

**ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ  
ОБРАЗОВАТЕЛЬНОЕ  
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ  
"МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ  
(НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ)"**

**ЖУРНАЛ ПРАКТИКИ**

Студента 2 курса

Гординского Дмитрия Михайловича

Институт №8 «Информационные технологии и прикладная математика»

Кафедра №804 «Теория вероятностей и компьютерное моделирование»

Учебная группа М8О-204Б-20

Направление 01.03.04

Прикладная математика

Вид практики Учебная (вычислительная) в Московском Авиационном Институте(НИУ)

Руководитель практики от МАИ Зайцева О.Б.

\_\_\_\_\_

Гординский Д.М /

/ 12 июля 2022 г.

**Москва, 2022**

### **1. Место и сроки проведения практики**

Дата начала практики 29 июня 2022 г.

Дата окончания практики 12 июня 2022 г.

Наименование предприятия МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ(НИИУ)

Название структурного подразделения Кафедра 804

### **2. Инструктаж по технике безопасности**

Платонов Е. Н. / / 29 июня 2022 г.

### **3. Индивидуальное задание студенту**

1. Тема "Анализ выживаемости"
2. Разобраться с теорией.
3. Привести пример решения задачи.
4. Написать отчет.

#### **4. План выполнения индивидуального задания**

1. Изучить теорию по Анализу выживаемости.
2. Ознакомиться с необходимыми библиотеками для работы с данными и их графическим представлением.
3. Решить задачу по анализу данных с применением методов анализа выживаемости.

*Руководитель практики от МАИ: Зайцева О.Б. / \_\_\_\_\_ /*

Платонов Е. Н. / \_\_\_\_\_ / 29 июня 2022 г.

### **5. Отзыв руководителя практики**

*Задание на практику выполнено в полном объеме. Материалы, изложенные в отчете студента, полностью соответствуют индивидуальному заданию.*

Руководитель

Платонов Е. Н. / / 12 июля 2022 г.

# Отчет студента

## Содержание

<b>1</b>	<b>Что такое “Анализ выживаемости”?</b>	<b>5</b>
1.1	Основные понятия . . . . .	5
1.1.1	Функция выживания (Survival function) . . . . .	5
1.1.2	Функция риска (Hazard function) . . . . .	5
1.1.3	Цензурирование (censoring) . . . . .	6
1.1.4	Медиана ожидаемого времени жизни (median number of survival days) . . . . .	6
1.1.5	Доверительный интервал (confidence interval) . . . . .	6
1.1.6	Усечение (truncation) . . . . .	6
1.2	Непараметрические методы оценивания распределения длительностей . . . . .	6
1.2.1	Оценка Каплана — Мейера . . . . .	6
1.2.2	Оценка Нельсона — Аалена . . . . .	7
1.2.3	Модель пропорциональных рисков ( <i>регрессионный анализ пропорциональных рисков Кокса</i> ) . . . . .	8
<b>2</b>	<b>Пример решения задачи</b>	<b>9</b>
2.1	Проанализируем половое соотношение . . . . .	9
2.2	Применяем <i>Оценку Каплана — Мейера</i> . . . . .	10
2.2.1	Сделаем таблицу событий . . . . .	10
2.2.2	Найдем вероятность выживания для каждого момента времени и вероятность с доверительным интервалом . . . . .	11
2.2.3	Найдем медиану времени выживания . . . . .	12
2.2.4	Найдем вероятность смерти для $\forall t$ . . . . .	13
2.3	Применяем <i>Оценку Нельсона — Аалена</i> . . . . .	14
2.3.1	Найдем риск для каждого момента времени и риск с доверительным интервалом . . . . .	14
2.3.2	Сравним кумулятивную функцию риска и кумулятивную плотность ( <i>вероятность смерти</i> ): . . . . .	15
2.4	Анализ выживания для групп . . . . .	15
2.4.1	Найдем вероятность выживания среди мужчин и женщин . . . . .	16
2.4.2	Сравним кумулятивную плотность выживания с кумулятивной функцией риска . . . . .	18
2.4.3	Найдем зависимость вероятности выживания от возрастной группы . . . . .	19
<b>3</b>	<b>Итоги</b>	<b>23</b>
3.1	Список материалов . . . . .	23

# 1 Что такое “Анализ выживаемости”?

*Анализ выживаемости* — набор статистических моделей, благодаря которым можно оценить вероятность наступления того или иного события. Анализ занимается моделированием процессов наступления *интересующих* нас (критических) событий для элементов той или иной совокупности (изначально — «смерти» для элементов совокупности живых существ).

*Интересным* событием может быть что угодно. Это может быть фактическая смерть, рождение, выход на пенсию и т. д.

Название “*survival analysis*” взято из медицины, т.к. цель анализа заключается в изучении продолжительности жизни пациента после приема препарата или других факторов влияния на здоровье.

## 1.1 Основные понятия

### 1.1.1 Функция выживания (Survival function)

Пусть  $T$  — неотрицательная случайная величина, представляющая собой время ожидания до наступления некоторого события. Для простоты будем использовать терминологию анализа выживаемости, называя исследуемое событие «смертью», а время ожидания — временем «выживания»

*Функция выживания* сопоставляет некоторому числу  $t$  вероятность того, что случайная величина  $T$  примет значение, не меньшее  $t$ . Иначе говоря, это вероятность того, что некоторое состояние «проживет» как минимум  $t$  единиц времени:

$$S(t) = \mathbb{P}\{T > t\} = 1 - \mathbb{P}\{T \leq t\}$$

Например, если мы хотим знать, какова вероятность того, что безработный индивид не сможет найти работу в течение полугода после начала поиска, то достаточно рассмотреть функцию выживания для  $t = 6$  месяцев.

### 1.1.2 Функция риска (Hazard function)

*Функцию риска* можно охарактеризовать как вероятность того, что событие произойдет за бесконечно малый интервал времени, при условии, что оно не произошло к моменту времени  $t$ .

$$h(t) = \lim_{dt \rightarrow 0} \frac{\mathbb{P}(t \leq T < t + dt | T \geq t)}{dt}$$

Числитель этого выражения — условная вероятность того, что событие произойдет в интервале  $(t, t + dt)$ , если оно не произошло ранее, а знаменатель — ширина интервала. Разделив одно на другое, получаем интенсивность осуществления события в единицу времени. Устремляя ширину интервала к нулю и переходя к пределу, получаем *мгновенную интенсивность осуществления события*.

Т. к. вышеперечисленные функции связаны друг с другом, можно показать, что:

$$S(t) = \exp\left(-\int_0^t h(x)dx\right)$$

Интеграл в фигурных скобках в этом уравнении называют *кумулятивным риском* и обозначают как:

$$H(t) = \int_0^t h(x)dx$$

Можно рассматривать  $H(t)$  как сумму всех рисков при переходе от момента времени 0 к  $t$ .

### 1.1.3 Цензурирование (censoring)

*Цензурирование* — вид неполноты информации, при котором наблюдения не содержат точной длительности изучаемого состояния. Различают цензурирование справа, слева и интервальное:

1. Цензурировано справа — о наблюдаемом состоянии известно лишь, что оно продлилось не менее определенного времени.
2. Цензурировано слева — о состоянии известно лишь, что оно продлилось не более определенного времени.
3. На интервале — известны только границы длительности.

