

Social Network Analysis Home Assignment 1

Vladislav Chagaev

15.05.2015 23:59

Contents

| | |
|--|----------|
| Power law. Descriptive network analysis | 1 |
| Problem 1 | 1 |
| Problem 2 | 4 |
| Problem 3. | 6 |

Power law. Descriptive network analysis

```
##
## Attaching package: 'igraph'

## The following objects are masked from 'package:stats':
##
##      decompose, spectrum

## The following object is masked from 'package:base':
##
##      union
```

Problem 1

Recall from the lecture that probability density function (PDF) for power law distributed variable is:

$$p(x) = Cx^{-\alpha}$$

Take logarithm of both sides:

$$\log p(x) = \log C - \alpha \log x$$

Now you can use standard R functions like `lm()` to calculate α coefficient via linear regression. However you might find that it is a bad idea.

Alternatively, you can compute cumulative density function (CDF)

$$f(x) = Pr(x < X)$$

of power law distribution. Good things about CDF of the power law are:

- It still has form of power law
- On log-log plot it looks more like a line

1. Derive the formula for CDF function of power law $F(x)$. It means you should calculate cumulative distribution function by integration of PDF function.

$$p(x) = Cx^{-\alpha}, \quad x > x_{min}$$

$$F(x) = \int_{x_{min}}^x p(x)dx = \int_{x_{min}}^x Ct^{-\alpha}dt$$

$$F(x) = C \frac{t^{1-\alpha}}{1-\alpha} \Big|_{x_{min}}^x = C \left(-\frac{x^{1-\alpha}}{\alpha-1} + \frac{x_{min}^{1-\alpha}}{\alpha-1} \right)$$

The antiderivative exists only if $\alpha > 1$. We know that $F(+\infty) = Pr(X < +\infty) = 1$, then

$$F(+\infty) = \int_{x_{min}}^{+\infty} p(x)dx = \int_{x_{min}}^{+\infty} Ct^{-\alpha}dt$$

$$F(x) = C \frac{t^{1-\alpha}}{1-\alpha} \Big|_{x_{min}}^{+\infty} = C \left(0 + \frac{x_{min}^{1-\alpha}}{\alpha-1} \right)$$

