

Practical Work

Lifetime Data Analysis

Rodrigo Arriaza, Alexander J Ohrt

14 desember, 2021

Introduction

We are given a data set on sexually transmitted diseases (STDs). This is data from a study about gonorrhea and chlamydia in 877 women. The objective with this practical work is to study possible risk factors for a reinfection with gonorrhea or chlamydia in women who have suffered one of both infections previously. The variables of interest are sociodemographic variables or those related to sexual practice. We have a lot of variables at our disposal, but have chosen to use the following, some for statistical reasons and some for medical reasons:

- Age: The age of the woman.
- NumPartners: The number of partners during the last 30 days.
- CondomUse: Use of condoms (1: always, 2: once in a while, 3: never)
- YearsSchool: Years of schooling.
- InitInfect: Initial infection (1: Gonorrhea, 2: Chlamydia, 3: both)
- InvVagAtExam: Involvement vagina at exam (1: yes; 0: no).
- DischargeExam: Discharge at exam (1: yes; 0: no)

The first three were chosen based on results from a [study](#) on gonorrhea reinfection in heterosexual STD clinic attendees. The study concluded that increased reinfection risk (of gonorrhea) was associated with younger age and a greater number of recent sex partners, among other risk factors. Moreover, the authors concluded that any type of condom use was a risk factor for reinfection with gonorrhea in women.

Another [publication](#) reports that, on average, 14% of women with clamydia and 12% of women with gonorrhea get reinfected, with younger women at higher risk. Moreover, they state that many adolescents treated for infection of one of the two STDs are reinfected within three to six months, usually because of resumed sexual contact with an untreated partner. Thus, the marital status might be interesting to analyse. However, this is not added, because, the ages in the data set are low, which most likely means that the amount in each level of **MaritalStatus** is very skewed towards “single”. This can be seen in the descriptive analysis below.

This [meta-analysis](#) reports that the relationship between race, socioeconomic status (SES) and chlamydial infection is not clear. It concludes that SES was not associated with chlamydia infection, where they tested for several variables, where level of parent’s education was one of them. Either way, we think it might be interesting to see if the years of schooling of the women (**YearsSchool**) have any impact on chlamydia reinfection and as is shown below it showed to be statistically significant during the exploratory analysis.

Moreover, we chose to use the initial infection (**InitInfect**) as an explanatory variable, because several of the studies above are only done on one of the two diseases, not on both at the same time. Because of this we wanted to investigate if the initial infection type is a risk factor and, if this is the case, if the risk differs based on which infection was suffered initially.

Naturally, the categorical variable which states if the woman is reinfected or not (**Reinfection**) will be used as a dependent variable in the analysis and the time until reinfection since the more time a subject is under study, the greater the risk of the event reoccurring.

Table 1: Statistical Significance of the Variables

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-4.4766626	0.6791361	-6.5917022	0.0000000
EthnicityW	-0.0786114	0.1576156	-0.4987540	0.6179527
MaritalStatusM	0.1142920	0.4681139	0.2441542	0.8071114
MaritalStatusS	0.5011754	0.3203698	1.5643657	0.1177317
Age	0.0188481	0.0156143	1.2071011	0.2273932
YearsSchool	-0.1689015	0.0442657	-3.8156308	0.0001358
InitInfect2	-0.3302518	0.1740868	-1.8970524	0.0578210
InitInfect3	-0.3318821	0.1755787	-1.8902183	0.0587288
NumPartners	0.1164568	0.0598373	1.9462257	0.0516276
OralSex12m1	-0.3703474	0.2387666	-1.5510855	0.1208812
OralSex30d1	-0.3246975	0.2643311	-1.2283739	0.2193066
RectalSex12m1	0.0669703	0.4881503	0.1371920	0.8908790
RectalSex30d1	-0.1627456	0.6172379	-0.2636675	0.7920361
AbPain1	0.2969178	0.1771403	1.6761734	0.0937042
SignDischarge1	0.1330009	0.1306664	1.0178660	0.3087416
SignDysuria1	0.1954606	0.1812469	1.0784219	0.2808455
CondomUse2	-0.1553543	0.2725108	-0.5700849	0.5686201
CondomUse3	-0.4582270	0.2819913	-1.6249684	0.1041693
SignItch1	-0.2209724	0.1750560	-1.2622958	0.2068424
SignLesion1	-0.2541307	0.3787052	-0.6710513	0.5021878
SignRash1	-0.0638066	0.4592994	-0.1389215	0.8895122
SignLymph1	0.2368538	0.5922357	0.3999317	0.6892069
InvVagAtExam1	0.5726933	0.2003764	2.8580874	0.0042620
DischargeExam1	-0.5805191	0.2691414	-2.1569301	0.0310111
AbnormNodeExam1	0.0801562	0.5157541	0.1554155	0.8764938

Statistical Variable Selection

As noted, in addition to medical criteria for selecting variables, we have used the following statistical model to select variables based on statistical criteria. Shown below. RODRI: EXPLAIN!