### 1.1.4 Медиана ожидаемого времени жизни (median number of survival days)

Это точка на временной оси, в которой кумулятивная функция выживания равна 0,5. Другими словами, *медиана* — время, когда ожидается, что половина пациентов будет жива. Это означает, что шанс выжить после этого времени составляет 50%.

### 1.1.5 Доверительный интервал (confidence interval)

Доверительный интервал — интервал, который покрывает неизвестный параметр с заданной надёжностью. Вероятность, с которой в условиях данного эксперимента полученные экспериментальные данные можно считать надёжными (достоверными), называют *доверительной вероятностью* или надёжностью. Величина доверительной вероятности определяется характером производимых измерений. Мы будем считать доверительную вероятность равной 95 %.

### 1.1.6 Усечение (truncation)

*Усечением*, или урезанием, называется вид неполноты информации, при котором какая-то область возможных значений длительности оказывается недостаточно представленной в выборке: состояния, длительность которых слишком велика или, наоборот, слишком мала, просто не включаются в анализируемые данные. В нашей задаче мы будем называть их (removed) — пациенты, которые больше не являются частью нашего эксперимента. Если человек умирает или подвергается цензуре, то он попадает в эту категорию.

## 1.2 Непараметрические методы оценивания распределения длительностей

При отсутствии цензурирования и усечения для оценивания закона распределения вероятностей может использоваться эмпирическая функция распределения, из которой легко получить оценки для других характеристик случайной величины: survival function etc. Но в нашем случае это невозможно, т. к. мы имеем дело с неполнотой данных. Эту проблему решают непараметрические методы оценки.

### 1.2.1 Оценка Каплана — Мейера

*Оценка Каплана-Мейера* — это непараметрическая статистика, используемая для оценки функции выживания на основе данных о жизни. В медицинских исследованиях она часто используется для измерения доли пациентов, живущих в течение

определенного времени после лечения или постановки диагноза. Например: подсчет количества времени, которое прожил конкретный пациент после того, как у него был диагностирован рак или началось его лечение.

$$\hat{S}(t) = \prod_{t_j \leq t} \frac{n_j - d_j}{n_j}$$

где

$\hat{S}(t)$  = Вероятность того, что испытуемый жив в момент времени  $t$

$n_j$  = Количество испытуемых, оставшихся в живых непосредственно перед моментом времени  $t_j$

$d_j$  = Количество событий в момент времени  $t_j$

Можем переписать формулу выше так:

$$S(t_j) = S(t_{j-1}) \left(1 - \frac{d_j}{n_j}\right)$$

где

$S(t_j)$  = Вероятность того, что испытуемый жив в момент времени  $t_j$

$n_j$  = Количество испытуемых, оставшихся в живых непосредственно перед моментом времени  $t_j$

$d_j$  = Количество событий в момент времени  $t_j$

$S(0) = 1$

$t_0 = 0$

### 1.2.2 Оценка Нельсона — Аалена

Мы можем визуализировать совокупную информацию о выживании, используя функцию риска *Нельсона-Аалена*  $h(t)$ . Функция риска  $h(t)$  дает нам вероятность того, что субъект, находящийся под наблюдением в момент времени  $t$ , имеет интересующее событие (смерть) в это время. Чтобы получить информацию о функции риска, мы не можем преобразовать оценку Каплана-Мейера. Для этого существует соответствующая непараметрическая оценка кумулятивной функции риска:

$$\hat{H}(t) = \sum_{t_j \leq t} \frac{d_j}{n_j}$$

где

$\hat{H}(t)$  = Кумулятивная вероятность риска

$n_j$  = Количество испытуемых, оставшихся в живых непосредственно перед моментом времени  $t_j$

$d_j$  = Количество событий в момент времени  $t_j$



### 1.2.3 Модель пропорциональных рисков (регрессионный анализ пропорциональных рисков Кокса)

*Модели пропорциональных рисков* — прогнозирование риска наступления события для рассматриваемого объекта и оценка влияния заранее определенных независимых переменных (предикторов) на этот риск.

В качестве решения для этого мы используем *регрессионный анализ пропорциональных рисков Кокса*, который работает как для количественных предикторов некатегориальных переменных, так и для категориальных переменных.

*Регрессионная модель Кокса (Cox regression)* — в анализе выживаемости математическая модель зависимости функции риска от независимых переменных-факторов. В анализе выживаемости решается задача оценки функции выживания или функций, производных от нее.

В нашей задаче мы попытаемся рассмотреть зависимость вероятности выживания от возрастной группы.

Целью метода пропорциональной риска Кокса является определение того, как различные факторы в нашем наборе данных влияют на интересующее нас событие.

$$h(t) = h_0(t) * \exp(b_1x_1 + b_2x_2 + \dots + b_nx_n)$$

где

$t$  = время выживания

$h(t)$  = функция риска

$x_1, x_2, \dots, x_n$  = ковариации

$b_1, b_2, \dots, b_n$  = влияния параметров ковариаций

$\exp(b_i)$  = коэффициент риска (Hazard Ratio [HR]),

если:

$b_i = 1 \Rightarrow \exp(b_i) = 0 \Rightarrow$  ковариат не оказывает влияния на риск.

$b_i < 1 \Rightarrow \exp(b_i) = 0 \Rightarrow$  ковариат оказывает отрицательное влияние на риск  $\Rightarrow$  положительное на время выживания.

$b_i > 1 \Rightarrow \exp(b_i) = 0 \Rightarrow$  ковариат оказывает положительное влияние на риск  $\Rightarrow$  отрицательное на время выживания.

## 2 Пример решения задачи

В качестве примера для анализа выживаемости возьмем заболевание *Chronic Granulomatous Disease*

*Хроническая гранулематозная болезнь*(ХГБ)[CGD] — это наследственное заболевание, которое возникает, когда тип лейкоцитов (фагоцитов), которые обычно помогают организму бороться с инфекциями, не работает должным образом. В результате фагоциты не могут защитить организм от бактериальных и грибковых инфекций. У людей с хронической гранулематозной болезнью могут развиваться инфекции в легких, коже, лимфатических узлах, печени, желудке и кишечнике или других областях. У них также могут образовываться скопления лейкоцитов в зараженных областях. У большинства людей ХГБ диагностируется в детстве, но у некоторых людей диагноз может не ставиться до зрелого возраста.

В датасете нас интересует:

- период, в течение которого наблюдали за пациентом ( $tstop - tstart$ )
- пол пациента ( $sex$ )
- $status = \{status == 0 = \text{alive}, status == 1 = \text{dead}\}$

В дальнейшем может пригодиться возраст ( $age$ ) для анализа выживаемости по группам.

