

# • DPCデータ分析演習

横浜市立大学  
データサイエンス研究科  
ヘルステータサイエンス専攻  
清水沙友里

## Q. PCが固まってしまいます

データの前処理(必要なデータに絞る、分割)、数式の最適化(ループ処理をしない)、データのリンク(外部データベースやCSVファイルからデータを直接読み込むことで、Excelシートにすべてをロードする必要を減らす)、マクロとVBA(繰り返し行作業はマクロやVBA(Visual Basic for Applications)を使用して自動化してもよい。動作が重くなりやすい関数は以下のようなものがあります)。

参照関数: VLOOKUP、HLOOKUP、INDEX、MATCH、OFFSETなど。大量のセルにわたって使用されると、計算に時間がかかる。

条件付き集計関数: SUMIFS、COUNTIFS、AVERAGEIFSなどは、多くの条件や大量のデータが関与する場合には、計算に時間の間隔がかかる。

再帰的な参照: 関数が他のセルを参照し、そのセルがさらに別の関数によって計算される場面では、計算量が急速に増加する可能性がある。

自作関数(VBA): VBAで作成されたユーザー定義関数も、処理が複雑であると時間がかかる。

外部データ参照: 外部データベースやWebサービスへのクエリを含む関数は、データベースのパフォーマンスに依存するため、遅くなる可能性がある。

関数のネスト: 関数の関数をネスト(入れ子)にして使用すると、計算の複雑性が増し、パフォーマンスに影響を与える。

パフォーマンスの問題を軽減するためには、計算の範囲を限定する、不要な計算を削除する、Excelの計算オプションを手動に設定する(自動計算をオフにする)などの方法があります。

## Q. データをグラフ化する際のポイント①

データをグラフ化する際には、何を伝えたいのか、また誰に伝えたいのかというコンテキストが非常に重要です。明確な目的: グラフには明確な目的が必要です。何を伝えたいのかを考え、それに最も適したグラフの種類を選びます。

シンブル: シンブルは複雑なグラフは混亂を招きます。必要な情報だけを表示し、無駄を排除します。

色の使用: 色は強力な視覚的要素ですが、多用すると逆効果になる場合もあります。基本的に2~3色に抑え、強調したい点に色を使います。

ラベルと凡例: グラフの各要素が何を意味するのかを明確にするために、適切なラベルと凡例を使用します。

スケールと軸: スケールが不適切だと、データの解釈が誤りやすくなります。0から始まる軸や、等間隔のスケールを使用すると良いでしょう。

## データ分析の前に…よくある質問



## Q. データをグラフ化する際のポイント①

データをグラフ化する際には、何を伝えたいのか、また誰に伝えたいのかというコンテキストが非常に重要です。明確な目的: グラフには明確な目的が必要です。何を伝えたいのかを考え、それに最も適したグラフの種類を選びます。

シンブル: シンブルは複雑なグラフは混亂を招きます。必要な情報だけを表示し、無駄を排除します。

色の使用: 色は強力な視覚的要素ですが、多用すると逆効果になる場合もあります。基本的に2~3色に抑え、強調したい点に色を使います。

ラベルと凡例: グラフの各要素が何を意味するのかを明確にするために、適切なラベルと凡例を使用します。

スケールと軸: スケールが不適切だと、データの解釈が誤りやすくなります。0から始まる軸や、等間隔のスケールを使用すると良いでしょう。

## Q. データをグラフ化する際のポイント②

データをグラフ化する際には、何を伝えたいのか、また誰に伝えたいのかというコンテキストが非常に重要

棒グラフ / 垂直または水平:

用途: カテゴリ別に数値を比較する。月別売上、部門別利益など。

折れ線グラフ:

用途: 時間に沿ったデータのトレンドを表示する。株価の時間経過による変動、年度ごとの売上推移など。

円グラフ:

用途: 全体に対する各部分の割合を示す。分類が多くなると読みにくくなる。通常は5~6個以下のカテゴリで使用する。

ヒストグラム:

用途: 一変数のデータセットの頻度分布を表示する。データを bin(範囲) に分け、各 bin に含まれるデータ点の数を棒グラフで表示

散布図:

用途: 二つの数値の関係性を探る。X軸とY軸にそれぞれ異なる変数を取り、各データ点をプロットします。

箱ひげ図:

