

EDA-lab

Janos Binder¹

European Molecular Biology Laboratory (EMBL),
Heidelberg, Germany

¹janos.binder@embl.de

June 3, 2014

Contents

1	Acknowledgements	1
2	Required Packages and other Preparations	2
3	Introduction	2
4	Base graphics and ggplot2	2
4.1	Building a Plot Layer by Layer	3
4.2	Some biological data	6
4.3	Setting axis limits	8
4.4	Faceting	10
4.5	Melting and casting data frames	12
5	Univariate Data Display	13
5.1	Frequency Table and Barplot	13
5.2	Scatterplots and Stripcharts	18
5.3	Histograms	22
5.4	Kernel density estimates	25
5.5	Boxplots	26
5.6	Violin plots	29
5.7	Empirical Cumulative Distribution Function	31
5.8	Quantile–Quantile (QQ) Plot	33
6	Descriptive Statistics	37
6.1	Measures of Central Tendency	38
6.2	Measures of Spread	38
7	Answers to Exercises	40

1 Acknowledgements

This tutorial is based on the material of Bernd Klaus (<http://www-huber.embl.de/users/klaus/main.html>).

2 Required Packages and other Preparations

```
load(url("http://www-huber.embl.de/users/klaus/BasicR/seqZyx.rda"))
library("TeachingDemos")
data(golub, package = "multtest")
library(biomaRt)
library(reshape2)
library(ggplot2)
library(plyr)
library(xlsx)
library(vioplot)
```

Should you get error messages, you can also install them quickly:

```
source("http://bioconductor.org/biocLite.R")
packs <- c("TeachingDemos", "multtest", "biomaRt",
           "reshape2", "ggplot2", "plyr", "xlsx", "vioplot")
biocLite(packs)
```

3 Introduction

In this lab, a few essential methods are given to display and visualize data. It quickly answers questions like: How are my data distributed? How can the frequencies of nucleotides from a gene be visualized? Are there outliers in my data? Does the distribution of my data resemble that of a bell-shaped curve? Are there differences between gene expression values taken from two groups of patients?

The most important central tendencies (mean, median) are defined and illustrated together with the most important measures of spread (standard deviation, variance, inter quartile range, and median absolute deviation).

We also introduce *ggplot2* a package to produce elegant graphics for data analysis.

4 Base graphics and ggplot2

The package *ggplot2* is one of the most commonly used graphics packages for *R*. It is an alternative to the basic graphic system of *R*, which is limited in various aspects. In this lab we will use basic graphics in *R* and corresponding *ggplot2* plots side by side.

ggplot2 is meant to be an implementation of the Grammar of Graphics, developed by L. Wilkinson, hence gg–plot. The basic notion is that there is a grammar to the composition of graphical components in statistical graphics, and by directly controlling that grammar, you can generate a large set of carefully constructed graphics tailored to your particular needs.

The central concept of the approach is that plots convey information through various aspects of their aesthetics. Aesthetics are mappings from the data to something you can visually perceive. Some aesthetics that plots use are:

- x position
- y position
- size of elements
- shape of elements
- color of elements

The elements in a plot are geometric shapes, like

- points

- lines
- line segments
- bars
- text

Some of these geometries have their own particular aesthetics. For instance:

- points
 - point shape
 - point size
- lines
 - line type
 - line weight
- bars
 - y minimum
 - y maximum
 - fill color
 - outline color
- text
 - label value

Each component is added to the plot as a layer, hence you might start with a simple mapping of the raw data to the x- and y-axes, creating a scatterplot. A second layer may be added by coloring the points according to a group they belong to and so on.

There are other basics of these graphics that you can adjust, like the scaling of the aesthetics, and the positions of the geometries.

The values represented in the plot are the product of various statistics. If you just plot the raw data, you can think of each point representing the identity statistic. Many bar charts represent the mean or the median statistic. Histograms are bar charts where the bars represent the binned count or density statistics and so on.

4.1 Building a Plot Layer by Layer

There's a quick plotting function in *ggplot2* called `qplot()` which is meant to be similar to the `plot()` function from base graphics. You can do a lot with `qplot()`, but it can be better to approach the package from the layering syntax.

All *ggplot2* plots begin with the function `ggplot()`. `ggplot()` takes two primary arguments, data is the data frame containing the data to be plotted and `aes()` are the aesthetic mappings to pass on to the plot elements.

As you can see, the second argument, `aes()`, isn't a normal argument, but another function. Since we'll never use `aes()` as a separate function, it might be best to think of it as a special way to pass a list of arguments to the plot.

The first step in creating a plot is to add one or more layers. Let's start with the iris data set as an example. Note that *ggplot2* always requires the specification of the data frame from which the variables used in the plot are drawn.

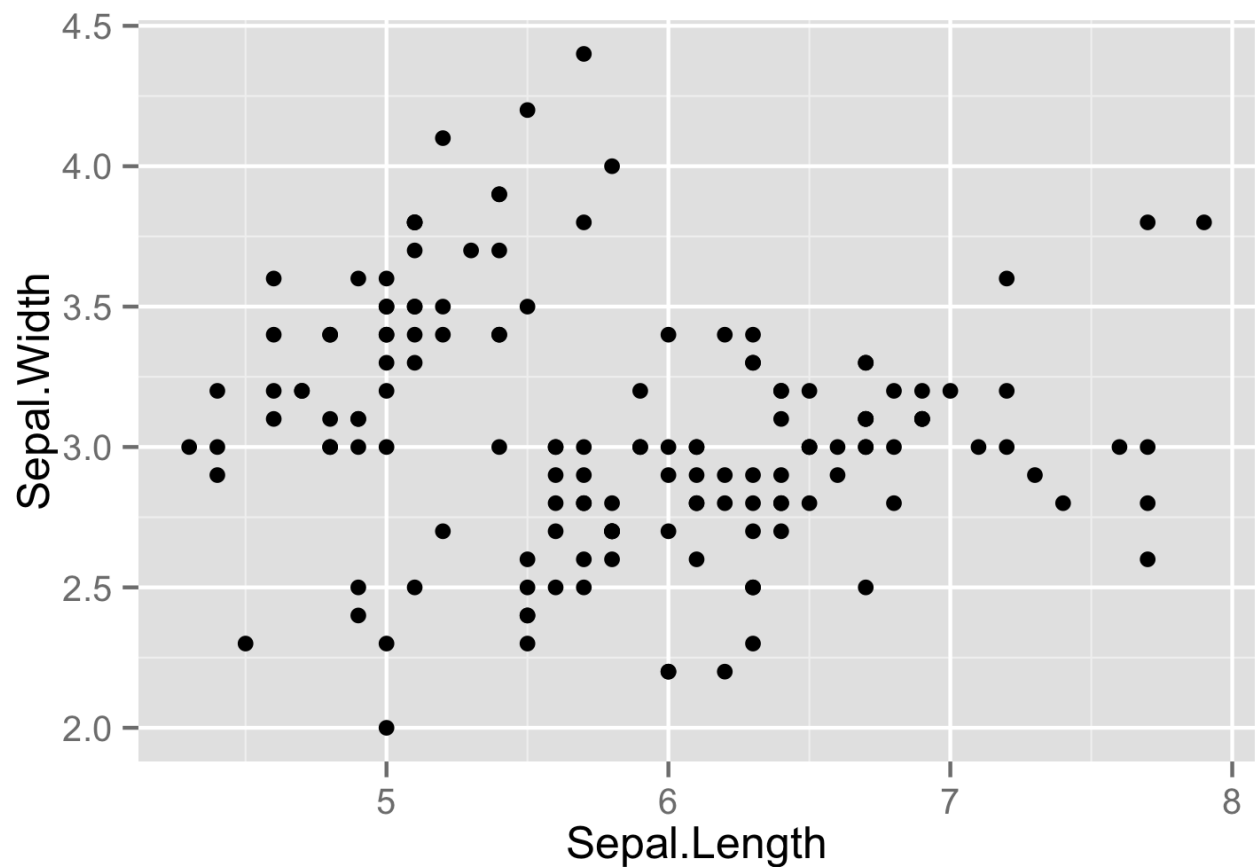
```
summary(iris)

#>      Sepal.Length  Sepal.Width  Petal.Length  Petal.Width      Species
#>      Min.       :4.30   Min.       :2.00   Min.       :1.00   Min.       :0.1   setosa       :50
#>      1st Qu.:5.10   1st Qu.:2.80   1st Qu.:1.60   1st Qu.:0.3   versicolor:50
#>      Median :5.80   Median :3.00   Median :4.35   Median :1.3   virginica  :50
#>      Mean    :5.84   Mean    :3.06   Mean    :3.76   Mean    :1.2
#>      3rd Qu.:6.40   3rd Qu.:3.30   3rd Qu.:5.10   3rd Qu.:1.8
#>      Max.    :7.90   Max.    :4.40   Max.    :6.90   Max.    :2.5

p <- ggplot(iris, aes(Sepal.Length, Sepal.Width))
```

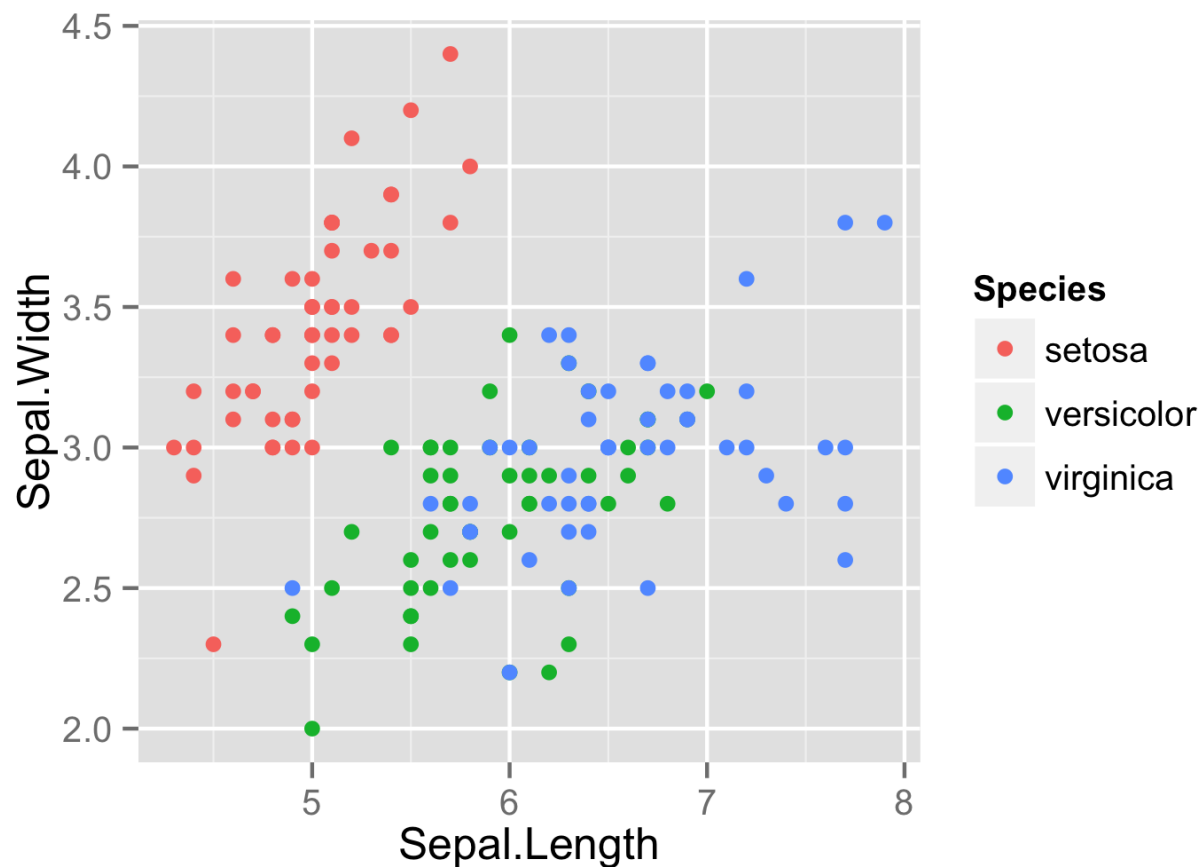
If you just type `p` or `print(p)`, you'll get back a warning saying that the plot lacks any layers. With the `ggplot()` function, we've set up a plot which is going to draw from the `iris` data, the `Sepal.length` variable will be mapped to the x-axis, and the `Sepal.width` variable is going to be mapped to the y-axis. However, we have not determined which kind of geometric object will represent the data. Let's add points, for a scatterplot.

```
p + geom_point()
```



Alternatively, this plot could have been produced with `qplot`. Additionally, you can map color to the species.

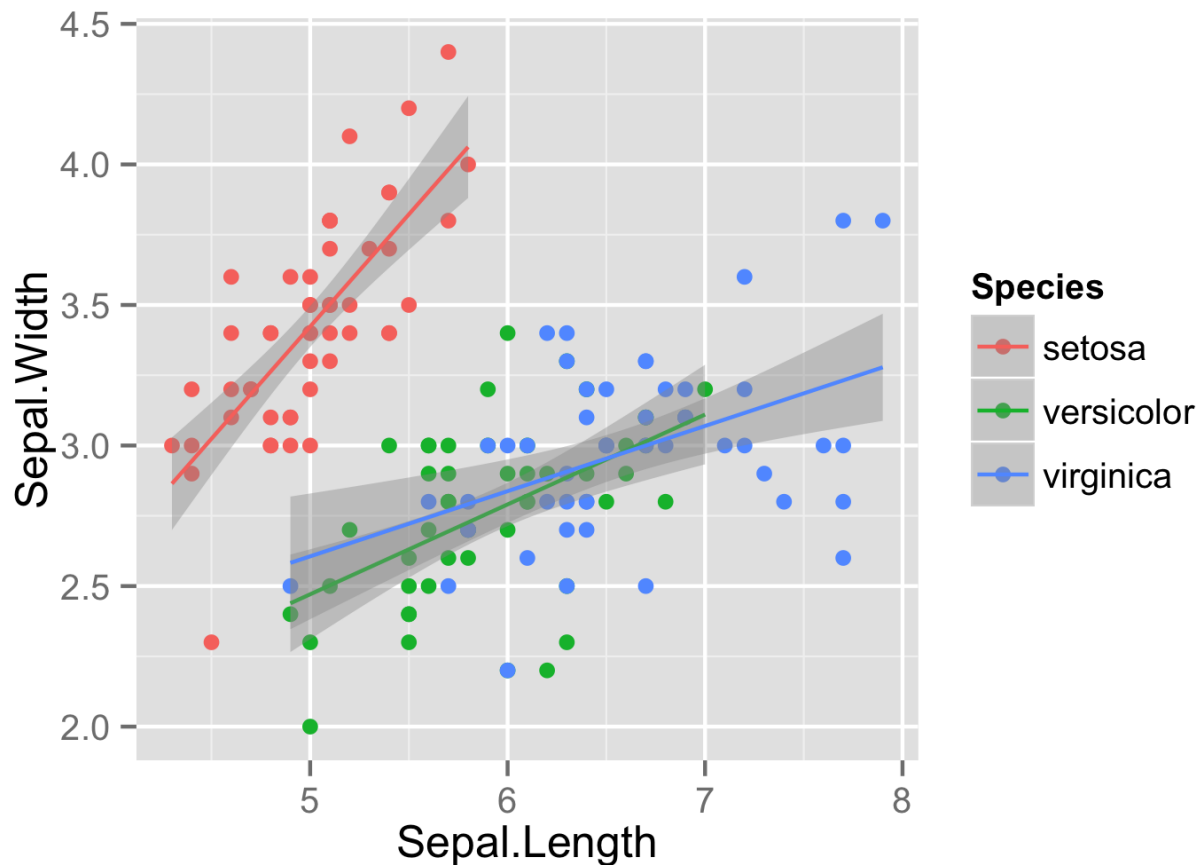
```
qplot(Sepal.Length, Sepal.Width, data = iris, color = Species)
```



We clearly see that the setosa plants have different Sepal.Length/Sepal.Width relationship compared to the other two species. The full documentation for *ggplot2* can be found at <http://docs.ggplot2.org/current/>. Apart from mapping data to aesthetics, *ggplot2* can handle statistical transformations of the data, i.e. easily create all the nice exploratory graphics we will look at below.

We will explore one of these transformations by adding a regression line to the data of each of the three plant species as a third layer.

```
ggsmooth <- (qplot(Sepal.Length, Sepal.Width, data = iris, color = Species)
+ stat_smooth(method = "lm"))
ggsmooth
```



The command `stat_smooth` first adds a statistical transformation to the existing data and then plots it using a certain geometry, in this case a special "smooth" geometry that is tailored to the plotting of regression fits. You can obtain the statistical transformations by looking at the saved plot and extracting the appropriate sublist.

```
transformed.data <- as.list(print(ggsmooth))$data[[2]]
```

Thus, you could also map the transformed data differently than the default geometry does it. This however does not make much sense in this case.