### 2.1 Проанализируем половое соотношение

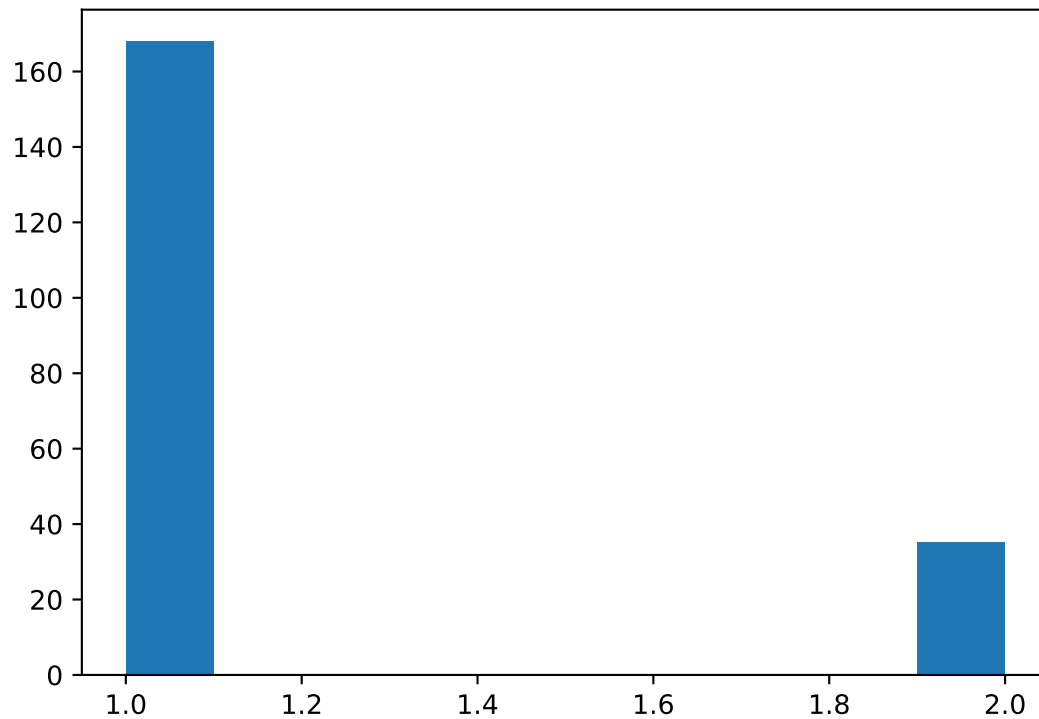
для начала подключим необходимые библиотеки...

*lifelines* — содержит необходимые нам методы для исследования вероятностей и времени жизни.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from lifelines import KaplanMeierFitter
from lifelines import NelsonAalenFitter
# считываем данные
data = pd.read_csv("src/cgd.csv")
head = data.head()
```

Unnamed: 0	id	center	random	treat	sex	age	...	steroids	propylac	hos.cat	tstart	enum	tstop	status	
0	1	1	Scripps Institute	1989-06-07	rIFN-g	2	12	...	0	0	US:other	0	1	219	1
1	2	1	Scripps Institute	1989-06-07	rIFN-g	2	12	...	0	0	US:other	219	2	373	1
2	3	1	Scripps Institute	1989-06-07	rIFN-g	2	12	...	0	0	US:other	373	3	414	0
3	4	2	Scripps Institute	1989-06-07	placebo	1	15	...	0	1	US:other	0	1	8	1
4	5	2	Scripps Institute	1989-06-07	placebo	1	15	...	0	1	US:other	8	2	26	1

```
data.loc[data.sex == "male", "sex"] = 1
data.loc[data.sex == "female", "sex"] = 2
plt.hist(data["sex"]) #гистограмма "пол"
plt.show()
```



## 2.2 Применяем *Оценку Каплана — Мейера*

```
kmf = KaplanMeierFitter()
# в нашем случае "status" === "dead"
data.loc[:, "time"] = data.loc[:, "tstop"] - data.loc[:, "tstart"] # time=tstop-tstart
kmf.fit(durations = data["time"], event_observed = data["status"])
```

```
## <lifelines.KaplanMeierFitter:"KM_estimate", fitted with 203 total observations, 127 right-
censored observations>
```

### 2.2.1 Сделаем таблицу событий

Нам это нужно, чтобы отделить цензурированные данные, получить необходимые временные данные для применения методов оценки.

```
print(kmf.event_table)
```

```
##          removed  observed  censored  entrance  at_risk
## event_at
## 0.0             0         0         0         203      203
## 2.0             1         1         0          0      203
## 4.0             3         2         1          0      202
## 5.0             1         1         0          0      199
## 6.0             1         1         0          0      198
```

```
## ...      ...      ...      ...      ...
## 371.0      1      0      1      0      7
## 373.0      2      1      1      0      6
## 376.0      1      0      1      0      4
## 382.0      1      0      1      0      3
## 388.0      2      0      2      0      2
##
## [154 rows x 5 columns]
```

где

- *event\_at* — хранит значение временной шкалы для нашего набора данных. т. е. когда пациент наблюдался в нашем эксперименте или когда был проведен эксперимент, хранит значение дней выживания для субъектов.
- *at\_risk* — хранит количество текущих пациентов, находящихся под наблюдением.

$$at\_risk = current\ patients\ at\_risk + entrance - removed$$

- *entrance* — хранит значение новопришедших пациентов. Т. е. во время проведения эксперимента появлялись новые больные.
- *censored* — если человек все еще жив по окончании эксперимента, то мы добавляем его в эту категорию.
- *observed* — содержит количество умерших пациентов во время эксперимента.
- *removed* — содержит количество пациентов, которые “выпадают” из эксперимента  
 $removed = observed + censored$

## 2.2.2 Найдем вероятность выживания для каждого момента времени и вероятность с доверительным интервалом

Для начала найдем вероятность выживания за время  $t$   
(см. раздел 1.2.1)

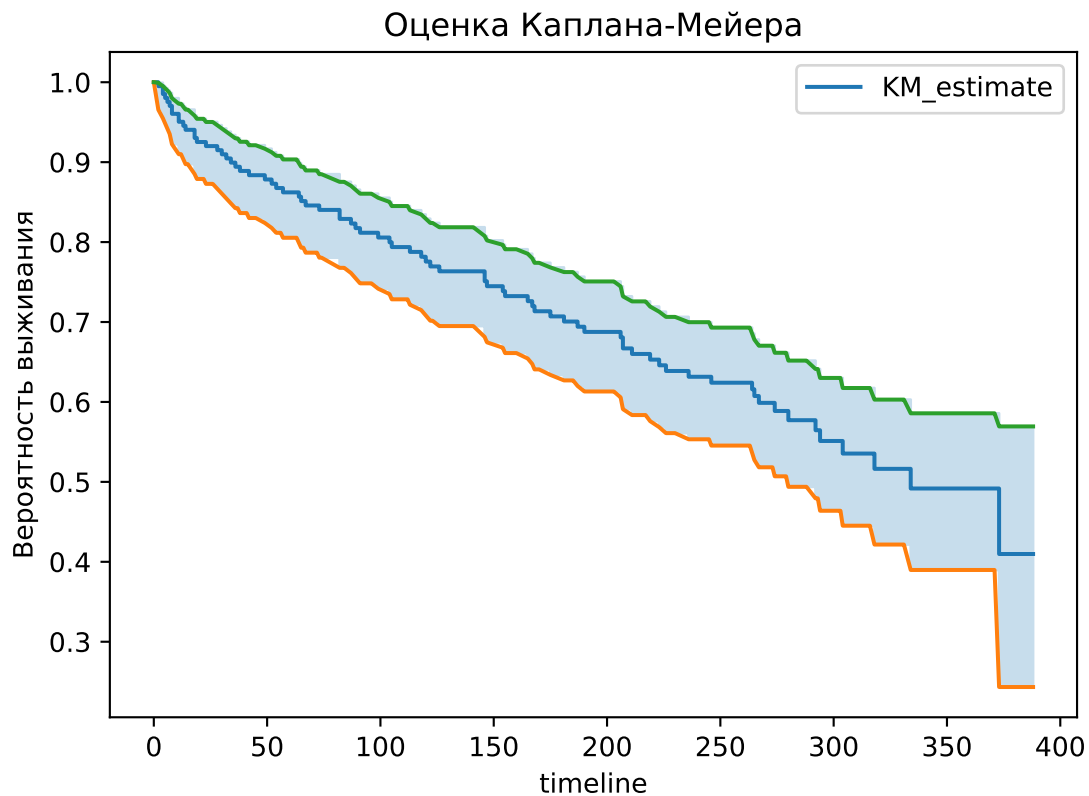
Не ограничивая общности,  $\forall t$ : Пусть  $t = 6 \Rightarrow$