$$C = \frac{\alpha-1}{x_{min}^{1-\alpha}}$$

2. Download Internet Network and plot PDF and CDF of the degree distribution in log-log scale

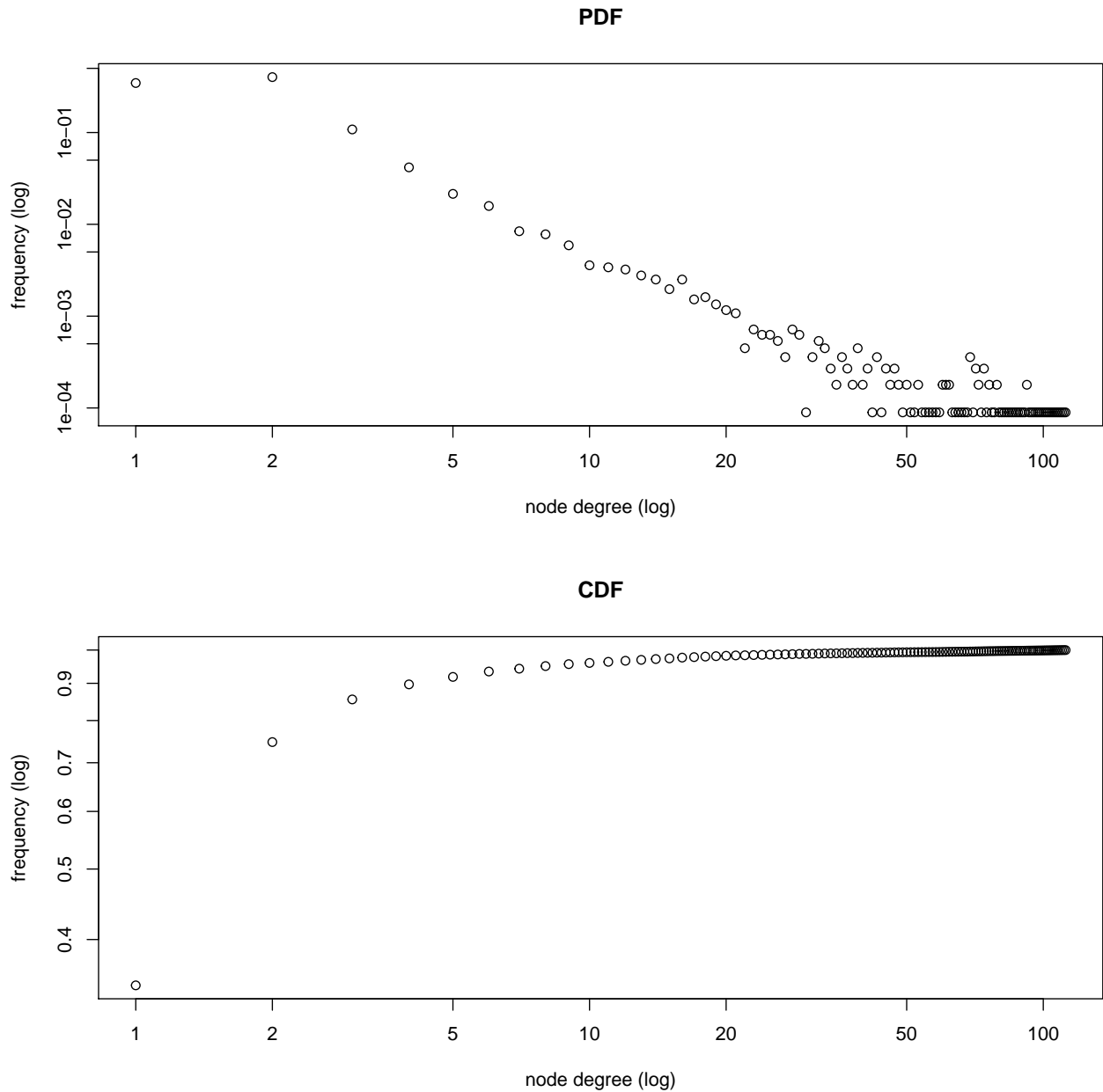
```
data <- read.table("Internet_AS.dat")
data <- t(as.matrix(data))
graph <- graph(data)
## get PDF and CDF
xmax=max(degree(graph))
degree_range <- as.numeric(1:xmax)
degree_dist <- degree.distribution(graph, cumulative = FALSE)[degree_range]

## delete zero values
nonzero.position <- which(degree_dist > 0)
degree_dist <- degree_dist[nonzero.position]
degree_range <- degree_range[nonzero.position]
cum_degree_range <-degree_range+1
compl_cum_degree_range <- cum_degree_range

cum_degree_dist <- cumsum(degree_dist)
compl_cum_degree_dist <- degree.distribution(graph, cumulative = TRUE)[cum_degree_range]
compl_cum_degree_dist <- compl_cum_degree_dist[cum_degree_range]

## plot two graphs side-by-side
par(mfrow=c(2,1))
plot(degree_dist, log = "xy", main = "PDF", xlab = "node degree (log)", ylab = "frequency (log)")

plot(cum_degree_dist, log = "xy", main = "CDF", xlab = "node degree (log)", ylab = "frequency (log)")
```



