# Q1

## HTTP：

HTTP日志结构:
root
 |-- host: string (nullable = true)
 |-- id.orig_h: string (nullable = true)
 |-- id.orig_p: long (nullable = true)
 |-- id.resp_h: string (nullable = true)
 |-- id.resp_p: long (nullable = true)
 |-- method: string (nullable = true)
 |-- orig_filenames: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- orig_fuids: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- orig_mime_types: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- origin: string (nullable = true)
 |-- proxied: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- referrer: string (nullable = true)
 |-- request_body_len: long (nullable = true)
 |-- resp_filenames: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- resp_fuids: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- resp_mime_types: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- response_body_len: long (nullable = true)
 |-- status_code: long (nullable = true)
 |-- status_msg: string (nullable = true)
 |-- tags: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- trans_depth: long (nullable = true)
 |-- ts: timestamp (nullable = true)
 |-- uid: string (nullable = true)
 |-- uri: string (nullable = true)
 |-- user_agent: string (nullable = true)
 |-- username: string (nullable = true)
 |-- version: string (nullable = true)

## DNS：

DNS日志结构:
root
 |-- AA: boolean (nullable = true)
 |-- RA: boolean (nullable = true)
 |-- RD: boolean (nullable = true)
 |-- TC: boolean (nullable = true)
 |-- TTLs: array (nullable = true)
 |    |-- element: double (containsNull = true)
 |-- Z: long (nullable = true)
 |-- answers: array (nullable = true)
 |    |-- element: string (containsNull = true)
 |-- id.orig_h: string (nullable = true)
 |-- id.orig_p: long (nullable = true)
 |-- id.resp_h: string (nullable = true)
 |-- id.resp_p: long (nullable = true)
 |-- proto: string (nullable = true)
 |-- qclass: long (nullable = true)
 |-- qclass_name: string (nullable = true)
 |-- qtype: long (nullable = true)
 |-- qtype_name: string (nullable = true)
 |-- query: string (nullable = true)
 |-- rcode: long (nullable = true)
 |-- rcode_name: string (nullable = true)
 |-- rejected: boolean (nullable = true)
 |-- rtt: double (nullable = true)
 |-- trans_id: long (nullable = true)
 |-- ts: timestamp (nullable = true)
 |-- uid: string (nullable = true)

## 使用 spark SQL API

```
+-------------------+------------+
|                uri|access_count|
+-------------------+------------+
|                  /|        9475|
|/admin/config.php...|        556|
|    /main.php?logout=1|      194|
|/top.php?stuff=15...|        191|
|            /top.php|        179|
|/main.php?stuff=1...|        172|
|    /get_latest_id.php|      159|
|/admin/config.php...|        138|
|    /cacti/index.php|        129|
|/en-US/api/messag...|        118|
|          /index.php|        105|
|/phpmyadmin/index...|         77|
|             /cacti/|         68|
|         /phpmyadmin/|        56|
|         /favicon.ico|        55|
|              /admin|         42|
|   /scripts/index.php|         40|
|             /icons/|         39|
|/en-US/api/search...|         39|
|    /cgi-bin/index.php|        37|
+-------------------+------------+
only showing top 20 rows
```

## 使用 spark DataFram

```
+----------------------------------------------------------------------------------+-----+
|uri                                                                               |count|
+----------------------------------------------------------------------------------+-----+
|/                                                                                 |9475 |
|/admin/config.php?type=tool&display=index&quietmode=1&info=stats&restrictmods=core/dashboard|556  |
|/main.php?logout=1                                                                |194  |
|/top.php?stuff=1583574484                                                         |191  |
|/top.php                                                                          |179  |
|/main.php?stuff=1583574484                                                        |172  |
|/get_latest_id.php                                                                |159  |
|/admin/config.php?type=tool&display=index&quietmode=1&info=info&restrictmods=core/dashboard|138  |
|/cacti/index.php                                                                  |129  |
|/en-US/api/messages/index                                                         |118  |
|/index.php                                                                        |105  |
|/phpmyadmin/index.php                                                             |77   |
|/cacti/                                                                           |68   |
|/phpmyadmin/                                                                      |56   |
|/favicon.ico                                                                      |55   |
|/admin                                                                            |42   |
|/scripts/index.php                                                                |40   |
|/icons/                                                                           |39   |
|/en-US/api/search/jobs?s=1331892438.21&s=1331892438.24                            |39   |
|/cgi-bin/index.php                                                                |37   |
+----------------------------------------------------------------------------------+-----+
only showing top 20 rows
```