```
# Вероятность выживания после 6 дней
e0 = kmf.event_table.iloc[0, :]
e2 = kmf.event_table.iloc[1, :]
e4 = kmf.event_table.iloc[2, :]
e6 = kmf.event_table.iloc[3, :]
s0 = (e0.at_risk - e0.observed)/e0.at_risk
s2 = (e2.at_risk - e2.observed)/e2.at_risk
s4 = (e4.at_risk - e4.observed)/e4.at_risk
s6 = (e6.at_risk - e6.observed)/e6.at_risk
s6 = s0 * s2 * s4 * s6
print(s6)
```

```
## 0.9802708121890239
```

Найдем вероятность выживания  $\forall t$ :

```
kmf.survival_function_  
plt.title("Оценка Каплана-Мейера")  
plt.ylabel("Вероятность выживания")  
kmf.plot()  
csf = kmf.confidence_interval_survival_function_  
plt.plot(csf["KM_estimate_lower_0.95"], label="lower")  
plt.plot(csf["KM_estimate_upper_0.95"], label="upper")  
plt.show()
```



По графику видно, что с течением времени вероятность выживания уменьшается.

### 2.2.3 Найдем медиану времени выживания

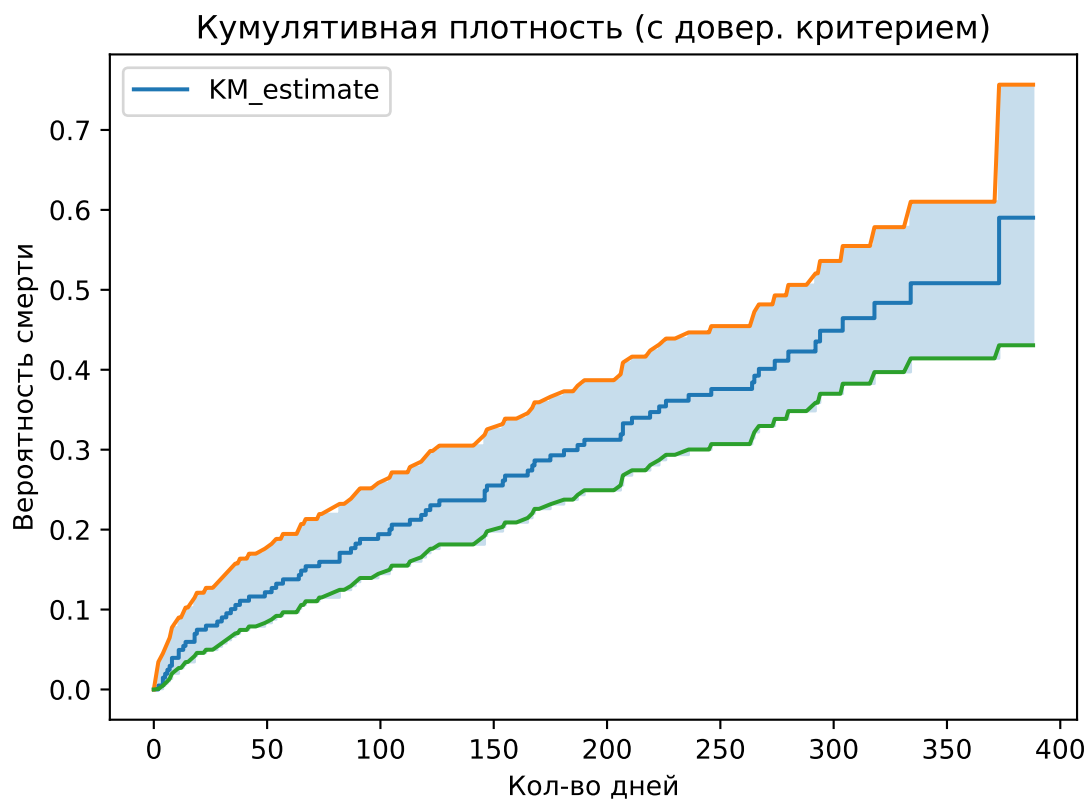
```
print("Медиана времени выживания", kmf.median_survival_time_)
```

```
## Медиана времени выживания 334.0
```

### 2.2.4 Найдем вероятность смерти для $\forall t$

Сделаем график кумулятивной функции плотности и кумулятивной плотности с доверительным критерием

```
kmf.plot_cumulative_density()
ccf = kmf.confidence_interval_cumulative_density_
plt.plot(ccf["KM_estimate_lower_0.95"], label="lower")
plt.plot(ccf["KM_estimate_upper_0.95"], label="upper")
plt.title("Кумулятивная плотность (с довер. критерием)")
plt.xlabel("Кол-во дней")
plt.ylabel("Вероятность смерти")
plt.show()
```



Видим полностью обратный график к вероятности выживания, что неудивительно, ведь кумул. ф-я плотности:

$$F(t) = 1 - S(t)$$

## 2.3 Применяем *Оценку Нельсона — Аалена*

```
naf = NelsonAalenFitter()
naf.fit(durations = data["time"], event_observed = data["status"])
```

