# AKKODIS INTERVIEW

Analysis of an oncological dataset

Erik De Luca

# SUMMARY

- Assignment description

- Setting the enviroment

- Descriptive Analysis

- Survival Analysis

- Deeper Analysis

# DATASET DESCRIPTION

The response variable, `SurvTime`, is the survival time in days of a lung cancer patient.

The covariates are:

- `Cell` (type of cancer cell),

- `Therapy` (type of therapy: standard or test),

- `Prior` (prior therapy: 0=no, 10=yes),

- `Age` (age in years),

- `DiagTime`(time in months from diagnosis to entry into the trial)

- `Kps`(performance status).

A censoring indicator variable Censor is created from the data, with the value 1 indicating a censored time and the value 0 indicating an event time. Since there are only two types of therapy, an indicator variable, Treatment, is constructed for therapy type, with: value 0 for standard therapy and value 1 for test therapy.

# EXERCISES

1. what was the maximum survival time for the cell type adeno?

2. what is the average age of subjects in this study?

3. which cell type appeared the most during this study?

4. Calculate descriptive statistics for all numeric variables within this dataset?

5. Perform a survival analysis to assess the survival time (variable SurvTime)? based on the cancerous cells (var Cell)? Consider applying survival functions/kaplan meier quartiles/cumulative incidence function/cox regression etc.

6. Perform an appropriate multivariable analysis to analyze the effect of independent variables age on the hazard ratio between the different cancerous cells (var Cell)?

# LOAD LIBRARIES

```r
 1  pacman::p_load(
 2    tidyverse,  # A set of many useful libraries
 3    readxl,     # To import the dataset from Excel
 4    here,       # To avoid problems with file directories
 5    janitor,    # To clean data in a fast way
 6    gt,         # Output tables
 7    gtsummary,  # Output tables for models and survival data
 8    survival,   # To manage survival data
 9    ggsurvfit,  # To plot survival analysis
10    tidycmprsk  # To fit survival models
11  )
```

# IMPORT DATA

```r
1  data_imported <- read_xlsx(
2    here("Akkodis interview", "Oncology_dataset_for_R.xlsx")
3    )
4
5  data_cleaned <- data_imported |>
6    clean_names() |>
7    remove_empty() |>
8    remove_constant()
9
10 data_cleaned |>
11   slice_sample(n = 1, by = c(therapy, cell)) |>
12   gt() |>
13   tab_header("Sample of the dataset cleaned",
14             "The sample was stratified by the variables therapy and cell")
```

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Sample of the dataset cleaned** | | | | | | | | | | |
| The sample was stratified by the variables therapy and cell | | | | | | | | | | |
| obs | therapy | cell | surv_time | kps | diag_time | age | prior | treatment | censor | event |
| 6 | Standard | Squamous | 10 | 20 | 5 | 49 | 10 | 0 | 0 | 1 |
| 16 | Standard | Small | 30 | 60 | 3 | 61 | 10 | 0 | 0 | 1 |
| 49 | Standard | adeno | 117 | 80 | 2 | 38 | 0 | 0 | 0 | 1 |
| 60 | Standard | large | 12 | 40 | 12 | 68 | 0 | 0 | 0 | 1 |
| 72 | Test | Squamous | 87 | 80 | 3 | 48 | 0 | 1 | 1 | 0 |
| 97 | Test | Small | 7 | 20 | 11 | 66 | 10 | 1 | 0 | 1 |

6

# WHAT WAS THE MAXIMUM SURVIVAL TIME FOR THE CELL TYPE ADENO?

```r
1  data_cleaned |>
2    filter(
3      surv_time == max(surv_time),
4      .by = cell
5      ) |>
6    select(obs, cell, surv_time, kps, age) |>
7    gt() |>
8    tab_header("Maximum Survival Time by Cell",
9               "Adeno's row is bolded") |>
10   tab_style(
11     style = cell_text(weight = "bold"),
12     locations = list(
13       cells_body(rows = cell == "adeno")
14       )
15     )
```

**Maximum Survival Time by Cell**

Adeno's row is bolded

| obs | cell | surv_time | kps | age |
|-----|------|-----------|-----|-----|
| 44 | Small | 392 | 40 | 68 |
| **52** | **adeno** | **162** | **80** | **64** |
| 58 | large | 553 | 70 | 47 |
| 70 | Squamous | 999 | 90 | 54 |

