

Tidyverse を活用した データ前処理の実践

静岡市立清水病院 呼吸器内科

森 和貴

MORI Kazutaka, M.D., Ph.D.

- 医師 21年め（臨床研修制度 1期生）

- 専門：呼吸器内科

- R遍歴：15年くらい

- 2009年頃～

- 2011年 RStudio が登場、飛びつく

- 2012年 大学院の後輩から EZR を知る

- 2017年頃～ Windows版 RStudioでの日本語対応に疲れ、仮想化に走る

- ✓ そこから、解析の再現可能性に関心を持つようになりました

- ✓ 2023年頃からRStudio Desktop版での日本語対応が改善し、Desktop版もまた使うようになりました

best subset selection procedure using Akaike's information criteria (AIC). All statistical analyses were performed using R version 2.11.1 (The R Foundation for Statistical Computing, Vienna, Austria, 2010). A value of $p < 0.05$ was considered significant for the results of all statistical analyses, and all tests were 2-sided.

はじめての論文
Mori K, *et al.* COPD. 2011

- 集計や分析に用いる生データを整えて加工すること全般をさす
- 「分析にかかる時間の8割は前処理の時間」といわれる
- 主な工程：
 - 収集した生データを **コンピュータ／Rで処理できる状態** に整形する
 - データクリーニング(クレンジング)
 - ✓ 欠損値、外れ値などの確認・処理
 - データ変換
 - ✓ データの「型」を整える (数値、順序変数、あり・なしの2値、など)
 - ✓ データの正規化・標準化 (表記のブレを揃える、など)
 - ✓ 匿名化 (年齢などを生の値→階級値に変換する操作も含む)

- Tidyverse とパイプ演算子のおさらい
- データ整形の指針
 - ～本当はデータファイルを作る前に確認しておきたいこと～
- 初級編：最新の癌種別死亡数のグラフを作りたい
 - ～人口動態統計のCSVファイルを整形してみる～
- 中級編：臨床研究のまとめファイルがやってきた
 - ～“ネ申エクセル” と戦うための基礎知識～

- Tidyverse とパイプ演算子のおさらい
- データ整形の指針
 - ～本当はデータファイルを作る前に確認しておきたいこと～
- 初級編：最新の癌種別死亡数のグラフを作りたい
 - ～人口動態統計のCSVファイルを整形してみる～
- 中級編：臨床研究のまとめファイルがやってきた
 - ～“ネ申エクセル” と戦うための基礎知識～

- Posit (旧 RStudio) 社の Hadley Wickham 氏らを中心に開発されているデータ処理のための  パッケージ群

```
> library(tidyverse)
— Attaching core tidyverse packages — tidyverse 2.0.0 —
✓ dplyr      1.1.4      ✓ readr      2.1.5
✓ forcats    1.0.0      ✓ stringr    1.5.1
✓ ggplot2     3.5.1      ✓ tibble     3.2.1
✓ lubridate  1.9.3      ✓ tidyr      1.3.1
✓ purrr       1.0.2
```

This package is designed to make it easy to install and load multiple 'tidyverse' packages in a single step.

- `library(tidyverse)` で 9つの “core” パッケージ が読み込まれる。
更に、インストールされているがこの時点では読み込まれていない多数の関連パッケージも控えている状態

Tidyverse の構成パッケージ

7

Core packages



| パッケージ | 用途 |
|-----------|--------------------|
| ggplot2 | データ可視化 |
| dplyr | データの加工、整理 |
| tidyr | データフレームの整形 |
| readr | データの読み書き |
| purrr | 繰り返し処理 |
| tibble | 改良版のデータフレーム |
| stringr | 文字列の処理 |
| forcats | factor型変数(順序尺度)の処理 |
| lubridate | 日時データの処理 |



readxl
Excelファイルの読み込み



googlesheets4
Googleスプレッドシートの読み込み



haven
SPSS, Stata, SAS のデータファイルの読み込み



magrittr
特殊なパイプ

⋮

パイプの基本: 左側を右側の関数の1つ目のパラメーター(引数)に挿入

data %>% **function(param)** ⇒ **function(data, param)**

data %>% **f(p1)** %>% **g(p2)** ⇒ **g(f(data, p1), p2)**

```
R 4.3.3 · C:/Home/GitLocal/mJ_Rpeer_2025Mar/
> f <- function(x, y = 1, z, ...) {
+   if (!is.numeric(y)) message("ERROR: parameter 'y' must be numeric.")
+   match.call()
+ }
> data <- data.frame()
>
> # 引数名を指定せずパイプで受け渡す場合
> data |> f("param", 12)
ERROR: parameter 'y' must be numeric.
f(x = data, y = "param", z = 12)
>
> # ズレないためには引数名を指定する
> data |> f(y = 123)
f(x = data, y = 123)
> # パイプの左側を反映する変数を指定する場合は placeholder を使う
> data |> f("param", z = _)
f(x = "param", z = data)
```

‘y’ のつもりが...