REMOVE THIS AFTER: I made a list of things that can be used in the explanation: * Despite the fact that the age of the woman is found to not be statistically significant in the method above, we have added it because of the mentioned studies (this is thus added based on medical criteria)

```
nb.model <- MASS::glm.nb(Reinfection ~ Ethnicity + MaritalStatus + Age + YearsSchool
  + InitInfect + NumPartners + OralSex12m + OralSex30d + RectalSex12m
  + RectalSex30d + AbPain + SignDischarge + SignDysuria + CondomUse
  + SignItch + SignLesion + SignRash + SignLymph + InvVagAtExam
  + DischargeExam + AbnormNodeExam + offset(log(TimeUntilReinf)),
  data=std_data)
s <- summary(nb.model)
k <- knitr::kable(s$coefficients, caption = 'Statistical Significance of the Variables')
kableExtra::row_spec(k, c(6,23,24), color='white', background = 'blue')
```

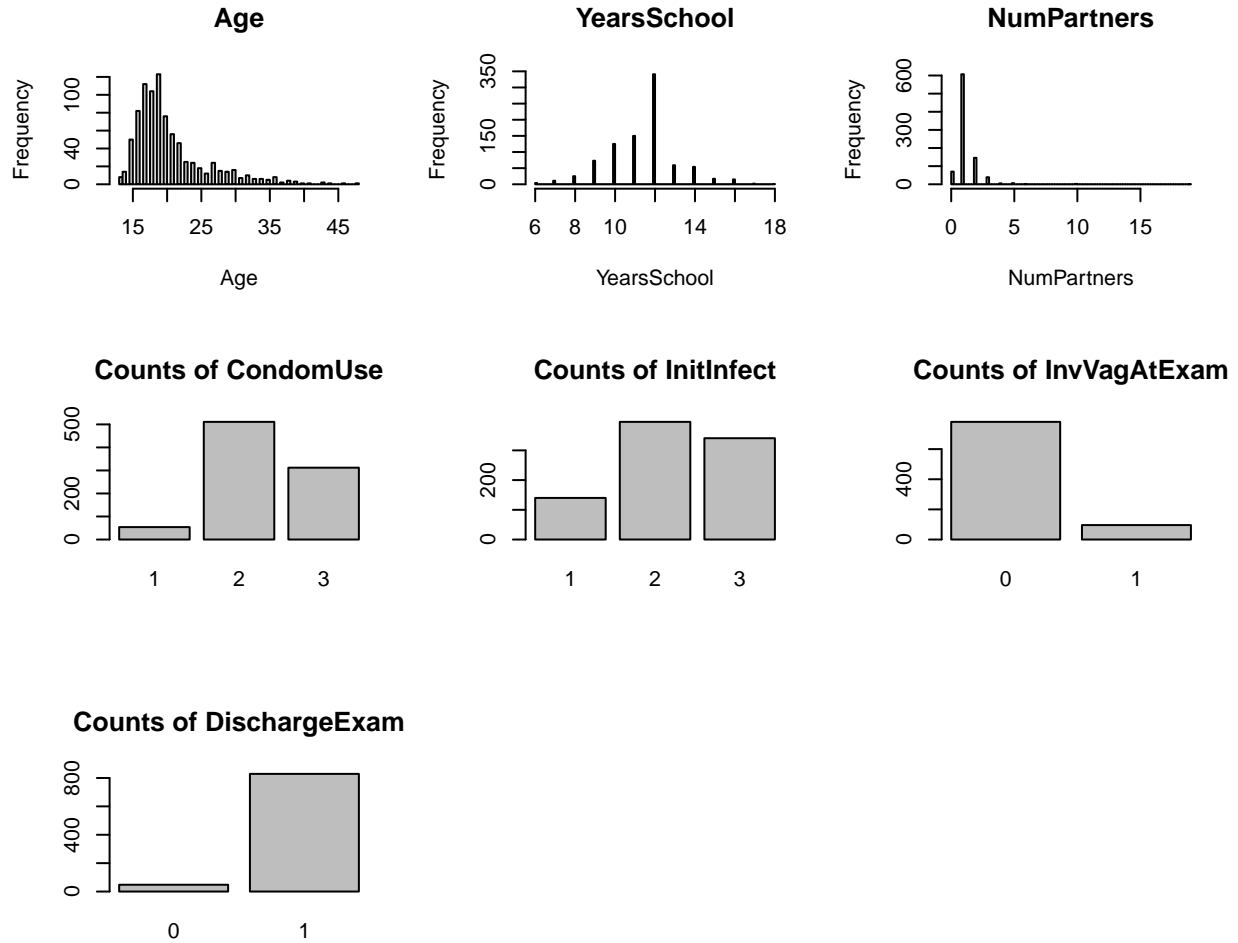
Finally, the vaginal involvement at exam (InvVagAtExam) and the discharge at exam (DischargeExam) are selected as variables in our analysis, since they are shown as statistically significant in the variable selection above.

Table 2: Corr. Between Continuous Variables

	Age	YearsSchool	NumPartners
Age	1.0000000	0.4316163	0.1348591
YearsSchool	0.4316163	1.0000000	0.0155090
NumPartners	0.1348591	0.0155090	1.0000000

Descriptive Analysis

In total, the data set contains 24 variables, but, as noted, we have selected only 7 of them in our analysis. Recall that the data set has 877 women. The percentage of right-censored data in the data set is 60.4, which is a relatively large part of the data set. The women were followed for 1529 days, then the study was stopped.