## 合并并计算验证

```
[Stage 12:===============================>
+---+--------------+-------------+
|uri|tcp_percentage|request_count|
+---+--------------+-------------+
+---+--------------+-------------+
```

验证统计：

```
[Stage 17:===============================>
+---------------------+-----------+
|total_matched_requests|unique_uris|
+---------------------+-----------+
|                    0|          0|
+---------------------+-----------+
```

Q2

实现参考：

https://notebooks.databricks.com/notebooks/CME/Survival_Analysis/index.html#Survival_Analysis_1.html

首先从以下链接下载 csv 文件至服务器

raw.githubusercontent.com/IBM/telco-customer-churn-on-icp4d/master/data/Telco-Customer-Churn.csv

该 csv 数据文件来自 IBM，旨在模拟一个虚构的电信公司的用户数据。每一行数据都代表着一个电信公司的订阅用户的个人信息，包括但不仅限于各自人口统计、服务计划、媒体使用情况、订阅状态以及在网时长和是否已经流失的信息。

我们要对这组数据进行生存分析，最重要的内容即是在网时长(Tenure)和是否流失(Churn)，从对这两组信息的分析，我们可以后续进行估计判断客户留存与在网时长的关系或某时间点内客户未流失的概率，进而对客户的市场需求和改良电信产品做出应对措施等。

我们会根据用户的订阅模式，将用户数据分为 bronze 和 silver 两部分，主要分析 silver 高级用户的留存率，便于指定针对化的措施。

具体执行方案：

前期配置，创建 spark、文件路径和表名等。然后从 csv 数据文件的表头获取 schema。

首先将全体数据都存在 bronze 中，进行预存储：

```
+----------+------+-------------+-------+----------+------+------------+----------------+----------------+----------------+----------------+----------------+-------------------+----------------+----------------+----------------+------------+------------+-----+
|customerID|gender|seniorCitizen|partner|dependents|tenure|phoneService|   multipleLines|internetService|     onlineSecurity|        onlineBackup|   deviceProtection|      techSupport|     streamingTV|    streamingMovies|    contract|paperlessBilling|     paymentMethod|monthlyCharges|totalCharges|Churn|
+----------+------+-------------+-------+----------+------+------------+----------------+----------------+----------------+----------------+----------------+-------------------+----------------+----------------+----------------+------------+------------+-----+
|7590-VHVEG|Female|          0.0|    Yes|        No|   1.0|          No|No phone service|             DSL|             No|            Yes|             No|                 No|              No|              No|Month-to-month|         Yes|    Electronic check|         29.85|       29.85|   No|
|5575-GNVDE|  Male|          0.0|     No|        No|  34.0|         Yes|              No|             DSL|            Yes|             No|            Yes|                 No|              No|              No|      One year|          No|        Mailed check|         56.95|      1889.5|   No|
|3668-QPYBK|  Male|          0.0|     No|        No|   2.0|         Yes|              No|             DSL|            Yes|            Yes|             No|                 No|              No|              No|Month-to-month|         Yes|        Mailed check|         53.85|      108.15|  Yes|
|7795-CFOCW|  Male|          0.0|     No|        No|  45.0|          No|No phone service|             DSL|            Yes|             No|            Yes|                Yes|              No|              No|      One year|          No|Bank transfer (au...|          42.3|     1840.75|   No|
|9237-HQITU|Female|          0.0|     No|        No|   2.0|         Yes|              No|     Fiber optic|             No|             No|             No|                 No|              No|              No|Month-to-month|         Yes|    Electronic check|          70.7|      151.65|  Yes|
|9305-CDSKC|Female|          0.0|     No|        No|   8.0|         Yes|             Yes|     Fiber optic|             No|             No|            Yes|                Yes|             Yes|             Yes|Month-to-month|         Yes|    Electronic check|         99.65|       820.5|  Yes|
|1452-KIOVK|  Male|          0.0|     No|       Yes|  22.0|         Yes|             Yes|     Fiber optic|             No|            Yes|             No|                 No|             Yes|              No|Month-to-month|         Yes|Credit card (auto...|          89.1|      1949.4|   No|
|6713-OKOMC|Female|          0.0|     No|        No|  10.0|          No|No phone service|             DSL|            Yes|             No|             No|                 No|              No|              No|Month-to-month|          No|        Mailed check|
```