右側の関数に元々指定されていた引数がひとつずつ後ろにずれることに注意。

エラーの原因となりやすいので、特にパイプを使うときは**できるだけ引数名を省略しない**ようにした方が安全。

2つのパイプ演算子

- 現在、主に使われているパイプ演算子は2種類存在する
- 基本的な挙動はほぼ同様なので、他の用途で tidyverse を読み込んでいるかやコーディング規約などで選択、読み替えて使用可能

| Magritter pipe | 名称 | Base pipe |
|---|-----------------------------------|---|
| <code>%>%</code> | 使用する記号 | <code> ></code> |
| 必要 (tidyverse / magritter) | パッケージ読み込み | 不要 (R 4.1以降) |
| 関数以外も可 (例: <code>x %>% .[1]</code>) | パイプの右側 | 関数のみ |
| <code>.</code> (複数可、引数名の省略可) | 第1引数以外 への受け渡し (placeholder) | <code>param = _</code> (1個だけ、引数名の省略不可) |
| 可能 (<code>%>% return()</code> でパイプを終了) | <code>return()</code> の使用 | <code> > return()</code> は使用できない |

Tee pipe : %T>%

- 2つの関数にデータを受け渡す
- 途中経過の確認に有効

```
R 4.3.3 ~ /  
> library(tidyverse)  
> library(magrittr)  
>  
> tibble(var1 = rnorm(3), var2 = rnorm(3)) %T>%  
+   print() %>%  
+   summarise(var1 = sum(var1), var2 = sum(var2))  
# A tibble: 3 × 2  
  var1    var2  
  <dbl> <dbl>  
1  0.0891 -0.377  
2  1.20   -0.485  
3 -0.558 -0.771  
# A tibble: 1 × 2  
  var1    var2  
  <dbl> <dbl>  
1  0.731 -1.63
```

tibble(...) %>%
print() の結果

tibble(...) %>%
summarise(...) の結果

Exposition pipe : %\$%

- 目的の関数にデータフレームを受け取る機能がない場合に有効

```
R 4.3.3 ~ /  
> diamonds %$% cor.test(price, carat)  
  
Pearson's product-moment correlation  
  
data: price and carat  
t = 551.41, df = 53938, p-value < 2.2e-16  
alternative hypothesis: true correlation is not equal to 0  
95 percent confidence interval:  
 0.9203098 0.9228530  
sample estimates:  
      cor  
0.9215913
```

cor.test() 関数には data= がないので、通常は
cor.test(diamonds\$price, diamonds\$carat)

- Posit社提供の **cheat sheet 集**

<https://posit.co/resources/cheatsheets/>

または、RStudio で Help – Cheat Sheets (上記の一部)

- **tidylog パッケージ**

- dplyr, tidyr の主要な関数に簡単な処理結果の報告を追加
- library(tidylog) でロードせず、`tidylog::function()` で確認後に必要に応じて `dplyr::function()` に戻す方が使い勝手が良い(私見)

```
R 4.3.3 ~|
> library(tidyverse)
> library(magrittr)
> diamonds %>%
+   tidylog::filter(cut == "Premium") %>%
+   head(5)
filter: removed 40,149 rows (74%), 13,791 rows remaining
# A tibble: 5 × 10
  carat cut      color clarity depth table price      x      y      z
  <dbl> <ord>    <ord> <ord>    <dbl> <dbl> <int> <dbl> <dbl> <dbl>
1  0.21 Premium E      SI1     59.8   61    326  3.89  3.84  2.31
2  0.29 Premium I      VS2     62.4   58    334  4.2   4.23  2.63
3  0.22 Premium F      SI1     60.4   61    342  3.88  3.84  2.33
4  0.2   Premium E      SI2     60.2   62    345  3.79  3.75  2.27
5  0.32 Premium E      I1      60.9   58    345  4.38  4.42  2.68
>
```

- magritter パッケージの **Tee pipe** で `print()` や `View()` を活用

- Tidyverse とパイプ演算子のおさらい
- データ整形の指針
 - ～本当はデータファイルを作る前に確認しておきたいこと～
- 初級編：最新の癌種別死亡数のグラフを作りたい
 - ～人口動態統計のCSVファイルを整形してみる～
- 中級編：臨床研究のまとめファイルがやってきた
 - ～“ネ申エクセル” と戦うための基礎知識～

- 「見る」あるいは「見せる」ための表や雑然とした表を、コンピュータで認識できて(**機械判読可能**)、解析に使用できる表に変形・変換する
- 大まかな目標となる代表的なルール
 - 総務省「統計表における機械判読可能なデータの表記方法の統一ルール」
https://www.soumu.go.jp/menu_news/s-news/01toukatsu01_02000186.html
 - 整然データ Tidy data
 - ✓ Wickham H. Journal of Statistical Software. 59: 1-23, 2014
 - ✓ Tidyverse はこの思想に基づいて設計されており、これらのパッケージをフル活用するためには tidy data を意識することが大切
- 元ファイルは極力編集せず、編集の内容や経過が後から確認できるようコードに残すようにする

- 2020年に、政府統計の総合窓口(e-Stat)に掲載する統計表におけるデータ表記方法の統一ルールとして総務省が策定したもの
- Excelなど特定のソフトや政府統計に依存しない部分を抜粋すると
 - ✓ 1セル1データとなっている
 - ✓ 数値データは数値属性とし、文字列（注:単位、注釈など）を含まない
 - ✓ スペースや改行等で体裁を整えていない
 - ✓ 項目名等を省略していない（注:「薬剤A」「B」「C」は「薬剤A」「薬剤B」「薬剤C」とする）
 - ✓ データの単位を記載している（注:数値とは別のセルに記載している）
 - ✓ 機種依存文字を使用していない
 - ✓ データが分断されていない・1シートに複数の表が掲載されていない

- Tidyverseの作者が提唱した概念で、「データの構造 (structure)」と「意味 (semantic)」を一致させることを目指している
 - ✓ 個々の変数 (variable) が1つの列 (column) をなす
 - ✓ 個々の観測 (observation) が1つの行 (row) をなす
 - ✓ 個々の値 (value) が1つのセル (cell) をなす
 - ✓ 個々の観測の構成単位の類型 (type of observational unit) が1つの表 (table) をなす