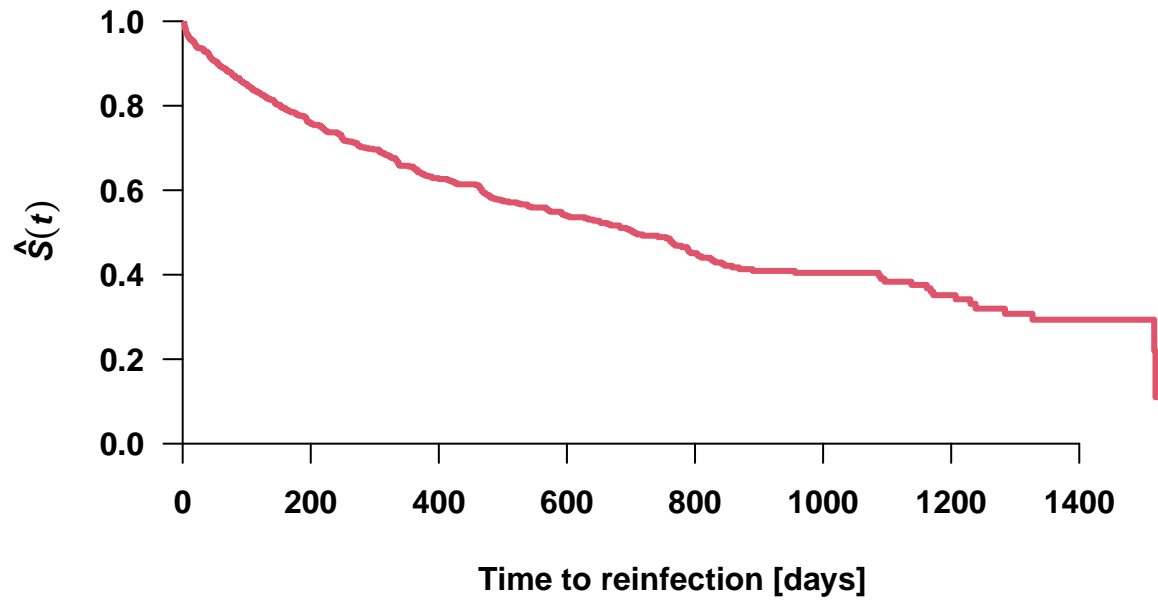


The three continuous variables we have chosen to use in the analysis are **Age**, **YearsSchool** and **NumPartners**. The correlations between the variables are shown in table 2. Note that the correlation between **Age** and **YearsSchool** is 0.43, which means that they are somewhat correlated. This could be interesting to have in mind in the following.

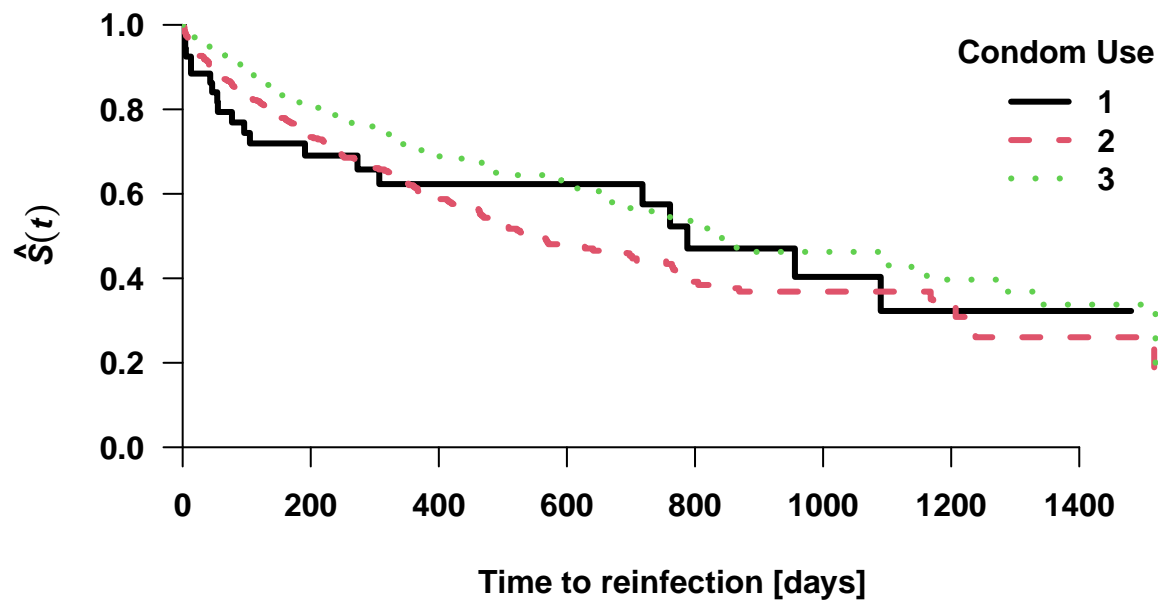
Nonparametric Analysis

The survival curve is estimated by means of Kaplan-Meier and plotted below. The curve below shows the general survival in the data set.