# WHAT IS THE AVERAGE AGE OF SUBJECTS IN THIS STUDY?

```
1  data_cleaned |>
2    summarise(
3      across(age, list(
4        mean = mean,
5        median = median,
6        sd = ~ round(sd(.)),
7        Q1 = ~ quantile(., .25),
8        Q3 = ~ quantile(., .75),
9        min = min,
10       max = max
11     ),
12     .names = "{fn}"
13     )
14   ) |>
15   pivot_longer(everything(), names_to = "statistics") |
16   gt() |>
17   tab_header(
18     "Age",
19     "Mean and other position and variance indicators"
20     ) |>
21   tab_style(
22     style = cell_text(weight = "bold"),
23     locations = list(
24       cells_body(rows = statistics == "mean")
25       )
```

| Age | |
|---|---|
| Mean and other position and variance indicators | |
| statistics | value |
| **mean** | **57.60** |
| median | 62.00 |
| sd | 11.00 |
| Q1 | 50.00 |
| Q3 | 65.25 |
| min | 34.00 |
| max | 81.00 |

# WHICH CELL TYPE APPEARED THE MOST DURING THIS STUDY?

```r
 1  data_cleaned |>
 2    tabyl(cell) |>
 3    adorn_pct_formatting(digits = 0) |>
 4    gt() |>
 5    tab_header("Cell's frequency") |>
 6    tab_style(
 7      style = cell_text(weight = "bold"),
 8      locations = list(
 9        cells_body(rows = n == max(n))
10        )
11    )
```

| Cell's frequency | | |
|---|---|---|
| cell | n | percent |
| **Small** | **41** | **41%** |
| Squamous | 35 | 35% |
| adeno | 9 | 9% |
| large | 15 | 15% |