```
## <lifelines.NelsonAalenFitter:"NA_estimate", fitted with 203 total observations, 127 right-
censored observations>
```

### 2.3.1 Найдем риск для каждого момента времени и риск с доверительным интервалом

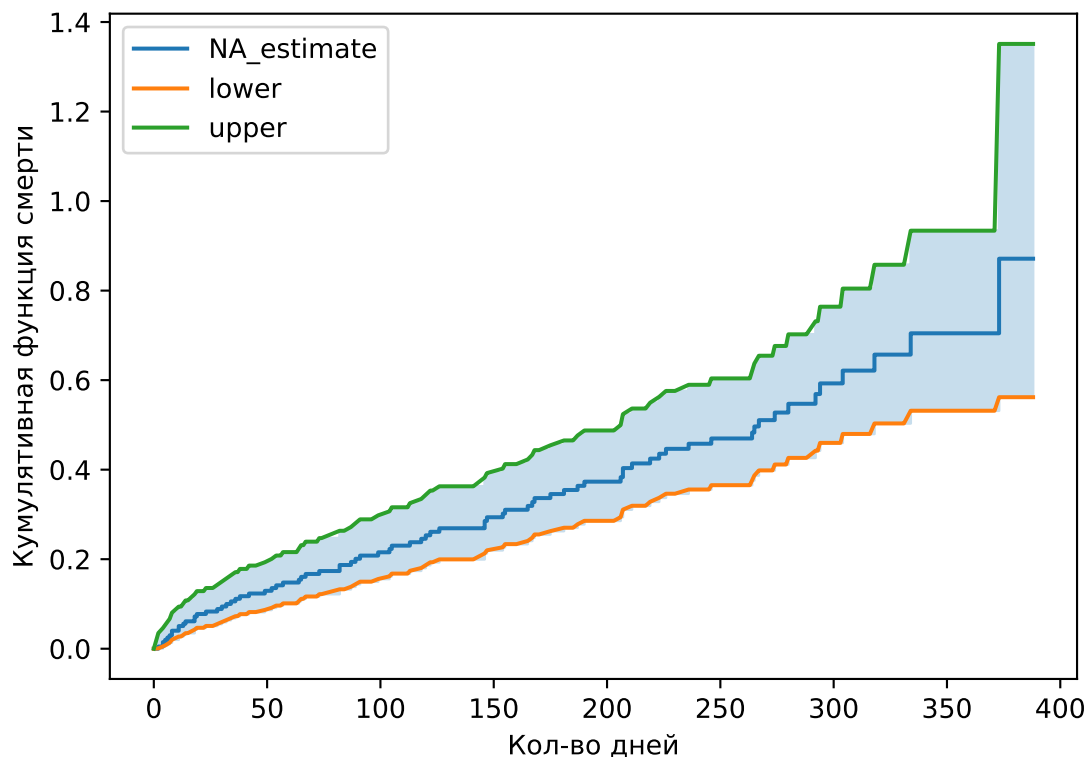
По нашей таблице событий (см. раздел 2.2.1) мы считаем функцию риска (см. раздел 1.2.2)

Также можем предсказать значение функции риска для  $\forall t$

```
print("300 дней: ", naf.predict(300))
```

```
## 300 дней: 0.5927191310947535
```

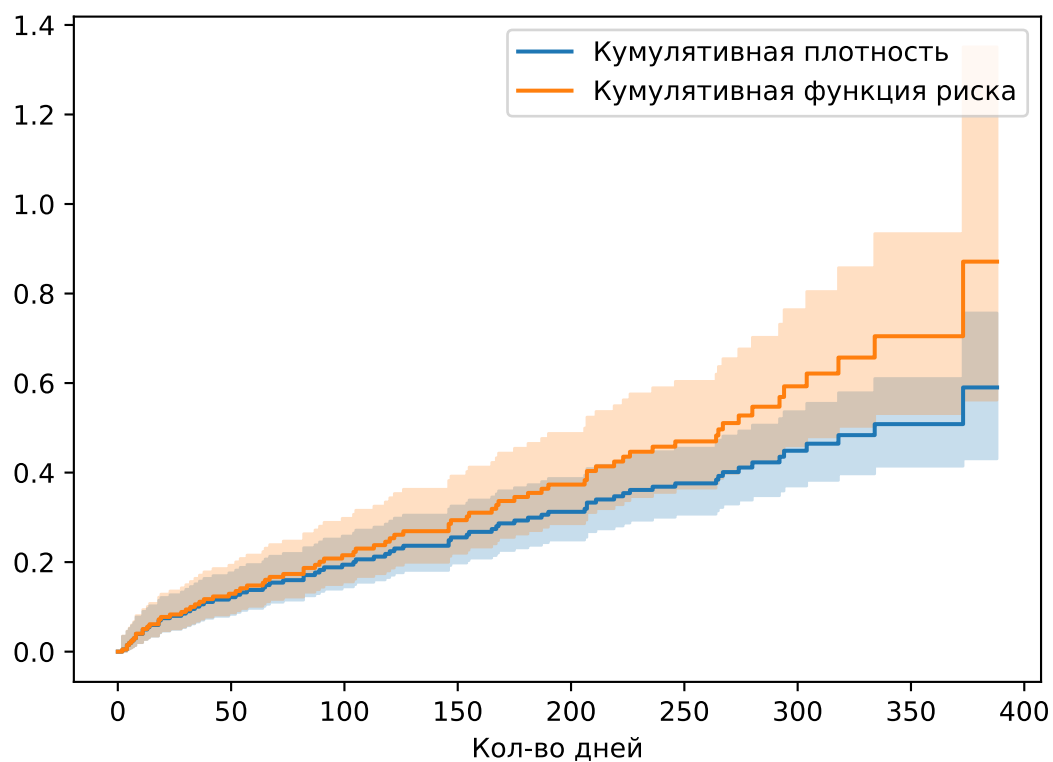
```
naf.plot_cumulative_hazard()
ci = naf.confidence_interval_
plt.plot(ci["NA_estimate_lower_0.95"], label="lower")
plt.plot(ci["NA_estimate_upper_0.95"], label="upper")
plt.xlabel("Кол-во дней")
plt.ylabel("Кумулятивная функция смерти")
plt.legend()
```



Другими словами, функция риска измеряет *общую сумму риска*, накопленного к моменту времени  $t$

### 2.3.2 Сравним кумулятивную функцию риска и кумулятивную плотность (вероятность смерти):

```
kmf.plot_cumulative_density(label="Кумулятивная плотность")
naf.plot_cumulative_hazard(label="Кумулятивная функция риска")
plt.xlabel("Кол-во дней")
plt.show()
```



## 2.4 Анализ выживания для групп

Сначала сравним выживаемость для мужчин и женщин:

```
kmfm = KaplanMeierFitter() # мужчины
kmff = KaplanMeierFitter() # женщины

data.loc[data.sex == "male", "sex"] = 1
data.loc[data.sex == "female", "sex"] = 2

male = data.query("sex == 1")
female = data.query("sex == 2")
```



```

kmfm.fit(durations=male["time"], event_observed=male["status"], label="male")
kmff.fit(durations=female["time"], event_observed=female["status"], label="female")

# сделаем таблицы событий отдельно для м. и ж.
kmfm.event_table
kmff.event_table
print(kmfm.event_table.head())
print(kmff.event_table.head())

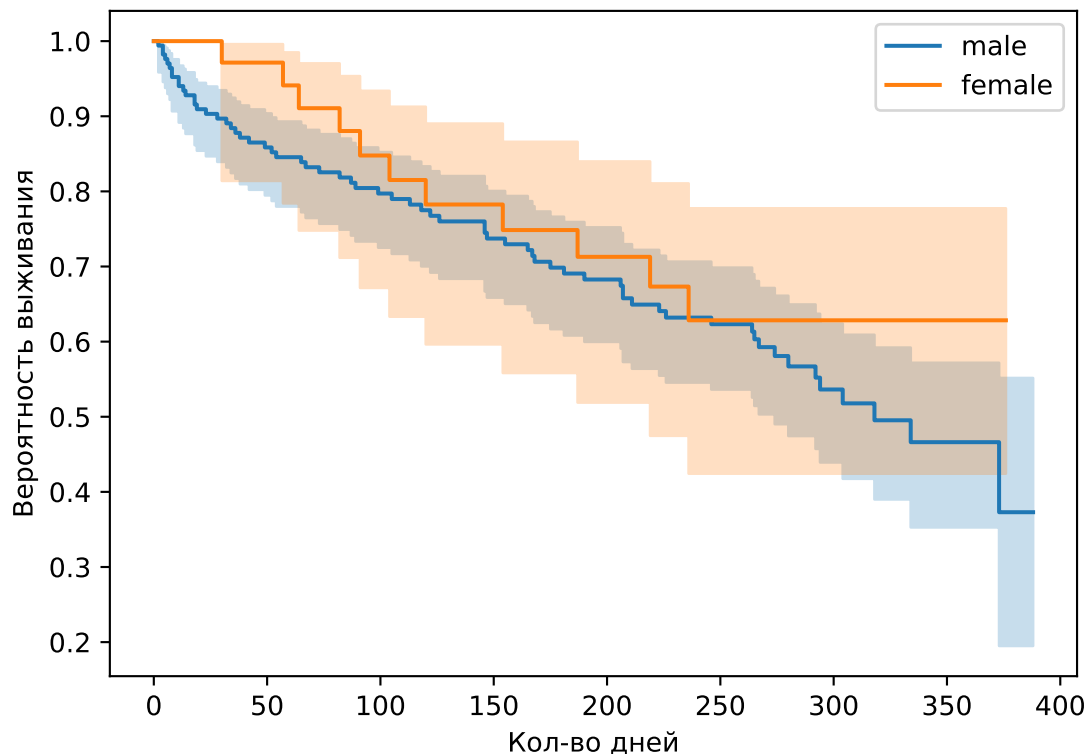
```

### 2.4.1 Найдем вероятность выживания среди мужчин и женщин