用途: データの四分位数と外れ値を視覚化する。特長、中央値、第一四分位数、第三四分位数などを一つのグラフで表示し、データのばらつきを理解しやすくします。

5

6

## Q. 知つていると便利なExcel関数は?

VLOOKUP / HLOOKUP: 一方のシートのデータをもとに、別のシートのデータを検索します。変数のコード化によく使います

=VLOOKUP(A1, \$D\$1:\$E\$2, 2, FALSE)  
=VLOOKUP(A1, \$D\$1:\$E\$2, 2, FALSE)  
D1:E2の範囲に変換テーブルがあり、A1の値(1または2)に基づいて"女"または"男"を返します。

XLOOKUP: 一つの列(または行)から特定の値を検索し、対応する値を別の列(または行)から返します

INDEX / MATCH: VLOOKUPよりも柔軟な検索が可能です。

SUMIF / COUNTIF: 特定の条件に一致するデータを集計します。これは、条件付きの集計にも使えます。

SUMPRODUCT: 複数の配列の要素同士を掛け合わせた後、その総和を計算します。これは、条件付きの集計にも使えます。

CONCATENATE / &: 文字列を結合します。

LEFT / RIGHT / MID: 文字列から特定の文字を抽出します。

IF / IFERROR: 条件に基づいて異なる計算を行い、エラーが発生した場合にはデフォルト値を返します。

=IF(A1=1, "女", IF(A1=2, "男", "不明"))

セルA1が1であれば"女"、2であれば"男"、それ以外であれば"不明"と表示されます。

AND / OR: 値数の条件を組み合わせて評価します。

LEN: 文字列の長さを返します。

TRIM: 文字列から余分なスペースを取り除きます。

## Q. 知つていると便利な関数 (Excel 365以降)

ダイナミック配列関数はExcel 365以降で使用できる新しいカテゴリの関数で、一つのセルに式を入力するだけで、複数の出力値を隣接するセルに自動的に「スピルオーバー」させることができます。

FILTER関数: 条件に合った行や列を抽出します。例えば、A1:C10の範囲に「ID」「名前」「年齢」という列がある場合、年齢が20歳以上の人だけを抽出したいとします。この場合、以下のようないFILITER関数を使用します。

=FILTER(A1:C10, C1:C10 >= 20)  
=SORT(A1:C10, 3, 1) この式は、A1:C10の範囲を3項目(年齢)に基づいて昇順(1)でソートします。

SEQUENCE関数: 重複を削除・名前列から重複を削除する場合、以下のようにします。

=SEQUENCE(10)  
=UNIQUE(B1:B10)

関数を組み合わせる: SORT + FILTER 年齢が20以上の人物を年齢でソートする場合、以下のように関数を組み合わせます。

=SORT(FILTER(A1:C10, C1:C10 >= 20), 3, 1)

## Q. データ量が少ない時に注意すべきことは?

■ まずは自院の課題が何であるか、分析を目的を明確に持つことが重要です。やみくもにデータを分析をしても得られるものが少なくなくなってしまいます。分析を行った結果、どのような行動変容に繋がれるかということを考えて分析します。

■ 記述統計(人数、平均、中央値など)を維持的に行なうことで比較可能になります。

■ データを可視化することで、データの傾向や問題点を直感的に把握することができます。

■ 小数サンプルの場合、外れ値があると解釈に影響を与える可能性が高くなります。外れ値が発生した要因を検討しましょう。

■ 別のグループや時間帯と比較することによって、見落としていた傾向やパターンを発見することができます。

■ ドメイン知識を活かして、数名で分析結果を話し合うと新たな視点が生まれます。

## Q. 分析でよく使われる指標は？

病院の理念がなにより大切。

- QI 国立病院機構を参考に
- 患者満足度調査。ケアの質を評価します

患者の状態

- 年齢構成、疾病構造、重症度(介護度)、

業務指標(KPI)

- LOS、病床利用率、紹介・逆紹介率、新患数と再来

競合と市場分析

- 対象領域における市場占有率
- 患者の地域分布、地域別に患者数を分析し、マーケティング戦略を考えます。

職員のパフォーマンスと満足度

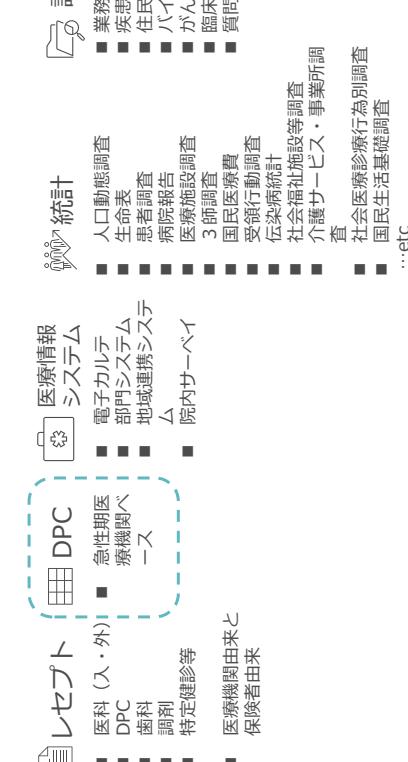
- 離職率：高い離職率は職員の不満や疲労、効率の低下を示す可能性があります。
- 職員満足度調査：モチベーションや職場環境に問題がないかを定期的に確認します。