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 1845 | 1980071 |
| Afghanistan | 2000 | 2666 | 20095360 |
| Brazil | 1999 | 31737 | 172006362 |
| Brazil | 2000 | 80488 | 174004898 |
| China | 1999 | 210258 | 127201272 |
| China | 2000 | 210766 | 128042583 |

variables

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 1845 | 1980071 |
| Afghanistan | 2000 | 2666 | 20095360 |
| Brazil | 1999 | 31737 | 172006362 |
| Brazil | 2000 | 80488 | 174004898 |
| China | 1999 | 210258 | 127201272 |
| China | 2000 | 210766 | 128042583 |

observations

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 1845 | 1980071 |
| Afghanistan | 2000 | 2666 | 20095360 |
| Brazil | 1999 | 31737 | 172006362 |
| Brazil | 2000 | 80488 | 174004898 |
| China | 1999 | 210258 | 127201272 |
| China | 2000 | 210766 | 128042583 |

values

Wickham H., et al.
R for Data Science (2e) .
<https://r4ds.hadley.nz/>

- Tidyverse とパイプ演算子のおさらい

- データ整形の指針

～本当はデータファイルを作る前に確認しておきたいこと～

- **初級編：最新の癌種別死亡数のグラフを作りたい**

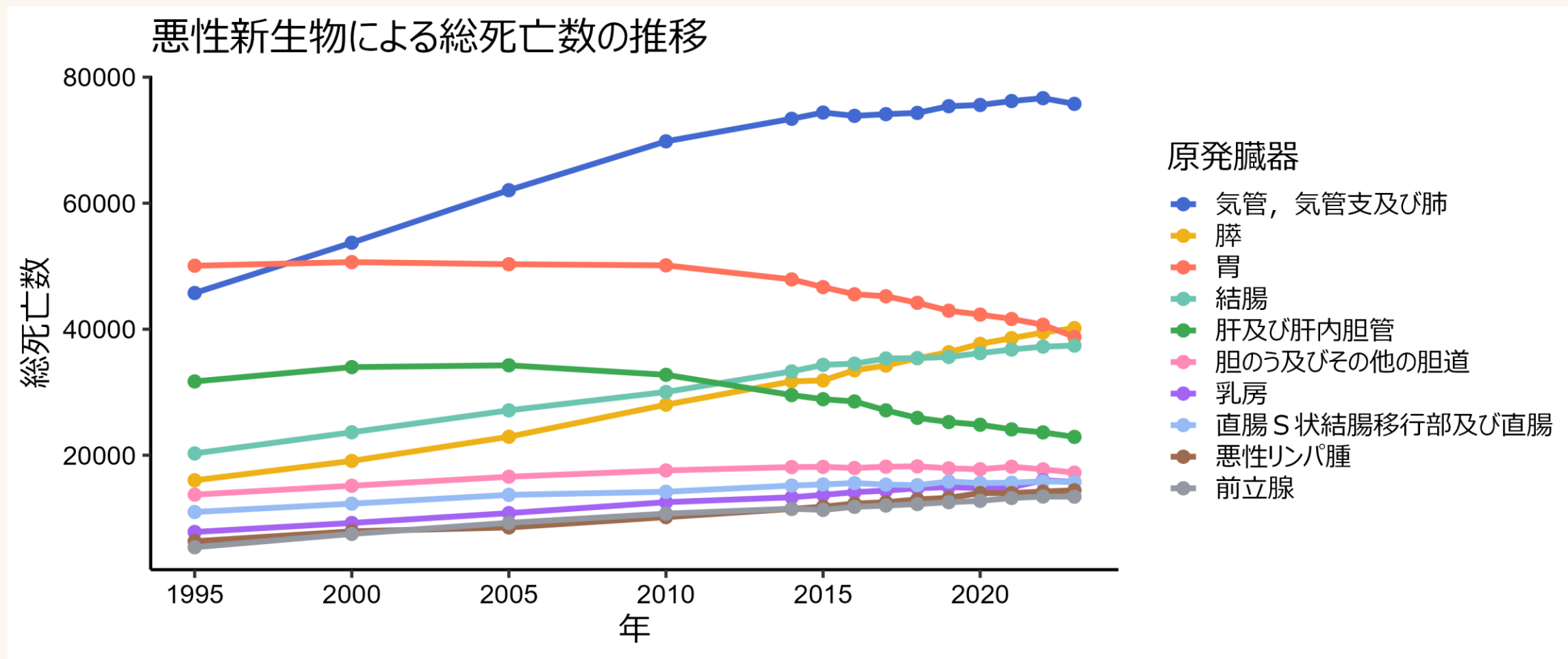
～人口動態統計のCSVファイルを整形してみる～



- 中級編：臨床研究のまとめファイルがやってきた

～“ネ申エクセル” と戦うための基礎知識～

厚労省 人口動態統計のCSVファイルを整理して、下の図を作成する



- 厚労省 人口動態統計より、2023年の死因別死亡数のデータを使用
- 政府統計の総合窓口(e-Stat)の該当ページ
<https://www.e-stat.go.jp/stat-search/files?tclass=000001041646&cycle=7&year=20230>
より、「5-13 死因(死因簡単分類)別にみた性・年次別死亡数及び死亡率(人口10万対)/ 2023年」の CSVファイルを事前にご自身の Working Directory にダウンロードしておいて下さい
 - ダウンロードされるファイル名: mc130000.csv
 - 直接リンク: <https://www.e-stat.go.jp/stat-search/file-download?statInfId=000040206118&fileKind=1>