4.2 Some biological data

Here we're looking at some real biological data: different doses of HGF (a cytokine) were applied to cells and the downstream effect to the phosphorylation of target proteins were assessed recording a time course-signal in different conditions. This data and the exercise ideas were provided by Lars Velten (Steinmetz lab). We first load the dataset.

```
proteins<-read.csv("http://www-huber.embl.de/users/klaus/BasicR/proteins.csv")[, -1]
head(proteins)
```

```
#>      Condition min Target   Signal   Sigma
#> 1 10ng/mL HGF  0  pAKT -6.71e+08 96749276
#> 2 10ng/mL HGF  5  pAKT  5.68e+08 97144057
#> 3 10ng/mL HGF 10  pAKT  1.05e+09 97215659
#> 4 10ng/mL HGF 20  pAKT -3.23e+08 97055787
#> 5 10ng/mL HGF 30  pAKT -9.00e+08 96908309
#> 6 10ng/mL HGF 60  pAKT  3.17e+08 96731346
```

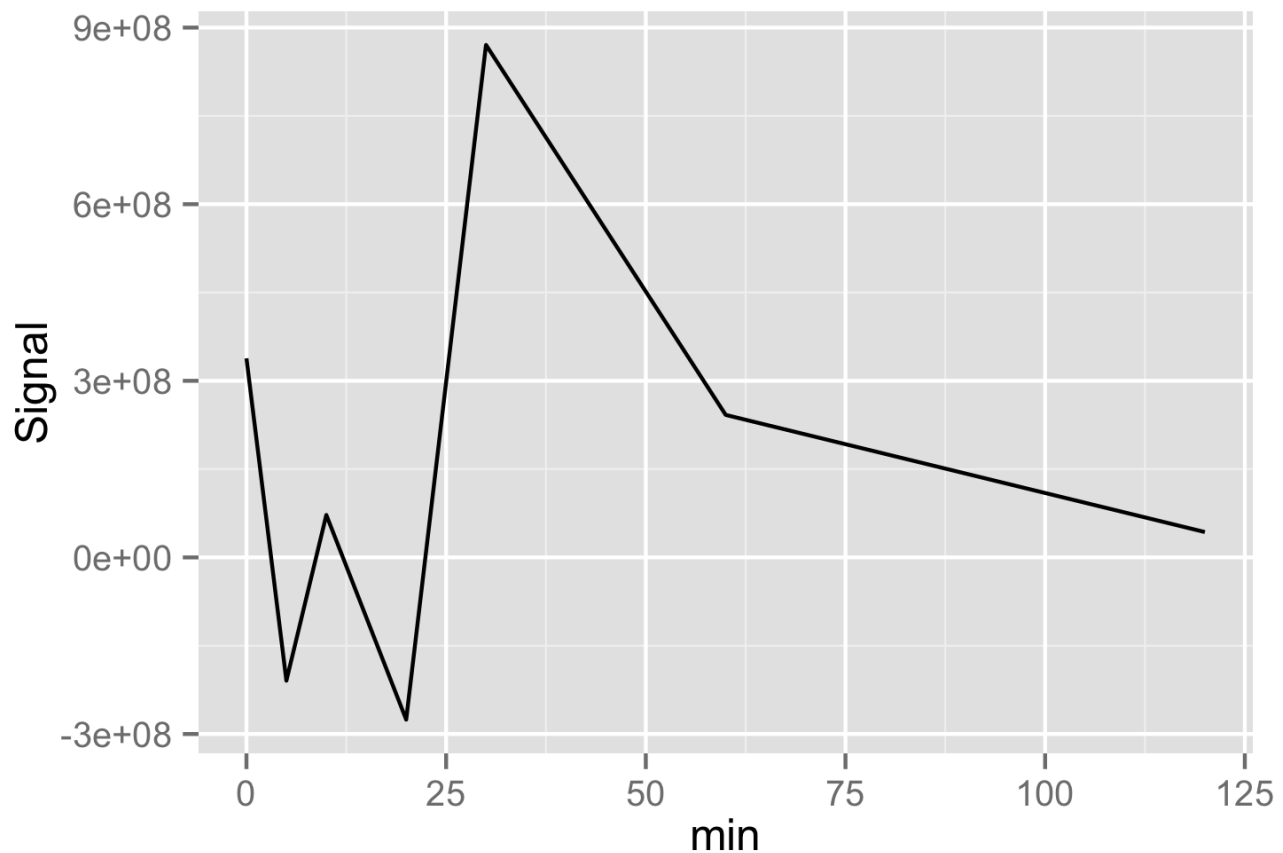
```
proteins_pMek <- subset(proteins, proteins$Target == "pMEK")
proteins_pMek_sub <- subset(proteins_pMek, proteins_pMek$Condition == "10ng/mL HGF")
```

We can start simple by only looking at the first condition of the "pMEK" protein target for now. We simply produce a line plot of the signal across time. In this plot we use the data `proteins_pMek_sub`, map `min` to the x-axis and `Signal` to the y-axis and use a line as a geometry. Additional modifications are added to the plot in the exercise.

```
proteins_pMek_sub
```

```
#>      Condition min Target   Signal   Sigma
#> 103 10ng/mL HGF   0  pMEK  3.38e+08 31005696
#> 104 10ng/mL HGF   5  pMEK -2.09e+08 31400418
#> 105 10ng/mL HGF  10  pMEK  7.20e+07 31199120
#> 106 10ng/mL HGF  20  pMEK -2.76e+08 31015194
#> 107 10ng/mL HGF  30  pMEK  8.70e+08 31140886
#> 108 10ng/mL HGF  60  pMEK  2.42e+08 31040783
#> 109 10ng/mL HGF 120  pMEK  4.31e+07 31072343
```

```
qplot(min, Signal, data = proteins_pMek_sub, geom = "line")
```



```
# or
#ggplot(proteins_pMek_sub, aes(min, Signal)) + geom_line()
```

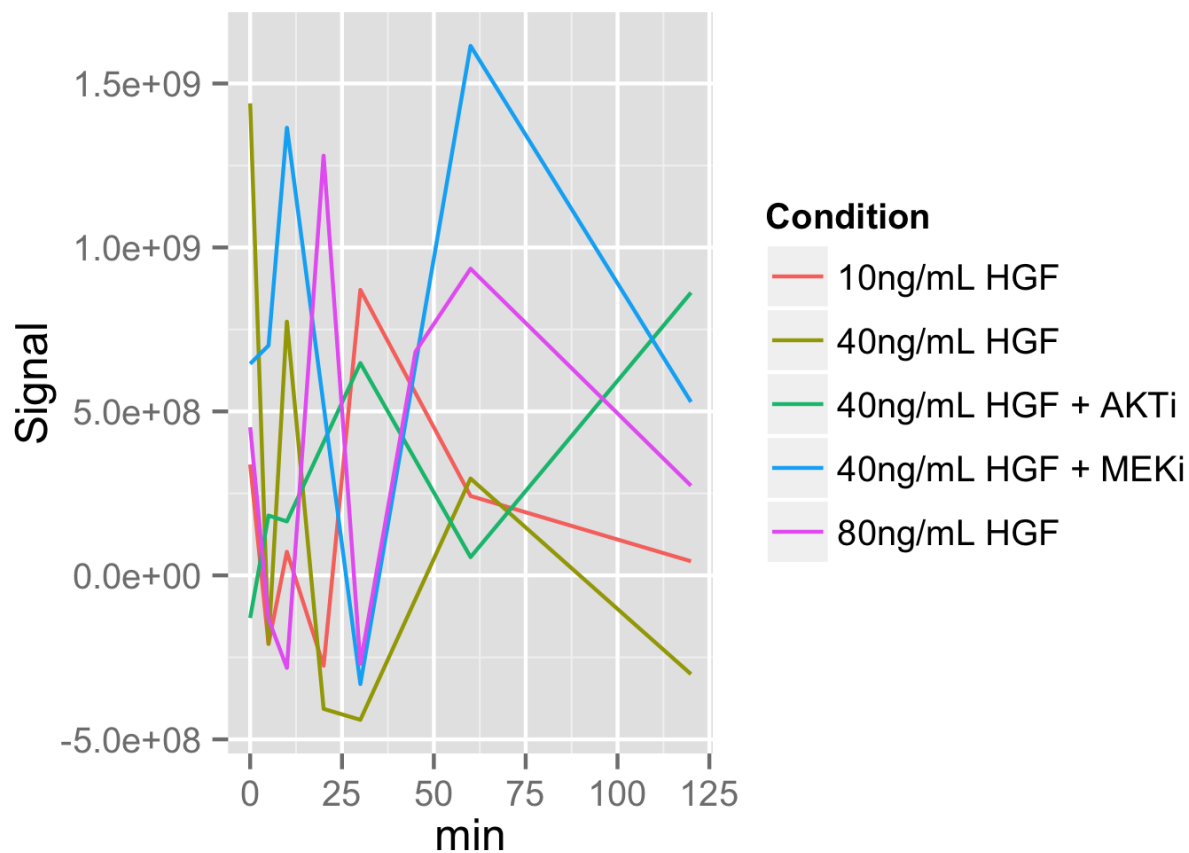
Exercise: Simple ggplot usage

Using the data `proteins_pMek_sub`, do the following

- Use points as a geometry instead of lines
- Use both lines and points
- Add errorbars `geom_errorbar` to the plot. This requires further aesthetics: `ymin` and `ymax`. The estimated error is stored in the variable `Sigma`.

We can also easily plot all conditions for the protein pMEK by mapping color to the experimental condition of pMEK.

```
qplot(min, Signal, data = proteins_pMek, geom = "line", color = Condition)
```

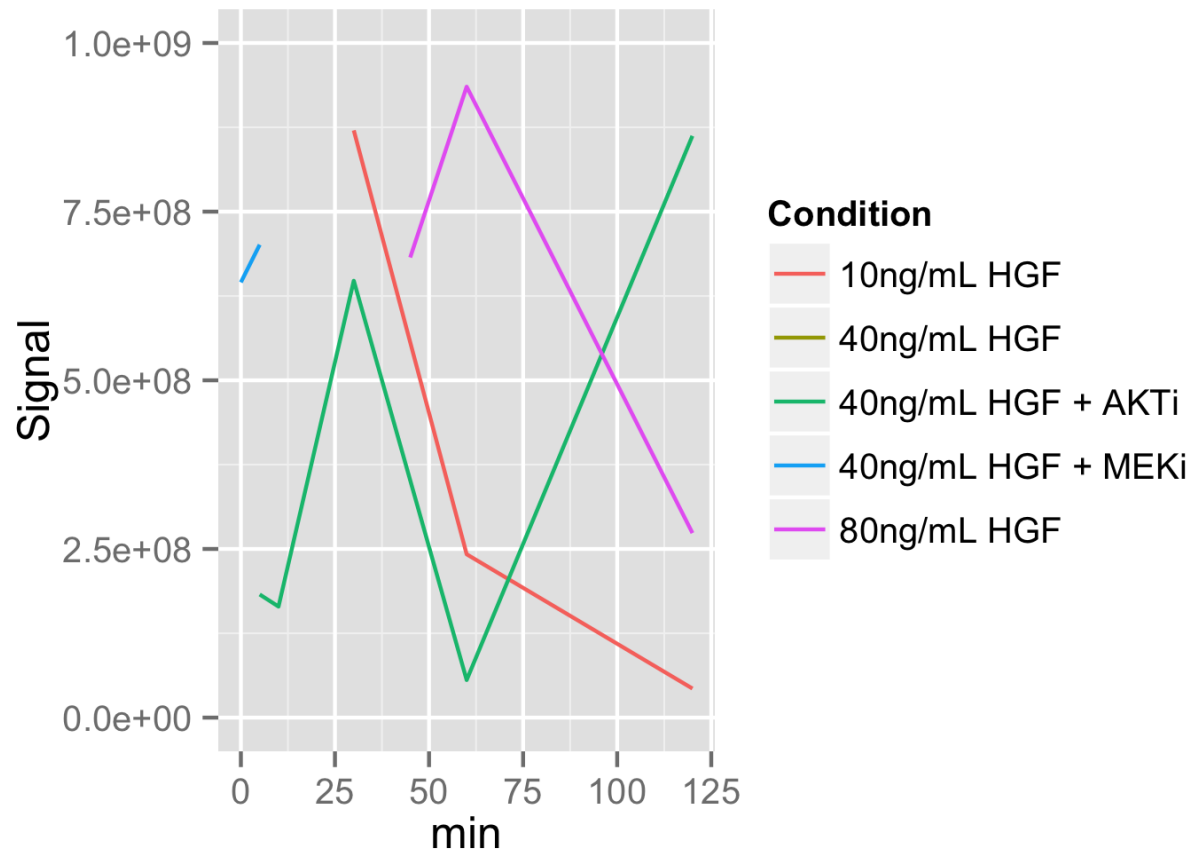


```
# or
#ggplot(proteins_pMek, aes(min, Signal, color = Condition)) + geom_line()
```

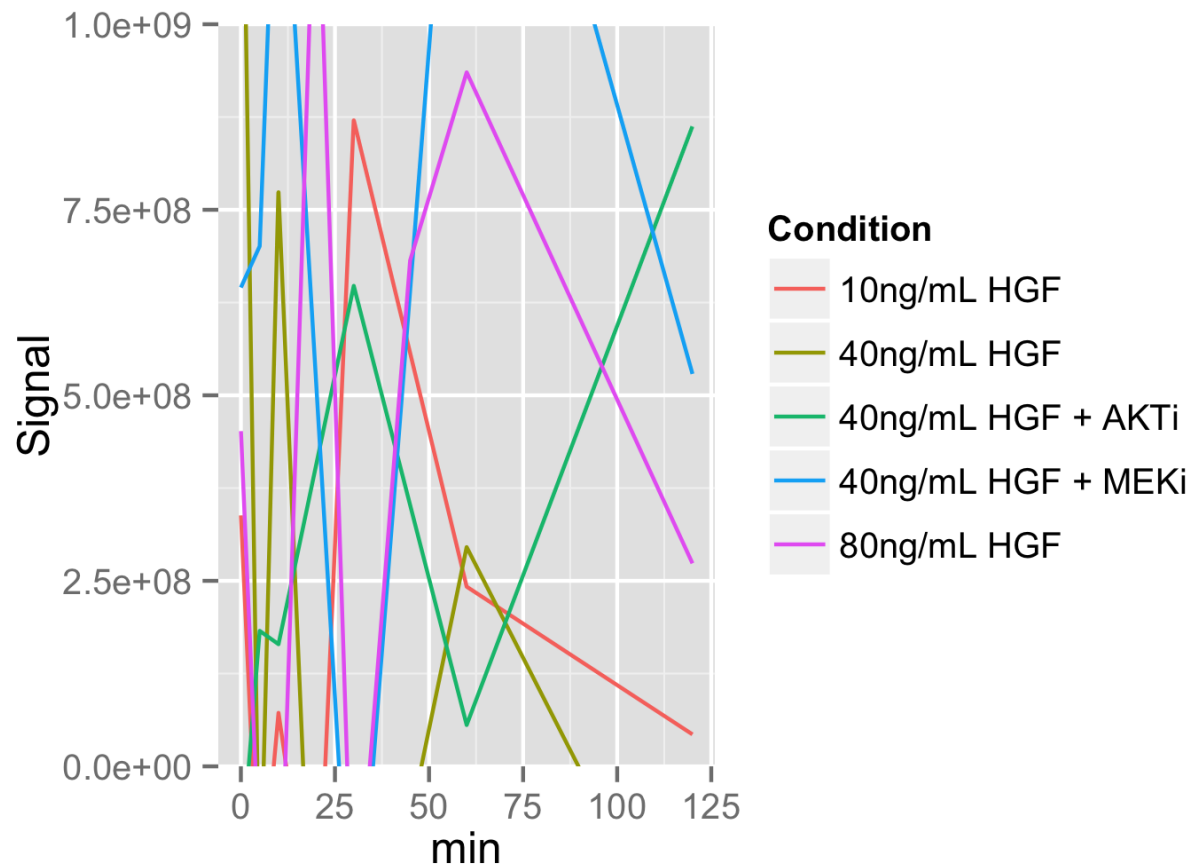
4.3 Setting axis limits

Another important aspect of data display are the limits of the x- and y-axis. You can change these limits by using `xlim()` and `ylim()`. However, the overall scale will not be "retrained" and only be limited. This essentially results in plotting only a subset of the data, excluding values outside of the limits. In order to actually "zoom-in" without excluding data points one has to use the `coord_cartesian()` command. In the code below, we first limit the y-axis to `1e9`, excluding all other data points, then we retrain the scale.


```
## limit scale
(qplot(min, Signal, data = proteins_pMek, geom = "line", color = Condition)
+ ylim(c(0,1e9)))
#> Warning: Removed 4 rows containing missing values (geom.path).
```



