unknown", text)
```

# Gas Prices in the Midwest

```r
# read a csv file published as a webfile
gp <- read.csv("http://www.hofroe.net/R workshops/03-r-format/data/midwest.csv")

# read (and find) a local csv file
gp <- read.csv(file.choose())

# reveals awful format
head(gp)
```

|   | Year.Month | Week.1 End Date | X Value | Week.2 End Date | X.1 Value | Week.3 End Date | X.2 Value | Week.4 End Date | X.3 Value | Week.5 End Date | X.4 Value |
|---|------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| 1 |            | End Date | Value | End Date | Value | End Date | Value | End Date | Value | End Date | Value |
| 2 | 1994-Nov   |          |       |          |       |          |       | 28-Nov | 1.122 |        |       |
| 3 | 1994-Dec   | 5-Dec    | 1.086 | 12-Dec   | 1.057 | 19-Dec   | 1.039 | 26-Dec | 1.027 |        |       |
| 4 | 1995-Jan   | 2-Jan    | 1.025 | 9-Jan    | 1.046 | 16-Jan   | 1.031 | 23-Jan | 1.054 | 30-Jan | 1.055 |
| 5 | 1995-Feb   | 6-Feb    | 1.045 | 13-Feb   | 1.04  | 20-Feb   | 1.031 | 27-Feb | 1.052 |        |       |
| 6 | 1995-Mar   | 6-Mar    | 1.053 | 13-Mar   | 1.042 | 20-Mar   | 1.048 | 27-Mar | 1.065 |        |       |

# Gas Prices in the Midwest

```
str(gp)

'data.frame':  212 obs. of  11 variables:
 $ Year.Month: Factor w/ 212 levels "","  1994-Dec",..: 1 3 2 8 7 11 4 12 10 9 ...
 $ Week.1    : Factor w/ 86 levels "","1-Apr","1-Aug",..: 86 1 52 18 65 69 26 10 56 31 ...
 $ X         : Factor w/ 197 levels "","0.905","0.918",..: 197 1 19 7 12 13 21 29 42 31 ...
 $ Week.2    : Factor w/ 86 levels "","10-Apr","10-Aug",..: 86 1 28 78 41 45 2 70 32 7 ...
 $ X.1       : Factor w/ 206 levels "","0.919","0.921",..: 206 1 17 14 12 13 27 39 45 34 ...
 $ Week.3    : Factor w/ 86 levels "","15-Apr","15-Aug",..: 86 1 52 18 65 69 26 10 56 31 ...
 $ X.2       : Factor w/ 199 levels "","0.91","0.929",..: 199 1 11 9 9 15 28 40 38 29 ...
 $ Week.4    : Factor w/ 85 levels "22-Apr","22-Aug",..: 85 82 51 17 64 68 25 9 55 30 ...
 $ X.3       : Factor w/ 201 levels "0.883","0.921",..: 201 29 9 14 13 15 32 44 34 27 ...
 $ Week.5    : Factor w/ 31 levels "","29-Apr","29-Aug",..: 31 1 1 16 1 1 1 9 1 27 ...
 $ X.4       : Factor w/ 74 levels "","0.955","1.023",..: 74 1 1 5 1 1 1 18 1 11 ...
```

needs some more work before we can analyze
(or even visualize the data)

# Gas Prices in the Midwest

Issues with the data:

- two lines of header information

- all variables are factor variables - but we know, that some are dates, some are numeric

# Your Turn

- Have a look at the parameters of read.table in the help (Hint: try `?read.table` to view the help) to solve the following problems:

  - Read the first two lines of the file into an object called 'gp_names'

  - Read everything but the first two lines into an object called 'gp_data'

# Reading Excel Data

- Need another package: `gdata`

```
read.xls(xls, sheet = 1, verbose = FALSE, pattern, ..., method = c("csv",
     "tsv", "tab"), perl = "perl")
```

```
library("gdata")
# get html page with an overview of the package functionality
help(package="gdata")

gp2 <- read.xls(file.choose(), sheet=1)
head(gp2)
```

|   | Year.Month | Week.1 | X | Week.2 | X.1 | Week.3 | X.2 | Week.4 | X.3 | Week.5 | X.4 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 |  | End Date | Value | End Date | Value | End Date | Value | End Date | Value | End Date | Value |
| 2 | 1994-Nov |  |  |  |  |  |  | 28-Nov | 1.122 |  |  |
| 3 | 1994-Dec | 5-Dec | 1.086 | 12-Dec | 1.057 | 19-Dec | 1.039 | 26-Dec | 1.027 |  |  |
| 4 | 1995-Jan | 2-Jan | 1.025 | 9-Jan | 1.046 | 16-Jan | 1.031 | 23-Jan | 1.054 | 30-Jan | 1.055 |
| 5 | 1995-Feb | 6-Feb | 1.045 | 13-Feb | 1.04 | 20-Feb | 1.031 | 27-Feb | 1.052 |  |  |
| 6 | 1995-Mar | 6-Mar | 1.053 | 13-Mar | 1.042 | 20-Mar | 1.048 | 27-Mar | 1.065 |  |  |

# Your Turn

- Read the file gasprices.xls into R and inspect it.

- What might be potential problems when analyzing the data?

# Package foreign

- Other file formats can be read using functions from package `foreign`

- SPSS: `read.spss`
  SAS:  `read.xport` (xport format)
       `read.ssd`  (permanent SAS data)
  Minitab: `read.mtp`
  Systat:  `read.systat`

# Your Turn

- The NHANES (National Health and Nutrition Survey) publishes data in SAS export format: http://wwwn.cdc.gov/nchs/nhanes/search/nhanes09_10.aspx

- Download one of the files and load into R