ファイル内容を確認

19

いきなり R で触る前に、まずは Excel などを開いてファイルの中身を俯瞰する

自動保存 オフ mc130000.csv

ファイル ホーム 挿入 ページレイアウト 数式 データ 校閲 表示 開発 ヘルプ JMP Acrobat

A1 : X Y Z 令和 5 年

1 令和 5 年 人口動態統計

2 上巻 死亡 第 5. 13 表 死因 (死因簡単分類) 別にみた性・年次別死亡数及び死亡率 (人口10万対)

3 注: 1) 「子宮の悪性新生物<腫瘍>」、「卵巣の悪性新生物<腫瘍>」及び「妊娠、分娩及び産じょく」の率については、女性人口10万対の死亡率である。

4 2) 「前立腺の悪性新生物<腫瘍>」の率については、男性人口10万対の死亡率である。

5 3) 「誤嚥性肺炎」の2016年 (平成28年) 以前死亡数は、死因基本分類コード「J69 固形物及び液状物による肺炎」の数値である。

6 4) 「間質性肺疾患」の2016年 (平成28年) 以前死亡数は、死因基本分類コード「J84 その他の間質性肺疾患」の数値である。

5) 「その他の呼吸器系の疾患 (J40-J47及びJ60-J69を除く)」の2016年 (平成28年) 以前死亡数は、死因簡単分類コード10600から、死因簡単分類コード10601及び10602を除いた数値である。

「A92.8A ジカ< Zika> ウイルス病」の数値である。

2019年3月29日公表) による再集計を行ったことにより、2017年 (平成29年) 以前の報告書とは数値が一致しない。

コード「U07.1 コロナウイルス感染症2019、ウイルスが同定されたもの」、

コロナウイルス感染症2019に関連する多系統炎症性症候群、詳細不明」の数値である。

死因基本分類コード「U12 エマージェンシーコードU12」の数値である。

10) 「その他の特殊目的用コード (22201及び22202を除く)」の2022年 (令和 4 年) 以前は、死因簡単分類コード22200から、死因簡単分類コード22201及び22202を除いた数値である。

13 死因 死亡数 死亡率

14 1995 2000 2005 2010 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 1995 2000 2005 2010 2014 2015 2016 2017

15 死亡総数 922139 961653 1083796 1029.8 1046.4 1075.9

16 01000 感染症及び寄生虫症 18925 19858 23538 20.1 20.1 19.9

17 1752 1.9 2 1.9

18 2296 1.6 1.5 1.9

19 30553 18.5 16.8 16.4

20 600 437 357 314 290 288 273 299 251 256 213 229 0.6 0.6 0.5 0.3 0.3 0.2 0.2

21 4364 5242 5748 6133 6237 7498 7554 7806 7607 7901 8298 8667 1.9 2.6 3.5 4.1 4.6 4.9 5

121 22000 特殊目的用コード ... 3466 16784 47661 38120 ...

122 22100 重症急性呼吸器症候群 [S A R] ... 3466 16784 47661 38120 ...

123 22200 特殊目的用コード ... 3466 16784 47661 38120 ...

124 22201 新型コロナウイルス感染症 ... 3466 16766 47661 38086 ...

125 22202 新型コロナウイルス感染症 ... 18 23 34 ...

126 22203 その他の特殊目的用コード ...

127 男 死亡総数 501276 525903 584970 633701 660340 666728 674946 690770 699138 707421 706834 738141 799420 802536

128 01000 感染症及び寄生虫症 10671 10907 12211 12795 12321 12307 12135 12021 11733 11531 10905 10889 11714 11910

129 01100 腸管感染症 476 524 733 999 1048 1036 1135 1012 1039 1009 995 900 938 979

準備完了 アクセシビリティ: 利用不可

総数、男性、女性の表が縦に積み重ねられ、その見出しは「死因」列内にある

ファイルの説明があり
実際のデータで始まらない

死亡総数だけコードがない
(他の表では 00000)

各年の死亡数、死亡率について
見出しが複数行にわたり、左端の列にしか項目名がない

(全角)スペースによる整形

数値がない部分には色々な記号
(このファイル内に凡例なし)

- 1ページの表に粗死亡数と人口10万人対の死亡率が一緒に掲載されているが、今回は「死亡数」の方を見ていくこととする
- 必要な処理としては以下のようなものが考えられる：
 1. CSVファイルを読み込むときに、先頭の解説行を読み飛ばす
 2. 列名(変数名)をRで扱いやすいように修正する
 3. 死因欄に性別の見出しも入っているので、性別と死因を分離する
 4. 死因欄の整形のためのスペースを取り除く
 5. 数値が入るべきセルにある記号を除去し、数値として扱えるように変換する
 6. (オプション)Tidy data に整形する

str_replace_all() 関数

21

`stringr::str_replace_all(string, pattern, replacement)`

- `string` 中に登場する `pattern` を `replacement` に置き換える
- 主に `dp1yr::mutate()` の中で使用する

R 4.3.3 · ~/

```
> stringr::str_replace_all("ABCabc123", "abc", "***")
[1] "ABC***123"
> stringr::str_replace_all("ABCabc123", "[:lower:]", "*")
[1] "ABC***123"
> stringr::str_replace_all("ABCabc123", "\\d", "*")
[1] "ABCabc***"
> stringr::str_replace_all("ABCabc123", "([:lower:]+)", "_\\1_")
[1] "ABC_abc_123"
> stringr::str_replace_all("ABCabc123", "^[:graph:]", "*")
[1] "*BCabc123"
> stringr::str_replace_all("ABCabc123", "[:graph:]+$", "*")
[1] "ABCabc12*"
>
```