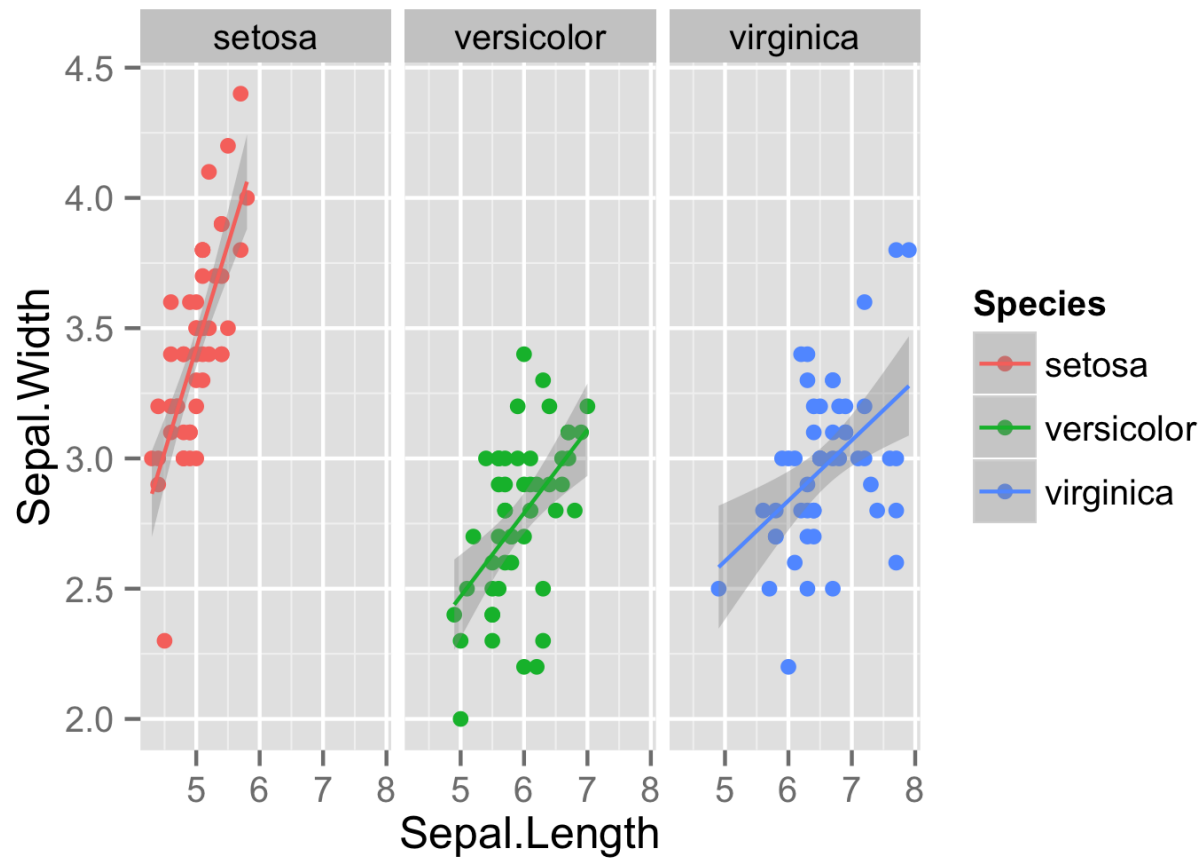
```
## retrain scale
(qplot(min, Signal, data = proteins_pMek, geom = "line", color = Condition)
+ coord_cartesian(ylim = c(0, 1e9)))
```



4.4 Faceting

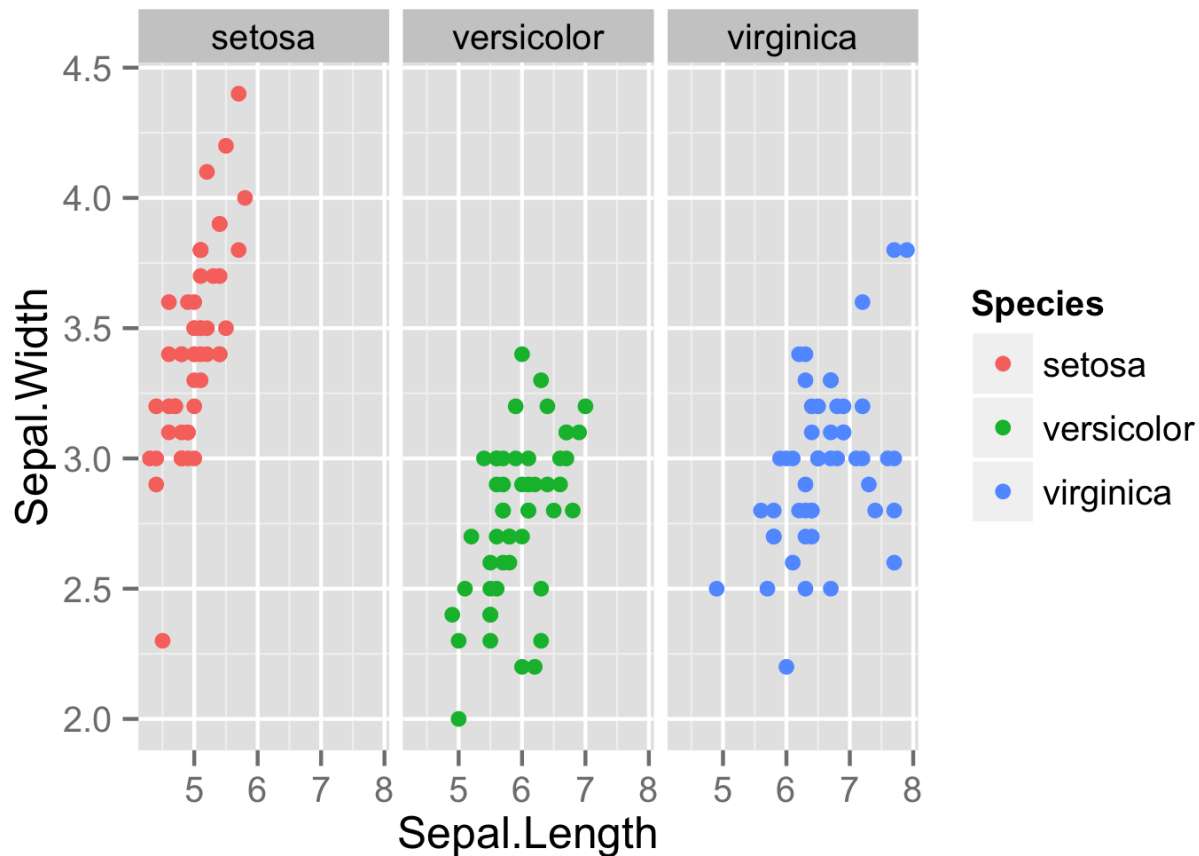
ggplot2 also allows you to easily split plots according to a factor variable, plotting parts of the data in different panels. Returning to the iris data, we can easily plot different species in different panels using `facet_wrap()`. This is the simplest faceting function, splitting according to one factor only.

```
ggsmooth + facet_wrap(~ Species)
```



The splitting is defined in a formula notation: `yfactor ~ xfactor`. Factors can also be combined, e.g. `yfactor ~ xfactor_1 + xfactor_2`. Faceting can also be used with `qplot()`, there you always have to specify a two-sided formula. If you only want to use a single splitting factor, you can use the dot notation:

```
(qplot(Sepal.Length, Sepal.Width, data = iris,  
color = Species, facets = . ~ Species))
```



Exercise: ggplot faceting

Using the data `proteins_pMek`, produce a plot split by the experimental condition factor using `facet_wrap()`.

4.5 Melting and casting data frames

The data table we loaded was already suitable for the plots we wanted to produce since every line represents exactly one observation. However, this is not necessarily the case and we might want to represent the different time points by different columns, not just a single one.

In time series analysis parlance our current data would be in "long" format, but we might want to transform it into a "wide" format, with a separate column for every time point. The package [reshape2](#) allows you to do this. `ggplot2` usually requires "long" formats, which can be obtained by using the function `melt`. Thus a "molten" data frame in `ggplot2` corresponds to a "long" format. "wide" formats can be computed using the function `dcast`. As an example we will now represent every time point of our data frame as a single column.

Note that the wide format is only compatible with a single numerical target variable, so we only include `signal` as a variable here. The casting grid is defined by a formula with the variables that will be represented by a single column on the left (x-variables) and the wide variables on the right side of the tilde (y-variables).

```
proteins_cast <- dcast(proteins, Target + Condition ~ min ,
                       value.var = "Signal")
head(proteins_cast)
```

```
#>      Target      Condition      0      5      10      20      30      45      60
#> 1   pAKT      10ng/mL HGF -6.71e+08 5.68e+08 1.05e+09 -3.23e+08 -9.00e+08      NA 3.17e+08
#> 2   pAKT      40ng/mL HGF -7.59e+08 3.32e+08 1.95e+08 1.29e+08 7.43e+07      NA 1.01e+09
#> 3   pAKT 40ng/mL HGF + AKTi -2.03e+08 1.62e+08 -2.41e+08      NA -4.69e+08      NA -3.31e+08
#> 4   pAKT 40ng/mL HGF + MEKi 7.30e+08 6.57e+08 8.73e+08      NA 1.39e+09      NA 7.10e+08
#> 5   pAKT      80ng/mL HGF -3.09e+08 8.36e+08 8.91e+08 1.12e+09 9.08e+08 3.8e+08 -6.95e+08
#> 6 pERK1      10ng/mL HGF 2.77e+08 2.37e+08 -3.36e+08 -1.48e+08 7.06e+08      NA 2.85e+07
#>      120
#> 1 -4.69e+08
#> 2 -7.95e+08
#> 3 5.92e+08
#> 4 -2.21e+08
#> 5 -2.16e+08
#> 6 5.52e+08
```

We can now melt the data frame again. For melting, you have to specify id variable and/or measurement variables. id variables identify a single line in the molten data frame and measured variables represent the measurements. If you only supply one of them, it is assumed that all the other variables belong to the other group. Columns that are not in either of the groups are discarded.

A factor column is then added to indicate to which former column a measurement belongs to. In our case this is the time point. Another useful function is `arrange` which allows you to reorder the data frame according to certain columns.

```
proteins_molten <- melt(proteins_cast, id.vars = c("Target", "Condition"),
  variable.name="Time", value.name= "Signal")
proteins_molten <-arrange(proteins_molten, Condition)
head(proteins_molten)

#>      Target Condition Time      Signal
#> 1   pAKT 10ng/mL HGF      0 -6.71e+08
#> 2 pERK1 10ng/mL HGF      0 2.77e+08
#> 3 pERK2 10ng/mL HGF      0 -1.91e+08
#> 4   pMEK 10ng/mL HGF      0 3.38e+08
#> 5   pAKT 10ng/mL HGF      5 5.68e+08
#> 6 pERK1 10ng/mL HGF      5 2.37e+08
```

Exercise: complex ggplot example

Use the data frame `proteins` to produce a plot of the time courses split by the experimental target and colored according to the experimental conditions. Add error bars to your plot.

5 Univariate Data Display

In order to study the distribution of data, various visualization methods are available. We will look at some of them in the following.

5.1 Frequency Table and Barplot

Discrete data occur when the values naturally fall into categories. A frequency table simply gives the number of occurrences within a category.

Example: Nucleotide Frequencies

A gene consists of a sequence of nucleotides $\{A, C, G, T\}$. The number of each nucleotide can be displayed in a frequency table. This will be illustrated by an Exon of the Zyxin gene, which plays an important role in cell adhesion. You can use the Bioconductor package *biomaRt* to query various genomic databases. The commented code below was originally used to retrieve the exon sequence. The command `strsplit` splits the sequence into its single letter parts. We already loaded the sequence at the beginning of the lab.

```
ensembl= useMart("ensembl",dataset="hsapiens_gene_ensembl")
seqZyx <- getSequence(id = "ENSG00000159840", type = "ensembl_gene_id",
mart = ensembl, seqType = "gene_exon")[1,]
seqZyx <- strsplit(as.character(seqZyx[1,1]), split = character(0))
seqZyx <- seqZyx[[1]]
save(seqZyx, file = "seqZyx.rda")
```

A (frequency) table, a corresponding pie plot and a barplot can be produced by following commands

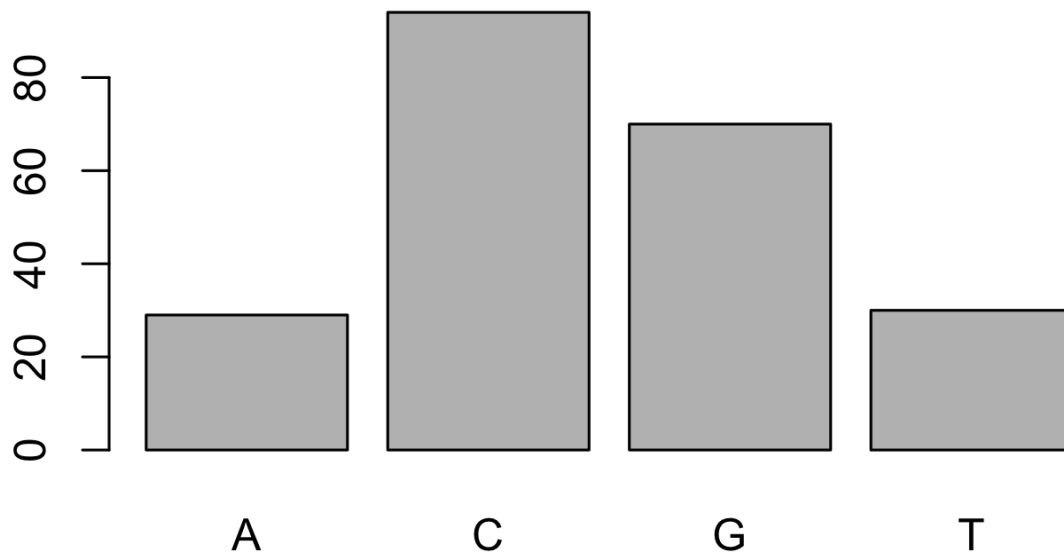
```
table(seqZyx) ## table

#> seqZyx
#>   A  C  G  T
#> 29 94 70 30

prop.table(table(seqZyx)) ## frequency table

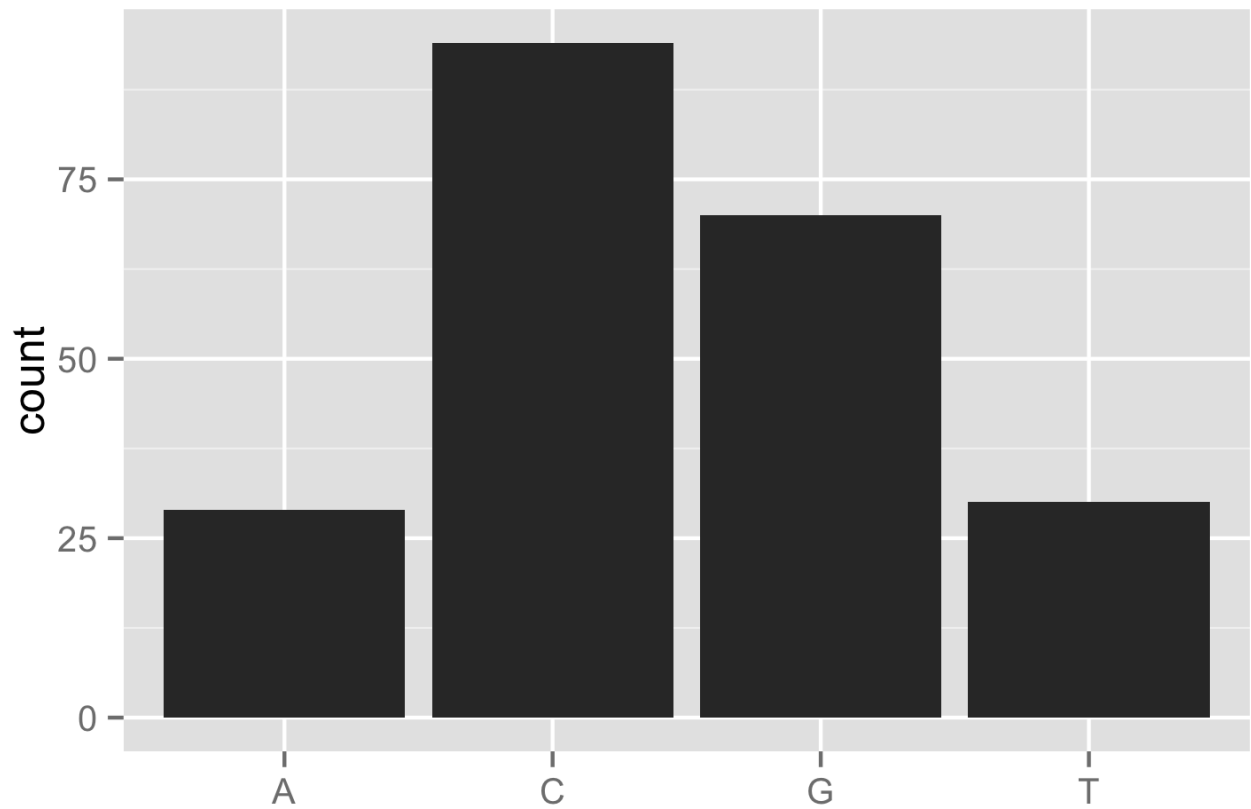
#> seqZyx
#>      A      C      G      T
#> 0.130 0.422 0.314 0.135

barplot(table(seqZyx))
```