3. Fit linear regression model to PDF and CDF to estimate α . Plot fitted models along with data

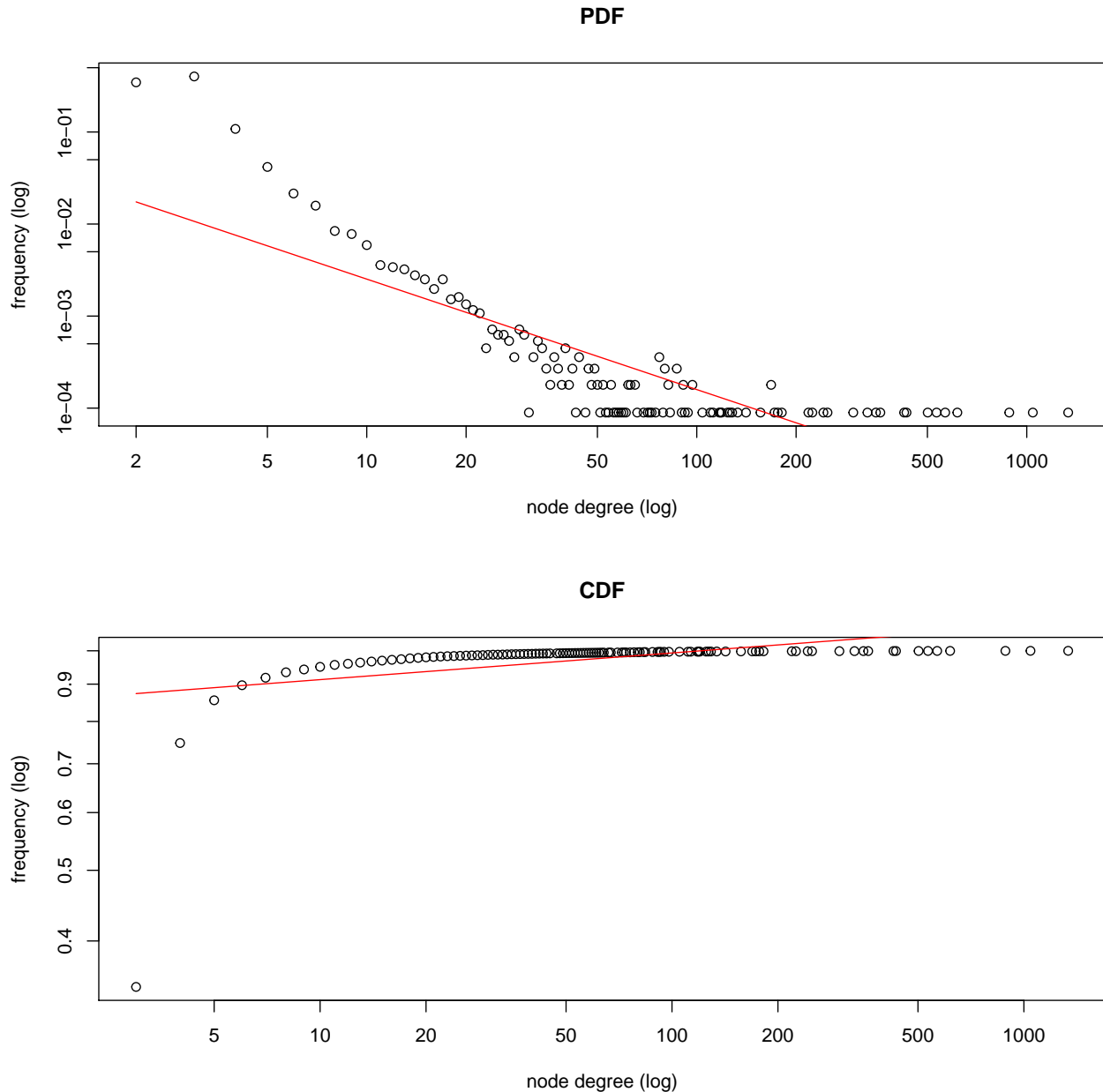
```
par(mfrow=c(2,1))
# fit linear regression
reg <- lm(log(degree_dist) ~ log(degree_range))
coefficients <- coef(reg)
func_powerlaw <- function(x) exp(coefficients[[1]] + coefficients[[2]] * log(x))
alpha <- -coefficients[[2]]
# plot PDF and its estimation
plot(degree_dist ~ degree_range, log = "xy", main = "PDF", xlab = "node degree (log)", ylab = "frequency",
     curve(func_powerlaw, col = "red", add = TRUE, n = length(degree_range))
# fit linear regression
```

```

reg <- lm(log(cum_degree_dist) ~ log(cum_degree_range))
coefficients <- coef(reg)
func_powerlaw <- function(x) exp(coefficients[[1]] + coefficients[[2]] * log(x))
alpha <- -coefficients[[2]]

plot(cum_degree_dist ~ cum_degree_range, log = "xy", main = "CDF", xlab = "node degree (log)", ylab = "frequency (log)",
     curve(func_powerlaw, col = "red", add = TRUE, n = length(cum_degree_range)))

