# 東京都データプラットフォーム 行政データ整備モデル事業 データ整備マニュアル

# 第1版

発行元:東京都

委託先:日本総合研究所、株式会社自動処理

2022年 3月

# 改訂履歴

版	日時	作成者/更新者	主な改訂内容
1.0	2022. 3. 31	株式会社自動処理/高木祐介	初版作成

# 目次

改訂履歴	
目次	2
手順別目次	
データ活用までのフローを知る	
データ整備モデル事業について	4
作業に使用するアプリケーション/アドインを準備する	5
アプリケーション	
アドイン	5
データを変換する前に準備する(事前準備)	6
1. データ構造の検討	
2. データ構造の調査	8
<ol> <li>CSV データ構造の決定</li></ol>	
4. CSV テーラルスの次足	18
6. ケース一覧	20
データ構造を統一する	
データ形式別の作業フロー(2次元形式データ変換)	22
データ形式別の作業フロー(2次元形式にならないデータ)	23
Power Query によるデータ変換	
AI-OCR によるデータ変換	
データを適切に配置し、データの精度を高める	
データの格納	
データのマッピングデータのクレンジング	
データの正規化	
・ データの保存(エクスポート)	
<b>ナータの1未仔(エクスパート)</b> データベース接続の設定	
データのエクスポート	113
FAQ	
FAV	

# 手順別目次

参照先 データ整備モデル事 データ活用までのフローを知る **p.4** 業の概要を理解する p.4 データ整備モデル事業について 作業に使用するアプリケーション/アドインを準備する アプリケーション **p.5** アドイン p.5 データ化対象のデー データを変換する前の事前準備 p.6 夕をどのような構造、 1. データ構造の検討 p.7 形式にするのかを検 2. データ構造の調査 p.8 討する 3. CSV データ構造の決定 p.14 4. CSV データ形式の決定 p.16 5. 機械判別可能なデータ作成のチェックポイント p.18 6. ケース一覧 p.20 Excel や PDF からデ データ構造を統一する p.22 ータ変換を行う Power Query によるデータ変換 (データ内容によりケ ケース[A] 保存形式 XLS/縦横多段クロス集計 p.27 ース[A]~ケース[E] ケース[B] 保存形式 XLS/単純クロス集計 p.48 p.55 に分けられるので、最 ケース[C] 保存形式 XLS/単純表形式または複数ファイル 適な手法を選択しま ケース[D] 保存形式 XLS/1 シートに単純表形式が種類別 p.78 す。) に縦に3つ並んでいる、等 AI-OCR によるデータ変換 p.88 ケース[E] 保存形式 PDF 外部調達によって実施想定API、もしくは p.100 データを適切に配置し、データの精度を高める 利活用しやすいデータ p.101 データの格納 に近づける p.103 データのマッピング p.105 データのクレンジング p.108 データの正規化 データベースを作る データの保存(エクスポート) p.114 4



# データ活用までのフローを知る

# データ整備モデル事業について

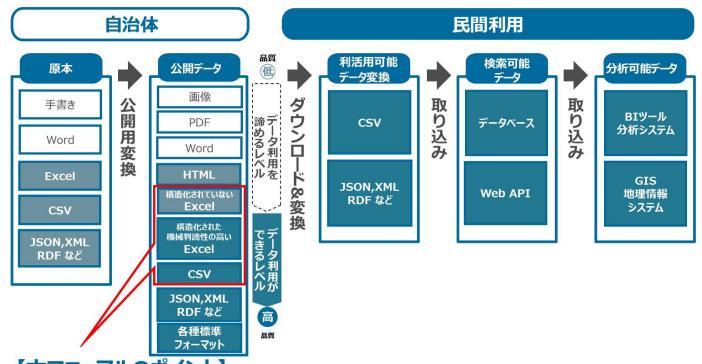
#### 民間利活用を目的にしたデータ化事業では

取り込まれたデータが、「業務システム用のマスタ」、「データ分析用の基礎データ」等として利用されます。企画の際の一過性な分析を行うケースを除き、民間利活用を考える場合、**簡単かつ安定的にデータベースに取り込む**ことができなければ、営業活動や、市場調査など実務に使う事が出来ません。その為、**データを公開する際には、データ利用者が利活用しやすい形(機械判読可能な形式)で提供する必要**があります。

データ化作業は、データベースを取り扱う**専門的な知識とスキルが必要**ですが、まずは**データ整備から活用までの一連の流れを把握することが重要**です。

例えば、手順 4 「データを適切に配置し、データの精度を高める」で、データベースのクレンジング作業は、データベースを精度良く取り扱うことのできる外注先(専門家)への発注や API の調達を前提にマニュアルを作成しています。

(マッピング、正規化する場合や、工程を自動化する場合も同様です。)



# 【本マニュアルのポイント】

# データ利活用を実現するために、まずは使いにくいExcelから使いやすいCSVへ

#### データ整備事業に使用するツール

本事業でのデータ化作業は、区市町村担当向けを想定している為、マニュアルや目視作業が必要な作業については、Excel で出来る範囲は極力 Excel にて対応を行うことを想定しています。

Excel 以外のツールについては「作業に使用するアプリケーション/アドインを準備する」(P.5)に記載しています。

またデータクレンジングなど作業に膨大な手間がかかるような処理については、API や外注など調達により作業を自動化する事を推奨します。

#### データ精度の扱い

Excel 等に格納されたデータは、安定的にデータベースに取り込むまでにさまざまな処理が必要になります。

正規化の過程では、人力でデータを加工したり、正規化ツールを使ったりする場面があります。東京都は 100%の精度で正規化することを必須とは考えていません。

# 作業に使用するアプリケーション/アドインを準備する

# アプリケーション

本マニュアルに記載された手順では、以下のアプリケーションを使用します。一部、有償アプリケーションの入手が必要になります。また、指定されたアプリケーションをお持ちでなくても、同等の機能を持つ互換アプリケーションを利用しても差し支えありません。

# Microsoft Excel(互換アプリケーション含む)

バージョン指定:なし

CSV ファイルを取り扱える互換アプリケーションであれば、本書記載の操作は可能です。

Excel や互換アプリケーションで CSV ファイルを取り扱う場合は、電話番号等先頭に 0(ゼロ)を含むデータに気をつける必要があります。 XLS 形式で保存すると、ゼロが消えてしまう場合があります。