任意の文字列を置き換える(基本形)

正規表現(ワイルドカード)

() で囲んだ pattern に合致した部分は replacement で引用できる

+ 1文字以上の繰り返し

^ 先頭

\$ 末尾



- Tidyverse とパイプ演算子のおさらい

- データ整形の指針

～本当はデータファイルを作る前に確認しておきたいこと～

- 初級編：最新の癌種別死亡数のグラフを作りたい

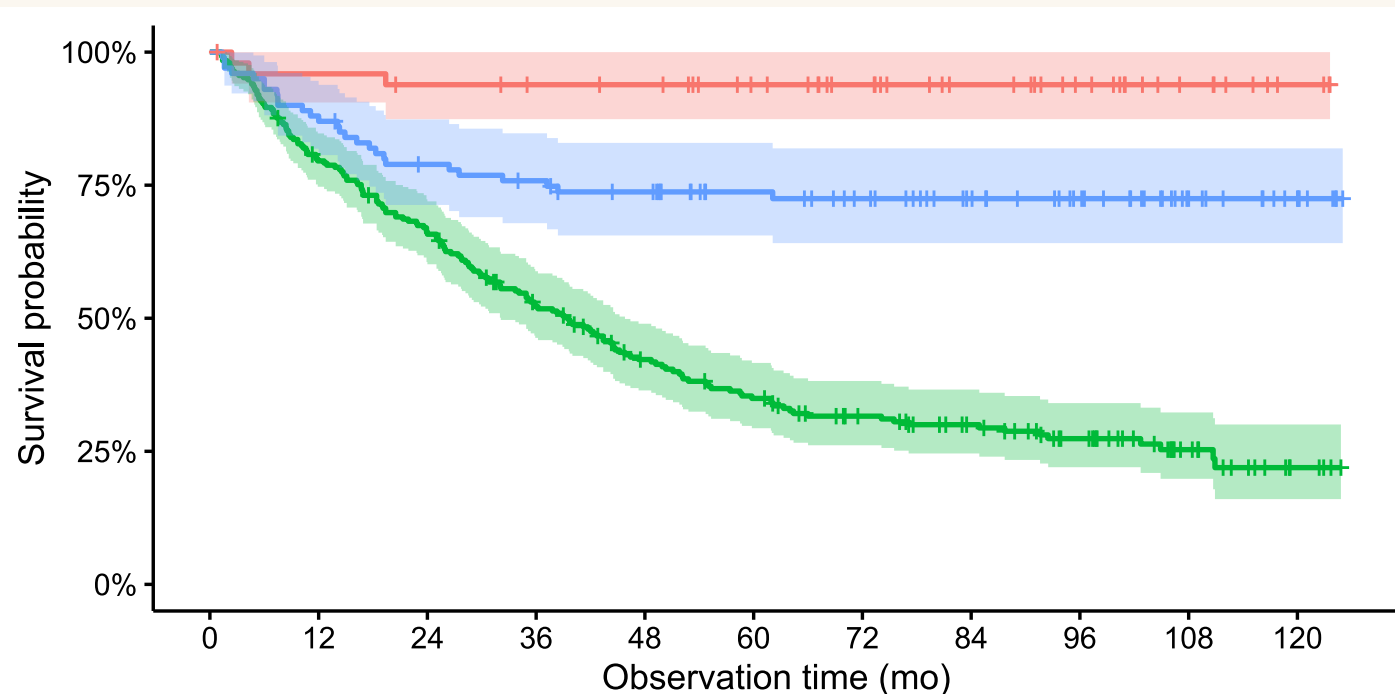
～人口動態統計のCSVファイルを整形してみる～

- 中級編：臨床研究のまとめファイルがやってきた

～“ネ申エクセル” と戦うための基礎知識～



臨床研究データを解析に使用できる形に整形し、以下の図表を作成する



Number at risk

| | | | | | | | | | | | |
|---------|-----|-----|-----|-----|----|----|----|----|----|----|----|
| dx=COP | 50 | 47 | 45 | 43 | 42 | 36 | 30 | 23 | 17 | 9 | 3 |
| dx=IPF | 250 | 197 | 163 | 124 | 93 | 76 | 60 | 49 | 36 | 18 | 5 |
| dx=NSIP | 100 | 88 | 77 | 73 | 68 | 58 | 51 | 42 | 32 | 19 | 10 |

| Characteristic | HR ¹ | 95% CI ¹ | p-value ² |
|----------------|-----------------|---------------------|----------------------|
| dx | | | |
| COP | — | — | |
| IPF | 23.7 | 3.29, 171 | 0.002 ** |
| NSIP | 10.4 | 1.37, 79.0 | 0.024 * |
| sex | | | |
| 女 | — | — | |
| 男 | 1.53 | 1.07, 2.19 | 0.020 * |
| age_enroll | 1.05 | 1.02, 1.08 | 0.002 ** |
| delta_FVC_1y | 0.17 | 0.10, 0.30 | <0.001 *** |

¹ HR = Hazard Ratio, CI = Confidence Interval

² *p<0.05; **p<0.01; ***p<0.001



データが集まりました

From: 共同研究者 <eraihito@daigaku.ac.jp>

To: 自分, ボス ▼



お疲れ様です。症例登録担当の xxx です。
臨床研究のデータが集まりました。相談した解析をお願いします。
シート「症例登録票」が症例登録のFileMakerからエクセル出力したものです。
不足データは各施設に確認して分かる範囲で記入しました。1年後以降の結果
はシート「アウトカム」で一緒に入れてあります。あと、登録時と経過のデータが
分かりにくかったので見出しを追加してあります。
個人情報保護のため、患者氏名は削除してありますので安心して下さい。
パスワードはいつも通りです。それではよろしく。