```

kmfm.survival_function_
kmff.survival_function_
kmfm.plot()
kmff.plot()
plt.xlabel("Кол-во дней")
plt.ylabel("Вероятность выживания")
plt.show()

```



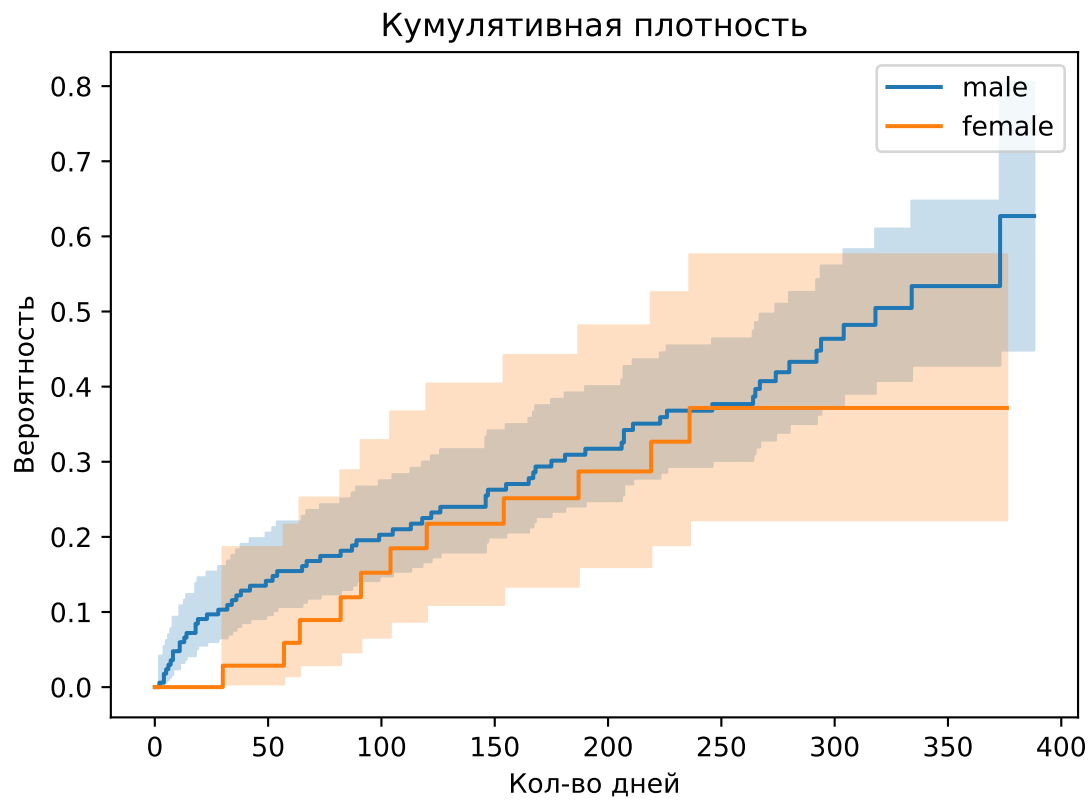
Несложно заметить, что вероятность выживания женщин выше, чем у мужчин. Однако такие выводы не совсем точны, т.к. мужчин много больше, чем женщин. Но до ~250 дня все равно вероятность выживания у женщин выше.

Получается, что кумулятивная плотность будет ниже у женщин

```

kmfm.plot_cumulative_density()
kmff.plot_cumulative_density()
plt.title("Кумулятивная плотность")
plt.xlabel("Кол-во дней")
plt.ylabel("Вероятность")
plt.show()