データの誤変換を避けるためには、より本書の作業に適した CSV エディターでデータを変換いただくことをおすすめします。

例: CSV エディター

https://www.asukaze.net/soft/cassava/

# CLOVA OCR Reader / AI-OCR 用

バージョン指定:なし

数字や文字をデータとして取り込めない場合は、OCR機能を利用してデータを生成させることが必要になります。

本書では CLOVA OCR (Reader)の AI-OCR 機能を利用して、精度よくデータを生成する手順を案内しています。

詳細はこちら

https://clova.line.me/clova-ocr/

# Adobe Acrobat DC(互換アプリケーション含む)/AI-OCR用

利用の目的: AI-OCR に取り込むファイルを生成

想定作業:保存、ページの回転、書き出し、の圧縮、アクションリストの利用

変換後のバージョン: Acrobat10.0

本書では Adobe Acrobat DC の書き出し機能を利用し、AI-OCR に文字が読み込まれやすいデータを生成する手順を案内しています。

詳細はこちら

https://www.adobe.com/jp/acrobat/features.html

#### アドイン

# Power Query (Microsoft Excel の標準またはアドイン)

バージョン: **Excel2016 では標準機能**(追加インストール不要)。 Excel2010、Excel2013 ではアドインとしてインストール必要。 Power Query (以前のバージョンの Excel でデータを取得& 変換と呼ばれる) を使用すると、外部データをインポートまたは接続し、 そのデータをニーズに合った方法で(列の削除、データ型の変更、テーブルの結合など) 整形できます。

本書では手順2で、Excelでデータを利用する際に使用します。

詳細はこちら

https://support.microsoft.com/ja-jp/office/excel-%E3%81%AE-power-

 $query-\%E3\%81\%AB\%E3\%81\%A4\%E3\%81\%84\%E3\%81\%A6-7104\\fbee-9e62-4cb9-a02e-5bfb1a6c536a$ 

# 手順 **1**

# データを変換する前に準備する(事前準備)

データ構造・形式統一前の準備として、データ形式別に変換処理を実施します。

#### 1. データ構造の検討 (参考 p. 7)

データの利用目的等によって、必要とされるデータ構造は異なります。具体的に求められているデータ形式が定まっていない場合、データ構造がどうあるべきかを検討します。

#### 2. データ構造の調査 (参考 p. 8)

検討したデータ構造に利用するデータ構造を調査します。

- ① **メタデータ調査** メタデータルールと利用イメージの検討等を参照します。メタデータとは、あるデータに関する属性情報などの付帯的なデータのことです。データのタイトルや説明、公開日などデータを検索する際などに使用します。
- ② **データ構造調査** 政府 CIO ポータルに配置された標準データセット等をはじめ、各種データ連携モデルやデータカタログサイト、 公開 API 等を参照します。

#### 3. データフォーマットの検討 (参考 p. 14)

推奨データセットや収集データを比較し、妥当なフォーマットを検討します。

- ① メタデータ検討
- ② 標準 CSV フォーマット検討
- 4. 出典元の収集 (参考 p. 18)
- 5. 機械判別可能なデータ作成のチェックポイント (参考 p. 18)
- 6. ケース一覧 (参考 p. 20)

# 1. データ構造の検討

データの利用目的等によって、必要とされるデータ構造は異なります。データ構造がどうあるべきかを検討します。

## 関連部署、団体から情報を収集

ここでは医療機関データをデータ化するための情報として、保健所や医師会からの診療所データを入手する例を示します。

入手元だけでなく、データの対象によっても、データの構造が異なっていることがあります。これからデータ化する Excel ファイルをどのデータ構造にするのかを検討します。

#### 保健所からの病院データ

病房	元(R2.12.1 開設届)				
No	施股名称	施設所在地	施設方會	診療科目名称	施設TEL
1	1 医療法人社団實理会 東京国際大堀病院	三處布下遷從4丁目8番40号		内科、麻酔科、泌尿器科、放射線科、循環器内科、病理診断科、婦 人科	0422-47-1000
2	2 医療法人財団拡友会 三處病院	三歲市下運從5丁目1番12号		内科、外科、整形外科、競科、リハビリテーション科、泌尿器科	0422-47-6101
3	3 被原病院	三城市下運出6丁目13番10号		内科、消化器科	0422-46-2251
4	4 医療法人財団站生会 野村開院	三城市下運從8丁目3番6号		外科、内科、聖光外科、リハビリテーション科、形成外科、指尿検 重料、神経内料、起神経分科、超反應内料、循環機内料、呼吸器内 料、放射線計制料、運方内料、機和ケア内料、内視機内料、加用内 料	0422-47-4848
E	5 公益財団法人 井之頭病院	三成市上運出4丁目14番1号			0422-44-5331
6	<b>第</b> 年上十年も素金 □ □ □ 中年   東京 → □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □		内料、消化器内料、消化器内料、清理器内料、外料、整形外料、加 神経外料、競科、リハビリテーション料、放射線料、皮膚料、泌尿 器料、呼吸器内料、呼吸器分料、神経内料、肛門外料、血液内料、 乳脂外料、脂肪中核、内分泌料	0422-44-616	
ī	7 杏林大学医学部付属病院	三嶋市新川 6 丁目 2 0 番 2 号		内核、小型核、皮膚核、維神核、外核、紅神道の核、腫形が核、心 関血管外核、形成外核、小児外核、次児養核、酸核、耳角いんこう 核、運核、構入核、改計機核、海除核、血核仁能分核、少にリ テーション核、美電分核、リウマチ核、呼吸機分核、神経内核、故 急核、原理整核、中枢接体、溶性核、水性核	0422-47-5511
9	8 医療法人計団組水会 長盆川病院	三志市大沢2丁目20番36号		護地科 曲利 内利 六爾内利	0422-31-8600