症例登録票_3 250220 Final-2 生存データつき Rev.1 解析用.xlsx

※フィクションです！！

承知しました……（とりあえずパスワード外して別名保存しよう）

ファイルを開く前に頭をよぎる可能性

25



データが集まりました

From: 共同研究者 <eraihito@daigaku.ac.jp>

To: 自分, ボス ▼



チェックリスト ⇒ 1セルに複数データ
がありそう

複数行の見出しが
セル結合がありそう

お疲れ様です。症例登録担当の xxx です。

臨床研究のデータが集まりました。相談した解析をお願いします。

シート「症例登録票」が症例登録のFileMakerからエクセル出力したものです。

不足データは各施設に確認して分かる範囲で記入しました。1年後以降の結果

はシート「アウトカム」で一緒に入れています。あと、登録時と結

違う形式で入力
されているかも

分かりにくかったので見出しを追加してあります。

個人情報保護のため、患者氏名は削除してありますので安心して下さい。

パスワードはいつも通りです。それではよろしく。

そのままではR
で開けないな...



症例登録票_3 250220 Final-2 生存データ

2つのデータの結合が必要だけど、
連結に使える情報は残っているかな

先方は mac かな。こちらは Windows
だけど文字化けしないと良いな

※フィクションです！

（あえずパスワード外して別名保存しよう）

ファイル内容を確認 (1)

26

シート1「症例登録票」

111

fx

糖尿病

不整脈

高血圧

% はじまりは、そのままではRの変数名として

| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|----|------|-------|--------|------------|----|-------|-----------|------|-------------|------|------|------|-------|
| 1 | 登録番号 | 施設 | 施設ID | 登録日 | 性別 | 登録時年齢 | 生年月日 | 診断 | 合併症 | KL6 | FVC | FEV1 | %DLco |
| 2 | 1 | 大学病院 | 1 | 2013/10/1 | 女 | 77 | 1936/1/4 | IPF | なし | 1010 | 2.1 | 1.95 | 88.6 |
| 3 | 2 | 大学病院 | 2 | 2013/10/3 | 男 | 68 | 1945/5/4 | IPF | 糖尿病不整脈 | 782 | 2.3 | 1.56 | 93.1 |
| 4 | 3 | 総合病院B | 3103 | 2013/10/3 | 女 | 67 | 1946/6/2 | COP | 糖尿病不整脈 | 912 | 2.4 | 2.22 | 99.5 |
| 5 | 4 | 大学病院 | 3 | 2013/10/4 | 女 | 75 | 1938/8/25 | COP | | | | | |
| 6 | 5 | 大学病院 | 4 | 2013/10/5 | | | | | 高血圧 | | | | |
| 7 | 6 | 総合病院A | 320540 | 2013/10/7 | | | | | 高血圧 | 1100 | 2.51 | 2.44 | 85.5 |
| 8 | 7 | 大学病院 | 5 | 2013/10/8 | | | | | 糖尿病、不整脈、高血圧 | 1011 | 2.53 | 2.1 | 101.5 |
| 9 | 8 | 大学病院 | 6 | 2013/10/8 | 女 | 59 | 1954/4/11 | IPF | 糖尿病不整脈高血圧 | 584 | 2.17 | 1.84 | 79.5 |
| 10 | 9 | 総合病院A | 186209 | 2013/10/11 | 女 | 60 | 1953/8/4 | NSIP | 糖尿病不整脈高血圧 | 987 | 2.42 | 2.42 | 97.2 |
| 11 | 10 | 総合病院A | 10125 | 2013/10/14 | 女 | 77 | 1936/5/10 | IPF | 糖尿病不整脈高血圧 | 1088 | 3.46 | 2.51 | |
| 12 | 11 | 市立病院X | 1 | 2013/10/15 | 女 | 68 | 1944/11/8 | NSIP | 糖尿病 | 869 | 3.07 | 3.04 | 81.1 |
| 13 | 12 | 大学病院 | 7 | 2013/10/16 | 女 | 64 | 1949/7/11 | IPF | 糖尿病 | 880 | 2.14 | 1.86 | 69.8 |
| 14 | 13 | 総合病院A | 67743 | 2013/10/17 | 男 | 66 | 1947/3/11 | IPF | 不明 | 2.57 | 2.19 | 80.1 | |
| | | | | | | 68 | 1945/9/11 | NSIP | 糖尿病不整脈高血圧 | | | | |
| | | | | | | 63 | 1950/7/13 | NSIP | 糖尿病 | | | | |
| | | | | | | 71 | 1949/7/20 | IPF | 高血圧 | | | | |

「不明」と空欄が混在している

手入力したところは「、」区切りになっ

施設IDは形式がバラバラ
登録番号は両方にあれば使えそう

FileMaker で複数選択チェックリスト
項目は改行区切りで1セルに入っ

% はじまりは、そのままではRの変数名として不適當

「不明」と空欄が混在している

手入力したところは「、」区切りになっている

施設IDは形式がバラバラ
登録番号は両方にあれば使えそう

FileMaker で複数選択チェックリストだった項目は改行区切りで1セルに入っている

ファイル内容を確認 (2)