接着对于 silver 用户，筛选所有拥有月度订阅（month-to-month）的在网用户：

```
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+--------------+------------+------
----------+------------+-----------+------------------+------------------+------------+-------------+------------+------
|customerID|gender|seniorCitizen|partner|dependents|tenure|phoneService|   multipleLines|internetService|onlineSecurity|onlineBackup|device
Protection|techSupport|streamingTV|streamingMovies|        contract|paperlessBilling|       paymentMethod|monthlyCharges|totalCharges|
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+--------------+------------+------
----------+------------+-----------+------------------+------------------+------------+-------------+------------+------
|7590-VHVEG|Female|          0.0|    Yes|        No|   1.0|          No|No phone service|            DSL|            No|            Yes|
No|          No|         No|             No|    Month-to-month|             Yes|    Electronic check|         29.85|       29.85|
|3668-QPYBK|  Male|          0.0|     No|        No|   2.0|         Yes|              No|            DSL|           Yes|            Yes|
No|          No|         No|             No|    Month-to-month|             Yes|        Mailed check|         53.85|      108.15|
|9237-HQITU|Female|          0.0|     No|        No|   2.0|         Yes|              No|    Fiber optic|            No|             No|
No|          No|         No|             No|    Month-to-month|             Yes|    Electronic check|          70.7|      151.65|
|9305-CDSKC|Female|          0.0|     No|        No|   8.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         99.65|       820.5|
|1452-KIOVK|  Male|          0.0|     No|       Yes|  22.0|         Yes|             Yes|    Fiber optic|            No|             No|
No|          No|        Yes|             No|    Month-to-month|             Yes|Credit card (auto...|          89.1|      1949.4|
|6713-OKOMC|Female|          0.0|     No|        No|  10.0|          No|No phone service|            DSL|            No|             No|
No|          No|         No|             No|    Month-to-month|              No|        Mailed check|         29.75|       301.9|
|7892-POOKP|Female|          0.0|    Yes|        No|  28.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         104.8|     3046.05|
|9763-GRSKD|  Male|          0.0|    Yes|       Yes|  13.0|         Yes|              No|            DSL|            No|             No|
No|          No|         No|             No|    Month-to-month|             Yes|        Mailed check|         49.95|      587.45|
|0280-XJGEX|  Male|          0.0|     No|        No|  49.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|Bank transfer (au...|         103.7|      5036.3|
|5129-JLPIS|  Male|          0.0|     No|        No|  25.0|         Yes|              No|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         105.5|     2686.05|
|4190-MFLUW|Female|          0.0|    Yes|       Yes|  10.0|         Yes|              No|            DSL|            No|             No|
Yes|          Yes|         No|             No|    Month-to-month|              No|Credit card (auto...|          55.2|      528.35|
|4183-MYFRB|Female|          0.0|     No|        No|  21.0|         Yes|             Yes|    Fiber optic|            No|             No|
No|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         90.05|      1862.9|
|8779-QRDMV|  Male|          1.0|     No|        No|   1.0|          No|No phone service|            DSL|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         39.65|       39.65|
|6322-HRPFA|  Male|          0.0|    Yes|       Yes|  49.0|         Yes|              No|            DSL|           Yes|             Yes|
No|          Yes|         No|             No|    Month-to-month|              No|Credit card (auto...|          59.6|      2970.3|
```