#### データ構造①

#### 保健所からの診療所データ

診療所(R2.12.1 開設届)				
No 施設名称	施設所在地	施設方會	診療科目名称	施設TEL
1 医療法人社団健晶会 下川整形外科	三歳市井の頭1-24-14		内・茎・リハ、	0422-43-5933
2 浜田耳島咽喉科	三嶋市井の頭1-30-13	サトビル101	R	0422-76-8733
3 笹本医院	三歳市井の頭1-31-22		内・小・外	0422-44-5802
4 医療法人社団護祥会 废真整形外科	三歳市井の頭2-1-17	石伊ビル弐203	至・リウ・リハ、	0422-79-7600
5 三流台ヒルズクリニック	三点市井の頭2-1-17	石伊ビル弐202	耳・形・皮・美皮、	0422-76-7722
6 松川内科クリニック	三底市井の頭2-1-17	石伊ビル武201	内・消内	0422-70-5525
7 三底台眼科	三歳市井の頭2-1-17	石伊ビル武204	競	0422-40-3003
8 医療法人社団研道会 高水クリニック	三嶋市井の頭2-14-2	パークブリージェ井の頭101	内・小	0422-76-1232
9 みこしばクリニック	三鷹市井の頭2-19-25		小・精・神	0422-49-5300
10 医療法人社団慈昭会 石井医院	三底市井の頭2-32-37	1階、地下1階	内・胃・外・蓋・皮・放・消・小・肛、	0422-44-3090
11 牟礼の里駅前クリニック	三底市井の頭 2 - 7 - 9	栗原ビル1階	内・整・外・消内・復内・皮	0422-40-6054
12 藤林医院	三歳市井の頭3-12-15		内科、小児科、呼吸器内科、粧尿病内科、漢方内科	0422-43-4322
13 医律法人社团团庭医院	三歳市井の頭3-21-16	2 階	内・小・放	0422-43-8367
14 医療法人社団慈司会 若林医院	三鷹市井の頭4-16-10		内・小・消、	0422-43-0526
16 武蔵野みどり診療所	三底市井の頭5-7-36	SAKURA1 Na	呼内・内・外	0422-24-9933
16 ヨシコクリニック	三處市井口1-22-24		内·小	0422-32-5517
17 医療法人社団美々会 斉藤皮膚科	三歳市井口2-13-26	1 File	皮	0422-33-1030
18 医療法人社団加藤整形外科医院	三咸市井口2-3-1	1 Fik	至・リウ・リバ、	0422-31-3332
19 医療法人社団 かえでこどもクリニック	三鷹市井口3-6-16	アップルかえで通りビル 1F-A	小	0422-39-3306

#### データ構造②

# 2. データ構造の調査

検討したデータ構造に利用するデータ構造を調査します。

メタデータ調査 メタデータルールと利用イメージの検討等を参照します。

A: メタデータルールと利用イメージの検討等を参照 https://cio.go.jp/dp2021\_07 (p. 8)

データ構造調査 各種データ連係モデルやデータカタログサイト等を参照します。

B:標準データセット等を参照 https://cio.go.jp/policy-opendata

C: 政府・自治体のデータカタログサイトを参照 https://datasetsearch.research.google.com/

D: 政府・自治体のホームページを参照 https://www.google.com

## メタデータ調査(A: メタデータ調査ルールと利用イメージの検討等を参照)

政府 CIO ポータルサイトのメタデータルールと利用イメージの検討 (https://cio.go.jp/dp2021\_07) にアクセスします。

BregDCAT( https://joinup.ec.europa.eu/collection/access-base-registries/solution/abr-bregdcat-ap )を基本としてメタデータを検討します。

サンプルの含まれる日本語の資料として、メタデータ比較検討用コンセプトペーパー

( https://cio.go.jp/sites/default/files/uploads/documents/dp2021\_07\_att.pdf ) がある為、参照しながら検討します。

※BregDCAT 方式の場合、データカタログサイトのデファクトスタンダードである CKAN のマッピングを行う事が出来る為、データカタログサイトに取り込む際には以下を参照に取り込む事としてください(https://github.com/ckan/ckanext-dcat#rdf-dcat-to-ckan-dataset-mapping)。

『2.4 メタデータの記載例』に基づきメタデータについて検討します。

記載例	
カタログ	
管理 ID	CA00001
タイトル	文字情報基盤
サブタイトル	
説明	行政で用いられる人名漢字、法人用漢字、かな等の約6万文字の漢字です。
キーワード	漢字、戸籍統一文字、住民基本台帳ネットワーク統一文字、登記統一文字、変体仮名
テーマ分類	文字
対象地域	全国
提供者	文字情報技術促進協議会
公開日	2011-10-26
最終更新日	2020-8-26
更新頻度	不定期
言語	ja
公開範囲	公開
公開条件	
ライセンス	IPA フォントライセンス v1.0
利用規約	クリエイティブ・コモンズ 表示 - 継承 2.1 日本 ライセンス条件
ホームページ	https://moji.or.jp/mojikiban/
カタログレコード	
トピックス	文字
説明	行政で用いられる人名漢字、法人用漢字、かな等の約6万文字の漢字です。
公開日	2011-10-26

記載例	
最終更新日	2020-8-26
言語	ja
データセット	
管理 ID	CA00001-DST001
タイトル	MJ 文字情報一覧表
サブタイトル	
バージョン	Ver.006.01
説明	各文字に関するコード、読み、字母、画数等をまとめた情報。
キーワード	フォント、ヨミガナ、画数
対象地域	全国
対象期間	
分類	全ての業務
提供者	文字情報技術促進協議会
作成者	情報処理推進機構
	組織名:一般社団法人 文字情報技術促進協議会
連絡先情報	メールアドレス:info@moji.or.jp
	フォーム URL: https://moji.or.jp/about/contact/
タイプ	Strict Open XML
来歴情報	2020年10月に、情報処理推進機構から文字情報技術促進協議会に信託譲渡
品質評価	正確性、網羅性
品質測定結果	公務で使うのに十分な品質
公開日	2011-10-26
最終更新日	2020-8-26
更新頻度	不定期
言語	ја
公開範囲	公開
公開条件	
準拠する標準	
関連ドキュメント	
ランディングページ	
データサービス	
管理 ID	DSV00001
タイトル	文字情報基盤検索サービス
説明	
キーワード	
対象データセット	
提供者	
タイプ	web
公開範囲	公開
公開条件	
ライセンス	
準拠する標準	
関連ドキュメント	
エンドポイント URL	
ランディングページ	
配信	mii 00001 yley zin
タイトル	mji.00601-xlsx.zip

記載例	
説明	
アクセスサービス	
バイトサイズ	7.3M
圧縮形式	zip
メディアタイプ	
公開日	
最終更新日	2019-05-01
期間	
ステータス	配信中
言語	ja
ライセンス	
利用規約	
準拠する標準	
関連ドキュメント	
アクセス URL	https://moji.or.jp/mojikiban/mjlist/
ダウンロード URL	https://moji.or.jp/wp-content/mojikiban/oscdl/mji.00601-xlsx.zip