# CALCULATE DESCRIPTIVE STATISTICS FOR ALL NUMERIC VARIABLES WITHIN THIS DATASET?

```
 1  data_cleaned |>
 2    select(-obs) |>
 3    summarise(
 4      across(
 5        where(is.numeric),
 6            list(
 7        mean = mean,
 8        median = median,
 9        sd = ~ round(sd(.)),
10        Q1 = ~ quantile(., .25),
11        Q3 = ~ quantile(., .75),
12        min = min,
13        max = max
14      ),
15      .names = "{col}-{fn}"
16      )
17    ) |>
18    pivot_longer(everything(), names_to = "statistics") |>
19    separate(col = statistics, sep = "-", into = c("column", "statistics")) |>
20    pivot_wider(names_from = statistics, values_from = value) |>
21    gt() |>
22    tab_header("Descriptive Statistics", "All numeric variables")
```
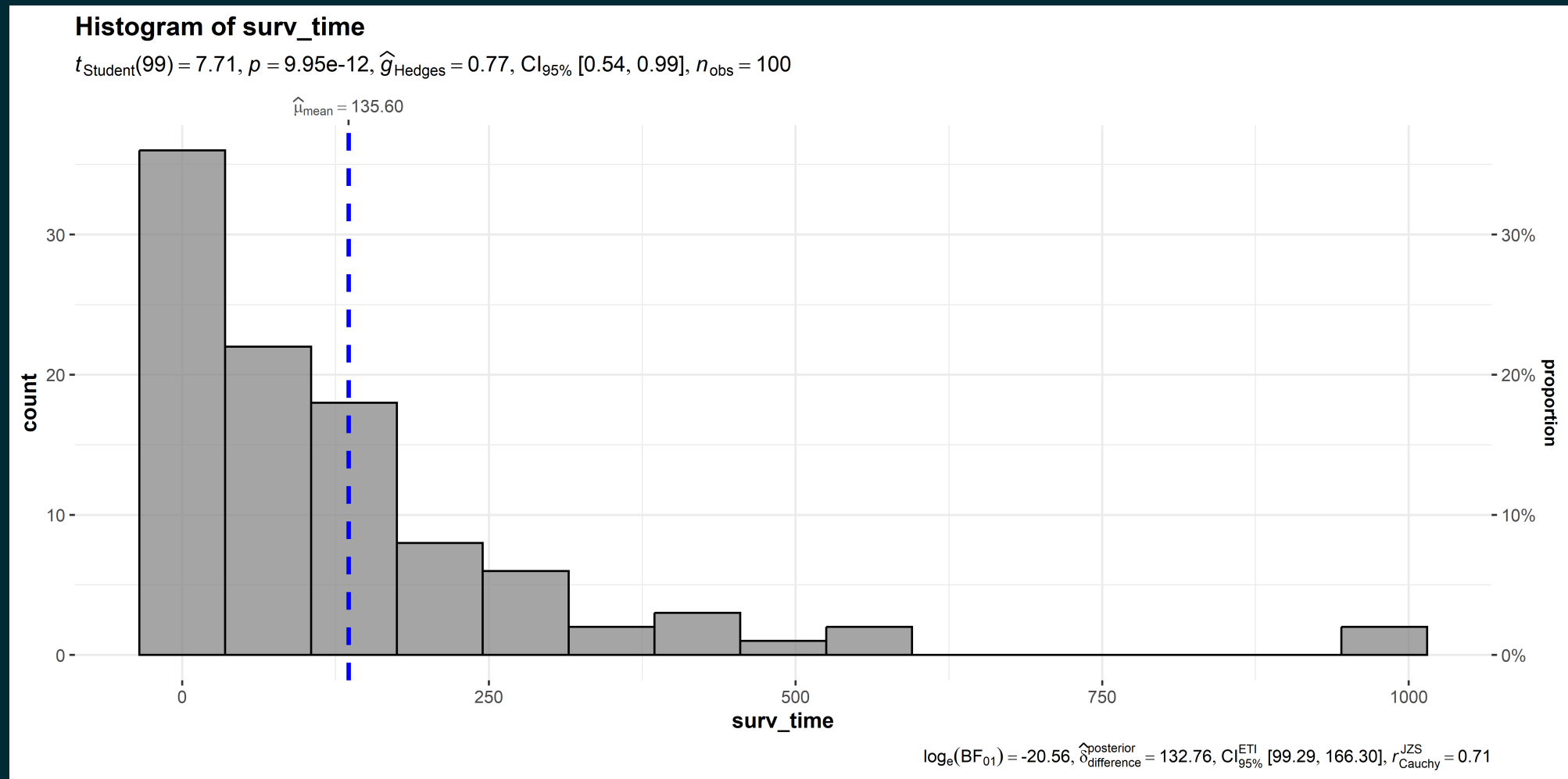
# CALCULATE DESCRIPTIVE STATISTICS FOR ALL NUMERIC VARIABLES WITHIN THIS DATASET?

### Descriptive Statistics
All numeric variables

| column | mean | median | sd | Q1 | Q3 | min | max |
|---|---|---|---|---|---|---|---|
| surv_time | 135.60 | 93.5 | 176 | 21.75 | 162.00 | 1 | 999 |
| kps | 58.65 | 60.0 | 20 | 40.00 | 80.00 | 20 | 90 |
| diag_time | 8.95 | 6.0 | 9 | 3.00 | 12.00 | 1 | 58 |
| age | 57.60 | 62.0 | 11 | 50.00 | 65.25 | 34 | 81 |
| prior | 3.10 | 0.0 | 5 | 0.00 | 10.00 | 0 | 10 |
| treatment | 0.31 | 0.0 | 0 | 0.00 | 1.00 | 0 | 1 |
| censor | 0.09 | 0.0 | 0 | 0.00 | 0.00 | 0 | 1 |
| event | 0.92 | 1.0 | 0 | 1.00 | 1.00 | 0 | 1 |

# A FOCUS ON SURVIVAL TIME

- We expect `surv_time` to follow an exponential distribution

```
1  ggstatsplot::gghistostats(data_cleaned, surv_time, binwidth = 70, title = "Histogram of surv_time")
```

**Histogram of surv_time**

$t_{Student}(99) = 7.71$, $p = 9.95e\text{-}12$, $\widehat{g}_{Hedges} = 0.77$, $CI_{95\%}$ [0.54, 0.99], $n_{obs} = 100$