```



Problem 2

Kolmogorov-Smirnov test describes how similar are two distributions. In our case, when we have fitted model and original data, we can calculate their CDFs and Kolmogorov-Smirnov test shows us how well model

approximates original data. In other words, it shows us the goodness-of-fit of our model.

$$D = \max_x \|f(x|\alpha, x_{min}) - f_{emp}(x)\|,$$

where $f(x|\alpha, x_{min})$ and $f_{emp}(x)$ are theoretical and empirical CDFs respectively. To estimate x_{min} of the fitted power-law model we can use KS test:

- Pick some x_{min} value
- Fit power-law distribution to data (that is estimation of α) – now we have $f(x|\alpha, x_{min})$
- Perform KS test – compute D statistic
- Finally, choose x_{min}^* that provides minimal value of D statistic among all KS tests run above.

In R all this stuff can be done in one line of code.

Again, use Internet Network

Properly load it into R and do following tasks:

1. Using `power.law.fit` find `xmin` value and corresponding `alpha`

```
d <- degree(graph)
fit <- power.law.fit(d, NULL, implementation = "plfit")
alpha <- fit$alpha
xmin <- fit$xmin
C <- (alpha-1)*xmin^(alpha-1)
alpha
```

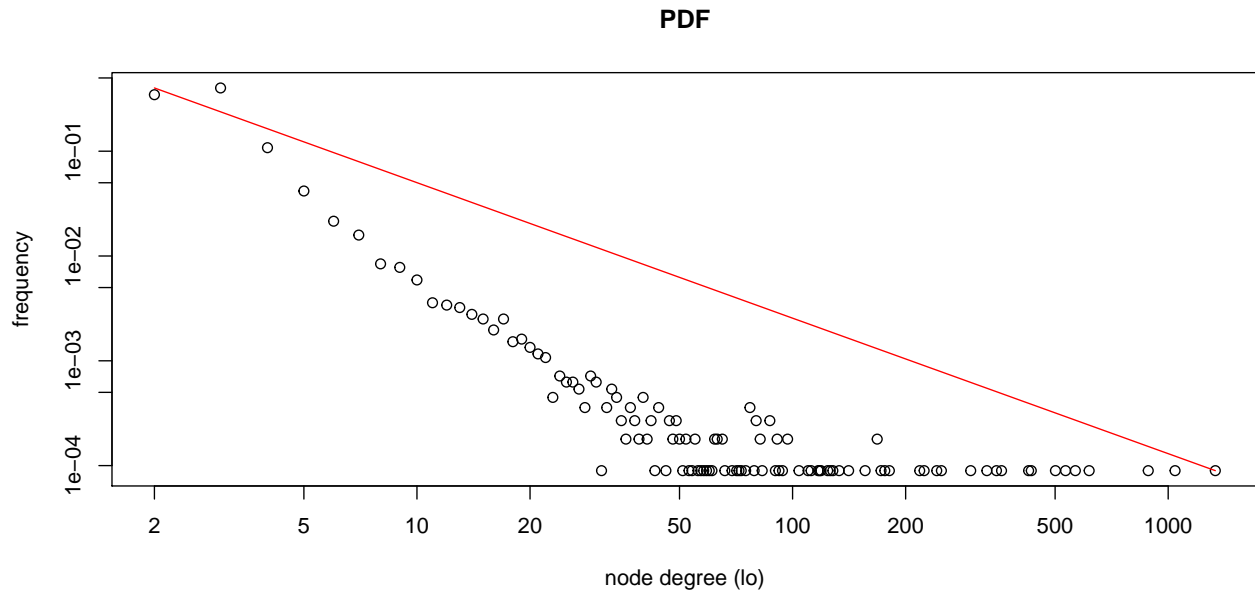
```
## [1] 2.08306
```

```
xmin
```

```
## [1] 7
```