メタデータの検討に当たっては、既に定義されているメタデータを参照します。

例えば DATA GO JP のサイト (https://www.data.go.jp/) にアクセスし、メタ情報を調査する場合には以下のように調査を行います。

拠点のデータセットをクリックします。

必要なメタデータを入手します。



メタデータ①

次に Google Dataset Search (https://www.data.go.jp/) を利用する方法もあります。

GoogleDataset Search では CKAN に登録されているデータが検索されるため、データカタログサイトに登録されているデータが全て検索されます。

例えば「医療機関」を調査する場合には、以下の通り検索を実施します。

医療機関データをクリックします。

必要なメタデータを入手します。



メタデータ②

## データ構造調査(B:推奨データセットを参照)

政府 CIO ポータルサイトのオープンデータ (https://cio.go.jp/policy-opendata) にアクセスします。

推奨データセットに含まれるデータに関しては、推奨データセットを基本として、拡張を行います。

推奨データセットは、1 データ 1 ファイルを基本としている為、推奨データセットが含まれる場合には、手順にしたがってデータ定義を行います。

推奨データセットに含まれないデータ構造については、別途他のデータ構造を確認して検討を進めます。

データ項目までは推奨データセットに記載されていますが、文字列長、データ形式などデータの表現についての記載は少ない為、別途検討する必要があります。

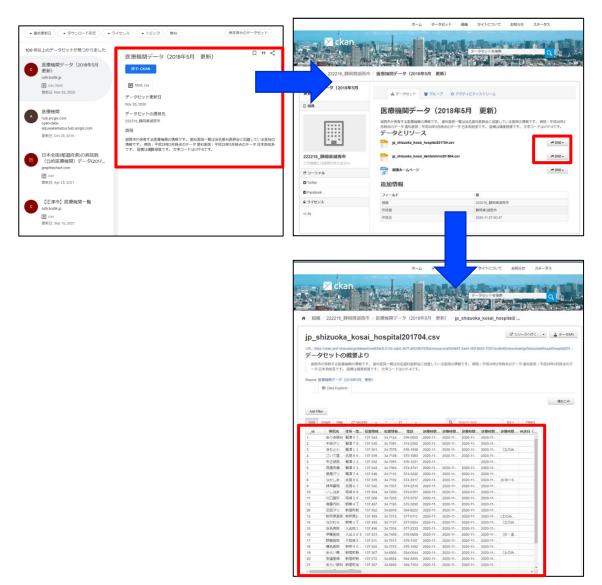


# データ構造調査(C:政府・自治体のデータカタログサイトを参照)

Google Dataset Search のサイト (https://www.data.go.jp/) にアクセスし、「医療機関」を調査します。
GoogleDataset Search では CKAN に登録されているデータが検索されるため、データカタログサイトに登録されているデータが全て検索されます。

医療機関データをクリックします。

必要なメタデータを入手します。



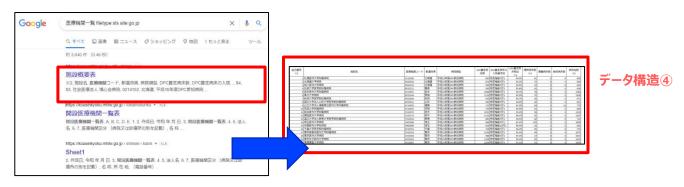
データ構造③

## データ構造調査(D:政府・自治体のホームページを参照)

Google 検索 (https://www.google.com/) にアクセスし、「**医療機関**一覧 filetype:xls site:go.jp」で検索します。 複数のデータ構造を閲覧し、1 行目が項目、2 行目以降がデータになっているデータを探します。

拠点の検索結果をクリックします。

必要な拠点データを入手します。



Google 検索 (https://www.google.com/) にアクセスし、「医療機関一覧 filetype:xls site:lq.jp」で検索します。

拠点の検索結果をクリックします。



必要と想われる拠点データを入手しますが・・・



データにより、1 行目が項目、2 行目以降がデータになっていない 場合があれば、 採用しません。

拠点の検索結果をクリックします。



必要な拠点データを入手します。

※他の都道府県・政令市・中核市に所在す				
	Office	医療機関	については当該自治体が指定・公示しています。	
医療機関名	v		医療機関所在地	T
医療法人 野瀬クリニック	4	tΣ	梅田1-2-2-200	
O Medical Clinic Osaka(ディー メディカル うりニック オオサカ)	4	ťΖ	梅田2-5-25 ハービスPLAZA 4F	
医療法人湖崎会 湖崎眼科梅田分院	1	tŒ	梅田3-1-1 サウスゲートビル17階	
医療法人伯肌会 大阪中央病院	4	区	梅田3-3-30	
医療法人知音会 堂島内科・消化器内科クリニック	4	ఠ	堂島2-4-27	
一般財団法人住友病院		ĽΚ	中之島5-3-20	
療法人緒秀会 インフュージョン(点滴)クリニック	4	区	大深町3-1 グランフロント大阪タワーB 9F	
法人凉您会 梅北眼科	4	r <u>ix</u>	大深町3-1 グランフロント大阪北館BIF	
法人思赐财団済生会支部 大阪府済生会中津病院	4	ĽŒ	芝田2-10-39	
太人順宵会 田中北樽田クリニック	1	ఠ포	芝田2-8-10 光栄ビル3階	
おさか内分泌診療所	4	βŒ	茶屋町4-6 タケムラビル4階	
医療法人 羅平胃腸科 放射線科	1	tΣ	角田町8-47 阪急グランドビル22F	
医療法人 社団皓歯会 阪急グランドビル診療所		区	角田町8-47 阪急グランドビル22階	
社会医療法人行同医学研究会 行同病院	1	tΣ	浮田2-2-3	
野町クリニック	1	区	神山町1-7 扇町メディックスモール4階	
ちこ歯科	4	ĽΖ	神山町1-7-3D	
米田内科胃腸科	4	ĽΚ	<b>管栄町5-17</b>	
医療法人 八杉クリニック	4	区	池田町1-75 ストークマンション101~103	
医療法人 天五診療所		ĽΚ	池田町6-10-201	
本出診療所		ĽΣ	同心1-8-3	
れんりこ歯科クリニック	1	区区	西天湯3-2-15 竹田ピル2F	

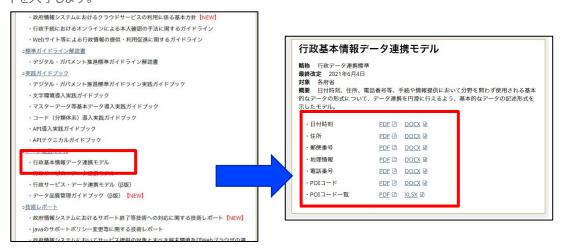
データ構造(5)

東京都のオープンデータカタログサイト( https://portal.data.metro.tokyo.lg.jp/ )にアクセスします。

標準ガイドライン群の中から、

データセットの中から、必要なデータセットを入手します。

「行政基本情報データ連携モデル」推奨データセットを入手します。



# 3. CSV データ構造の決定

「データ構造の調査」の結果を元にまとめます。

## CSV データフォーマット検討

#### 1. 推奨データセットと収集データを比較し、妥当なフォーマットを検討する

標準データセット データ構造① (P. 7)

(P. 11)

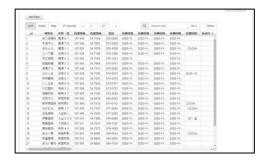
都道府県 コード又は 市区町村 コード NO 都道府県 名 市区町村 名 名称 名称\_カナ 医療機関 の種類 住所 方書 緯度 奴庇