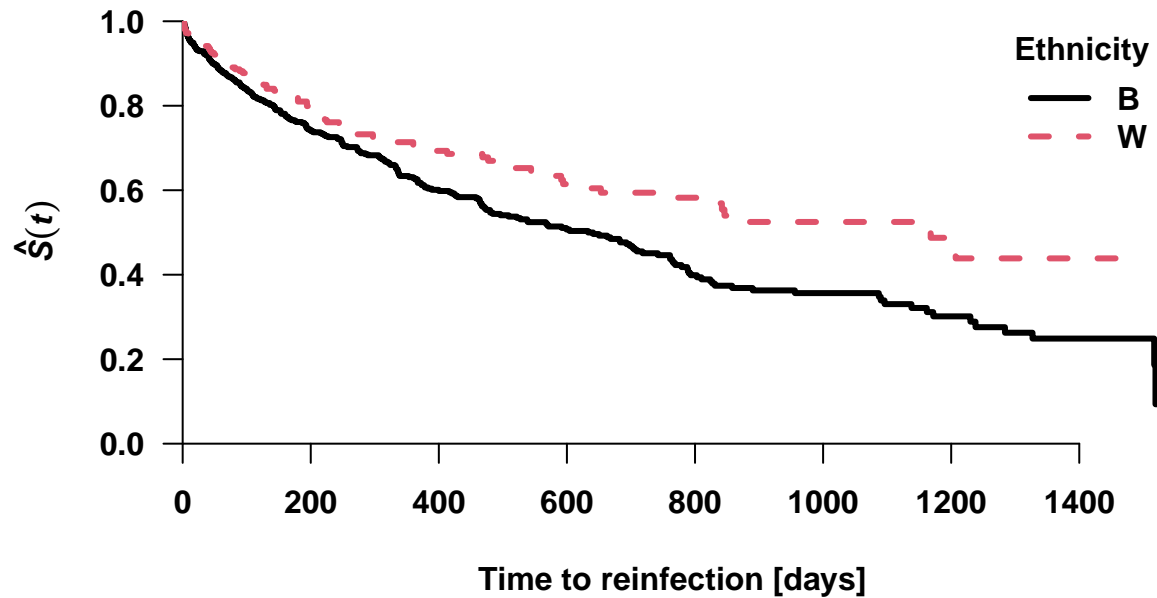
Survival Function



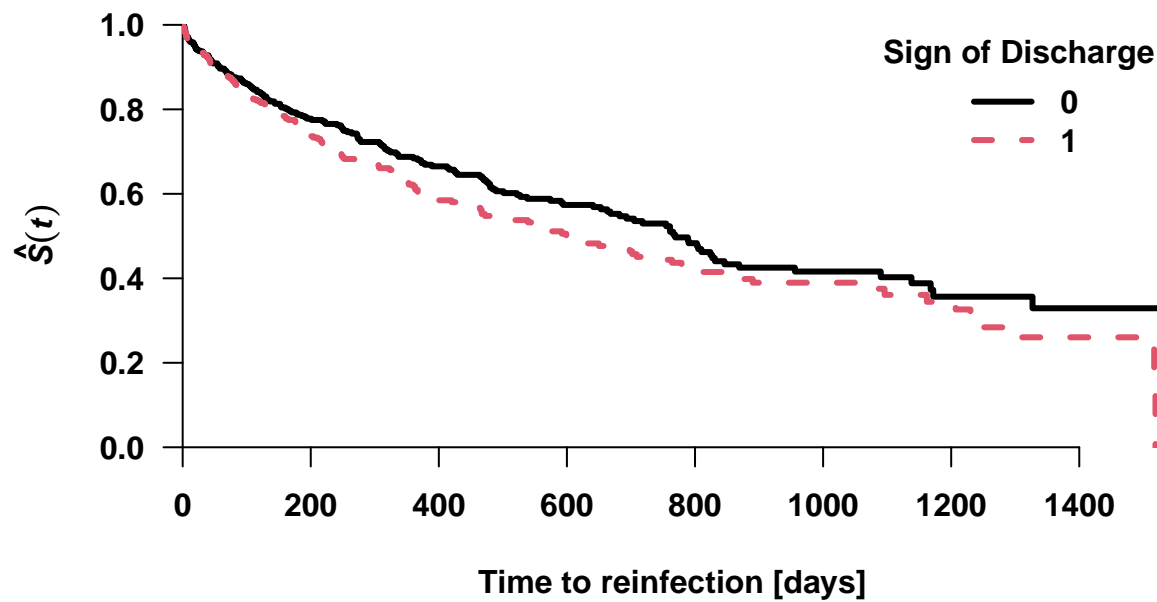
Survival Function as Function of Condom Use