对旧数据进行删除并创建新数据库，分别将 bronze 和 silver 用户数据写入。
通过 spark SQL API 显示用户数据：

```
Bronze Customers Data:
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+---------------+------------
----+------------+------------+----------------+------------------+------------+------------------+--------------+
--------------+------------+-----+
|customerID|gender|seniorCitizen|partner|dependents|tenure|phoneService|   multipleLines|internetService|   onlineSecurity|      onlineB
ackup|   deviceProtection|    techSupport|       streamingTV|   streamingMovies|        contract|paperlessBilling|       paymentMethod|
monthlyCharges|totalCharges|Churn|
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+---------------+------------
----+------------+------------+----------------+------------------+------------+------------------+--------------+
--------------+------------+-----+
|7590-VHVEG|Female|          0.0|    Yes|        No|   1.0|          No|No phone service|            DSL|             No|
Yes|          No|         No|             No|No|    Month-to-month|             Yes|    Electronic check|
29.85|       29.85|   No|
|5575-GNVDE|  Male|          0.0|     No|        No|  34.0|         Yes|              No|            DSL|            Yes|
No|          Yes|         No|             No|No|          One year|              No|        Mailed check|
56.95|      1889.5|   No|
|3668-QPYBK|  Male|          0.0|     No|        No|   2.0|         Yes|              No|            DSL|            Yes|
Yes|          No|         No|             No|No|    Month-to-month|             Yes|        Mailed check|
53.85|      108.15|  Yes|
|7795-CFOCW|  Male|          0.0|     No|        No|  45.0|          No|No phone service|            DSL|            Yes|
No|          Yes|        Yes|             No|No|          One year|              No|Bank transfer (au...|
42.3|     1840.75|   No|
|9237-HQITU|Female|          0.0|     No|        No|   2.0|         Yes|              No|    Fiber optic|             No|
No|          No|         No|             No|No|    Month-to-month|             Yes|    Electronic check|
70.7|      151.65|  Yes|
|9305-CDSKC|Female|          0.0|     No|        No|   8.0|         Yes|             Yes|    Fiber optic|             No|
No|          Yes|         No|            Yes|Yes|    Month-to-month|             Yes|    Electronic check|
99.65|       820.5|  Yes|
|1452-KIOVK|  Male|          0.0|     No|       Yes|  22.0|         Yes|             Yes|    Fiber optic|             No|
Yes|          No|         No|            Yes|No|    Month-to-month|             Yes|Credit card (auto...|
89.1|      1949.4|   No|
|6713-OKOMC|Female|          0.0|     No|        No|  10.0|          No|No phone service|            DSL|            Yes|
No|          No|         No|             No|No|    Month-to-month|              No|        Mailed check|
29.75|       301.9|   No|
```
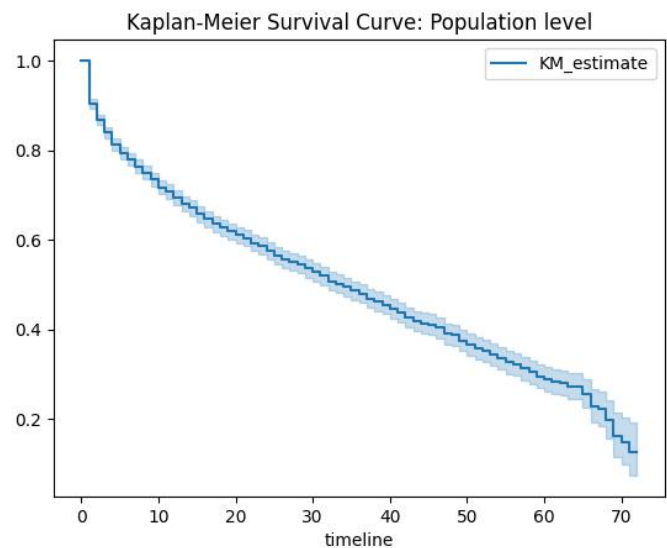
```
Silver Monthly Customers Data:
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+--------------+------------+------
----------+------------+-----------+------------------+------------------+------------+-------------+------------+-----+
|customerID|gender|seniorCitizen|partner|dependents|tenure|phoneService|   multipleLines|internetService|onlineSecurity|onlineBackup|device
Protection|techSupport|streamingTV|streamingMovies|        contract|paperlessBilling|       paymentMethod|monthlyCharges|totalCharges|churn|
+----------+------+-------------+-------+----------+------+------------+----------------+---------------+--------------+------------+------