2. Put fitted model along with empirical PDF (CDF)

```
par(mfrow=c(2,1))
fit_pdf <- function(x) return(C*x^(-alpha))
plot(degree_range,degree_dist, log = "xy", main = "PDF", xlab = "node degree (lo)", ylab = "frequency")
#lines(x, x^(-alpha), col="green")
par(new=TRUE)
curve(fit_pdf,from=xmin,to=xmax, log="xy", col = "red", add = FALSE, n = length(degree_range),main = ""
#x=xmin:xmax
#plot(log(x), log(C * x ^ (-alpha)), col="red", xlab="", ylab="", yaxt="n")
#lines(log(x), log(C * x ^ (-alpha)), col="brown")
```



Problem 3.

For Wikipedia vote network (clear up comments in the beginning of the file) derive the following characteristics:

1. The number of vertices and edges
2. The number of loops (edges that start and end at the same vertex)
3. The number of symmetrical edges
4. Degree distribution (without considering the direction)
5. The number of nodes with a degree greater than 1 and with a degree greater than 15
6. Find strongly connected components and their sizes.
7. Take subgraph of the original graph, which consists of the first 80 vertices and set color into red for those nodes in which the number of incoming edges is greater than the number of outgoing edges. Otherwise, set color in blue. For nodes with the same number of incoming and outgoing edges set color into green. Besides that, increase the size of vertices with a maximum value of transitivity (for example, you may set size into 10 for these nodes and 1 for others).
8. Take subgraph from the previous task and find maximal connected component. For this component highlight any way that corresponds to the diameter of the subgraph. How many such paths are in this graph?
9. Make average neighbor degree vs node degree scatter plot (one point on the plot per node) and aggregated plot, averaging over all nodes with the same degree (aggregated average vs degree, one value per degree). Explain your observations.
10. Make local clustering coefficient vs node degree scatter plot (one point on the plot per node) and aggregated, averaging over all nodes with the same degree (aggregated average vs degree, one value per degree). Explain your observations.

1. The number of vertices and edges.

```
initial <- read.table("Wiki-vote.txt")
initial <- t(as.matrix(initial))
graph <- graph(initial)
vcount(graph)
```

```
## [1] 8297
```

```
ecount(graph)
```

```
## [1] 103689
```

2. The number of loops (edges that start and end at the same vertex)

```
sum(is.loop(graph), na.rm=TRUE)
```

```
## [1] 0
```