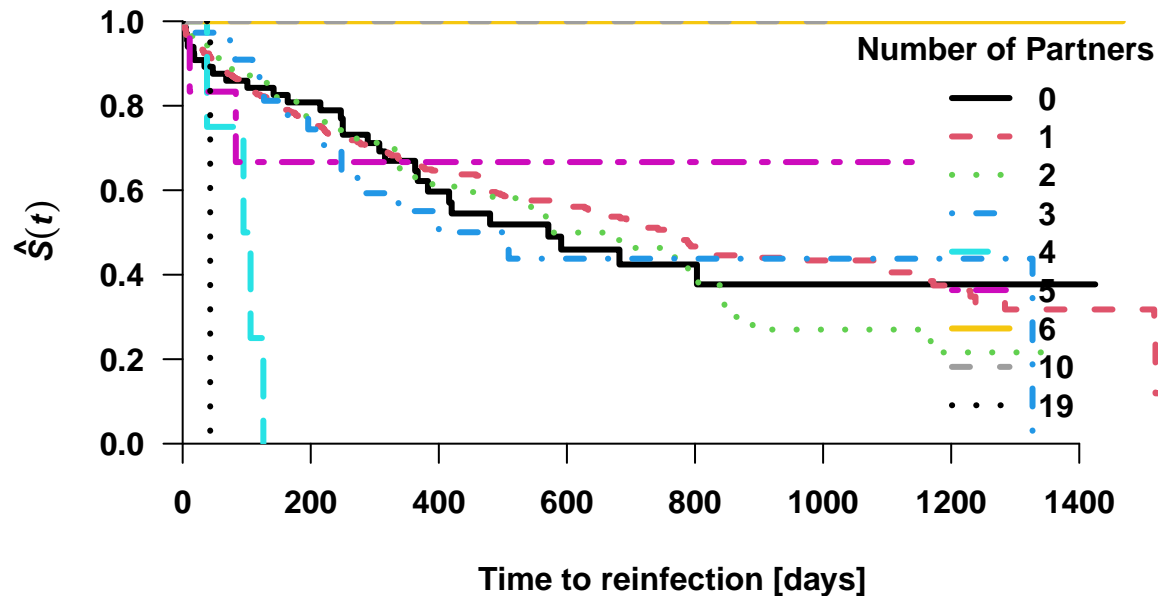
Survival Function as Function of Ethnicity



Survival Function as Function of Sign of Discharge



Survival Function as Function of Number of Partners



The median survival time is 705 days.

Comparison of survival functions by means of nonparametric tests, such as the logrank test.

Below, logrank test for all curves plotted above are done. Other types of tests can also be used (Fleming-Harrington) and test can be done on different variables in the data set.

```
#>
#> K-sample test for right-censored data
#>
#> Parameters: rho=0, lambda=0
#> Distribution: counting process approach
#>
#> Data: Surv(TimeUntilReinf, Reinfection) by CondomUse
#>
#>      N Observed Expected    O-E (O-E)^2/E (O-E)^2/V
#> CondomUse=1  54      21   19.3   1.72   0.154   0.164
#> CondomUse=2 511     210  185.0  24.97   3.369   7.338
#> CondomUse=3 312     116  142.7 -26.69   4.993   8.606
#>
#> Chisq= 8.6 on 2 degrees of freedom, p-value= 0.0133
#> Alternative hypothesis: survival functions not equal

#>
#> Two-sample test for right-censored data
#>
#> Parameters: rho=0, lambda=0
#> Distribution: counting process approach
#>
#> Data: Surv(TimeUntilReinf, Reinfection) by Ethnicity
#>
#>      N Observed Expected    O-E (O-E)^2/E (O-E)^2/V
#> Ethnicity=B 585     264   238   25.9   2.82   9.03
#> Ethnicity=W 292      83   109 -25.9   6.16   9.03
```

```

#>
#> Statistic Z= -3, p-value= 0.00266
#> Alternative hypothesis: survival functions not equal

#>
#> Two-sample test for right-censored data
#>
#> Parameters: rho=0, lambda=0
#> Distribution: counting process approach
#>
#> Data: Surv(TimeUntilReinf, Reinfection) by SignDischarge
#>
#>
#>      N Observed Expected    O-E (O-E)^2/E (O-E)^2/V
#> SignDischarge=0 472      171      187 -15.6      1.31      2.84
#> SignDischarge=1 405      176      160  15.6      1.53      2.84
#>
#> Statistic Z= 1.7, p-value= 0.0918
#> Alternative hypothesis: survival functions not equal

#>
#> Trend FH test for right-censored data
#>
#> Parameters: rho=0, lambda=0
#> Distribution: counting process approach
#>
#> Data: Surv(TimeUntilReinf, Reinfection) by NumPartners
#>
#>
#>      N Observed Expected    O-E
#> NumPartners=0   70      29 27.5754  1.425
#> NumPartners=1  607     234 244.7937 -10.794
#> NumPartners=2  146      61 55.5667  5.433
#> NumPartners=3   39      16 14.0435  1.957
#> NumPartners=4    5       4  0.6262  3.374
#> NumPartners=5    6       2  1.8913  0.109
#> NumPartners=6    1       1  1.2190 -0.219
#> NumPartners=10   2      29  1.1948 27.805
#> NumPartners=19   1     234  0.0894 233.911
#>
#> Statistic Z= 269, p-value= 0
#> Alternative hypothesis: survival functions not equal

```

Fit of a parametric survival model

After trying to fit Weibull, log-logistic and lognormal log-linear models, we concluded that the Weibull model is best suited to our data.