----------+------------+-----------+------------------+------------------+------------+-------------+------------+-----+
|7590-VHVEG|Female|          0.0|    Yes|        No|   1.0|          No|No phone service|            DSL|            No|            Yes|
No|          No|         No|             No|    Month-to-month|             Yes|    Electronic check|         29.85|       29.85|    0|
|3668-QPYBK|  Male|          0.0|     No|        No|   2.0|         Yes|              No|            DSL|           Yes|            Yes|
No|          No|         No|             No|    Month-to-month|             Yes|        Mailed check|         53.85|      108.15|    1|
|9237-HQITU|Female|          0.0|     No|        No|   2.0|         Yes|              No|    Fiber optic|            No|             No|
No|          No|         No|             No|    Month-to-month|             Yes|    Electronic check|          70.7|      151.65|    1|
|9305-CDSKC|Female|          0.0|     No|        No|   8.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         99.65|       820.5|    1|
|1452-KIOVK|  Male|          0.0|     No|       Yes|  22.0|         Yes|             Yes|    Fiber optic|            No|             No|
No|          No|        Yes|             No|    Month-to-month|             Yes|Credit card (auto...|          89.1|      1949.4|    0|
|6713-OKOMC|Female|          0.0|     No|        No|  10.0|          No|No phone service|            DSL|            No|             No|
No|          No|         No|             No|    Month-to-month|              No|        Mailed check|         29.75|       301.9|    0|
|7892-POOKP|Female|          0.0|    Yes|        No|  28.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          Yes|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         104.8|     3046.05|    1|
|9763-GRSKD|  Male|          0.0|    Yes|       Yes|  13.0|         Yes|              No|            DSL|            No|             No|
No|          No|         No|             No|    Month-to-month|             Yes|        Mailed check|         49.95|      587.45|    0|
|0280-XJGEX|  Male|          0.0|     No|        No|  49.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          Yes|        Yes|            Yes|    Month-to-month|             Yes|Bank transfer (au...|         103.7|      5036.3|    1|
|5129-JLPIS|  Male|          0.0|     No|        No|  25.0|         Yes|              No|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         105.5|     2686.05|    0|
|4190-MFLUW|Female|          0.0|    Yes|       Yes|  10.0|         Yes|              No|            DSL|            No|             No|
Yes|          Yes|         No|             No|    Month-to-month|              No|Credit card (auto...|          55.2|      528.35|    1|
|4183-MYFRB|Female|          0.0|     No|        No|  21.0|         Yes|             Yes|    Fiber optic|            No|             No|
Yes|          No|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         90.05|      1862.9|    0|
|8779-QRDMV|  Male|          1.0|     No|        No|   1.0|          No|No phone service|            DSL|            No|             No|
Yes|          Yes|        Yes|            Yes|    Month-to-month|             Yes|    Electronic check|         39.65|       39.65|    1|
|6322-HRPFA|  Male|          0.0|    Yes|       Yes|  49.0|         Yes|              No|            DSL|            No|             No|
No|          Yes|         No|             No|    Month-to-month|              No|Credit card (auto...|          59.6|      2970.3|    0|
```