$\widehat{\mu}_{mean} = 135.60$

$\log_e(BF_{01}) = -20.56$, $\widehat{\delta}_{difference}^{posterior} = 132.76$, $CI_{95\%}^{ETI}$ [99.29, 166.30], $r_{Cauchy}^{JZS} = 0.71$
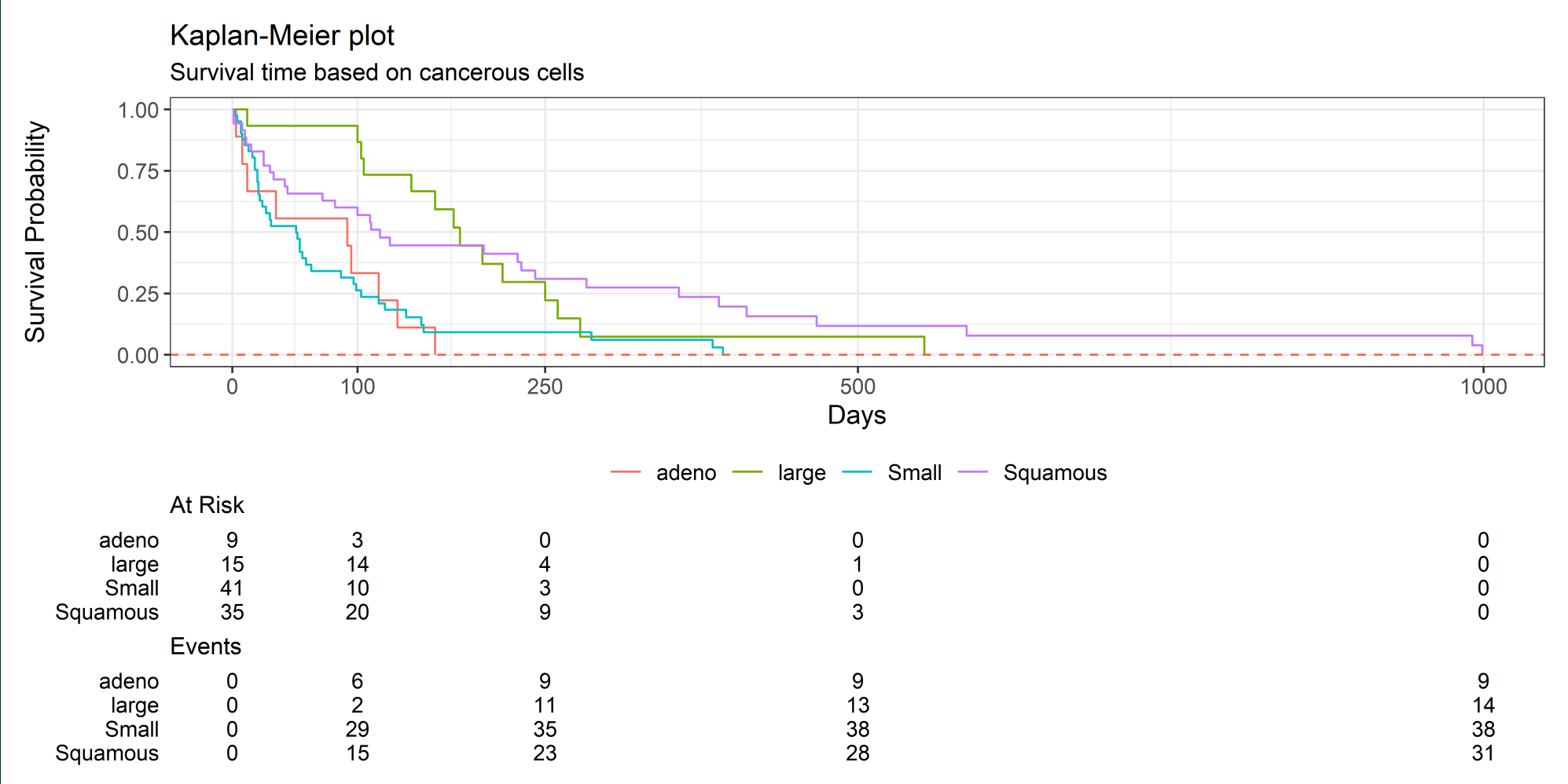
# SURVIVAL ANALYSIS

Perform a survival analysis to assess the survival time? based on the cancerous cells? Consider applying survival functions/kaplan meier quartiles/cumulative incidence function/cox regression etc.