用詞(#2.12.1 開設庫)				
No 無計名称	境别形在地	別形力量	影響がきなか	MRTEL
1. 医療主人た工力協会 東京国际大心共和	三次中下選上47日日日40号		内的、用杂种、以应量的、放射性的、推理器内的、用键数形的、排 人的	0422-47-1000
2. 医康夫人物学以实在 三共传统	立成者で連出る丁姓に会にこ号		門外、共和、盟和共和、兼料、テハビリテーション和、定定豊和	0422-47-6000
1.位于元年	三月万下進州4丁日13年10年		F6. 3680	0422-46-2251
43#1/hIS12 BINK	医成性下腺业化下面 3 整化等		内札、内札、福祉の利、リハビリチーション和、別点の利、日に除 重払、押払力払、以付出の利、明日用力払、海関目力払、予以前力 払、記述前別回利、直分内払、決力ケア内払、円式前り払、減乏内 払	0422-47-4548
5 公益計画は人 当に際内間	型用市上建区47日14第1号			0422-44-5331
《京學集》、位置水產金 三月中央開開	正成年上課出57首23巻10号		門外、現住的内容、現住者内容、確保能力な、内容、監視内容、 情報の対象、部外、ジハビリテーションは、配料制体、皮膚性、皮膚 新外、子供器内容、子供器内容、溶液内容、医内内容、直及内容、 別場があ、再合す、認定れた、対当分科	0422-84-6063
7. 四州大学区学副作業機構	三点の前別6丁目2の参2号		四外、の元利、成業外、精神は、カ州、松中電力所、整合外外、心 設定者が終、影成内外、小光内外、北方線体、燃料、原合外によっ か、変数、減入が、放射線が、水内の、水中では一分、 ケーションが、東京が利、フローカー、アイモデ州、同位内が、加 を外、電影型化外、上面製料、子の影響、水上銀子	6422-47-8611
京西東京人社団会会会 美谷川市院 ·	世界市大大大工会との乗りを受		BINS, BALL DALL CROSS	6422-31-5500

データ構造② (P. 7)

D #R7(R2.12.1 開致(編)				
io SRAM	BRRETI	お記させ	飲得料当名村	MRTEL
1 医療法人は国際基金 下川屋町内町	正成物式の第1-24-14		カ・草・サハ	0422-43-5033
2 HEROSEN	工用物体の第1-30-13	910101	×	6422-76-8733
5 性中进门	三点考井の領1-31-22		A-4-5	6422-44-5862
4 医奎耳人拉氏螺状体 交流放射的柱	出席市井の開2-1-17	89'040203	第・リテ・リハ、	0422-79-7660
1 三衣台にルズテリニック	工术市中の第2-1-17	世界にん見るのと	寒・形・皮・美女、	0422-76-7722
4 松川門科クリニック	正成年件の第2-1-17	程序にかけまり1	5-95	6422-10-5525
7 正成台號5	正成等計の第2-1-17	石伊ビル代294	et .	6422-40-3068
※原本人は言語連合 事水とリニック	工机中排均開立-14-2	バータブリージェ持の第101	rt - o.	6400-76-0000
<b>ま</b> みこしばクリニック	王京市井の第2-19-25		0.10.10	6422-49-5500
10 医家体人和国会社会 石井美物	三点布井の様2-32-37	19. NT19	内・胃・水・臓・皮・状・治・水・肌	0422-44-3060
11 作礼の意影的テリニッテ	三点市場の第2-7-9	京事でか1階	内・整・所・河内・保内・疣	0422-40-6064
12 無外区院	工具市井の領3-12-15		内科、小児科、牙奇菌内科、粒层类内科、重方内科	6422-43-4322
13 医甲基人社团共和国民	正成年十の領3-21-16	28	7 · 4 · 3	0422-43-6367
14 医療法人社会の司会 正性医院	工机中計の第4-16-10		PL-0-19.	0422-43-0506
15 女変がみどり目標所	工机中排の開5-7-36	DAKURAIM	96-6-8	6400-04-9933
16 ランコクリニック	三点申井口1-22-24		M-a	6422-32-5517
1) 医療法人心医療不会 光極大明和	三点市中口2-13-26	18	π	0422-39-0000
12 医蒙洛人拉耳拉等整形外科医院	五点市井口2-3-1	1.90	種・ジウ・ジハ、	6422-21-3332
19 選挙導入知道 かんでこどもケリニック	工机市井口3-6-16	アップルかんで送りたル コドース	0	6422-39-3366
William Co.				

データ構造③ (P. 11)



データ構造④ (P. 11)

医療機関名	v	医療機関所在地
医療法人 野瀬クリニック	#FIX	梅田1-2-2-200
D Medical Clinic Osaka(ディー メディカル クリニック オオザカ)	北区	梅田2-5-25 ハービスPLAZA 4F
医療法人凝畸会 湖崎県料積田分院	北区	梅田3-1-1 サウスゲートビル17階
医療法人伯里会 大阪中央病院	北区	梅田3-3-30
医療法人知音会 堂島内料・消化器内料クリニック	4比区	堂島2-4-27
一般附近法人住友病院	地区	中之島5-3-20
医療法人結秀会 インフュージョン(点滴)クリニック	北区	大塚町3-1 グランフロント大阪タワーB 9F
医療法人流悠会 梅北眼科	4FIE	大深町3-1 グランフロント大阪北館81F
社会福祉法人思赐财団済生会支部 大阪府済生会中津病院	北区	芝田2-10-39
医療法人順育会 田中北樽田クリニック	16区	芝田2-8-10 光栄ビル3階
おおきか内分泌診療所	1 <b>U</b> Z	茶屋町4-6 タケムラビル4階
医療法人 藤平雪鵑科 放射線科	11/1X	角田町8-47 阪急グランドビル22F
医療法人 社団結歯会 阪急グランドビル診療所	北区	角田町8-47 阪急グランドビル22階
社会医療法人行用医学研究会 行間病院	16区	淨田2-2-3
馬町クリニック	16K	神山町1-7 扇町メディックスモール4階
きちこ歯科	#UK	神山町1-7-3D
米田内科胃腸科	1UX	<b>管</b> 栄町5-17
医療法人 八杉クリニック	北区	池田町1-75 ストークマンション101~103
医療法人 天五診療所	地区	池田町6-10-201
本出診療所	4UE	同心1-8-3
れんりこ歯科クリニック	地区	西天湯1-2-15 竹田ビル2F

データ構造⑤ (P. 11)