以上就完成了模拟对电信公司的用户数据进行简单的预处理

接下来将采用 Kaplan-Meier 曲线来对数据进行生存分析。生存曲线以留存时间为横轴，其他特定数据指标为纵轴，绘制成连续型的阶梯形曲线，用以说明生存时间与生存率之间的关系。

从 silver 数据内获取 tenure 和 churn 并拟合 KM 模型。

`<lifelines.KaplanMeierFitter:"KM_estimate", fitted with 3351 total observations, 1795 right-censored observations>`
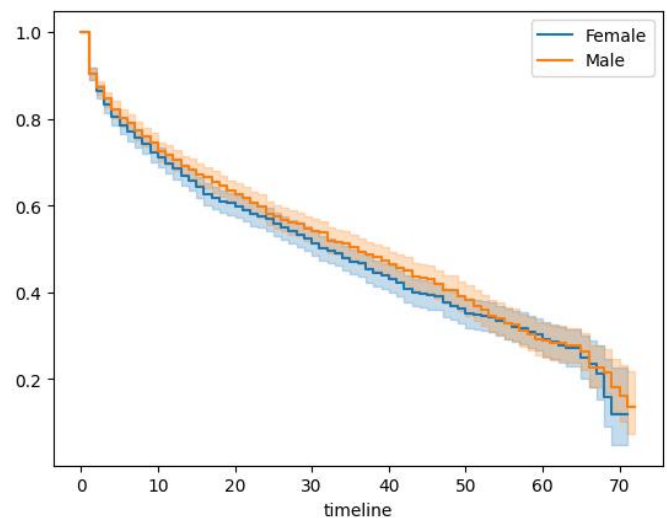
首先绘制曲线来观察整体的生存率：

`<Axes: title={'center': 'Kaplan-Meier Survival Curve: Population level'}, xlabel='timeline'>`
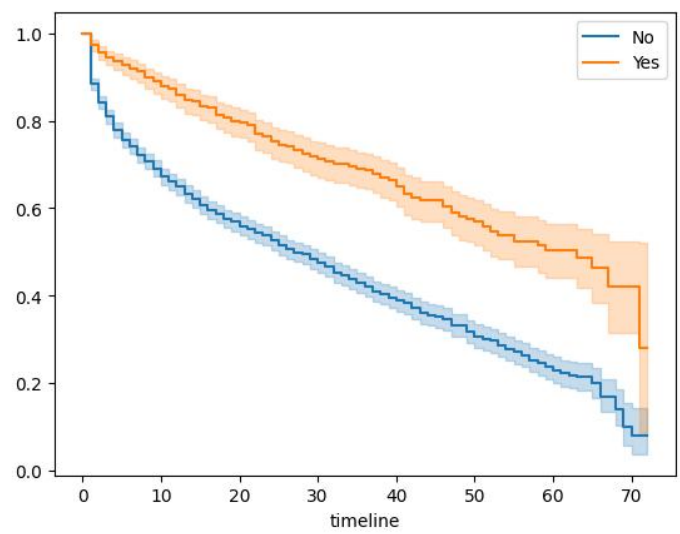