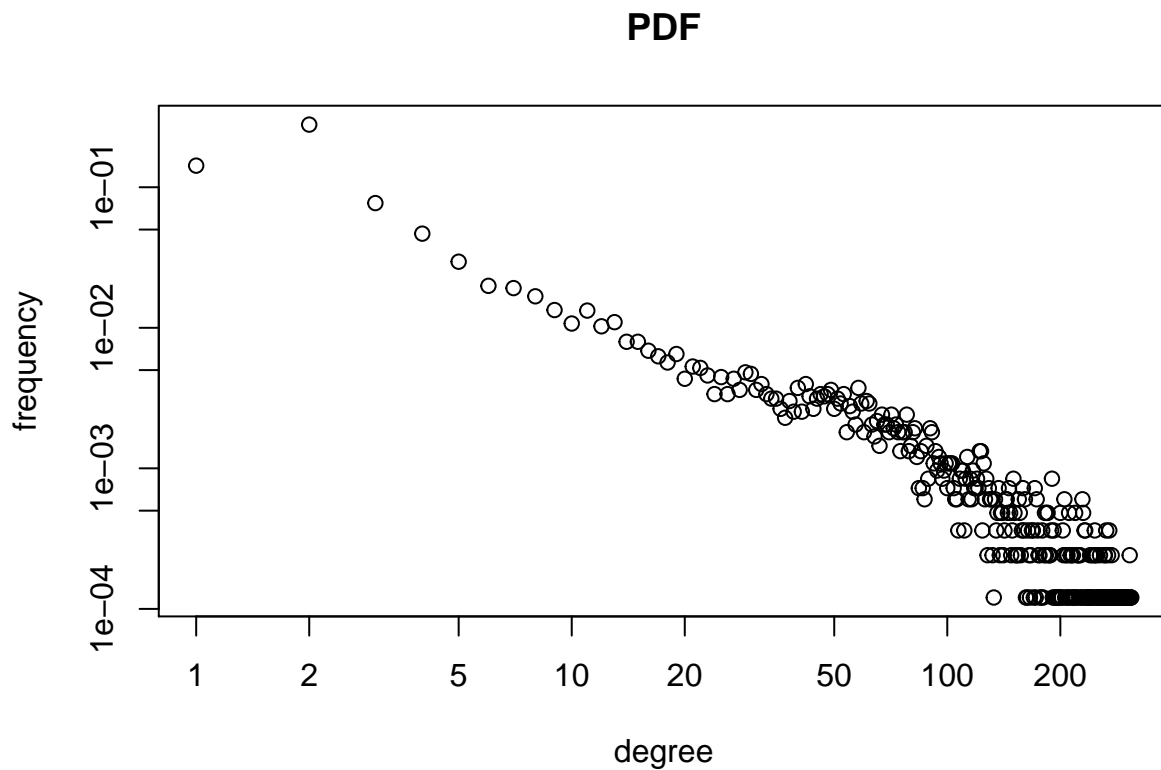
3. The number of symmetrical edges

```
reciprocity(graph)*ecount(graph)/2
```

```
## [1] 2927
```

4. Degree distribution

```
degree_dist <- degree.distribution(graph)  
plot(degree.distribution(graph, mode='all')[which(degree_dist > 0)], log = 'xy', main = 'PDF', xlab =
```



5. The number of nodes with a degree greater than 1 and with a degree greater than 15

```
sum(degree(graph) > 1, na.rm=TRUE)
```

```
## [1] 4800
```

6. Find strongly connected components and their sizes.

```
components <- clusters(graph, mode = "strong")
#print(components$size)
print(components$no)
```

```
## [1] 6998
```

There is some variation in the results of task number 7 and 8. You can compute the strongly and weakly connected components. Depending on the type of connectivity we obtain different results.