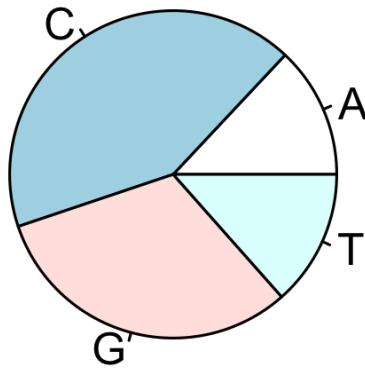
In *ggplot2* we first have to transform our data into an appropriate `data.frame`. We then map the bases to the x-axis and add binning by using the bar geometry.

```
ggBarplot = ggplot( data = data.frame(bases = factor(seqZyx)), aes(x = bases))  
ggBarplot = ggBarplot + xlab("") + geom_bar()  
ggBarplot
```

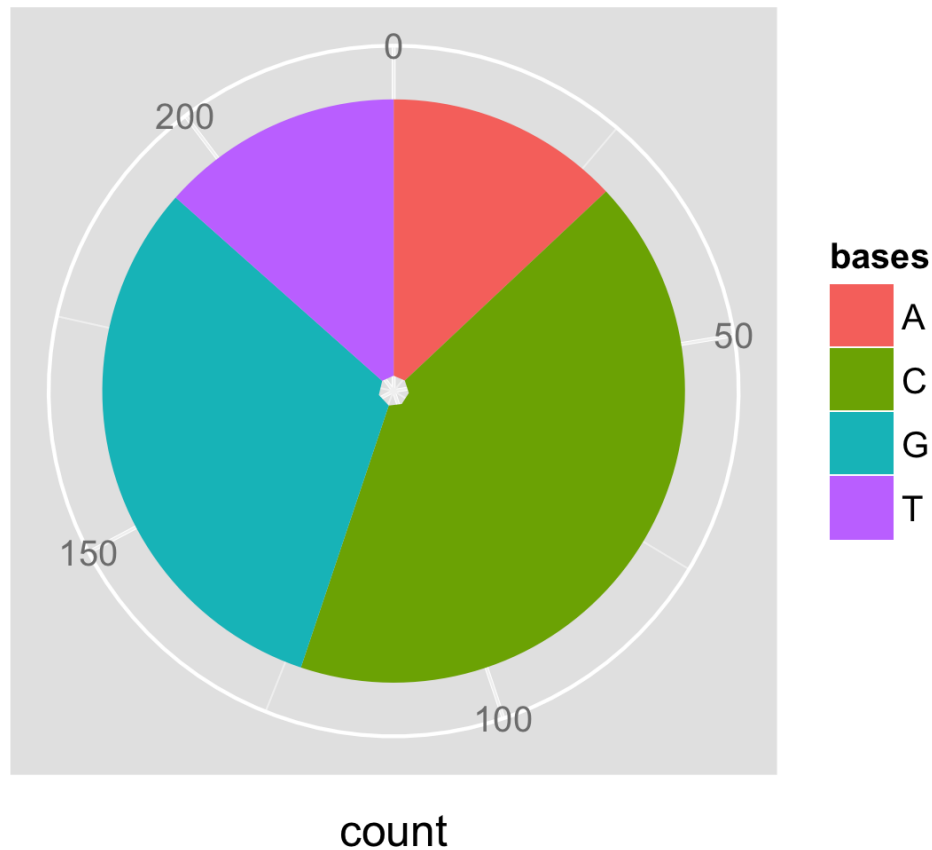


In order to plot a pie chart with *ggplot2*, we first have to create a stacked bar plot by mapping the base counts to the fill of the bar and then use polar coordinates.

```
pie(table(seqZyx))
```

```
ggPie = ggplot( data = data.frame(bases = factor(seqZyx)), aes(x = factor(1), fill = bases))
ggPie = ggPie + xlab(NULL) + geom_bar() + scale_x_discrete(breaks=NULL)
ggPie + coord_polar(theta = "y")
```



5.2 Scatterplots and Stripcharts

An elementary method to visualize data is the so-called scatterplot, which simply plots two sets of data against each other or, if applied to a single data set, plots each data point of the set according to its index.

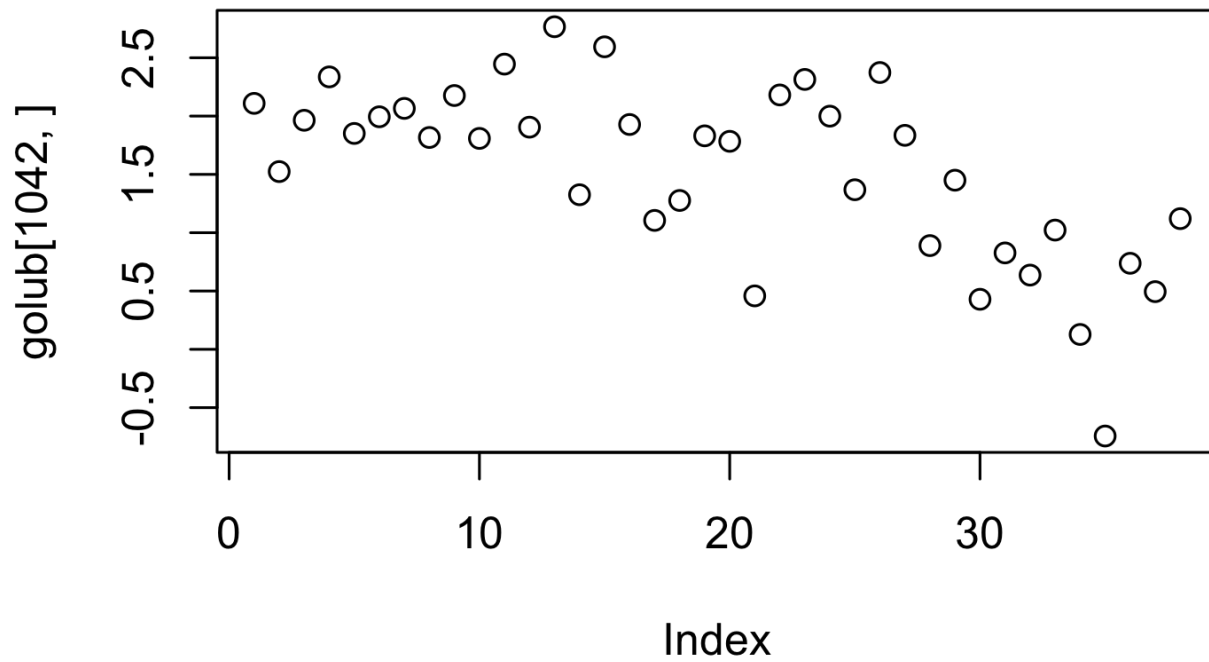
A stripchart is similar to a scatterplot, but plots all data points of a single data set horizontally or vertically. This is particularly useful in combination with a factor that distinguishes members from different experimental conditions or patients groups.

Example: Visualizing Gene CCND3

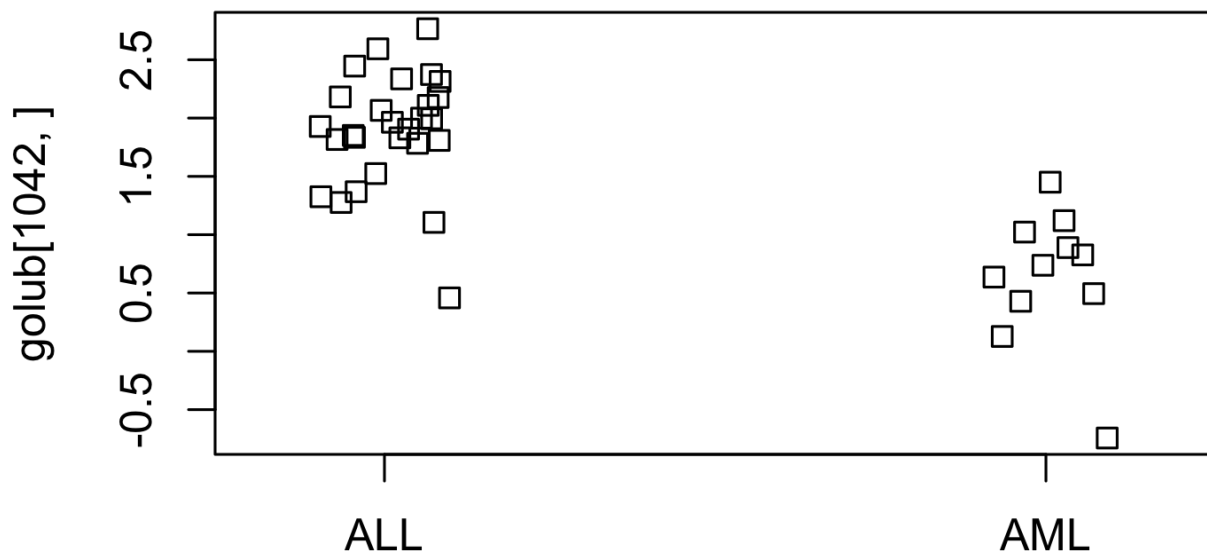
Many visualization methods will be illustrated by the Golub microarray data on two sub types of leukemia. We shall concentrate on the expression values of gene "CCND3 Cyclin D3", which are collected in row 1042 of the data matrix `golub`. To plot the data values one can simply use `plot(golub[1042,])`. In the resulting plot the vertical axis gives the size of the expression values and the horizontal axis the index of the patients. This is a scatterplot of a single data set.

It can be observed that the values for patient 28 to 38 are somewhat lower, but, indeed, the picture is not very clear because the groups are not plotted separately. To produce two adjacent stripcharts one for the ALL and one for the AML patients, we use a factor called `gol.fac` which separates the classes of patients.

```
plot(golub[1042,])
```



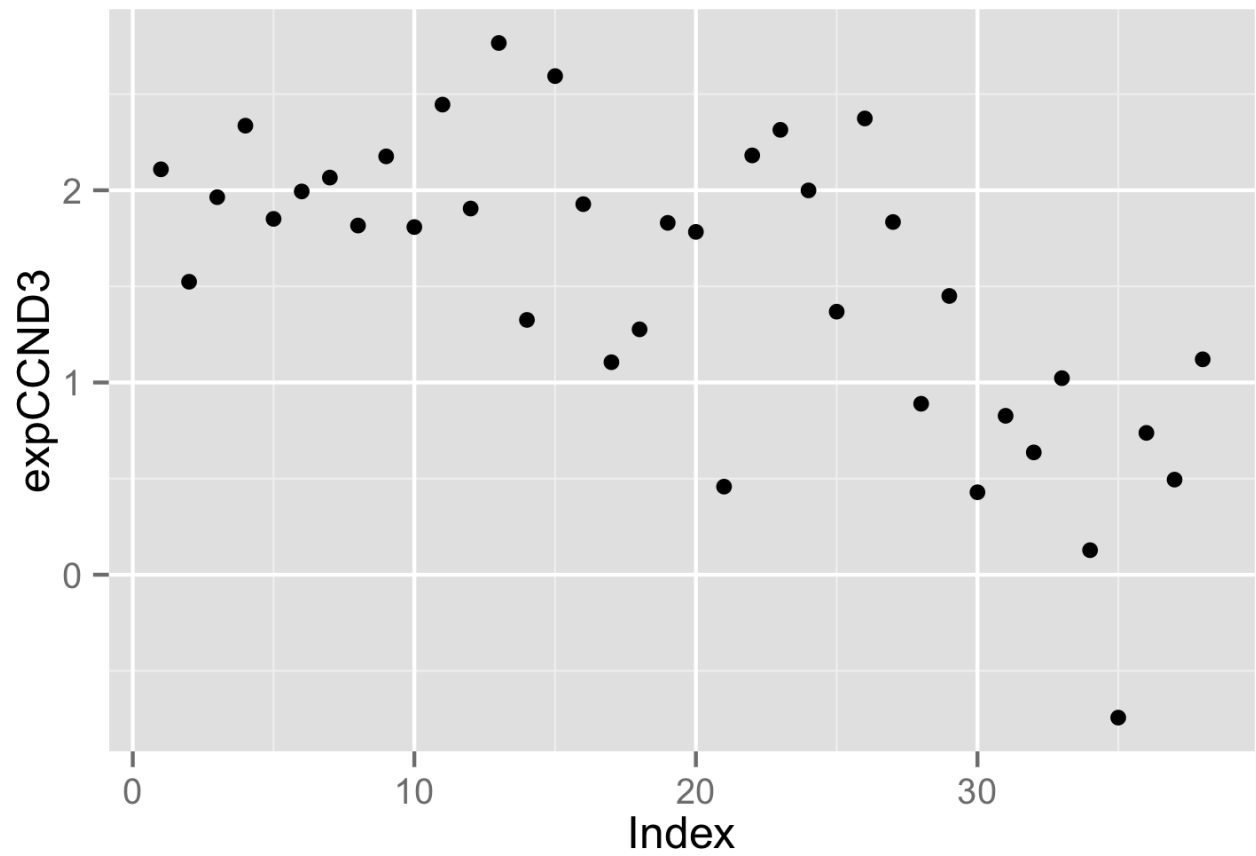
```
gol.fac <- factor(golub.cl, levels=0:1, labels= c("ALL", "AML"))  
stripchart(golub[1042,] ~ gol.fac, method="jitter", vertical = TRUE)
```



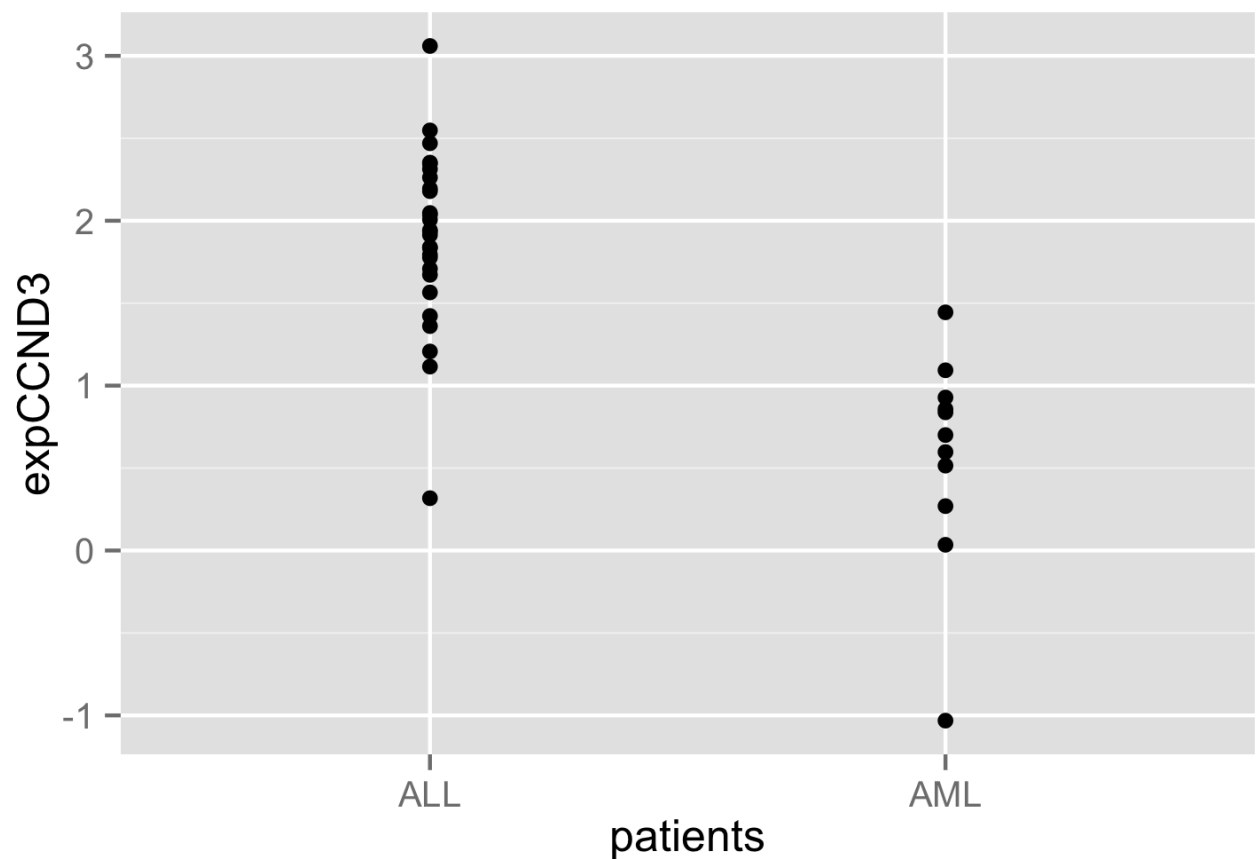
From the resulting figure it can be observed that the CCND3 expression values of the ALL patients tend to have larger expression values than those of the AML patients. The option `jitter` will “smear” the data points a little bit. In *ggplot2* these plots can be created like this:

```
CCND3 <- data.frame(expCCND3 = golub[1042,], patients = gol.fac)

ggScatter = qplot(seq_along(expCCND3), expCCND3, data = CCND3, geom = "point")
ggScatter + xlab("Index")
```



```
ggStrip = ggplot(CCND3, aes(x=patients, y=expCCND3))  
ggStrip + geom_point(position = position_jitter(w = .0, h = .30))
```



The “jittering” in *ggplot2* occurs after the plot has been created and the data has been mapped to aesthetics that is why it is called a “position adjustment”. Note that the degree of jittering can be controlled explicitly in *ggplot2*.