ATER TI	mos.	88983-7		ANNE	PCB2K 并数	の代数主向状力 入所書す料	NAIS.	BHICKS 14	PRILID	MERKE	MANUEL III
	元級的な対抗性の対	831730		THOUSE PORT	20	WETTER-771	34.5%	- 41		-	116
	と指摘とTVRDE	Jaconico.		するいな名 一年10年18		or transcrite	YLA'N				346
	NUMBER PROF.	900,004	1274	平成31年第200年10月日		TR.T. TRATE T-171					
	<b>公共・予及予約を集相</b> 性		85	<b>で出かる場合の表示的</b>		WENE-CO.					
			2018	<b>学成の在後の中の教育</b>		WE THAT : TT					
	BULL PAIR	(80) (1) (6)	201	学式の研究の大学の研究		TR.\$1640.1751					
	NG+YEYERENG		trill	子式の英葉から野な母は		TO E TRANSPORT I					
				で見いなないのかの利用		07540-75s					
	COCAPSA機能用交換NOC字形型的性		16.00	VICINES HISTORY		PRE-276-071					
	OLE PRIMER	301000	264	平式10年度200季10年後		TRIMETOIL					900
	DeENATHERS.		85.4	<b>V.Kook Work 新公司</b>		·特里特的-771					
	ROS HA THIX			中式いなない。中の内性		10 T TEACH (77)	26.4%				
	<b>第二大学会人的性大学会学的対象的状</b>		DIT.	学成別性集団の創設の機能		19.2 Nam (7)1	94.5%				
	地工资和大学相談			子がおるをからからのは		** T 14.80 - 1/2					
			143	VANAS MERK		VPT MACHOTTS					845
	子葉大学哲学的行業的技	SHOWA	TR	平成の作用(34) 新田田田		(175 MAT 2764E)	94.2%				716
	BERNSON FREME		819	平式10年第10年的 <b>第</b>		TOTAL PROPERTY.					3070
	RESILEYAN		261	<b>でだいは第四十二回日</b>		祖子協能で行る					
	重常女子包持人学病院	\$96,5642	30.0	平成10年第10年の開発		世史を見る	25.4%				1424
			81								

#### 2. CSV フォーマットのデータ項目を整理する

#### 東京都データプラットフォーム データ整備モデル事業

			対象		対象	対象	対象	対象	対象	対象		
				1	5	7	1 0	10				
	_											
	ゲータセット名	X # 10 X	X FAX	医療機関	医療機関	医療検禁	医療物質	医療機関	医療物質	2.792		
						<b>有模区</b>				正在布		
	形式					elsx				risk		
			2 机块区的推开一覧 维纳的推科目的									
	#HTURL	-	-	-	-	-	-	-	_	-		
	ALTON L											
									病院、診療所が別シート			
HR/No	注意点	10 TO 4. 57 W 19	[数據報告] (特殊文は無い)毎に数数	COMMAND ACTION OF THE	CROSSESS TO COMPANY OF THE PARTY THAT	10 10 A 10 W W	病院、動療所、歯科動療所が別シード			GFI開布開放?		
41110	2.0		- David Cold Cold Cold	TAXABLE STATE		変要の	THE DISTRICT STREET		00.000000	07,000,000		
						計算を						
		開設中	MIN+			0350						
テータス	α		<b>第点</b>			英点		_				
	_					女童 ウ						
						計可用						
		開設の	解設中		Willia	9350						
テータス	α		第立			第上						
He4	- a	0	0	0	0	0				_		
T#9	a			0	0	0						
人格						tr .	4	2	di di			
IAM IAM				Α	Δ	9	4		Δ	-		
IH-6	a	a .		Ô.	Ô	<u>a</u>	Δ.	1	<u>a</u>	<u>a</u>		
127		0	0	0	0	0	O .	0	0	0		
	Œ	-		-	0	-	-	-	-	-		
使曲号	Œ	0		0	0	0	-		-	-		
res:		0				0	-		-	-		
rane	\$7.16E	0		0		0	0		0	0		
在地	78	Δ	Δ	Δ		Δ	Δ		Δ	Δ		
fate	無地号	Δ	Δ	Δ		Δ	Δ		Δ	Δ		
16%	建物	Δ	Δ	Δ		Δ	Δ		Δ	Δ		
供電話音号	α	Δ	Δ	Δ		Δ	Δ		0	0		
<b>以</b> 権	表記揺れ	有	有	有	有	-	-		-	-		
人格	Œ	Δ					-					
技術氏名	Œ	0	0	0		-	-		-	-		
人格	世紀揺れ	-	-	-		tξ	-		-	-		
人格	Œ	-	-	-		Δ	-		-	-		
<b>英老氏名</b>	Œ	-	-	-	-	0	-	-	-	-		
理會氏名	Œ	Δ	Δ	Δ		Δ	-		-	-		
443	Œ	O .	0			0	-	-	-	-		
更事由名称	α	-	-		-	-	-		-	-		
更属出年月日	Œ	-	-	Δ	-	-	-	-	-	-		
更許可率月日	α	-	-	Δ		-	-		-	-		
放射可半月目	Œ	Δ	Δ	-		-	-	-	-	-		
国种可年月日	α	-	-	-		0	-		-	-		
批准出年月日	Œ	0	0	-	0	-	-	-	-	-		
可 <b>然</b> 似乎月日	Œ	-	-	-		Δ	-		-	-		
粉年月日	α	0	0	-	0	0	-	-	-	-		
止催出年月日	Œ	Λ.	Δ	-	^	Δ.	-	-	-	-		
可用了年月日	g g		-	-		۸	-	-	-	-		
心無 / 平月日  止薬出半月日	a a	Α	Δ.		^	Δ.			-			
止無出 中月日 止年月日	a	in the second	Α		Δ	Δ.			_	_		
原料日 原料日	a a	0	A. C.		-	0.0						
	a a	0	0		0	-			0	0		
考 イアウトバターン	18	-	-	-		-	-		-	-		
		4	4	4	4	4	4	4	4	4		

#### データ整備マニュアル

今回はデータ処理を行う 為のデータフォーマットを検 討した為に、複数市区町 村にてデータが格納されて いる共通項目をデータ項 目とする方針としていま す。

#### 3. CSV フォーマットのデータ項目を整理する

#### 推奨データセット 収集データマッピング 項目検討結果

		## ## ## ##   1
都道府県コード又は市区町村コード	市区町村コード	都道府県コード又は市区町村コード
NO		NO
都道府県名		都道府県名
	都道府県コード	都道府県コード
市区町村名		市区町村名
	市区町村コード	市区町村コード
名称	施設名	名称
名称_カナ	カタカナ	名称_カナ
医療機関の種類		医療機関の種類
	郵便番号	郵便番号
住所	住所	住所
方書		方書
緯度	緯度	緯度
経度	経度	経度
電話番号	代表電話番号	電話番号
内線番号		内線番号
FAX番号		FAX番号
	開設者法人格	開設者法人格
	開設者	開設者
法人番号	法人の場合、法人番号	法人番号
法人の名称	法人の場合、法人名	法人の名称
医療機関コード		医療機関コード
診療曜日		診療曜日
診療開始時間		診療開始時間
診療終了時間		診療終了時間
診療日時特記事項		診療日時特記事項
時間外における対応		時間外における対応
	業種種別	業種種別
診療科目	診療科目	診療科目
病床数		病床数
URL		URL
備考		備考