7. Take subgraph of the original graph, which consists of the first 80 vertices and set color into red for those nodes in which the number of incoming edges is greater than the number of outgoing edges. Otherwise, set color in blue. For nodes with the same number of incoming and outgoing edges set color into green. Besides that, increase the size of vertices with a maximum value of transitivity (for example, you may set size to 10 for these nodes and 1 for others).

```
subgraph <- induced.subgraph(graph, V(graph)[1:80])

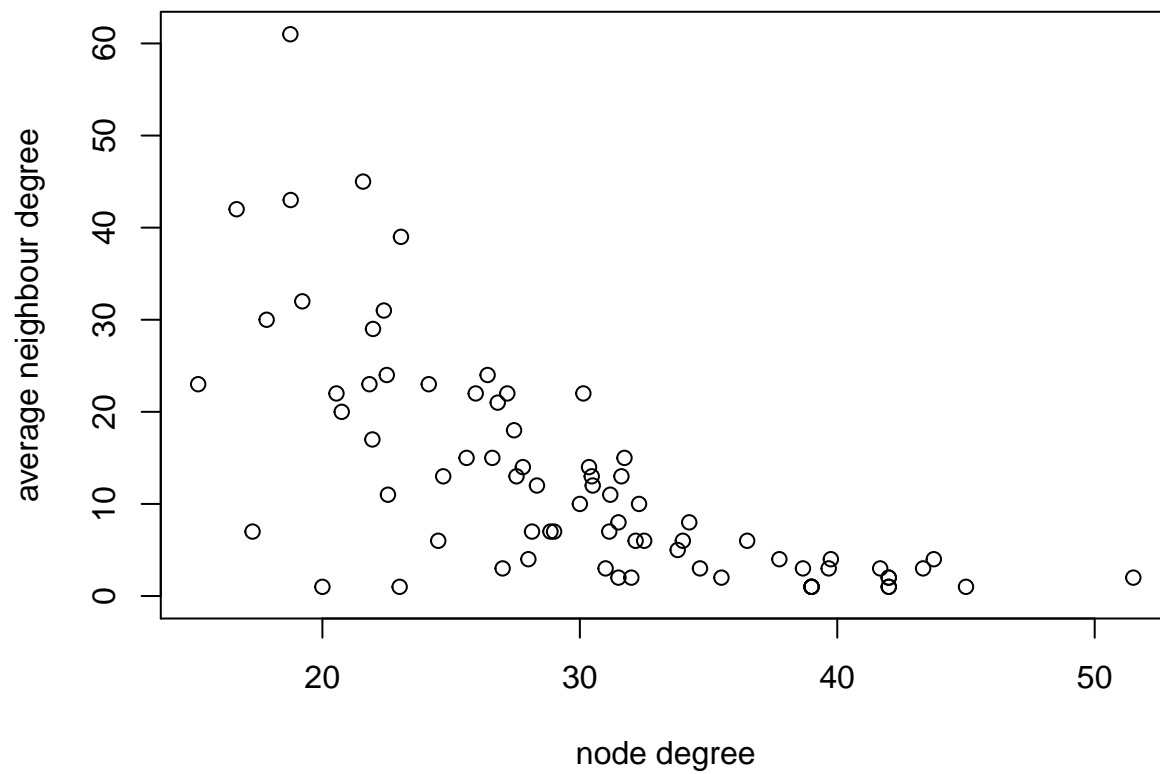
V(subgraph)$ins <- degree(subgraph, v = V(subgraph), mode = "in")
V(subgraph)$outs <- degree(subgraph, v = V(subgraph), mode = "out")
V(subgraph)$color <- "blue"
V(subgraph)[ins > outs]$color <- "red"
V(subgraph)[ins == outs]$color <- "green"
V(subgraph)$size <- 1
V(subgraph)[which(transitivity(subgraph, type="local") == 1)]$size <- 10
```

8. Take subgraph from the previous task and find maximal connected component. For this component highlight any way that corresponds to the diameter of the subgraph. How many such paths are in this graph?

```
m <- shortest.paths(subgraph) m <- max(m[m != Inf])
clust = V(subgraph)[which(cccmembership == which(ccccsize == max(ccc$size)))]
dis <- get_diameter(subgraph)
V(subgraph)[intersect(clust,dis)]$color <- "yellow"
```