# KAPLAN-MEIER PLOT

```r
1  survfit2(Surv(surv_time, event) ~ cell, data = data_cleaned) |>
2    ggsurvfit(type = "survival") +
3    labs(
4      title = "Kaplan-Meier plot",
5      subtitle = "Survival time based on cancerous cells",
6      x = "Days"
7    ) +
8    scale_x_continuous(breaks = c(0, 100, 250, 500, 1000)) +
9    add_risktable(times = c(0, 100, 250, 500, 1000)) +
10   geom_hline(yintercept = 0, color = "tomato", linetype = "dashed")
```

# KAPLAN-MEIER PLOT



Kaplan-Meier plot
Survival time based on cancerous cells

| | | adeno | large | Small | Squamous |

| At Risk | | | | |
|---|---|---|---|---|
| adeno | 9 | 3 | 0 | 0 | 0 |
| large | 15 | 14 | 4 | 1 | 0 |
| Small | 41 | 10 | 3 | 0 | 0 |
| Squamous | 35 | 20 | 9 | 3 | 0 |

| Events | | | | |
|---|---|---|---|---|
| adeno | 0 | 6 | 9 | 9 | 9 |
| large | 0 | 2 | 11 | 13 | 14 |
| Small | 0 | 29 | 35 | 38 | 38 |
| Squamous | 0 | 15 | 23 | 28 | 31 |

# COMPARISON OF SURVIVAL CURVES

- Calculate the Log Rank Test

```
1  survdiff(Surv(surv_time, event) ~ cell, data = data_cleaned)
```

```
Call:
survdiff(formula = Surv(surv_time, event) ~ cell, data = data_cleaned)

               N Observed Expected (O-E)^2/E (O-E)^2/V
cell=adeno     9        9      5.2      2.78      3.04
cell=large    15       14     19.8      1.70      2.26
cell=Small    41       38     23.8      8.54     12.44
cell=Squamous 35       31     43.2      3.46      7.27

 Chisq= 18.4  on 3 degrees of freedom, p= 4e-04
```