The correct variables should be inserted here after we decide on them!!

```

#>
#> Call:
#> survreg(formula = s2 ~ Age + NumPartners + CondomUse + YearsSchool +
#>   InitInfect + InvVagAtExam + DischargeExam, data = std_data,
#>   dist = "weibull")
#>
#>      Value Std. Error      z      p
#> (Intercept)   3.9516    0.6338  6.23 4.5e-10

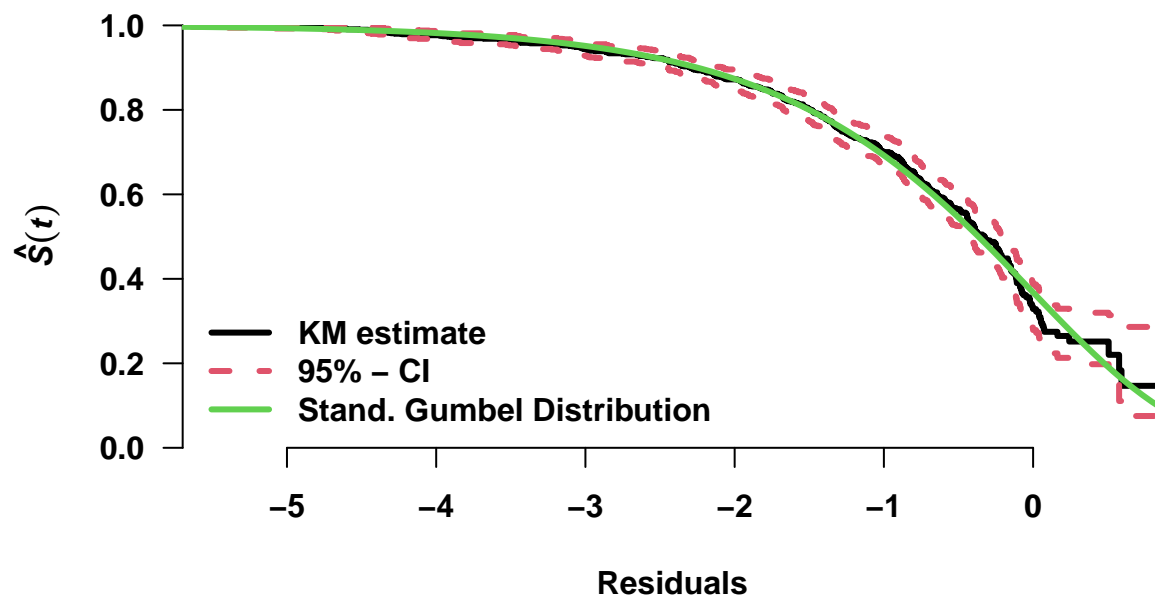
```

```

#> Age          0.0101      0.0161  0.63 0.52912
#> NumPartners  -0.0139     0.0681 -0.20 0.83835
#> CondomUse2    0.0788     0.2995  0.26 0.79244
#> CondomUse3    0.4228     0.3106  1.36 0.17350
#> YearsSchool   0.1704     0.0487  3.50 0.00047
#> InitInfect2   0.5110     0.1907  2.68 0.00738
#> InitInfect3   0.3149     0.1915  1.64 0.10011
#> InvVagAtExam1 -0.5092     0.2212 -2.30 0.02133
#> DischargeExam1 0.4601     0.2894  1.59 0.11184
#> Log(scale)    0.2606     0.0445  5.85 4.9e-09
#>
#> Scale= 1.3
#>
#> Weibull distribution
#> Loglik(model)= -2674.9   Loglik(intercept only)= -2697.1
#>  Chisq= 44.43 on 9 degrees of freedom, p= 1.2e-06
#> Number of Newton-Raphson Iterations: 7
#> n= 877

```

Residuals of the Weibull Regression Model



The residuals seem to fit relatively nicely to the Gumbel distribution, which indicates that the Weibull is reasonable to use **Check that the terminology I used here is correct!**

```

#>
#> Call:
#> survreg(formula = s2 ~ Age + NumPartners + CondomUse + YearsSchool +
#>   InitInfect + InvVagAtExam + DischargeExam, data = std_data,
#>   dist = "lognormal")
#>
#>              Value Std. Error      z      p
#> (Intercept)   2.881331   0.776808  3.71 0.00021
#> Age           -0.000895   0.018939 -0.05 0.96232
#> NumPartners   -0.019853   0.077571 -0.26 0.79800
#> CondomUse2     0.411402   0.360931  1.14 0.25435
#> CondomUse3     0.811166   0.375119  2.16 0.03059
#> YearsSchool    0.204027   0.057861  3.53 0.00042