观察留存时间中位数：

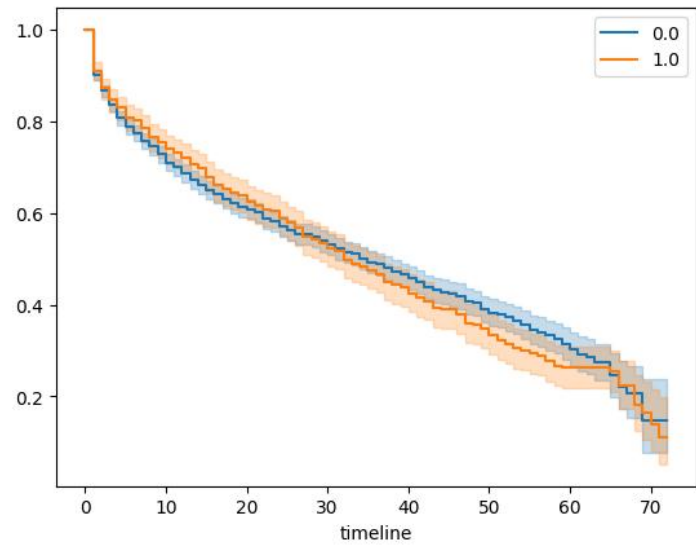`np.float64(34.0)`

定义函数绘制在协变量水平上的 Kaplan-Meier curve 和打印 Log-rank 测试结果有关性别和生存时间的关系：



| | test_statistic | p | -log2(p) |
|---|---|---|---|
| **Female  Male** | 2.038938 | 0.153317 | 2.705414 |

有关网络安全和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| **No** | **Yes** | 141.60316 | 1.187554e-32 | 106.053706 |

有关高龄人士和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| **0.0** | **1.0** | 0.125471 | 0.723174 | 0.467584 |

有关 partner 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 135.758896 | 2.252911e-31 | 101.807981 |

有关 dependents 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 35.031241 | 3.244576e-09 | 28.199323 |

有关 phoneService 和生存时间的关系：



| | test_statistic | p | -log2(p) |
|---|---|---|---|
| **No Yes** | 1.683709 | 0.194432 | 2.36266 |

有关 multipleLines 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| **No** | **No phone service** | 12.382712 | 4.333273e-04 | 11.172255 |
| | **Yes** | 72.358368 | 1.794602e-17 | 55.629114 |
| **No phone service** | **Yes** | 1.500291 | 2.206266e-01 | 2.180322 |

有关 internetService 和生存时间的关系：



| | test_statistic | p | -log2(p) |
|---|---|---|---|
| **DSL   Fiber optic** | 25.172866 | 5.241449e-07 | 20.863531 |

有关 streamingTV 和生存时间的关系：



| | test_statistic | p | -log2(p) |
|---|---|---|---|
| **No   Yes** | 12.93926 | 0.000322 | 11.601718 |

有关 streamingMovies 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 17.941685 | 0.000023 | 15.422016 |

有关 onlineBackup 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 189.482865 | 4.122979e-43 | 140.799221 |

有关 deviceProtection 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 71.496825 | 2.777047e-17 | 54.999226 |

有关 techSupport 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| No | Yes | 90.430334 | 1.916059e-21 | 68.822348 |

有关 paperlessBilling 和生存时间的关系：



| | test_statistic | p | -log2(p) |
|---|---|---|---|
| No Yes | 8.340802 | 0.003876 | 8.011049 |

有关 paymentMethod 和生存时间的关系：



| | | test_statistic | p | -log2(p) |
|---|---|---|---|---|
| Bank transfer (automatic) | Credit card (automatic) | 0.061543 | 8.040732e-01 | 0.314601 |
| | Electronic check | 91.191889 | 1.303937e-21 | 69.377616 |
| | Mailed check | 43.536998 | 4.160192e-11 | 34.484559 |
| Credit card (automatic) | Electronic check | 79.991082 | 3.761035e-19 | 61.205504 |
| | Mailed check | 39.684613 | 2.984678e-10 | 31.641706 |
| Electronic check | Mailed check | 0.898320 | 3.432326e-01 | 1.542741 |

有关 paperlessBilling 和生存时间的关系：

在完成分析之后，可以提取生存概率以用于其他程序进行预测分析等。
定义一个函数用于获取特定生存率（以 DSL 为例）：

| | DSL |
|---|---|
| 0 | 1.000000 |
| 1 | 0.902698 |
| 2 | 0.864380 |
| 3 | 0.834702 |
| 4 | 0.810522 |
| 5 | 0.794352 |
| 6 | 0.783900 |
| 7 | 0.776362 |
| 8 | 0.768486 |
| 9 | 0.750833 |