今回は推奨データセットと 収集データをマッピングした 結果を図の通りマッピング し、CSV フォーマットを検 討しています。

# 4. CSV データ形式の決定

ISO、JIS、WWC、RFCを検索

ISO、JIS、WWC、RFCを参照しデータ形式を

データ連携モデル

https://cio.go.jp/guides https://www.google.com

#### 1. 項目ごとにデータ形式を検討する(行政基本情報データ連携モデル)

行政基本情報データ連携モデルのデータフォーマットを確認し、一般的な項目については行政基本情報データ連携モデルを参照します。 https://cio.go.jp/guides

例えば、日付については ISO8601 及び JIS X 0301 (日付及び時刻の表記) に準拠するようにします。



#### 1.1 日付

日付のデータは以下の形式とする。半角を使用する。

#### YYYY-MM-DD

YYYY: 西暦年4桁

MM :月2桁 (1桁の場合には前に0をつける) DD :日2桁 (1桁の場合には前に0をつける)

#### YYYY-MM-DDTHH:MM:SS+hh:mm

hh : UTC に対して進んでいる「時」

mm : UTC に対して進んでいる「分」(通常は00)

例) 2017-09-01T09:30:00+09:00 (日本) 例) 2017-09-01T00:30:00=+ (英国)

#### 2. 項目ごとにデータ形式を検討する(ISO、JIS)

ISO、JISにて定義されているコードがあるか確認します。

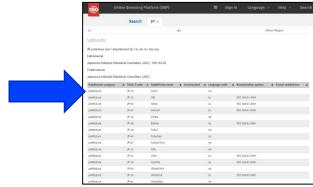
性別、都道府県、市区町村など既にコード化されているデータが存在します。その場合には出来る限り、既存のデータに追加してコードを付与して情報が公開される事が望ましいです。

例えば都道府県の標準コードを検索するには下記のように検索して調査を行います。

Google 検索 ( https://www.google.com/ ) にアクセスし、「"都道府県" JIS 標準 OR ISO 標準」で検索します。

調査する際には、出典元を確認し、最新のコードで あることを確認する事が必要です。





#### 3. 項目ごとにデータの表現を検討する(共通語彙基盤)

ISO、JISにて定義されているコードがあるか確認します。

性別、都道府県、市区町村など既にコード化されているデータが存在します。その場合には出来る限り、既存のデータに追加してコードを付与して情報が公開される事が望ましいです。

例えば都道府県の標準コードを検索するには下記のように検索して調査を行います。

コア語彙 (https://imi.go.jp/core/) や共通語彙基盤コア語彙 2.4.2 (https://imi.go.jp/ns/core/Core242.html) にアクセスし、表現を検索します。

人型	名称型	氏名型	住所型
連絡先型	電話番号型	組織型	業務組織型
法人型	人数型	構成員型	組織関連型
場所型	座標型	ID型	ID体系型
製品型	製品個品型	地物型	土地型
施設型	施設関連型	建物型	駐車場型
設備型	イベント型	活動型	関与型
測定值型	数量型	容量型	面積型
重量型	長さ型	時間型	金額型
価格型	実体型	状況型	日時型
日付型	期間型	期間スケジュール型	イベントスケジュール型
定期スケジュール型	詳細スケジュール型	詳細スケジュール規則型	コード型
単位コード型	コードリスト型	アクセス型	アクセス区間型
概念型	事物型	参照型	記述型
サービス型	文書型	対象型	制約型
コード制約型	範囲制約型	期間割約型	



#### 4. データ構造とデータ形式の検討結果をまとめる

これまでの結果をまとめ、データの表現についての検討結果を一元化します。

#### 推奨データセット 収集データマッピング 項目検討結果 データの表現についての検討結果

	I		
都道府県コード又は市区町村コード	市区町村コード	都道府県コード又は市区町村コード	, , , , ,
NO		NO	数字
都道府県名		都道府県名	総務省『全国地方公共団体コード』に準拠する
	都道府県コード	都道府県コード	JISX0401に準拠する
市区町村名		市区町村名	総務省『全国地方公共団体コード』に準拠する
	市区町村コード	市区町村コード	JISX0402に準拠する
名称	施設名	名称	英数字は半角、ひらがなカタカナ漢字は全角
名称_カナ	施設名_カタカナ	名称_カナ	名称からカタカナを作成する
医療機関の種類		医療機関の種類	
	郵便番号	郵便番号	行政基本情報データ連携モデルに準拠する
住所	住所	住所	行政基本情報データ連携モデルに準拠する
方書		方書	
緯度	緯度	緯度	行政基本情報データ連携モデルに準拠する
経度	経度	経度	行政基本情報データ連携モデルに準拠する
電話番号	代表電話番号	電話番号	行政基本情報データ連携モデルに準拠する
内線番号		内線番号	
FAX番号		FAX番号	
	開設者法人格	開設者法人格	法令に準拠する
	開設者	開設者	英数字は半角、ひらがなカタカナ漢字は全角
法人番号	法人の場合、法人番号	法人番号	国税庁法人番号公表サイトの公開情報に準拠する
法人の名称	法人の場合、法人名	法人の名称	国税庁法人番号公表サイトの公開情報に準拠する
医療機関コード		医療機関コード	
診療曜日		診療曜日	
診療開始時間		診療開始時間	
診療終了時間		診療終了時間	
診療日時特記事項		診療日時特記事項	
時間外における対応		時間外における対応	
	業種種別	業種種別	確認中
診療科目	診療科目	診療科目	診療科区分(厚生労働省 様式コード表)に準拠する
病床数		病床数	
URL		URL	
備考		備考	

# 5. 機械判読可能なデータ作成のチェックポイント

データが機械判読可能かどうかで、データ化の精度が左右されます。精度の高いデータ化を目指すには、目視も含むチェックが必要です。 下記を確認してください。

#### ファイル形式

- ファイル形式は Excel か CSV となっているか  $\Box 1$
- Word 形式のファイルが埋め込まれた場合は Excel ファイル等にコピーアンドペーストして保存したか  $\square 2$