27

シート2「アウトカム」

| | | | | | | | | | | | | | | | | | | | |
|---------------------------------|----|------|----|------------|-----|------|-------|------|-------|------|-------|------|------|------|-------|------|-------|--|--|
| F1 | | 1年後 | | | | | | | | | | | | | | | | | |
| 見出しが2行構成 | | | | セル結合！ | | | | | | | | | | | | | | | |
| | | | | 1年後 | | | | 3年後 | | | | 5年後 | | | | | | | |
| 登録番号 施設 施設ID 転帰日 転帰 打ち切り0死亡1 | | | | KL6 | FVC | FEV1 | %DLco | KL6 | FVC | FEV1 | %DLco | KL6 | FVC | FEV1 | %DLco | | | | |
| 3 | 1 | 大学病院 | 1 | 2023/8/24 | 0 | 1060 | 2.03 | 1.62 | 74.7 | 1142 | 1.55 | 1.45 | 61.9 | 949 | 1.04 | 0.94 | | | |
| 4 | 2 | 大学病院 | 2 | 2023/5/1 | 1 | 522 | 2.45 | 2.43 | 110.7 | 561 | 2.33 | 2.29 | | 569 | 2.28 | 2.1 | 105.6 | | |
| 5 | 4 | 大学病院 | 4 | 2024/1/27 | 0 | 1227 | 2.24 | 1.8 | 64.2 | 181 | 1.43 | 52.1 | 707 | 1.57 | 1.22 | 45.3 | | | |
| 6 | 5 | 大学病院 | 5 | 2024/2/13 | 0 | 919 | | | | | | | | | 1.59 | 74 | | | |
| 7 | 7 | 大学病院 | 6 | 2013/11/15 | 1 | 762 | 1.87 | 1.78 | 56.9 | 925 | 1.34 | 1.24 | | 1121 | 1.2 | 1.01 | | | |
| 8 | 8 | 大学病院 | 7 | 2024/2/7 | 0 | 625 | 2.75 | 2.32 | 85.1 | 714 | 2.52 | 2.08 | 68.1 | 250 | 2.36 | 1.87 | 66.9 | | |
| 9 | 12 | 大学病院 | 8 | 2024/3/19 | 0 | 692 | 3.58 | 2.78 | 94.1 | 499 | 3.25 | 2.67 | 81.1 | 541 | 2.97 | 2.48 | 67.8 | | |
| 10 | 15 | 大学病院 | 9 | 2024/1/22 | 0 | 1878 | 2.5 | 1.69 | 64.7 | 3059 | 1.68 | 1.68 | 58.9 | 2142 | 1.34 | 1.22 | | | |
| 11 | 17 | 大学病院 | 10 | 2023/5/18 | 0 | 804 | 2.52 | 2.41 | 95.8 | 781 | 2.43 | 2.16 | 82.9 | 811 | 2.32 | 1.86 | 76.1 | | |
| 12 | 19 | 大学病院 | | | | 1045 | 1.9 | 1.47 | 84.6 | 1050 | 1.58 | 1.19 | 73.2 | | 1.46 | 0.93 | | | |
| | | | | | | 1412 | 1.98 | 1.91 | | 1598 | | | | | | 1.01 | 33.7 | | |
| 16 | 28 | 大学病院 | 14 | 2020/8/18 | 0 | 623 | 2.94 | 1.94 | 69 | 520 | | | | | | 1.2 | 46.4 | | |
| 17 | 32 | 大学病院 | 15 | 2024/3/22 | 0 | 610 | 2.8 | 2.30 | 60.7 | 306 | 2.45 | 2.42 | 86.8 | 454 | 2.27 | 2.16 | 76.8 | | |

ファイル内容を確認 (3)

28

シート2「アウトカム」

```
readxl::read_xlsx("ip_registry_data.xlsx", sheet = 2) %>% View()
```

セル結合の範囲の左上端に読み込まれ、他は空白扱い

| | ...1 | ...2 | ...3 | ...4 | ...5 | 1年後 | ...7 | ...8 | ...9 | 3年後 | ...11 |
|----|------|------|------|-------|-------------|------|--------------------|--------------------|-------|------|--------------|
| 1 | 登録番号 | 施設 | 施設ID | 転帰日 | 転帰 打ち切り0死亡1 | KL6 | FVC | FEV1 | %DLco | KL6 | FVC |
| 2 | 1 | 大学病院 | 1 | 45162 | 0 | 1060 | 2.0299999999999998 | 1.62 | 74.7 | 1142 | 1.55 |
| 3 | 2 | 大学病院 | 2 | 41647 | 1 | NA | NA | NA | NA | NA | NA |
| 4 | 4 | 大学病院 | 3 | 45053 | 0 | 522 | 2.4500000000000002 | 2.4300000000000002 | 110.7 | 561 | 2.33 |
| 5 | 5 | 大学病院 | 4 | 45318 | 0 | 1227 | 2.2400000000000002 | 1.8 | 64.2 | NA | 1.81 |
| 6 | 7 | 大学病院 | 5 | 45335 | 0 | 919 | 2.39 | 1.96 | NA | 1005 | 2.11 |
| 7 | 8 | 大学病院 | 6 | 41593 | 1 | NA | NA | NA | NA | NA | NA |
| 8 | 12 | 大学病院 | 7 | 45329 | 0 | 762 | 1.87 | 1.78 | 56.9 | 925 | 1.34 |
| 9 | 15 | 大学病院 | 8 | 45370 | 0 | 625 | 2.75 | 2.3199999999999998 | 85.1 | 714 | 2.52 |
| 10 | 17 | 大学病院 | 9 | 45313 | 0 | 692 | 3.58 | 2.78 | 94.1 | 499 | 3.25 |
| 11 | 19 | 大学病院 | 10 | 45064 | 0 | 1878 | | | | 059 | 1.68 |
| 12 | 21 | 大学病院 | 11 | 43657 | 0 | 804 | | | | 81 | 2.4300000000 |
| 13 | 22 | 大学病院 | 12 | 44129 | 0 | 1045 | | | | 050 | 1.58 |
| 14 | 26 | 大学病院 | 13 | 44366 | 0 | 1412 | | | | 598 | 2.04 |
| 15 | 28 | 大学病院 | 14 | 44061 | 0 | 623 | 2.94 | 1.94 | 69 | 520 | 2.7 |