9



## DPCデータを分析しよう！

### 医療系のさまざまデータ



研究法に係わらず、実臨床のデータから  
強く実践を志向した研究(real-world evidence)

11

### DPCデータとは

ファイル名	内容
患者半数で把握	簡易診断情報表 (カルテのサブレーフのような情報)
病院3施設半数で把握	施設情報表 (施設名、届け出法人の名称、加算額の算定方法など)
病院4施設半数で把握	医療機関情報の収集の有無に関する情報表
病院5施設半数で把握	医療機関分点収集表により算定した患者にかかる診療報酬請求書
ロファイル	診療報酬請求書(DPCレポート提出) →DPC実績報告の提出
ヒューリック	患者半数で把握 人差患者医療料金表に基づく出来高算算用情報 出未記載セバ付情報 外来診療料の標準料金表
外来E-Hospital	外未診療料の標準料金表に基づく出未記載算定料金 外未診療料の標準料金表
Hファイル	一括請求用 医療・看護必要量に関する情報
患者半数で把握	保険者半数で把握 医療・看護必要量に関する情報

- 様式1は退院時サマリーのイメージ。性別、年齢、身長、体重、病名、入院目的や入院経路、手術術式、重症度などを含む。患者がどのような状態で入院し、主にどのような手術を受け、何日間入院して、どのような状態で退院していくのか、概略が把握できます。
- EFファイルは診療報酬明細書情報。医科点数表に基づく出来高による診療区分の算定分を範囲とする。手術、検査、方等の診療区分のほか、実施年月日や行為回数、診療明細名稱や薬剤の使用教書・基準単位などの情報が収載されています。
- Hファイルは重症度、看護必要度に係る調査票を書く評価項目の点数。様式1と異なり、1日ごとに情報を入力する。産科等の患者についても作成する必要がある

## DPCデータで医療の質を評価する

### アウトカム分析

- 診療の成果を分析して医療の質を評価する方法
- 手術成功率、合併症発生率、死亡率など直接的な結果を評価する指標なので、医療の質の良し悪しが分かりやすく判断しやすいが、公平な評価が難しい指標
- 診療行為の明確性が高いEFファイルなどの診療データを用いて、診療の内容をより細かく比較する
- 診療プロセス分析を行うことができる

### プロセス分析

- 診療内容の詳細を分析して医療の質を評価する方法
  - ダイドラインの遵守率などを評価
  - 特にアウトカムと密接に関連するようなプロセス評価が重視されている
- ストラクチャーナリティ分析**
- 設備や人員体制の視点が主。診療報酬上の施設基準、法令に規定される人員、委員会等の整備等を評価する
  - 医療機能の特性や目指している医療に合致する医療提供体制かどうか

13

14

国立病院機構 臨床評価指標 Ver.4.1 2022 <https://moho.hosp.go.jp/files/00018119.pdf>

## DPCデータを活用して臨床指標を分析する

### 手術ありの患者の肺血栓塞栓症の予防対策の実施率を分析する(リスクレベルが高いリスク)

#### 分子

- 分母のうち、肺血栓塞栓症の予防対策(弾性ストッキングの着用、間歇的空気圧迫装置の利用、抗凝固療法のいずれか、または2つ以上)を実施した患者数

### 肺血栓塞栓症発症のリスクレベルが「高」の手術を施行した退院患者数

	年齢	病院集計				
		病院数	平均	標準偏差	中央値	25パーセンタイル
2021	*	****○○○○	81	82	87	96.2
2020	*	* *	95.2	96.1	95.5	95.2
2019	**	*** ○○○○	8.9	9.4	5.7	5.7
	目標値	97.9	97.6	98.1	94.2	95.5

14

国立病院機構 臨床評価指標 Ver.4.1 2022 <https://moho.hosp.go.jp/files/00018119.pdf>

## DPCデータを活用して臨床指標を分析する

### 75歳以上退院患者の入院中の予期せぬ骨折発症率

#### 分子

- 75歳以上退院患者の入院日から数えて2日目以降退院日までに骨折を発症した患者数

### 75歳以上の退院患者数

#### 分子

- 分母のうち、当該入院の入院日から数えて2日目以降退院日までに骨折を発症した患者数

	年齢	病院集計				
		病院数	平均	標準偏差	中央値	25パーセンタイル
2021	0	101	102	102	0.1	0.1
2020	1	0.2	0.2	0.2	0.2	0.2
2019	2	0.2	0.2	0.2	0.1	0.1
	目標値	0.0	0.0	0.0	0.0	0.0