5.3 Histograms

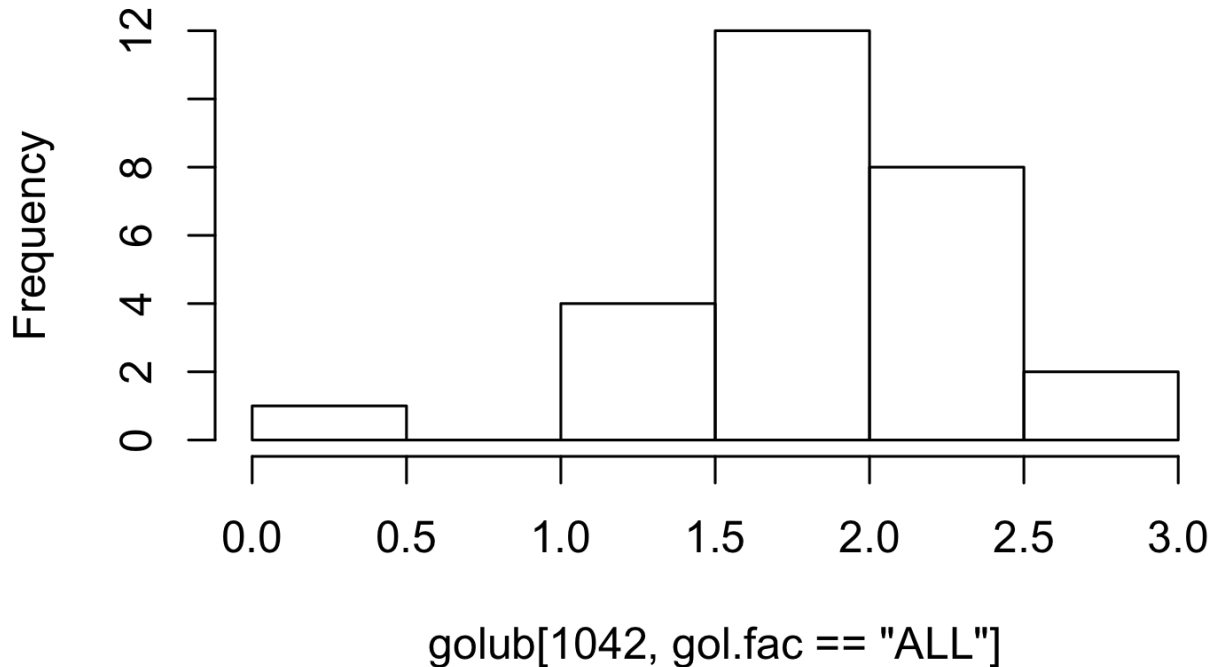
Another method to visualize data is by dividing the range of data values into a number of intervals and to plot the frequency per interval as a bar. Such a plot is called a histogram.

Example: Histogram of CCND3

A histogram of the expression values of gene CCND3 of the acute lymphoblastic leukemia patients can be produced as follows:

```
hist(golub[1042, gol.fac=="ALL"])
```

Histogram of golub[1042, gol.fac == "ALL"]



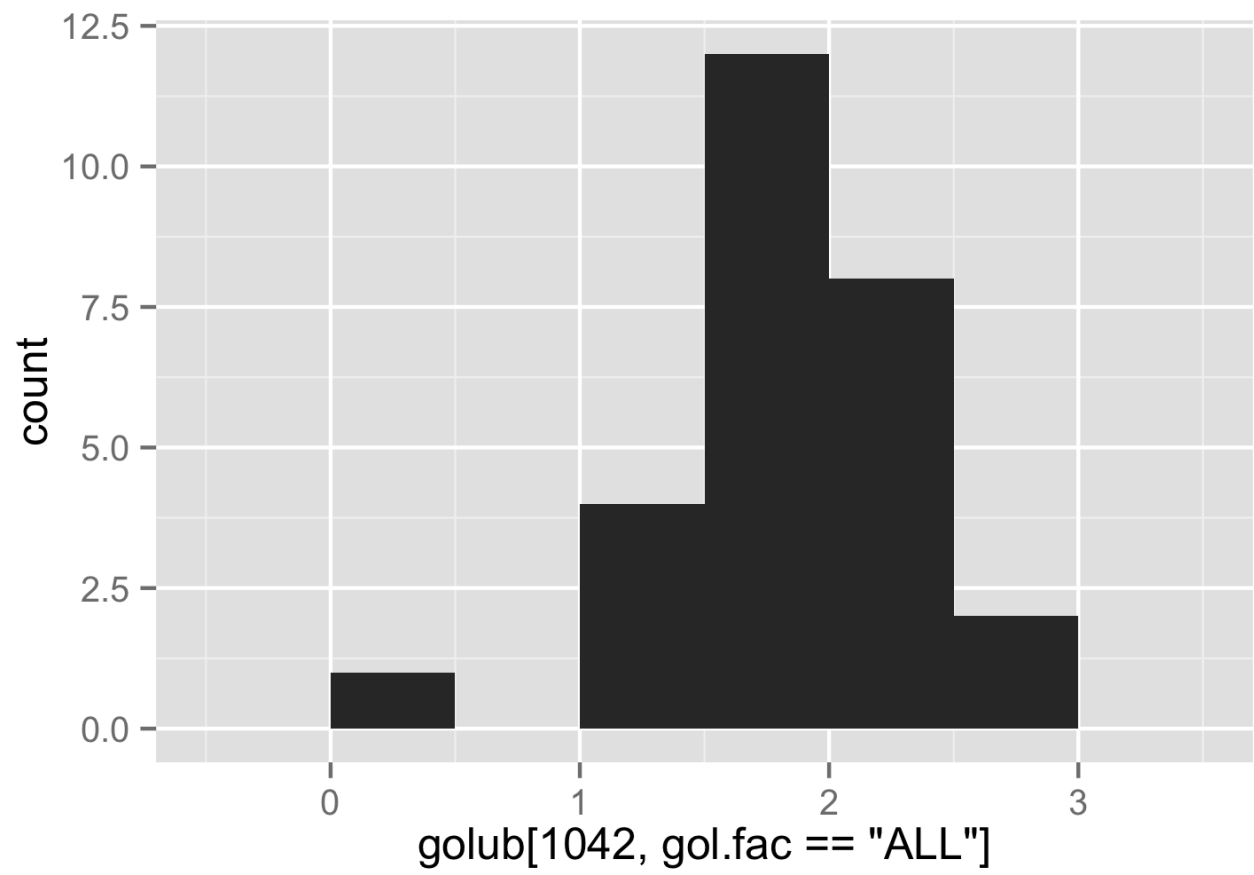
The function `hist` divides the data into 5 intervals having width equal to 0.5. Observe that one value is small and the other ones are more or less symmetrically distributed around the mean. The number of bins can be adjusted by the `breaks` option, choosing `prob = TRUE`, will give the relative frequency for each bin on the y -axis instead of the absolute frequencies.

ggplot2 uses binning as the default geometry if only x -values are supplied, so creating a histogram is really easy. By default it uses the range of the data divided by 30 as the binwidth. This can be easily adjusted. Here, we produce a barplot given the absolute counts per bin and then a classical histogram giving the relative frequencies.

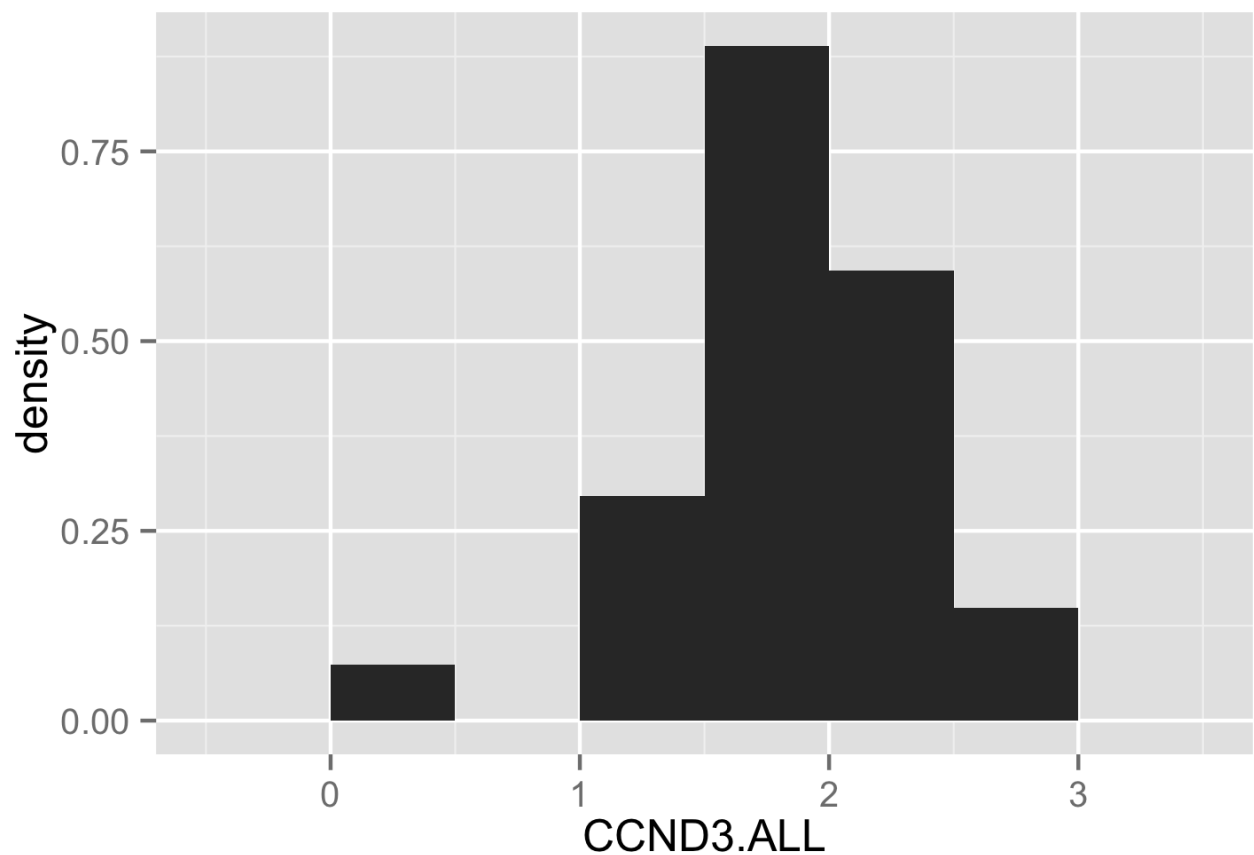
In the code below, `aes(y = ..density..)` means that we map the y -coordinate to the relative frequency estimate as returned by the stat `_bin` statistical transformation used by the “`geom_histogram`” function.

```
CCND3.ALL = data.frame(CCND3.ALL = golub[1042, gol.fac=="ALL"])
ggplotAbsFreq = qplot(golub[1042, gol.fac=="ALL"], binwidth = 0.5)
ggplotHist = ggplot(CCND3.ALL, aes(x=CCND3.ALL ))
ggplotHist = ggplotHist + geom_histogram(aes(y = ..density..), binwidth = 0.5)
```

```
ggplotAbsFreq
```



```
ggplotHist
```

5.4 Kernel density estimates

Kernel density estimates are a method to turn the histogram into a smooth density estimate. Given data x_1, \dots, x_n and a constant called "bandwidth" h , the kernel density estimate is given by:

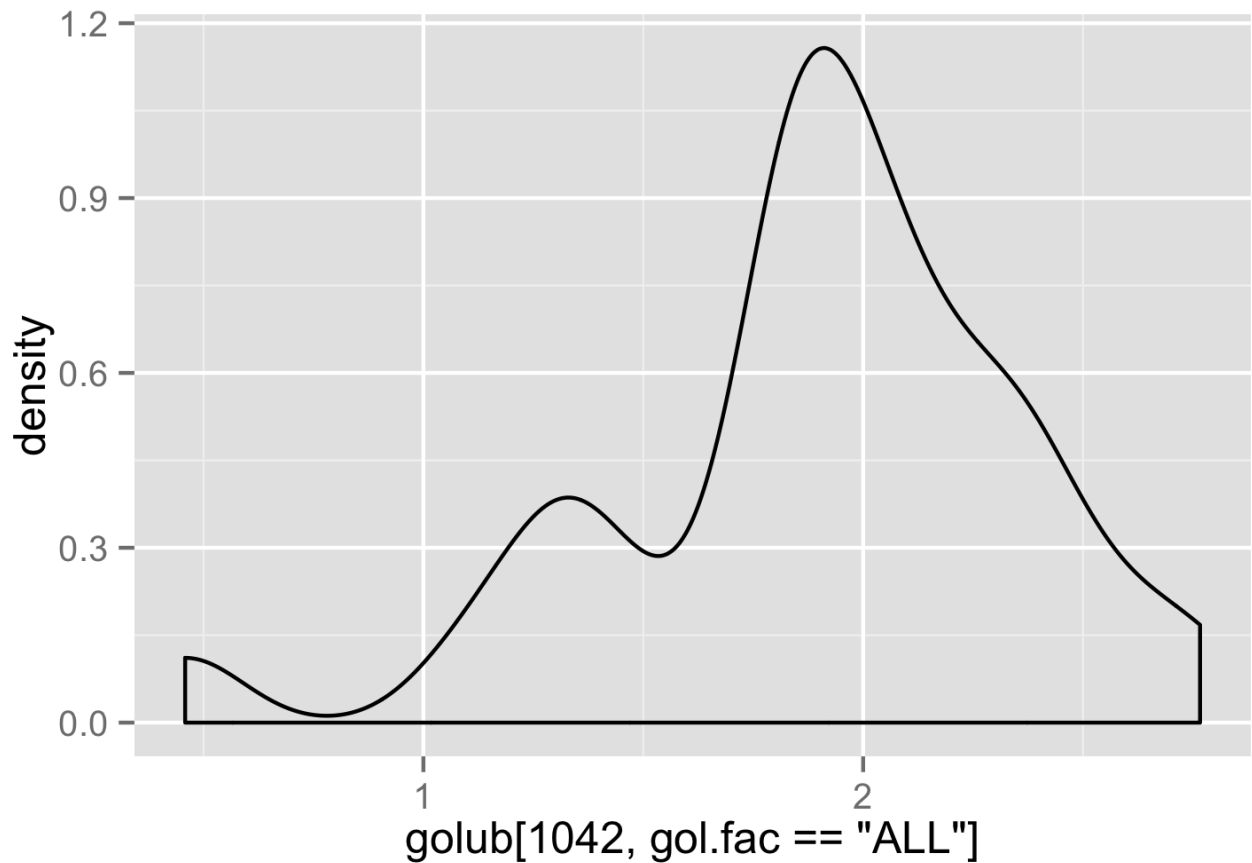
$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - x_i}{h}\right)$$

Typical kernels are:

- Bisquare kernel: $K(u) = \frac{15}{16}(1 - u^2)^2$ for $u \in [-1, 1]$ and 0 otherwise
- Gauß kernel: $K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right)$ for $u \in \mathbb{R}$

They readily computed using the *R* function `density` or the appropriate ggplot statistical transformation, which uses this function.

```
CCND3.ALL = data.frame(CCND3.ALL = golub[1042, gol.fac=="ALL"])
qplot(golub[1042, gol.fac=="ALL"], geom = "density")
```



Since we do not have many data points, the kernel density estimates gives a somewhat different impression than the histogram. Also note that kernel density estimates do not respect the boundaries of the data set by default.

5.5 Boxplots

It is always possible to sort n data values to have increasing order $x_1 \leq x_2 \leq \dots \leq x_n$, where x_1 is the smallest, x_2 is the first-to-the smallest, etc.

Let $x_{0.25}$ be a number for which it holds that 25% of the data values x_1, \dots, x_n are smaller. That is, 25% of the data values lay on the left side of the number $x_{0.25}$, the reason for which it is called the first quartile or the 25th percentile.