# COX REGRESSION

- **adeno** is the reference category, so the hazard ratios are relative to it

```r
1  cox_model <- coxph(Surv(surv_time, event) ~ cell, data = data_cleaned)
2  summary(cox_model)
```

```
Call:
coxph(formula = Surv(surv_time, event) ~ cell, data = data_cleaned)

  n= 100, number of events= 92


                coef exp(coef) se(coef)       z Pr(>|z|)
celllarge    -1.0003    0.3678   0.4358 -2.295  0.02172 *
cellSmall    -0.1062    0.8993   0.3741 -0.284  0.77653
cellSquamous -1.0385    0.3540   0.3990 -2.603  0.00925 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


             exp(coef) exp(-coef) lower .95 upper .95
celllarge       0.3678      2.719    0.1565    0.8640
cellSmall       0.8993      1.112    0.4320    1.8721
cellSquamous    0.3540      2.825    0.1619    0.7738


Concordance= 0.604  (se = 0.033 )
Likelihood ratio test= 17.32  on 3 df,    p=6e-04
Wald test         = 17.33  on 3 df,    p=6e-04
```
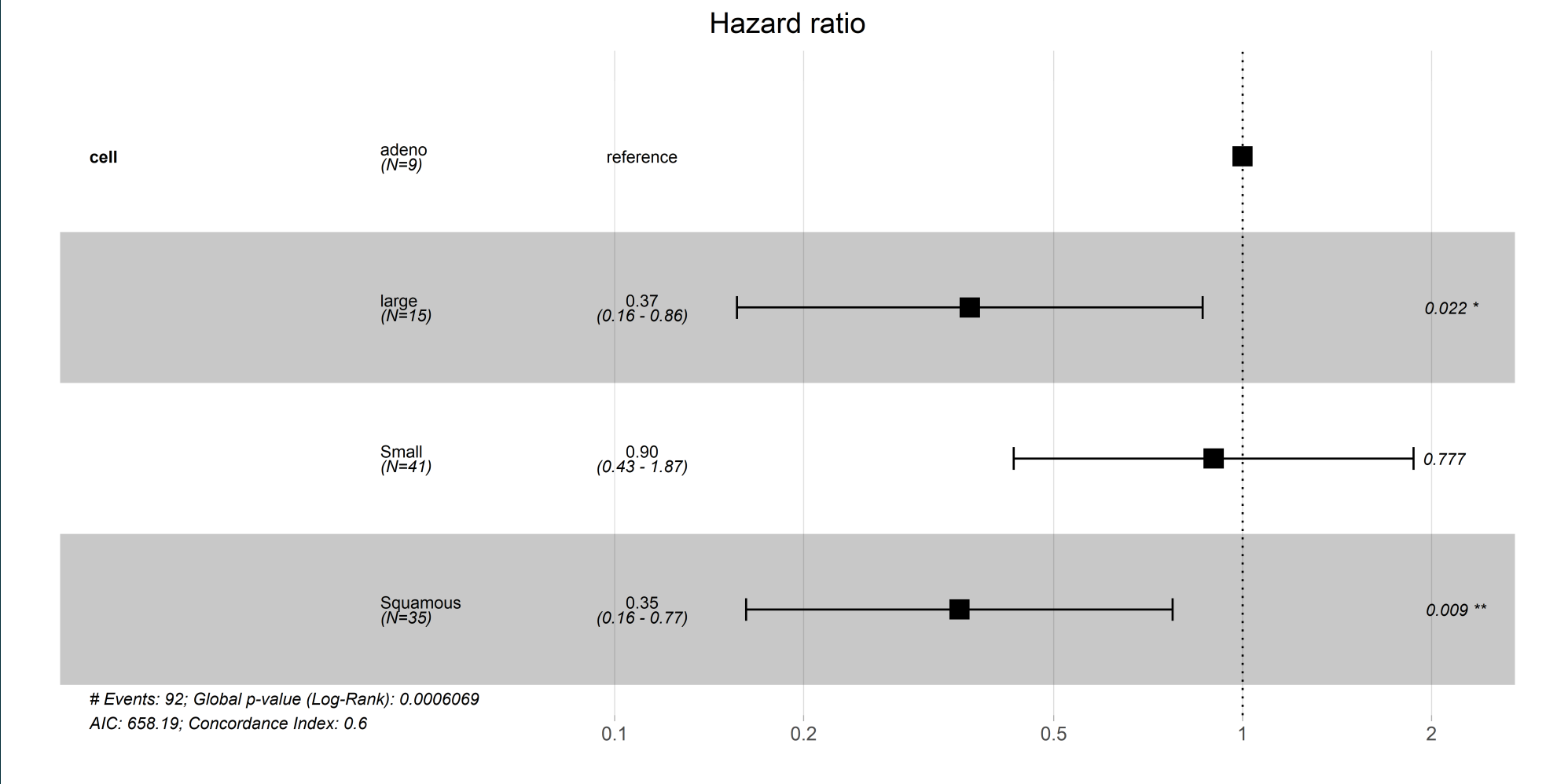
# VISUALIZE THE COEFFICIENTS

```
1  survminer::ggforest(cox_model)
```



Hazard ratio

| | | | |
|---|---|---|---|
| **cell** | adeno (N=9) | reference | |
| | large (N=15) | 0.37 (0.16 - 0.86) | 0.022 * |
| | Small (N=41) | 0.90 (0.43 - 1.87) | 0.777 |
| | Squamous (N=35) | 0.35 (0.16 - 0.77) | 0.009 ** |

# Events: 92; Global p-value (Log-Rank): 0.0006069
AIC: 658.19; Concordance Index: 0.6

# MULTIVARIABLE ANALYSIS

Perform an appropriate multivariable analysis to analyze the effect of independent variables age on the hazard ratio between the different cancerous cells (var Cell)?