```

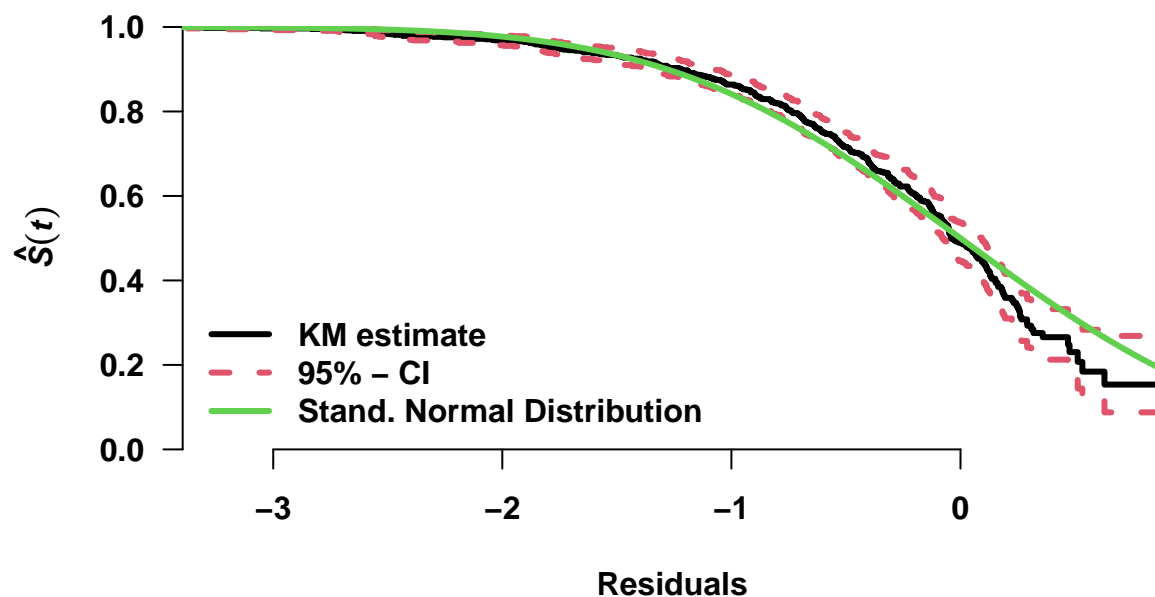


```

#> InitInfect2      0.373146    0.242880    1.54 0.12446
#> InitInfect3      0.318740    0.249560    1.28 0.20153
#> InvVagAtExam1    -0.777588    0.270837   -2.87 0.00409
#> DischargeExam1   0.725573    0.365444    1.99 0.04709
#> Log(scale)       0.724527    0.039725   18.24 < 2e-16
#>
#> Scale= 2.06
#>
#> Log Normal distribution
#> Loglik(model)= -2688.6    Loglik(intercept only)= -2709.4
#>  Chisq= 41.51 on 9 degrees of freedom, p= 4e-06
#> Number of Newton-Raphson Iterations: 4
#> n= 877

```

Residuals of the Lognormal Regression Model



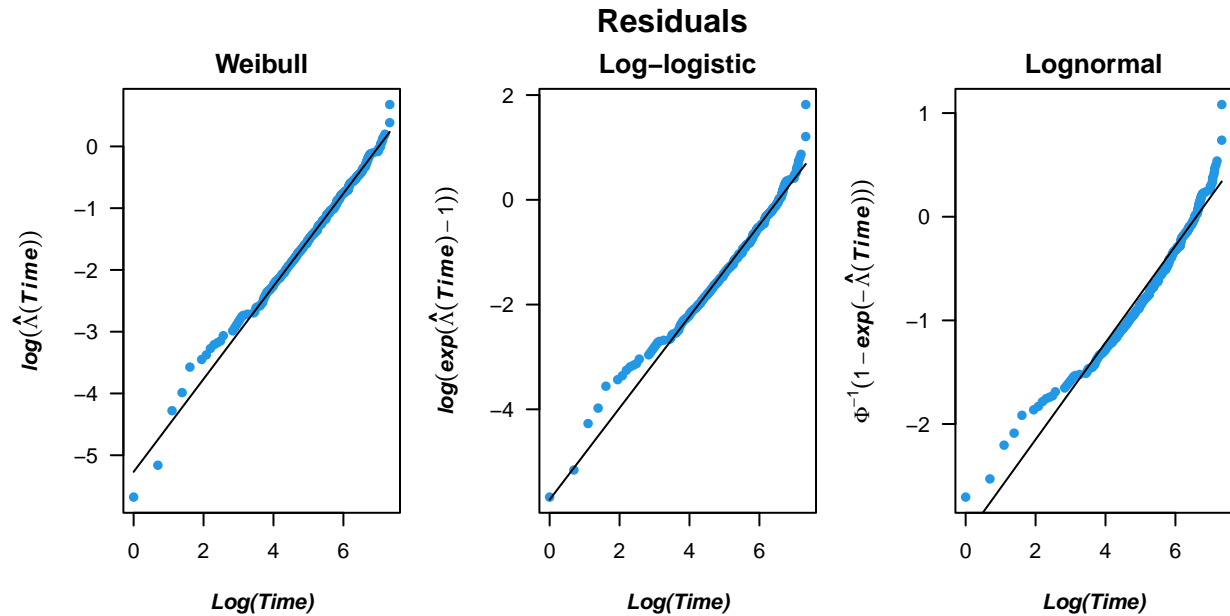
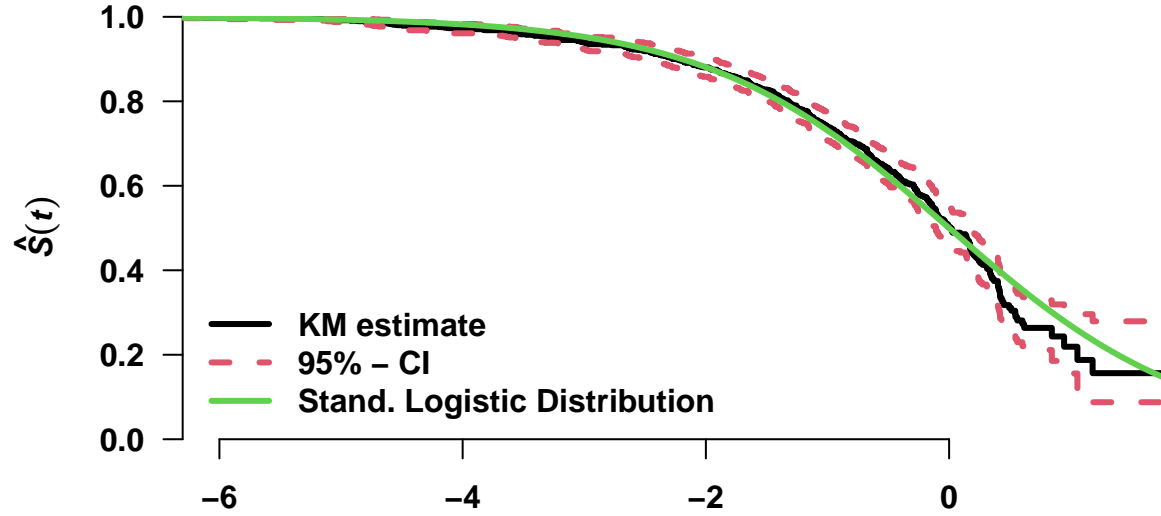
```

#>
#> Call:
#> survreg(formula = s2 ~ Age + NumPartners + CondomUse + YearsSchool +
#>   InitInfect + InvVagAtExam + DischargeExam, data = std_data,
#>   dist = "loglo")
#>
#>              Value Std. Error      z      p
#> (Intercept)   3.26979    0.72210  4.53 5.9e-06
#> Age           0.00152    0.01777  0.09 0.93188
#> NumPartners   -0.05643    0.07503 -0.75 0.45201
#> CondomUse2     0.26490    0.35206  0.75 0.45179
#> CondomUse3     0.63442    0.36424  1.74 0.08156
#> YearsSchool    0.19049    0.05365  3.55 0.00038
#> InitInfect2    0.42277    0.21873  1.93 0.05326
#> InitInfect3    0.32407    0.22067  1.47 0.14194
#> InvVagAtExam1 -0.70568    0.25677 -2.75 0.00599
#> DischargeExam1 0.54808    0.34697  1.58 0.11420
#> Log(scale)     0.09847    0.04479  2.20 0.02791
#>
#> Scale= 1.1