The second quartile is the value $x_{0.5}$ such that 50% of the data values are smaller. It is also called the median. Similarly, the third quartile or 75th percentile is the value $x_{0.75}$ such that 75% of the data is smaller.

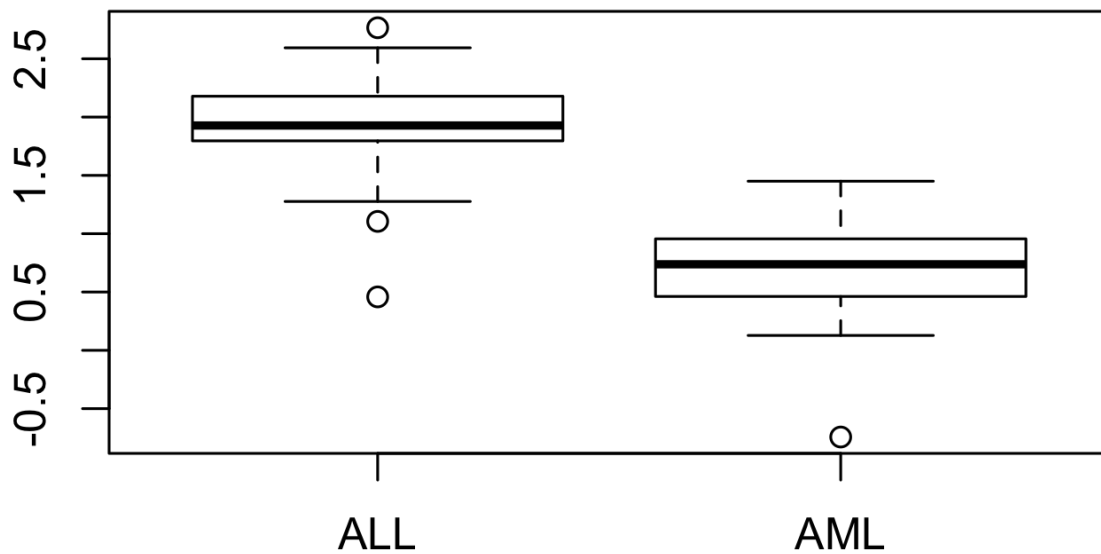
A popular method to display data is by drawing a box around the first and the third quartile, a bold line segment for the median, and the smaller line segments (whiskers) for the smallest and the largest data values. Such a data display is known as a box-and-whisker plot or simply boxplot.

Example: Boxplot of CCND3

A view on the distribution of the expression values of the ALL and the AML patients on gene CCND3 can be obtained by constructing two separate boxplots adjacent to one another. To produce such a plot the factor `gol.fac` is again very

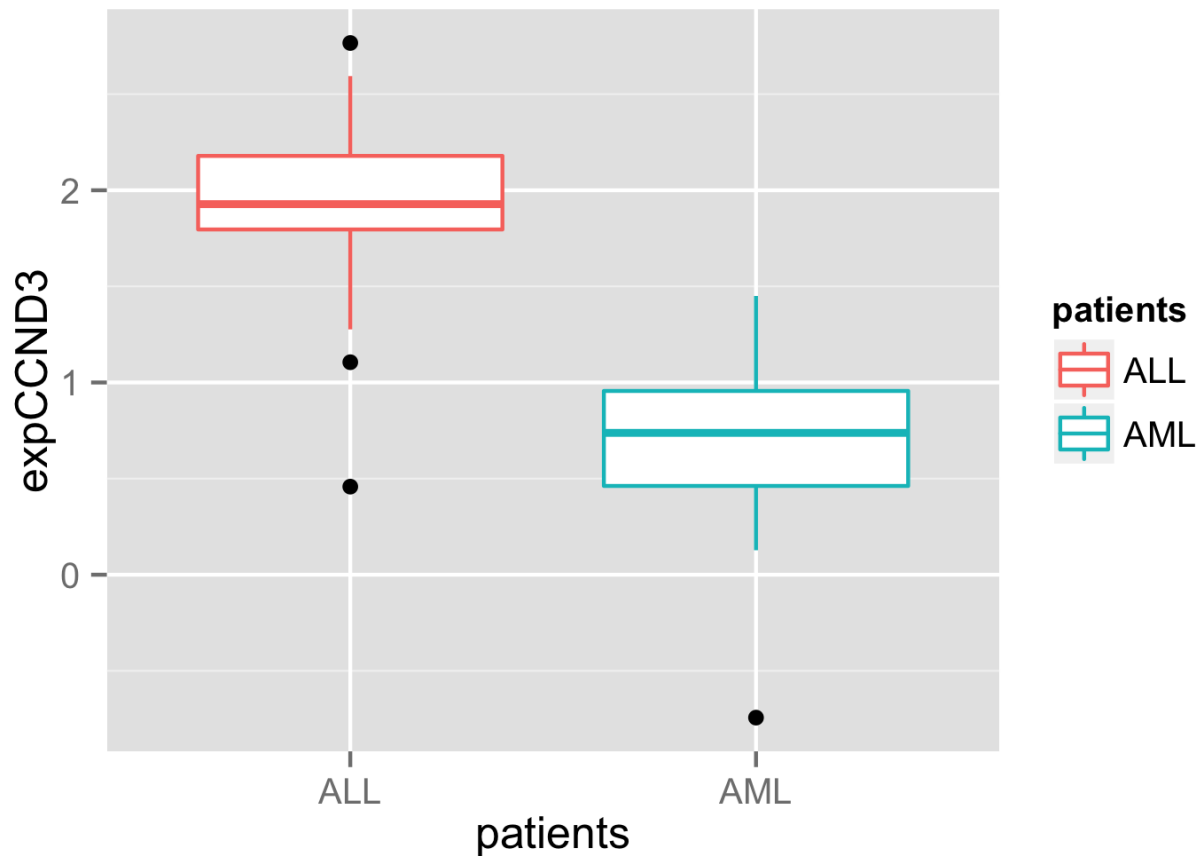
useful.

```
boxplot(golub[1042,] ~ gol.fac)
```



In *ggplot2* a boxplot can be exactly as a stripchart, we just have to change the geometry of the plot:

```
ggBox= ggplot(CCND3, aes(x=patients, y=expCCND3, color = patients))  
ggBox + geom_boxplot()
```



From the position of the boxes, it can be observed that the gene expression values for ALL are larger than those for AML. Furthermore, since the two sub-boxes around the median are more or less equally wide, the data are quite symmetrically distributed around the median. The quantiles can be computed with the function `quantile`. By default the four “quartiles” are returned but other quantiles can be obtained by specifying the `probs` parameter

```
quantile(golub[1042,])
#>      0%    25%    50%    75%   100%
#>  -0.743  1.043  1.813  2.049  2.766

## custom quantiles
quantile(golub[1042,], probs = c(0, 0.1, 0.5))
#>      0%    10%    50%
#>  -0.743  0.484  1.813
```

Outliers are data values laying far apart from the pattern set by the majority of the data values. The implementation in R of the (modified) boxplot draws such outlier points separately as small circles. A data point x is defined as an outlier point if

$$x < x_{0.25} - 1.5 \cdot (x_{0.75} - x_{0.25})$$

or

$$x > x_{0.75} + 1.5 \cdot (x_{0.75} - x_{0.25})$$

From the boxplot above it can be observed that there are outliers among the gene expression values of ALL patients. These are the smaller values 0.45827 and 1.10546, and the largest value 2.76610. The AML expression values have one

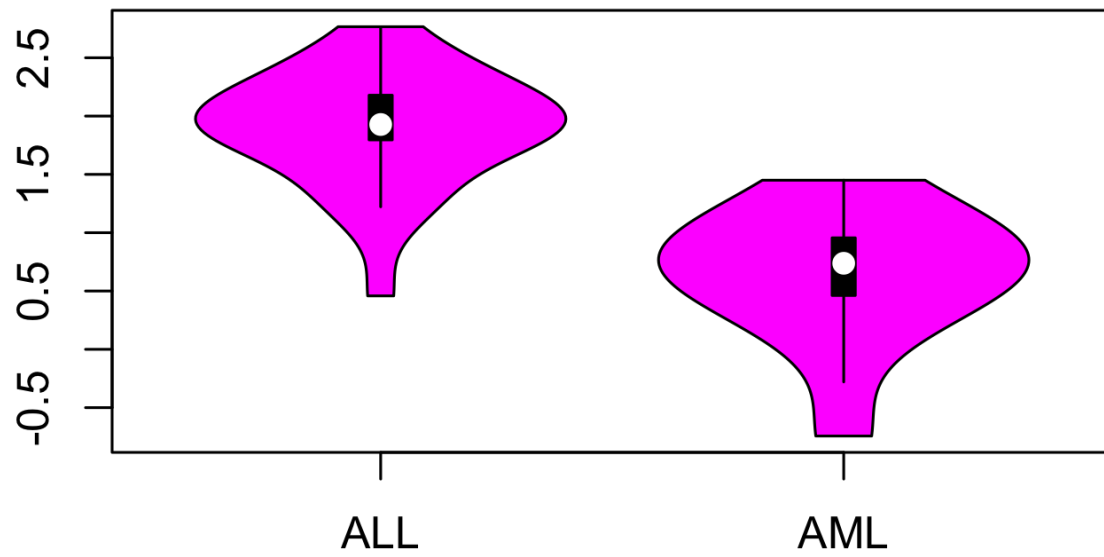
outlier with value -0.74333. To define extreme outliers, the factor 1.5 is raised to 3.0. These numbers can be conveniently extracted from the `boxplot.stats` function:

```
boxplot.stats(golub[1042, gol.fac == "ALL"])  
  
#> $stats  
#> [1] 1.28 1.80 1.93 2.18 2.59  
#>  
#> $n  
#> [1] 27  
#>  
#> $conf  
#> [1] 1.81 2.04  
#>  
#> $out  
#> [1] 2.766 1.105 0.458  
  
### outliers are in the "out" element of the list
```

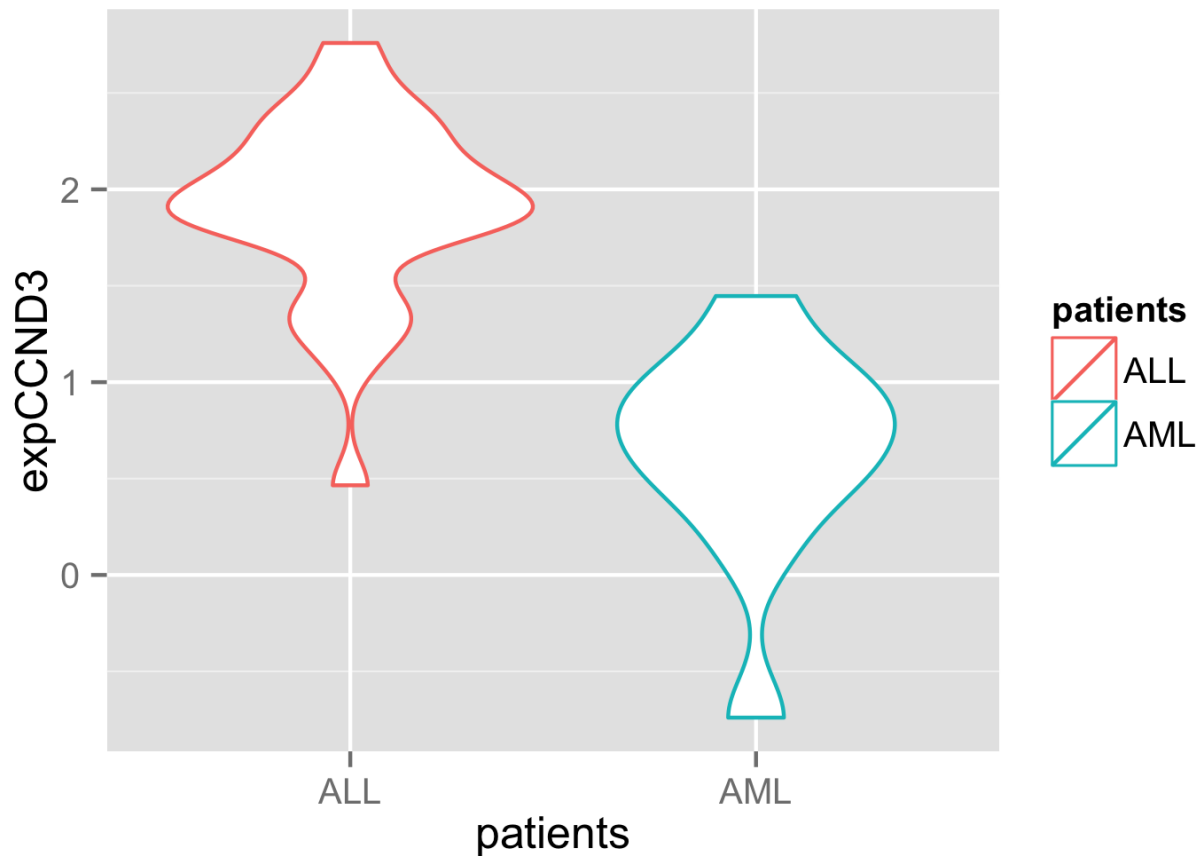
5.6 Violin plots

A violin plot is a combination of a boxplot and a kernel density estimate. Specifically, instead of straight borders for the boxes, a kernel density estimate is displayed.

```
vioplot(split(golub[1042,], gol.fac)[[1]], split(golub[1042,], gol.fac)[[2]],  
        names = c("ALL", "AML"))
```



```
ggBox + geom_violin()
```



5.7 Empirical Cumulative Distribution Function

Another way to visualize data is the empirical cumulative distribution function (ecdf). It is an estimator for the actual cumulative distribution function, which gives the probability of observing values less than a given value t . Especially if you have a large number of data points, it gives a useful visualization tool for your data. For every number t the ecdf $\hat{F}_n(t)$ is the fraction of data points that are smaller than t .

$$\hat{F}_n(t) = \frac{\text{number of elements in the sample} \leq t}{n}$$

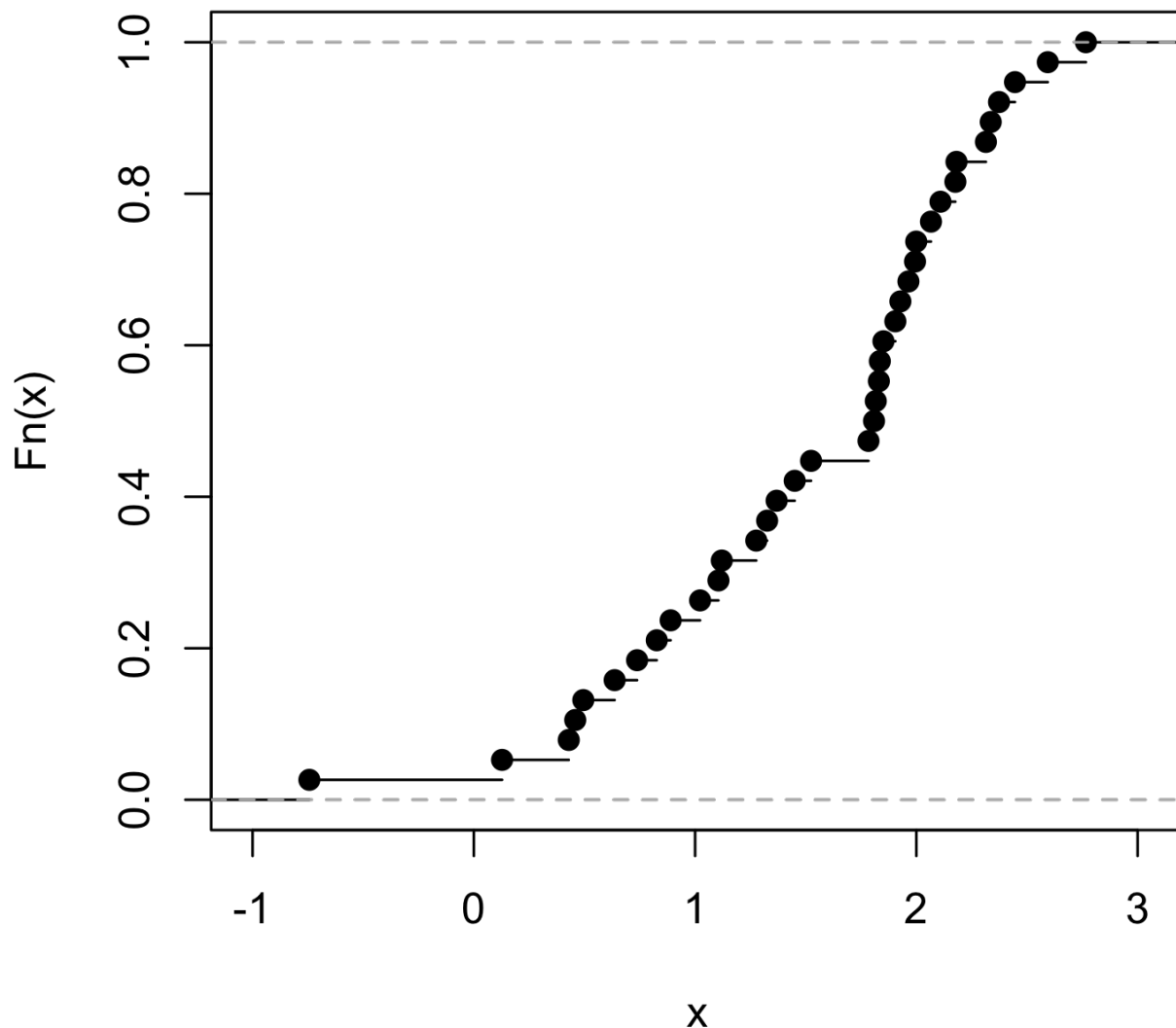
Obviously, just as the cdf the ecdf will always start 0 and end up at 1. You can easily read off quantiles from the ecdf, since the y -axis represents the quantiles of the data.