	年齢	病院集計				
		病院数	平均	標準偏差	中央値	25パーセンタイル
2021	0	101	102	102	0.1	0.1
2020	1	0.2	0.2	0.2	0.2	0.2
2019	2	0.2	0.2	0.2	0.1	0.1
	目標値	0.3	0.3	0.2	0.2	0.2%以下

PCデータの前処理(済)

- | MDC  | D | OPD<br>(手術の<br>年齢) |
|------|---|--------------------|
| 医療 1 |   |                    |
| 医療 2 |   | 65歳以上              |
| 医療 3 |   | 65歳以上              |
| 医療 4 |   | 65歳以上              |
| 医療 5 |   | 65歳以上              |
| 医療 6 |   | 在院日                |

様式1とEFファイル、Dファイルを連結した分析を行うため、1入院1患者による分析IDを作成します。様式1、EFファイル、Dファイルに分析用IDを付与します。複数のFFIデータを保持します

✓ 1人の患者が同年に複数回入院するとき、複数のFFIデータを保持します

✓ 様式1、EFファイル、Dファイル全てで同じ作業を実施

✓ DPCコード別の分析を行った後、DファイルからDPCコード情報を様式1に追加

✓ MDC(主要診断群分類)を追加

✓ DPC14桁から頭2桁

✓ DPC6桁を追加

✓ DPC14桁のうち上6桁が“DPC6”

DPCコードから定期OPD(手術の有無)を追加

誕生年月日から年齢を計算

“65歳以上高齢者”列を作成し、0と1でカテゴリー化

在院日数を計算

✓ 様式1の開始日～終了までの情報

実施説明資料をデータと照らし合せて確認する

2021年4月1日時点 2020年3月31日までの累計	2021年度 「DPC導入の影響評価に係る調査」 実施説明資料	2021年4月1日
2021年4月1日時点 2020年3月31日までの累計	2021年度 「DPC導入の影響評価に係る調査」 実施説明資料	2021年4月1日

EFファイルを使った分析

病院内で急性心筋梗塞の診療ガイドラインがどれだけ遵守されているのか確認したいと考えています。何がいい方法はありませんか?

「基本的な対応策として、診療カーデラインで推奨されている薬剤の投与と、薬剤の適正化や質の改善につなげましょう」として、診療率を調べてみると、手は1つです。薬剤の投与と、薬剤の適正化や質の改善につなげましょう

第三期

では、一方として、日本倦怠症学会がまとめた「ST上型急性心筋梗塞の診断に関するガイドライン」(2013年改訂版)では、記載されているHMG-CoA還元酵素阻害薬(スタチア)の投与実行率を「STアーチルームで見る」に記載しています。付録CD-ROMに記載した「スタチアマスク」というExcelファイルを用いて、STアーチルームでSTアーチルームで記載しているマスクを表示しているのです。その情報と「標準」のデータをリンクさせて、急性心筋梗塞者の投与実態について集計できます。

EFファイルを使った分析

診療ガイドラインが遵守されているか、あるいはどれくらい逸脱した例があるかを検証する作業は、診療の適正化を図る上で非常に重要な算筋と言える。心筋梗塞の重篤防止に指管管理を行うことは、高いエビデンスがあり、日本循環器学会がまとめて「ST上昇型心筋梗塞の診療に関するガイドライン（2013年改訂版）」では、「禁忌がないければ開設する」とが推奨されている。コレステロール値にかかわらずHMG-CoA還元酵素阻害薬（タチクマ）を開設するにどこが選択されるべきであるかなど、今回の例では、主病名が急性心筋梗塞の患者に対するタチクマ投与の施行率を集計するにこどめたが、タチクマが投与されなかつた患者について院内調査を行い、非投与となつた理由について問へ、以後、同様のケースで改善策が図られるかどうかなどもぜひ検討したい。

- ① スタチングマスターを使って、EFファイルでスタチングが投与されている患者を特定する
  - ② スタチングが投与された患者の分析用IDのリストを作成する
  - ③ リストを使って、同じ分析用IDの患者に1のログをたてる
  - ④ ピボット機能を使って、急性心筋梗塞でスタチングが投与された患者の数を集計する

## 分析時に参照する参考資料の例