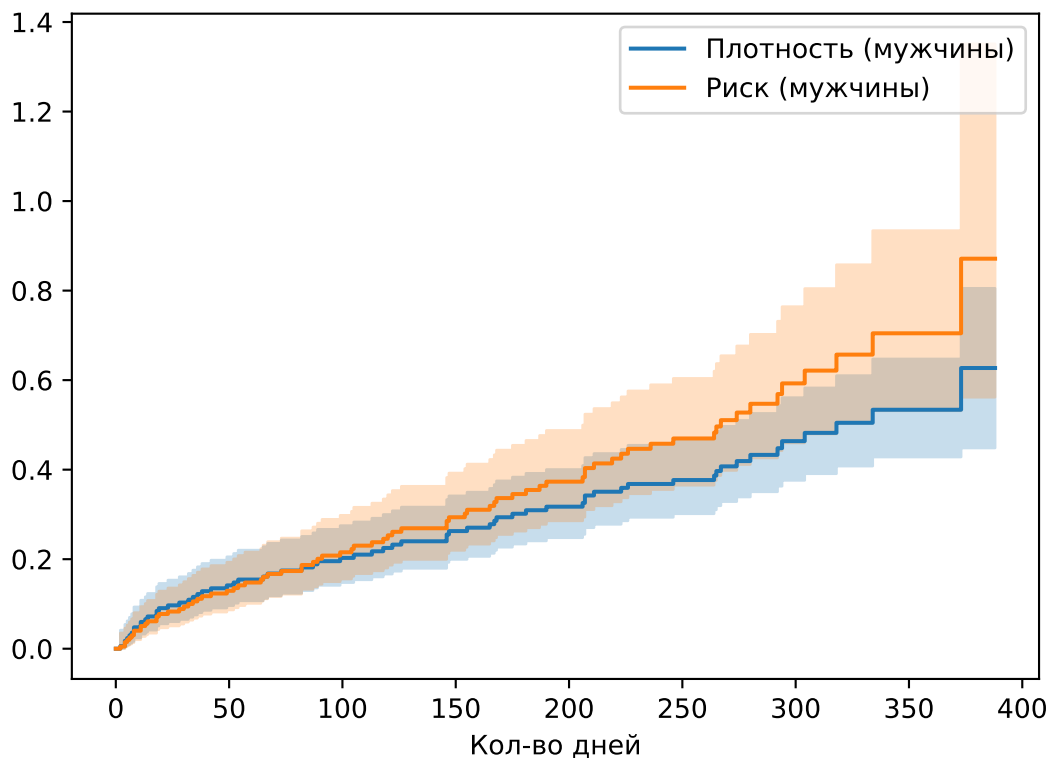
```



## 2.4.2 Сравним кумулятивную плотность выживания с кумулятивной функцией риска

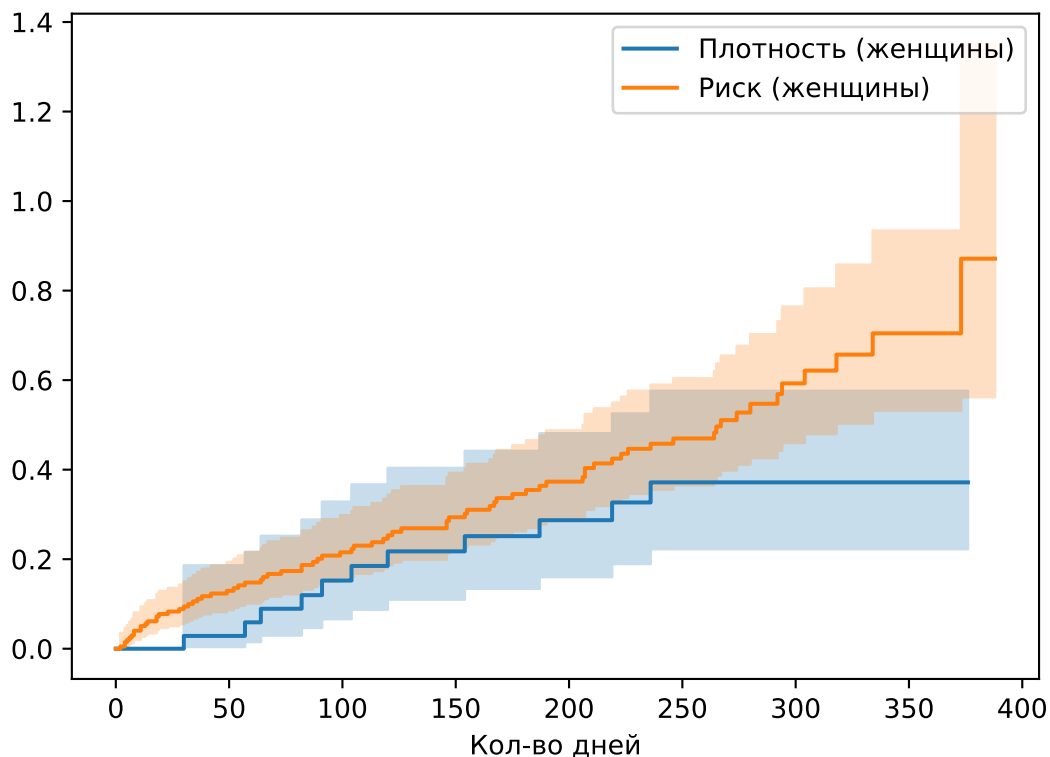
мужчины:

```
nafm = NelsonAalenFitter()
nafm.fit(durations = data["time"], event_observed = data["status"])
kmfm.plot_cumulative_density(label="Плотность (мужчины)")
nafm.plot_cumulative_hazard(label="Риск (мужчины)")
plt.xlabel("Кол-во дней")
plt.show()
```



женщины:

```
naff = NelsonAalenFitter()
naff.fit(durations = data["time"], event_observed = data["status"])
kmff.plot_cumulative_density(label="Плотность (женщины)")
naff.plot_cumulative_hazard(label="Риск (женщины)")
plt.xlabel("Кол-во дней")
plt.show()
```



⇒ с течением времени риск увеличивается.

### 2.4.3 Найдем зависимость вероятности выживания от возрастной группы

```
from lifelines import CoxPHFitter
data = pd.read_csv("src/cgd.csv")
data = pd.read_csv("src/cgd.csv")
data = data.drop("center", axis=1)
data = data.drop("id", axis=1)
data = data.drop("random", axis=1)
data = data.drop("steroids", axis=1)
data = data.drop("inherit", axis=1)
data = data.drop("hos.cat", axis=1)
data = data.drop("propylac", axis=1)
data = data.drop("enum", axis=1)
```

Удалим из наших данных строки с нулевыми значениями

```
data = data.dropna(subset=['sex', 'age', 'treat', 'height',
                           'weight', 'tstart', 'tstop', 'status'])
```

```
# возьмем объект kmf из (раздела 2.2)
data.loc[:, "time"] = data.loc[:, "tstop"] - data.loc[:, "tstart"] # time=tstop-tstart
data.loc[data.sex == "male", "sex"] = 1
data.loc[data.sex == "female", "sex"] = 2
```

```
data.loc[data.treat == "placebo", "treat"] = 1
data.loc[data.treat == "rIFN-g", "treat"] = 2
kmf.fit(durations = data["time"], event_observed = data["status"])
kmf.event_table
data = data[['time', 'treat', 'age', 'height', 'weight', 'sex', 'status']]
```

```
cph = CoxPHFitter()
cph.fit(data, "time", event_col="status")
```

```
<lifelines.CoxPHFitter: fitted with 203 total observations, 127 right-censored observations>
    duration col = 'time'
    event col = 'status'
    baseline estimation = breslow
    number of observations = 203
    number of events observed = 76
    partial log-likelihood = -348.68
    time fit was run = 2022-07-10 17:08:50 UTC

---
      coef  exp(coef)    se(coef)   coef lower 95%   coef upper 95%  exp(coef) lower 95%  exp(coef) upper 95%
covariate
treat    -1.15     0.32     0.27         -1.68         -0.61           0.19           0.54
age       -0.08     0.92     0.03         -0.15         -0.02           0.86           0.98
height     0.00     1.00     0.01         -0.02          0.02           0.98           1.02
weight     0.02     1.02     0.02         -0.01          0.05           0.99           1.05
sex       -0.05     0.95     0.33         -0.71          0.60           0.49           1.83