Example: ECDF of CCND3

From the plot of the ecdf of CCND3, it can be seen that the median of the overall data is roughly 1.8. The steep increase of the ecdf for values greater than the median shows that larger values are more common in the data set.

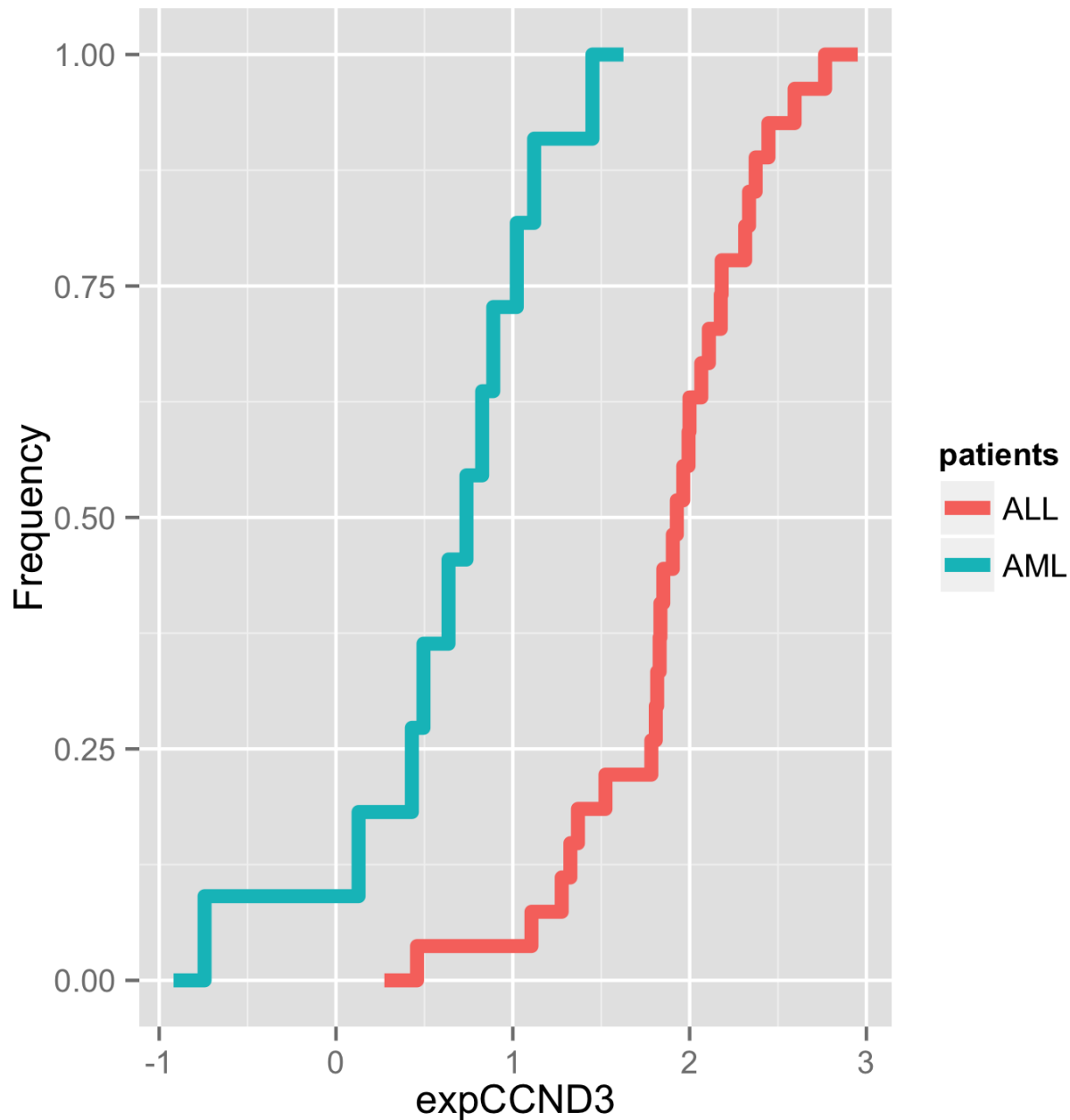
```
plot(ecdf(golub[1042, ]), main = "ecdf of CCND3")
```

ecdf of CCND3



In *ggplot2* there exists a specialized statistic which performs the calculation necessary to obtain the ecdf. It is also easy to split the ecdf calculation by a factor, here we clearly see that for CCND3, the expression for ALL patients is always higher than the expression for AML patients.

```
ggECDF = ggplot(CCND3, aes(x=expCCND3, color = patients))  
ggECDF + stat_ecdf(geom = "step", size = 2) + ylab("Frequency")
```

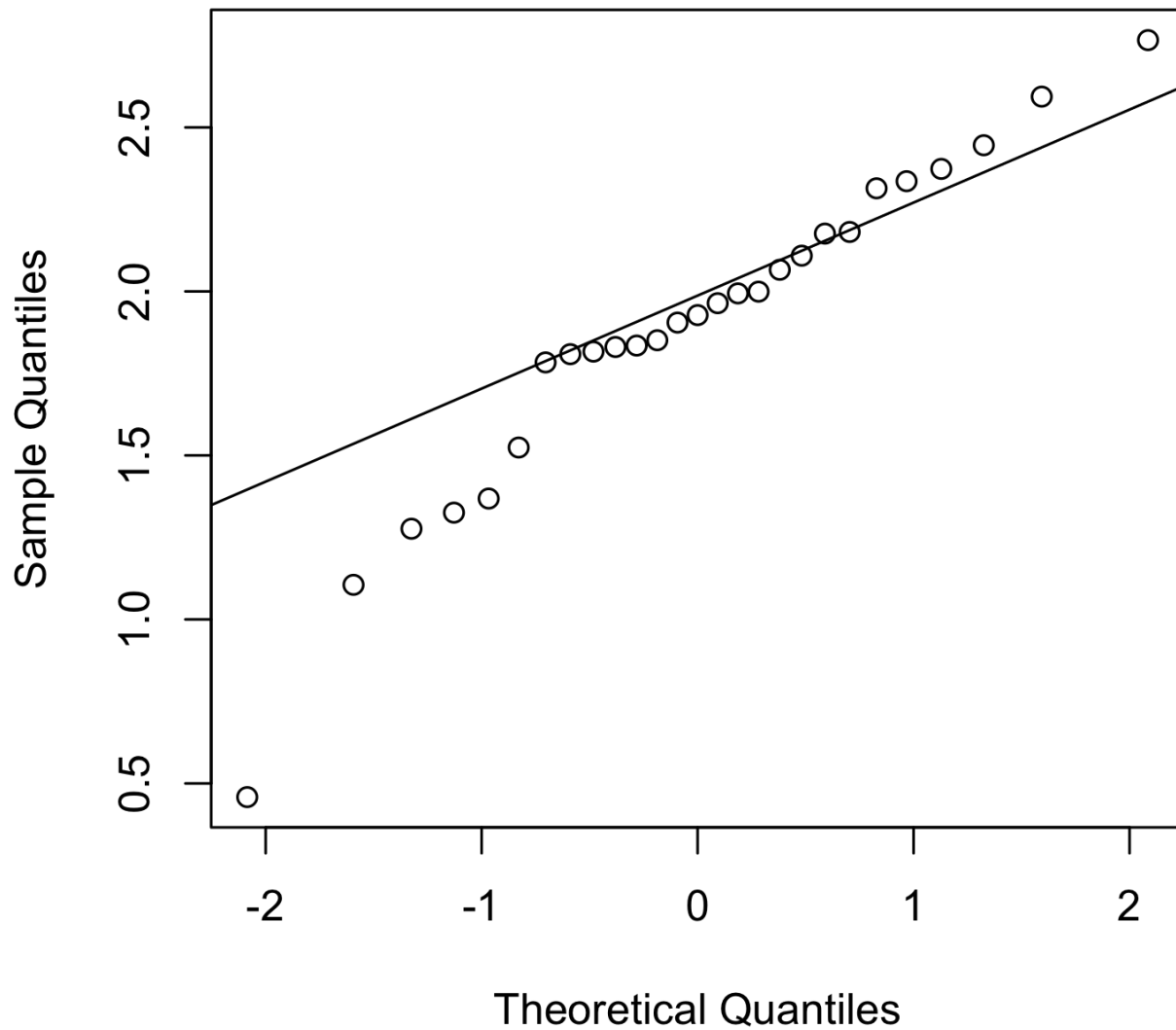
5.8 Quantile–Quantile (QQ) Plot

A method to compare the (e)cdf of a data set to another (e)cdf is so-called quantile-quantile (Q–Q) plot. In such a plot, usually the quantiles of the data set are displayed against the corresponding quantiles of the normal distribution (bell-shaped). Hence, the empirical cdf of the data are compared to the (corresponding) normal cdf.

A straight line is added representing points which correspond exactly to the quantiles of a corresponding normal distribution. By observing the extent in which the points appear on the line, it can be evaluated to what degree the data are normally distributed. That is, the closer the values appear to the line, the more likely it is that the data are normally distributed. To produce a Q–Q plot of the ALL gene expression values of CCND3 one may use the following code

```
qqnorm(golub[1042, gol.fac=="ALL"])  
qqline(golub[1042, gol.fac=="ALL"])
```

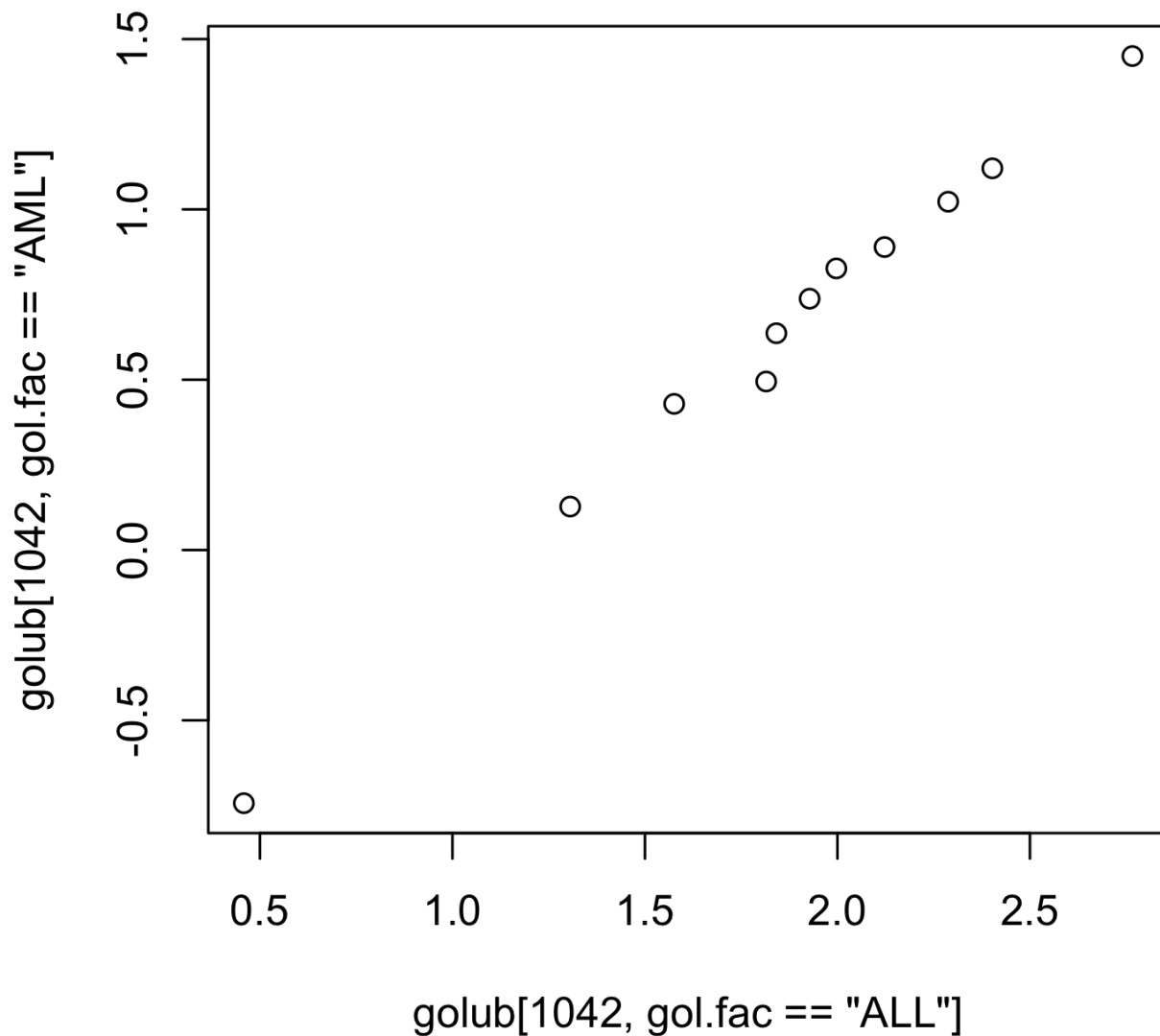
Normal Q-Q Plot



From the resulting figure it can be observed that most of the data points are on or near the straight line, while a few others are further away. This is an expected behavior for gene expression data. The above example illustrates a case where the degree of non-normality is moderate so that a clear conclusion cannot be drawn.

The QQ-plot can also be used to compare the distributions of two data sets to one another by plotting their quantiles against each other. In order to do this, you can call the function `qqplot`. A comparison of the distributions / ecdfs of AML and ALL gene expression values is given by:

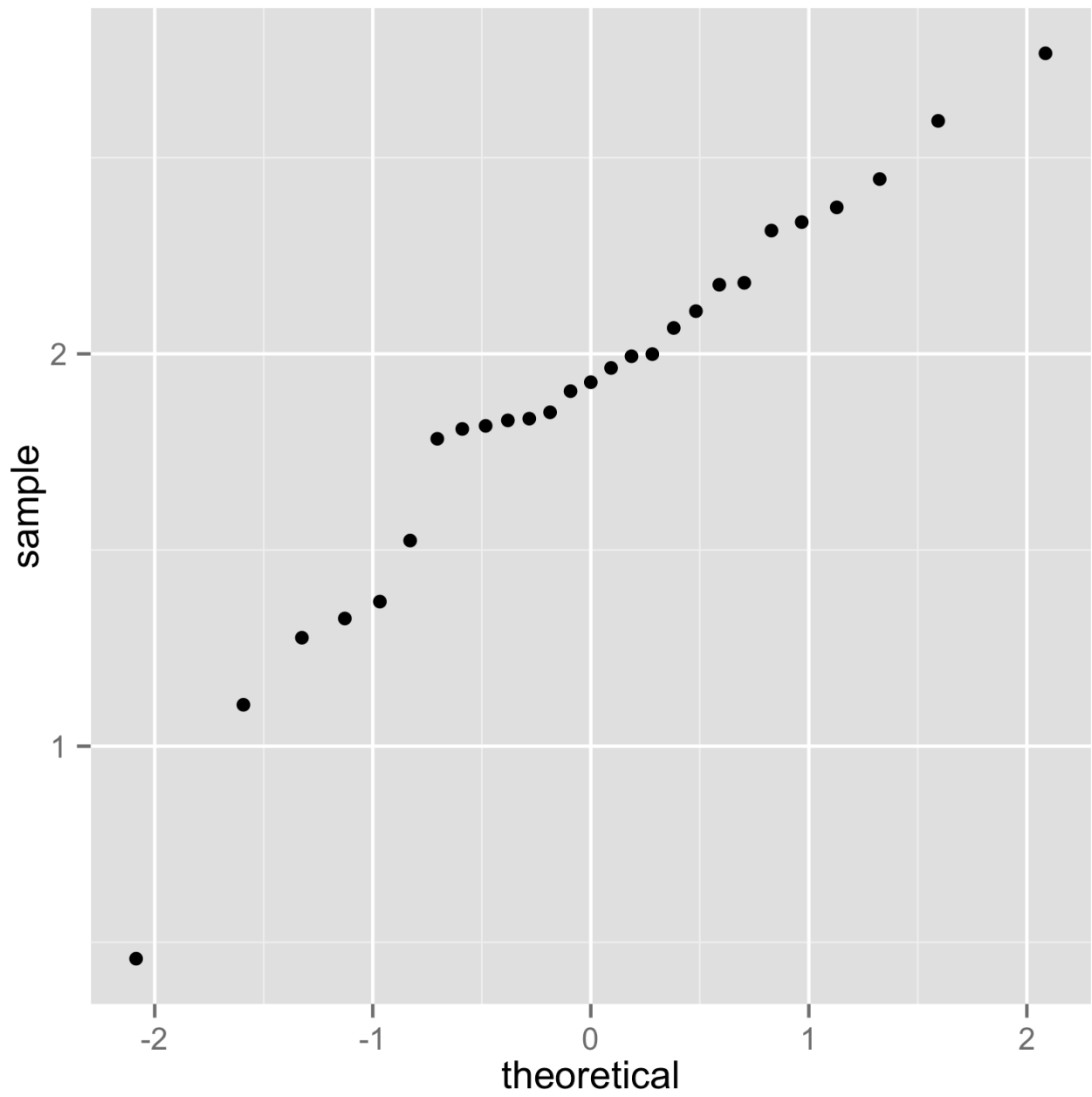
```
qqplot(golub[1042, gol.fac=="ALL"], golub[1042, gol.fac=="AML"])
```