2行目の見出しはデータの1行目扱い

RとExcelの数字の丸めの違いで
数値が変わっている
(臨床データとしてはほぼ誤差レベル)

日付データが整数値になってしまっている

生存分析に使用することを考え、観察期間や各種変数を整理する：

- 列名はRで扱いやすい英単語ベースに改名する
 - 4回登場する KL6, FVC, FEV1, %DLco は時期を区別できるよう工夫する
- 2つのシートは「登録番号」をキーにして連結する
- 日付データは整数値(Excelのシリアル値)から日付に再変換する
- 観察期間がないので「転帰日 - 登録日」で求める
- 合併症は、「糖尿病」「不整脈」「高血圧」それぞれの有無に分解する
- 誤差が生じている数値があれば、小数点以下の桁数を揃えて丸める
- 入力間違いによる外れ値がないか確認

- 間質性肺炎の多施設レジストリ研究をイメージした模擬データ (n = 400)

- 観察期間は10年、診断(病型)は以下の3つに簡略化

- IPF（特発性肺線維症）

- ✓ 特発性（他に原因がない）間質性肺炎の中で最多
- ✓ 難治性で予後不良
- ✓ 喫煙の影響あり、男性に多い

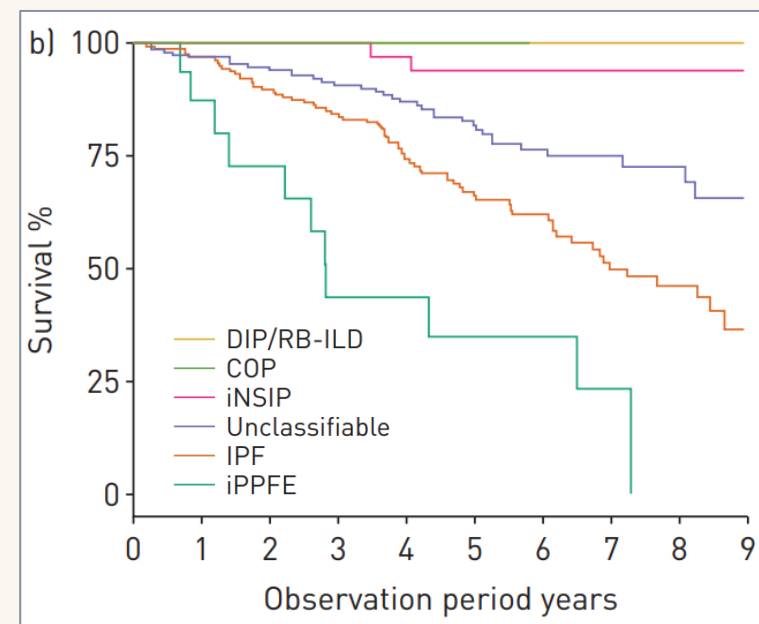
- NSIP（非特異性間質性肺炎）

- ✓ 治療反応性・予後などは雑多
- ✓ 膠原病に関連する例も多く、女性に多い

- COP（特発性器質化肺炎）

- ✓ 治療反応性良く、予後良好

- 努力肺活量の低下、高齢、男性などが間質性肺炎における既知の予後不良因子



Fujisawa T, *et al.* Eur Respir J. 2019
(実際のレジストリ研究) より引用改変

これらの値まで取り扱う時間がないですが、toy data として残しておきます。

- 努力肺活量

通常の呼吸から目一杯息を吸い込んだ後、力いっぱい一気に吐き出す検査

- 最初の1秒間で吐き出した空気の量が「1秒量 FEV_{1.0}」
- 最後まで吐き切った総量が「努力肺活量 FVC」
 - ✓ したがって必ず FEV_{1.0} < FVC となり、それなりの相関関係がみられる

- 肺拡散能 DL_{CO} (%DL_{CO} は年齢体格等からの予測値に対する値)

- 測定手技の関係で、最低 1.5L（基本的には 2.0L以上）の肺活量が必要
 - ✓ それ以下では実際の拡散能によらず計測できない ⇒ (FVC低値例では) MAR型の欠測が生じる

- KL-6

※もちろん病態の悪化により呼吸機能全体が実施不能となる場合もあるので、実際にはMNAR型の欠測も含まれます

- 血液検査項目。代表的な肺組織障害(≡間質性肺炎の病勢)のバイオマーカー

- データの前処理では、データクリーニングや変換操作の前に生データをコンピュータで認識できて解析に使用できる(機械判読可能な)表に変形・変換する過程が重要かつ時間がかかる部分
- Tidyverse の各関数を用いてデータを操作する際には、
 - 一息に操作を行わず、少しずつ結果を確認しながら進める
 - 作業のお供に
 - ✓ Posit社提供の各種チートシート
 - ✓ tidylog パッケージ
 - ✓ magritter の Tee pipe
 - ✓ AI アシスタント