# COX REGRESSION WITH `cell` AND `age`

```r
1  multivariable_cox_model <- coxph(Surv(surv_time, event) ~ cell + age, data = data_cleaned)
2
3  multivariable_cox_model |>
4    tbl_regression(exponentiate = T) |>
5    add_global_p() |>
6    add_n(location = "level") |>
7    add_nevent(location = "level") |>
8    bold_labels() |>
9    bold_p() |>
10   italicize_levels()
```

| Characteristic | N | Event N | HR | 95% CI | p-value |
|---|---|---|---|---|---|
| **cell** | | | | | **<0.001** |
| *adeno* | 9 | 9 | — | — | |
| *large* | 15 | 14 | 0.36 | 0.15, 0.85 | |
| *Small* | 41 | 38 | 0.85 | 0.40, 1.80 | |
| *Squamous* | 35 | 31 | 0.34 | 0.15, 0.74 | |
| **age** | 100 | 92 | 1.01 | 0.99, 1.03 | 0.5 |
| Abbreviations: CI = Confidence Interval, HR = Hazard Ratio | | | | | |

# DIAGNOSTIC OF THE MODEL

- Testing the proportional hazards assumption for the multivariable Cox regression model

```r
1  survminer::ggcoxzph(cox.zph(multivariable_cox_model))
```

# DEEPER ANALYSIS

```r
1  data_cleaned |>
2    select(surv_time, event, cell, age,
3           therapy, diag_time) |>
4    tbl_uvregression(
5      method = coxph,
6      y = Surv(surv_time, event),
7      exponentiate = T
8    ) |>
9    add_global_p() |>
10   add_n(location = "level") |>
11   add_nevent(location = "level") |>
12   bold_labels() |>
13   bold_p() |>
14   italicize_levels()
```

| Characteristic | N | Event N | HR | 95% CI | p-value |
|---|---|---|---|---|---|
| **cell** | | | | | **<0.001** |
| *adeno* | 9 | 9 | — | — | |
| *large* | 15 | 14 | 0.37 | 0.16, 0.86 | |
| *Small* | 41 | 38 | 0.90 | 0.43, 1.87 | |
| *Squamous* | 35 | 31 | 0.35 | 0.16, 0.77 | |
| **age** | 100 | 92 | 1.01 | 0.99, 1.03 | 0.4 |
| **therapy** | | | | | 0.2 |
| *Standard* | 69 | 64 | — | — | |
| *Test* | 31 | 28 | 0.72 | 0.45, 1.16 | |
| **diag_time** | 100 | 92 | 1.02 | 1.0, 1.05 | 0.14 |
| Abbreviations: CI = Confidence Interval, HR = Hazard Ratio | | | | | |

# JUST A CHECK

- kps (Key Performance Status) was not included because it's a measure of the patients that we observe during the therapy not before, as `age`, `diag_time` or `cell`