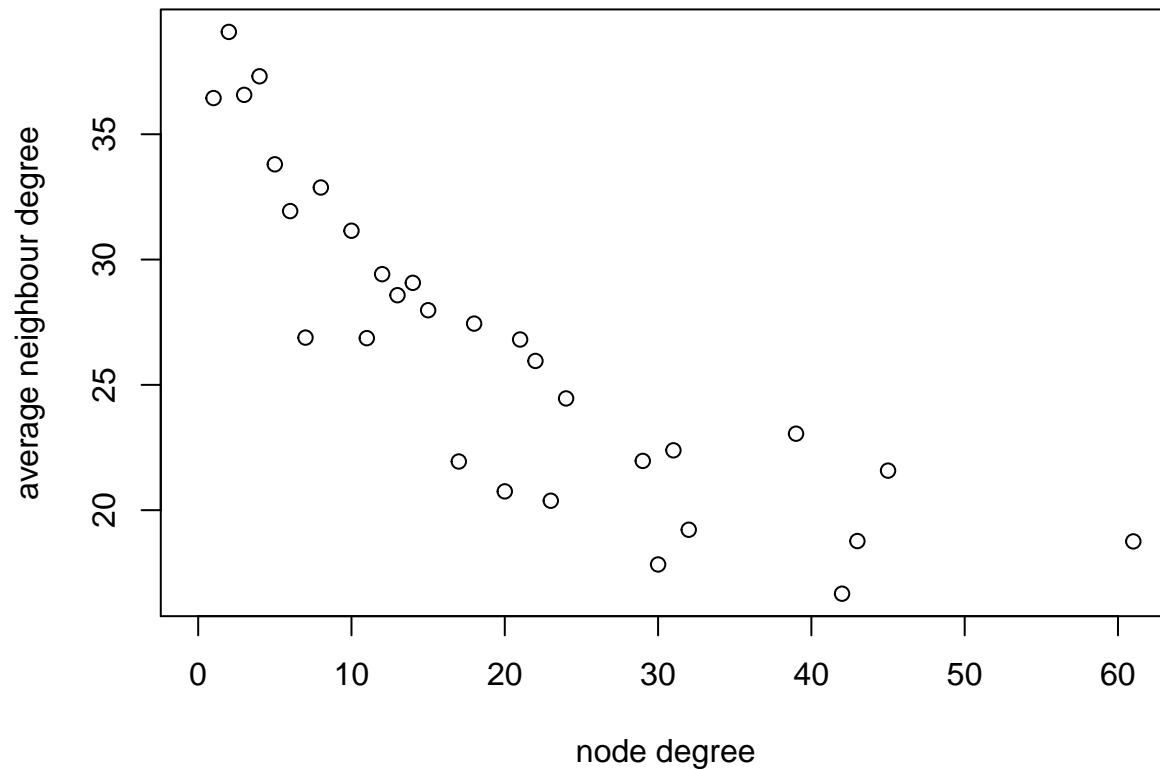
di <- diameter(subgraph, directed = TRUE) p <- path.length.hist(subgraph)\$res[di] p ##### 9. Make average neighbour degree vs node degree scatter plot (one point on the plot per node) and aggregated plot, averaging over all nodes with the same degree (aggregated average vs degree, one value per degree). Explain your observations.

```
# get average neighbour degree
avnd <- graph.knn(subgraph)$knn
# get node degree
nd <- degree(subgraph)
# scatter plot
plot(avnd, nd, xlab = "node degree", ylab = "average neighbour degree")
```

```
agg_avnd <- c()
index = 1
for (d in unique(nd)) {
  agg_avnd[index] <- sum(avnd[which(nd==d)])/length(which(nd==d))
  index = index + 1
}
plot(unique(nd), agg_avnd, main = "Aggregated plot", xlab = "node degree", ylab = "average neighbour degree")
```

Aggregated plot



10. Make local clustering coeff vs node degree scatter plot (one point on the plot per node) and aggregated, averaging over all nodes with the same degree (aggregated average vs degree, one value per degree). Explain your observations.

Hint: The code is similar to the previous one. Local clustering coefficients for each node could be found with the use of “transitivity” function.