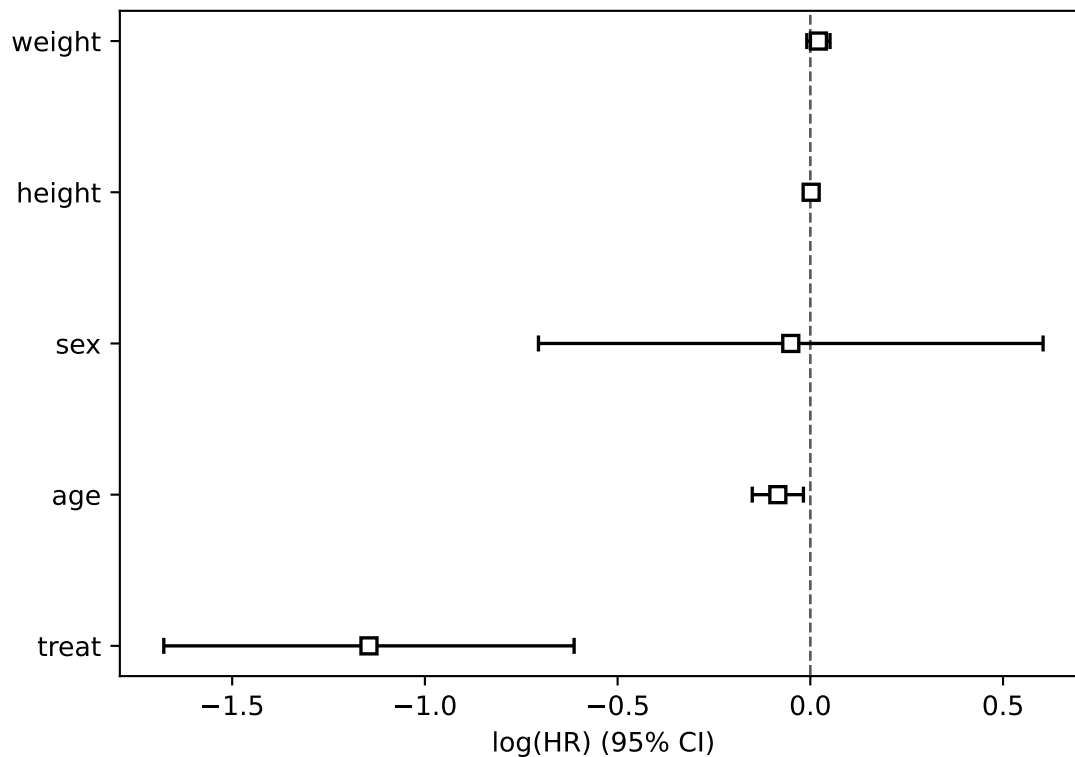
      cmp to      z      p   -log2(p)
covariate
treat      0.00 -4.22 <0.005    15.30
age         0.00 -2.49   0.01     6.28
height      0.00  0.22   0.82     0.28
weight      0.00  1.36   0.17     2.53
sex         0.00 -0.15   0.88     0.19
---
Concordance = 0.67
Partial AIC = 707.36
log-likelihood ratio test = 28.13 on 5 df
-log2(p) of ll-ratio test = 14.83
```

Мы знаем, что значение  $p < 0.05$  считается статистически значимым. Здесь мы видим, что *age* имеет 0.01, *treat* менее 0.005. Поэтому рассмотрим группировку “по возрасту” и “по лечению”.

Заметим, что значение  $HR = \exp(coef)$  для *age* равно 0.92,  $\Rightarrow$  слабая зависимость риска смерти от возраста пациента, т.к. если остальные ковариаты = const,  $\Rightarrow$  если вы относитесь к взрослой возрастной группе, то имеете на 8% меньше смертельный риск, чем, например, дети. Для *treat* ситуация такая: т.к.  $HR$  равна 0.32,  $\Rightarrow$  сильная зависимость риска смерти от вида лечения: пациенты, принимающие *rIFN-g(Interferon gamma)*, имеют на 68% меньше смертельный риск, нежели те, кому давали *placebo*.

Построим график, чтобы убедиться в наших предположениях:

```
cph.plot()
```

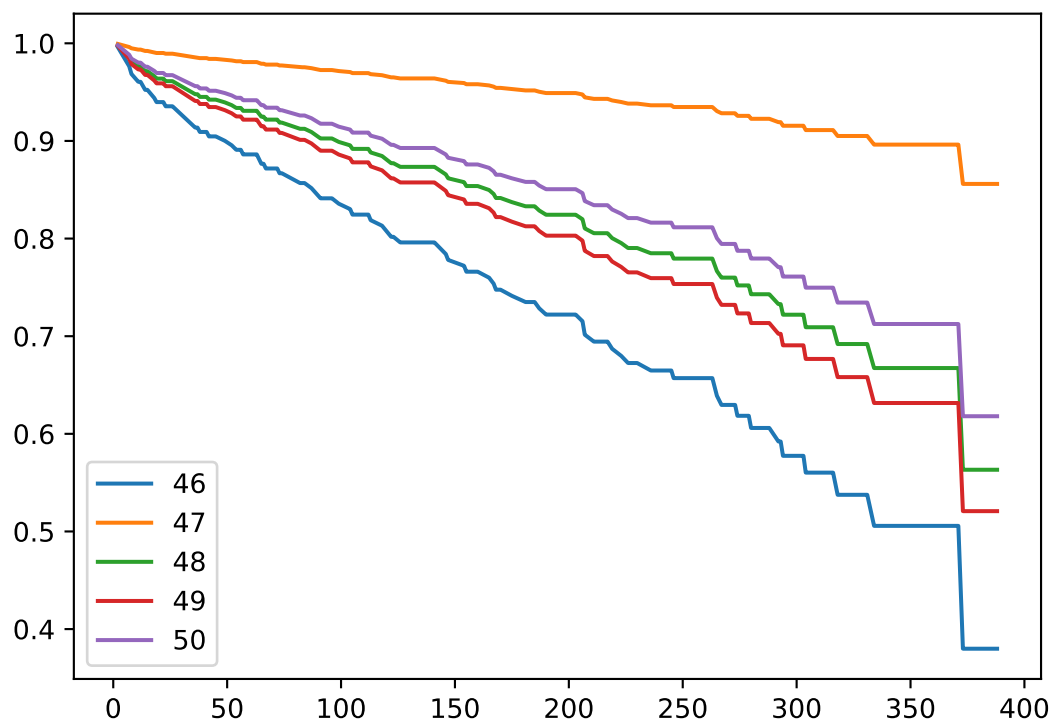


Рассмотрим некоторых пациентов, чтобы убедиться в правильности оценки наших параметров, влияющих на смертельный риск

```
print(data.iloc[46:51, :])
```

```
##      time treat  age height weight sex  status
## 46   105     1   29  175.0   73.1   1        0
## 47   376     2   31  167.0   51.8   2        0
## 48   360     2    7  121.0   19.9   1        0
## 49   306     1   26  153.0   46.9   2        0
## 50   160     2   12  136.5   30.0   1        0
```

```
ddata = data.iloc[46:51, :]
cph.predict_survival_function(ddata).plot()
plt.show()
```



Из-за того, что основным весовым параметром у нас получился *treat*, видно, что пациент (#46) имеет самую низкую вероятность выживания, а пациент (#47) мало того, что самый взрослый, так еще и принимал лекарство *rIFN-g*, поэтому имеет максимальную вероятность выживания среди данной выборки.

## 3 Итоги

Проанализировав выживаемость пациентов с заболеванием CGD, можно сказать, что вероятность выживания со временем падает  $\Rightarrow$  чем больше времени проходит, тем больше смертельный риск.

Анализ показал, что на нашей выборке у женщин немного больше шансов на выживание, нежели у мужчин. Выживаемость в старших возрастных группах на  $\sim 8\%$  выше, чем в младших. Также мы выяснили, что если пациент принимает лекарство, а не плацебо, то вероятность смертельного исхода значительно падает — на  $\sim 68\%$ .

### 3.1 Список материалов

- <https://towardsdatascience.com/force-of-mortality-in-bathtub-shaped-lifetimes-b0425cb66925>
- <https://pub.towardsai.net/survival-analysis-with-python-tutorial-how-what-when-and-why-19a5cfb3c312>
- “Анализ панельных данных и данных о длительности состояний”, Ратникова Т.А., Фурманов К.К.
- <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2394368/>
- <http://www-eio.upc.edu/~pau/cms/rdata/datasets.html>
- [https://wiki5.ru/wiki/Proportional\\_hazards\\_model](https://wiki5.ru/wiki/Proportional_hazards_model)
- <http://statsoft.ru/home/textbook/modules/stsurvan.html>