- If included, the model results are as follows:

```
1  coxph(Surv(surv_time, event) ~ kps, data = data_cleaned)
```

```
Call:
coxph(formula = Surv(surv_time, event) ~ kps, data = data_cleaned)

         coef exp(coef)  se(coef)     z        p
kps -0.029414  0.971014  0.006193 -4.75 2.04e-06

Likelihood ratio test=22.31  on 1 df, p=2.319e-06
n= 100, number of events= 92
```
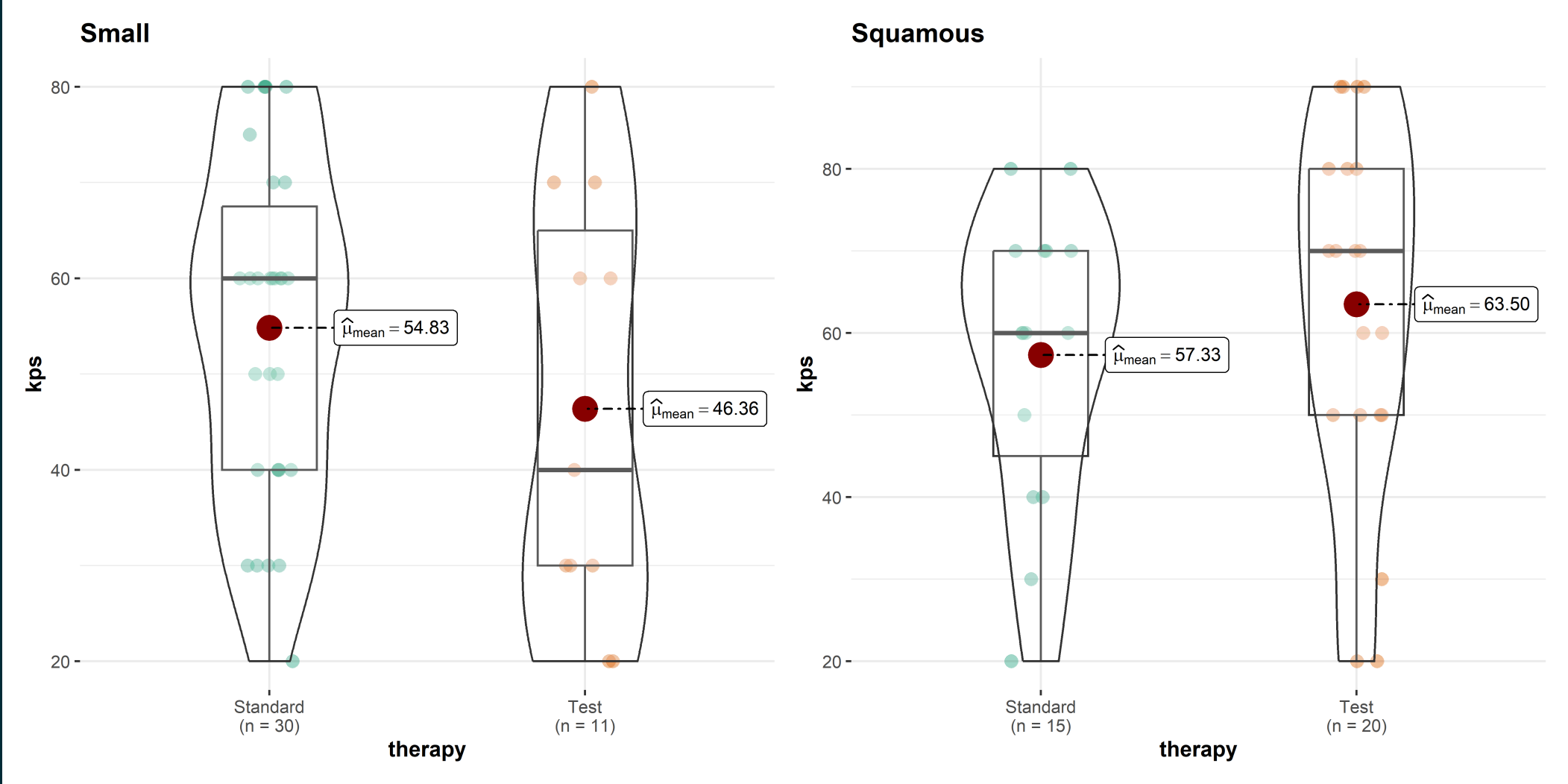
# WHAT ABOUT THE THERAPY PURPOSED?

```
1   data_cleaned |>
2     semi_join(data_cleaned |> filter(therapy == "Test"), by = "cell") |>
3     select(therapy, event, surv_time, cell, kps, diag_time) |>
4     tbl_strata(
5       strata = cell,
6       .tbl_fun =
7         \(x) x |>
8         tbl_summary(
9           by = therapy,
10          type = kps ~ "continuous"
11        ) |>
12        add_p()
13    )
```

| Characteristic | Small | | | Squamous | | |
|---|---|---|---|---|---|---|
| | Standard N = 30[1] | Test N = 11[1] | p-value[2] | Standard N = 15[1] | Test N = 20[1] | p-value[2] |
| event | 28 (93%) | 10 (91%) | >0.9 | 13 (87%) | 18 (90%) | >0.9 |
| surv_time | 53 (20, 122) | 21 (8, 87) | 0.080 | 100 (25, 144) | 157 (32, 373) | 0.3 |
| kps | 60 (40, 70) | 40 (30, 70) | 0.2 | 60 (40, 70) | 70 (50, 80) | 0.3 |
| diag_time | 4 (3, 11) | 4 (2, 11) | >0.9 | 9 (5, 11) | 7 (3, 13) | 0.5 |

[1] n (%); Median (Q1, Q3)

[2] Fisher's exact test; Wilcoxon rank sum test

# LET'S PLOT IT!

# THANKS