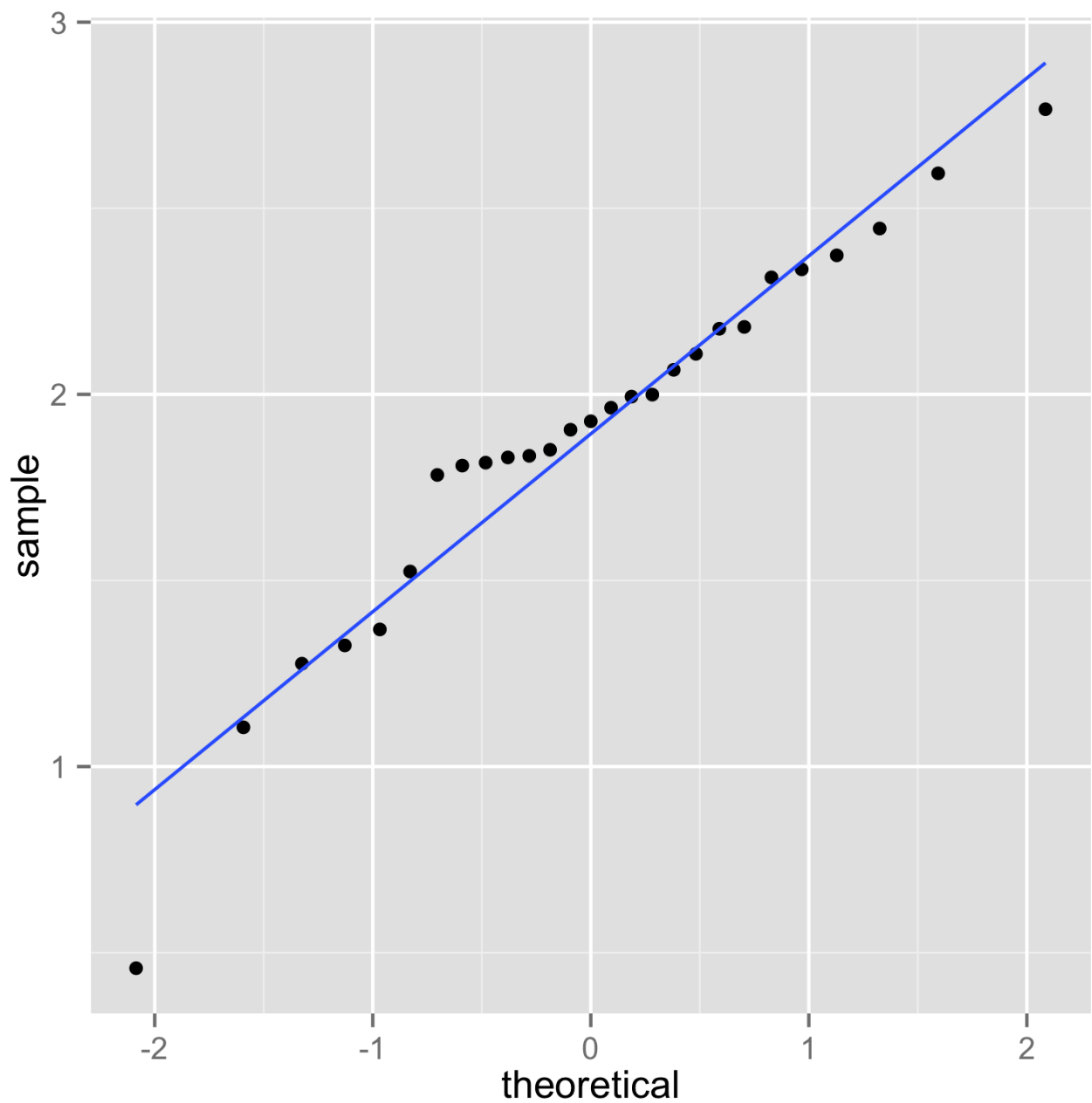
As we can see, the quantiles lie on a straight line, which indicates that their general shape agrees. However, they might have different means and standard deviations.

In *ggplot2* there exists a specialized statistic which performs the calculation necessary to obtain the QQ-plot. The `qnorm` equivalent is a bit harder to obtain, we have to extract the augmented data frame returned by the `stat_qq` function and add a linear smoother.

```
ggQQ = ggplot(CCND3.ALL, aes(sample=CCND3.ALL)) + stat_qq()
#ggQQ
tp = cbind(CCND3.ALL, as.data.frame((print(ggQQ)$data)[[1]][, c("sample", "theoretical"))])
```



```
ggQQ + stat_smooth(mapping = aes(x = theoretical, y = sample), data = tp, method = "lm", se = F )
```



6 Descriptive Statistics

There exist various ways to describe the central tendency as well as the spread of data. In particular, the central tendency can be described by the mean or the median, and the spread by the variance, standard deviation, interquartile range, or median absolute deviation. These will be defined and illustrated.

6.1 Measures of Central Tendency

The most important descriptive statistics for central tendency are the mean and the median. The sample mean of the data values x_1, \dots, x_n is defined as:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} (x_1 + \dots + x_n).$$

Thus the sample mean is simply the average of the n data values. Since it is the sum of all data values divided by the sample size, a few extreme data values may largely influence its size. In other words, the mean is not robust against outliers.

The median is defined as the second quartile or the 50th percentile, and is denoted by $x_{0.50}$. When the data are symmetrically distributed around the mean, then the mean and the median are equal. Since extreme data values do not influence the size of the median, it is very robust against outliers. Robustness is important in biological applications because data are frequently contaminated by extreme or otherwise influential data values.

Example: Mean and Median of CCND3

To compute the mean and median of the ALL expression values of gene CCND3 Cyclin D3 consider the following.

```
mean(golub[1042, gol.fac=="ALL"])
#> [1] 1.89
median(golub[1042, gol.fac=="ALL"])
#> [1] 1.93
```

Note that the mean and the median do not differ much so that the distribution seems to be quite symmetric.

6.2 Measures of Spread

The most important measures of spread are the standard deviation, the interquartile range, and the median absolute deviation. The standard deviation is the square root of the sample variance, which is defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right).$$

Hence, it is the average of the squared differences between the data values and the sample mean. The sample standard deviation s is the square root of the sample variance and may be interpreted as the distance of the data values to the mean. The variance and the standard deviation are not robust against outliers.

The interquartile range is defined as the difference between the third and the first quartile, that is $x_{0.75} - x_{0.25}$. It can be computed by the function `IQR(x)`. More specifically, the value `IQR(x)/1.349` is a robust estimator of the standard deviation.

The median absolute deviation (MAD) is defined as a constant times the median of the absolute deviations of the data from the median. In R it is computed by the function `mad` defined as the median of the sequence $|x_1 - x_{0.5}|, \dots, |x_n - x_{0.5}|$ multiplied by the constant 1.4826. It equals the standard deviation in case the data come from a bell-shaped (normal) distribution. Because the interquartile range and the median absolute deviation are based on quantiles, these are robust against outliers.

Example: Measures of Spread for CCND3

These measures of spread for the ALL expression values of gene CCND3 can be computed as follows.

```
sd(golub[1042, gol.fac=="ALL"])
#> [1] 0.491

IQR(golub[1042, gol.fac=="ALL"]) / 1.349
#> [1] 0.284

mad(golub[1042, gol.fac=="ALL"])
#> [1] 0.368
```

Due to the three outliers (cf. boxplot above) the standard deviation is larger than the interquartile range and the mean absolute deviation. That is, the absolute differences with respect to the median are somewhat smaller than the root of the squared differences.

Exercise: Illustration of Mean and Standard Deviation

- Compute the mean and the standard deviation for 1, 1.5, 2, 2.5, 3.
- Compute the mean and the standard deviation for 1, 1.5, 2, 2.5, 30.
- Comment on the differences.

Exercise: Plotting Gene Expressions of CCND3

Use the gene expressions from "CCND3" of Golub collected in row 1042 of the object golub.

- Produce a stripchart for the gene expressions separately for the ALL as well as for the AML patients. Hint: Use a factor for appropriate separation.
- Rotate the plot to a vertical position and keep it that way for the questions to come.
- Color the ALL expressions red and AML blue. Hint: Use the `col` parameter.
- Add a title to the plot. Hint: Use `title`.
- Change the boxes into stars. Hint: Use the `pch` parameter.

Exercise: Plotting Gene Expressions using ggplot2

In this exercise we will plot the gene expressions for a couple of genes in the golub data set in a single figure, using a separate panel for each gene.

- Turn the golub data set into a data frame with genes in the columns. Add the group descriptor `gol.fac` as a additional column and turn it into a factor.
- Select a random sample of 6 genes from the golub data. HINT: Use the function `sample` to do this.
- Melt the resulting data frame containing the selected genes in such a way that all the gene expression values are in a single column.
- Use this data frame and `ggplot2` to produce a pdf file containing boxplots separated per patient group for each of the randomly selected genes. Also add the raw data points to the plot.

Exercise: Comparing Normality for Two Genes

Consider the gene expression values in row 790 and 66 of the Golub data.

- Produce a boxplot for the expression values of the ALL patients and comment on the differences. Are there outliers?
- Produce a QQ-plot and formulate a hypothesis about the normality of the genes.

- (c) Compute the mean and the median for the expression values of the ALL patients and compare these. Do this for both genes.

7 Answers to Exercises

Exercise: Simple ggplot usage

- (a) Use points as a geometry instead of lines
 (b) Use both lines and points
 (c) Add errorbars `geom_errorbar` to the plot. This requires further aesthetics: `ymin` and `ymax`. The estimated error is stored in the variable `Sigma`.

Solution: Simple ggplot usage

```
#a
#####
plot_1 <- ggplot(aes( x = min, y = Signal ), data = proteins_pMek_sub)
plot_1 <- plot_1 + geom_point()
plot_1 <- plot_1 + xlab("Time [min]") + ylab("pMEK Signal")

#b
#####
plot_1 <- plot_1 + geom_line()

#c
#####
plot_1 <- plot_1 + geom_errorbar(aes(ymax = Signal+2*Sigma,
ymin = Signal-2*Sigma))
```

Exercise: ggplot faceting

Using the data `proteins_pMek_sub`, to produce a plot split by the experimental condition factor using `facet_wrap()`.

Solution: ggplot faceting

```
(qplot(min, Signal, data = proteins_pMek, geom = "line", color = Condition)
+ facet_wrap( ~ Condition) + geom_errorbar(aes(ymax = Signal+2*Sigma,
ymin = Signal-2*Sigma)))
```

Exercise: complex ggplot example

Use the data `proteins`, to produce a plot of the time courses split by the experimental target and colored according to the experimental conditions. Add error bars to your plot.

Solution: complex ggplot example

```
(qplot(min, Signal, data = proteins, geom = "line", color = Condition)
+ facet_wrap( ~ Target, ncol = 2) + geom_errorbar(aes(ymax = Signal+2*Sigma,
ymin = Signal-2*Sigma)))
```


Exercise: Illustration of Mean and Standard Deviation

- Compute the mean and the standard deviation for 1, 1.5, 2, 2.5, 3.
- Compute the mean and the standard deviation for 1, 1.5, 2, 2.5, 30.
- Comment on the differences.

Solution: Illustration of Mean and Standard Deviation

```
#Use x<- c(1,1.5,2,2.5,3) and mean(x) and sd(x) to obtain
#that the mean is 2 and the standard deviation is 0.79
(b) Now the mean is 7.4 and the standard deviation dramatically increased
(c) The outlier increased the mean as well as the standard deviation.
```

Exercise: Plotting Gene Expressions of CCND3

Use the gene expressions from "CCND3" of Golub collected in row 1042 of the object golub.

- Produce a so-called stripchart for the gene expressions separately for the ALL as well as for the AML patients. Hint: Use a factor for appropriate separation.
- Rotate the plot to a vertical position and keep it that way for the questions to come.
- Color the ALL expressions red and AML blue. Hint: Use the `col` parameter.
- Add a title to the plot. Hint: Use `title`.
- Change the boxes into stars. Hint: Use the `pch` parameter.

Solution: Plotting Gene Expressions of CCND3

```
gol.fac <- factor(golub.cl,levels=0:1, labels= c("ALL","AML"))
stripchart(golub[1042,] ~ gol.fac,method="jitter")
stripchart(golub[1042,] ~ gol.fac,method="jitter",vertical = TRUE)
stripchart(golub[1042,] ~ gol.fac,method="jitter",col=c("red", "blue"),
vertical = TRUE)
stripchart(golub[1042,] ~ gol.fac,method="jitter",col=c("red", "blue"), pch="*"
,vertical = TRUE)
title("CCND3 Cyclin D3 expression value for ALL and AML patients")
```

Exercise: Plotting Gene Expressions using ggplot2

In this exercise we will plot the gene expressions for a couple of genes in the golub data set in a single figure, using a separate panel for each gene.

- Turn the golub data set into a data frame with genes in the columns. Add the group descriptor `gol.fac` as a additional column and turn it into a factor.
- Select a random sample of 6 genes from the golub data. HINT: Use the function `sample` to do this.
- Melt the resulting data frame containing the selected genes in such a way that all the gene expression values are in a single column.
- Use this data frame and [ggplot2](#) to produce a pdf file containing boxplots separated per patient group for each of the randomly selected genes. Also add the raw data points to the plot.

Solution: Plotting Gene Expressions using ggplot2

```
#a
#####
golub.df = as.data.frame(cbind(t(golub), gol.fac))
golub.df$gol.fac = as.factor(golub.df$gol.fac)

#b
#####
rand.sample <- c(sample(dim(golub)[1],6),3052)

#c
#####
dataForPlot = melt(golub.df[, rand.sample],
  id = 'gol.fac')

#e
#####
pdf("boxplots.pdf", width = 14, height = 10)
p <- qplot(gol.fac, value, data = dataForPlot) +
  geom_boxplot(aes(fill = gol.fac)) +
  facet_wrap(~ variable) +
  geom_point(colour = 'black', alpha = 0.5)
p
dev.off()
```

Exercise: Comparing Normality for Two Genes

Consider the gene expression values in row 790 and 66 of the Golub data.

- Produce a boxplot for the expression values of the ALL patients and comment on the differences. Are there outliers?
- Produce a QQ-plot and formulate a hypothesis about the normality of the genes.
- Compute the mean and the median for the expression values of the ALL patients and compare these. Do this for both genes.

Solution: Comparing Normality for Two Genes

```
#Comparing two genes
#(a) Use
boxplot(golub[66,]~gol.fac)
dev.new()
boxplot(golub[790,]~gol.fac)
#to observe that 790 has three
#outliers and 66 has no outlier.
# (dev.new() opens a new graphical window)

#(b) Use
qqnorm(golub[66,gol.fac=="ALL"])
qqline(golub[66,gol.fac=="ALL"])
dev.new()
qqnorm(golub[790,gol.fac=="ALL"])
qqline(golub[790,gol.fac=="ALL"])
#to observe that nearly all values of 66 are on the line, where as for
```

```
#790 the three outliers are way of the normality line. Hypothesis:  
#The expression values of 66 are normally distributed, but those of  
#row 790 are not.  
  
 #(c) Use  
mean(golub[66,gol.fac=="ALL"])  
median(golub[790,gol.fac=="ALL"])  
#and#  
mean(golub[790,gol.fac=="ALL"])  
median(golub[790,gol.fac=="ALL"])  
#The mean (-1.174024) is larger than the median (-1.28137) due to  
#outliers on the right hand side. For the gene in row 66 the mean is  
#1.182503 and the median 1.23023. The differences are smaller.
```