```

```
#>
#> Log logistic distribution
#> Loglik(model)= -2680.2   Loglik(intercept only)= -2701.5
#>  Chisq= 42.58 on 9 degrees of freedom, p= 2.6e-06
#> Number of Newton-Raphson Iterations: 4
#> n= 877
```

Residuals of the log-logistic Regression Model



The probability plots above also show that the Weibull is the better parametric model for the data, because the log-logistic and lognormal models clearly do not fit the line in the tails.

But how do we interpret this model fit? First of all, the model we have fit has the expression

$$Y = \ln(T) = \mu + \gamma^T \mathbf{Z} + \sigma W,$$

where $W \sim EV(0, 1)$,

Table 3: Parameter Estimates, AF and HR for each Parameter Estimate

	Parameter.Estimate	AF	HR
(Intercept)	3.9516047	0.0192238	0.0475893
Age	0.0101018	0.9899491	0.9922457
NumPartners	-0.0138832	1.0139800	1.0107560
CondomUse2	0.0788112	0.9242144	0.9410747
CondomUse3	0.4227912	0.6552154	0.7219443
YearsSchool	0.1704370	0.8432962	0.8769191
InitInfect2	0.5109972	0.5998971	0.6745026
InitInfect3	0.3149121	0.7298530	0.7845268
InvVagAtExam1	-0.5091714	1.6639119	1.4804896
DischargeExam1	0.4601294	0.6312020	0.7014676

$$\gamma^T = (\gamma_{Age}, \gamma_{Ethn.}, \gamma_{NumPartn.}, \gamma_{Cond.}, \gamma_{YSchool}, \gamma_{SignDisch})$$

are the estimated parameters and

$$\mathbf{Z}^T = (Age, Ethn., NumPartn., Cond., YSchool, SignDisch),$$

is the vector of values. **These need to be changed according to the values used also!!** Thus, each of the quantities $\exp(\gamma_i)$ s can be interpreted as the unitary change in time until reinfection (when the covariate i is continuous), or the change in time until reinfection when changing level (when the covariate i is categorical with different levels), when the other explanatory variables are kept fixed. **Any other interpretation of the model fit they are looking for you think?**

In the Weibull model, the acceleration factor (AF) is calculated using the equation

$$AF = \exp(-\hat{\gamma}_i)$$

and the hazard ratio (HR) is calculated using the equation

$$HR = \exp(-\hat{\gamma}_i/\hat{\sigma}).$$

In this case, the model fit gives the scale $\hat{\sigma} \approx 1.298$. These values are calculated for each of the covariates below.

Consider an example using the covariate **CondomUse** when explaining the interpretation of the covariates in terms of the AF. From the table above it is apparent that the AF of **CondomUse3** versus **CondomUse1** is ≈ 0.655 . This means that the reinfection time for a person that never uses a condom is ≈ 0.655 times the reinfection time for a person that always uses a condom **Not sure that this makes sense!? I think it makes sense with the coefficient value given from the model above, but does not make sense in real life, as this suggests that not using a condom is protective!** The interpretation in terms of the AF is similar when considering the other covariates, except for when considering the **Age** and **NumPartners**, which is not categorical **Perhaps it indeed makes sense to considering the Age in this way also, even though it is weird to treat the Age this way?**

The relative hazards (RH) can be calculated using the equation ...

Is the relative hazards the same as the hazard ratio??

Fit of a semi-parametric survival model

The proportional hazards model is fit. **This is probably done in next lab!**

Conclusions