#### データ・項目

#### 構造 □3 □4 セルの結合は解除したか ※1 □5 オブジェクトを使用している場合、正しく変換されているか □6 改行コードが入っていないか ※2※3 ヘッダーは1行に収まっているか(可能な限りヘッダーは1行とする方が良い)※4※5※6 □7 同一列に内容の異なるデータが混在していないか ※6 □8 □9 月毎、年毎等でブックを分ける場合は、ブック名、シート名のフォーマットとヘッダー列名等レイアウトは揃っているか 数値 □10 数値データは(文字列を含まない)数値属性、かつ半角数字となっているか ※2 $\Box 11$ 数式を使用している場合は、数値に修正しているか データの単位を記載しているか $\square 12$ □13 年月日記載の際は、半角数字で記載しているか(全角数字、漢数字、元年は使用しない) ※2 $\Box 14$ 英数字記載の際は半角英数字で記載しているか ※2 文字·体裁 □15 スペースや改行等で体裁を整えている場合、対策を行ったか ※2※3 項目名等が省略されている場合、不足情報を補ったか □16 $\Box$ 17 機種依存文字を使用していないか

- e-Stat の時間軸コードの表記、西暦表記文は和暦に西暦の併記がされているか □18
- 地域コードまたは地域名称が表記されているか □19
- 数値データの同一列内に特殊記号(秘匿等)が含まれていないか □20
- $\square$ 21 取消線や装飾文字を使用していないか ※1
- データとして取り込むべき情報がコメントにて記載されていないか ※7  $\square$ 22

#### ※1~※6でエラーが発生する場合は下記のような整形処理が必要となります。

- ※1 フォーマットやレイアウトが異なる場合は個別に整形処理する
- ※2 値を置換する
- ※3 改行コードは半角スペースに変換する(改行コードが入っていると CSV 化する際、別行データとして扱われるため)
- ※4 列を結合/分離する
- ※5 行列を入れ替える
- ※6 同一列に内容の異なるデータが混在する場合、同一内容で列を分割するため、行列の入れ替えや列の分割等複雑な整形処 理が必要となる
- ※7 Excel マクロで処理する

## 表の構成

- □23 データが分断されていないか
- □24 1シートに複数の表が掲載されていないか
- □25 データが増える際に、列ではなく行が増える構成か(データ項目が縦に並ぶ「縦持ち」の表になっているか)

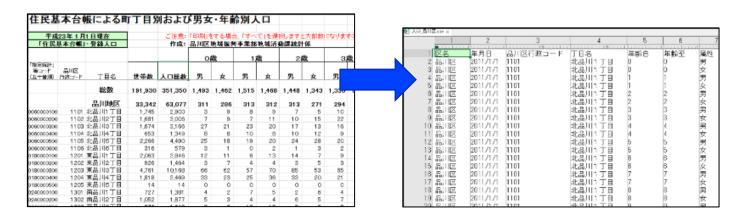
出典( $\Box$ 1、 $\Box$ 3~ $\Box$ 5、 $\Box$ 10~ $\Box$ 12、 $\Box$ 15~ $\Box$ 20、 $\Box$ 23、 $\Box$ 24):総務省統計局.統計表における機械判読可能なデータの表記方法の統一ルールの策定.報道資料,令和 2 年 12 月 18 日,別紙,p.1-22. (online), https://www.soumu.go.jp/menu\_news/s-news/01toukatsu01\_02000186.html, (参照 2022-01-31) / 出典をもとに、チェックリスト形式に改変。

#### 6. ケース一覧

データ変換作業において、本書ではデータ保持の状況ごとに 5 種類に分類し、各ケースにあわせてデータ化手順を掲載しています。 以下、5 種類のケースを示します。

# ケース[A]: 保存形式 XLS/縦横多段クロス集計

Excel 上でコピー、貼り付け等の操作や、結合セルを解除することで CSV データに変換できるデータのレイアウトの Excel データを、 CSV データに変換します。



# ケース[B]:保存形式 XLS/単純クロス集計

1 行目が項目行、2 行目以降がデータ行となり、1 レコードが 1 行で保持されたデータを、CSV データに変換します。



# ケース[C]:保存形式 XLS/単純表形式または複数ファイル

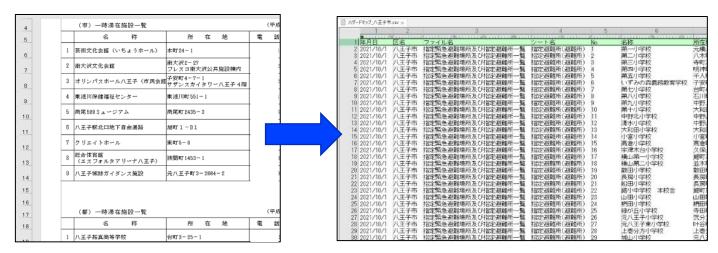
1 行目が項目行、2 行目以降がデータ行となった、表形式のデータを、CSV データに変換します。 項目の粒度や目的が混在している場合は、変換したい CSV ファイルにあわせて表を分割する必要があります。



# ケース[D]:保存形式 XLS/1シートに単純表形式が種類別に縦に3つ並んでいる、等

1シートに複数の単純表形式が混在しているレイアウトのデータを、CSV データに変換します。

CSV データにするにはファイルを分割する必要がある。例えば同一シート内に単純表形式が縦に3つあれば、1シートに1表、かつ3ファイルに分割する必要があります。



# ケース[E]: 保存形式 PDF

・AI-OCR を利用して読み取り、や PDF を CSV データに変換します。

区分		****		人口		前月比	増減	区分		世帯数	人口			######################################	
区分		世帯数	男	女	āt	世帯数	人口				男	女	āt	世帯数	시
	日本人	119,108	99,438	113,124	212,562	20	113	港区総数	日本人	119,108	99,438	113,124	212,562	20	1
港区総数	外国人	_	10,374	9,187	19,561	_	65		外国人	_	10,374	9,187	19,561	_	
	合計	_	109,812	122,311	232,123	_	178		合計	_	109,812	122,311	232,123	_	
	日本人	20,888	16,589	17,998	34,587	22	27	芝地区総合支所管内	日本人	20,888	16,589	17,998	34,587	22	
芝地区総合支所管内	外国人		1,551	1,240	2,791	_	10		外国人	_	1,551	1,240	2,791	_	
	合計	_	18,140	19,238	37,378	_	37		合計	_	18,140	19,238	37,378	_	
	日本人	27,638	21,556	24,974			19	麻布地区総合支所管内	日本人	27,638	21,556	24,974	46,530	△ 7	
麻布地区総合支所管内	外国人		3,980	3,577	7,557		74		外国人	_	3,980	3,577	7,557	_	
0_1002071013	合計	_	25,536	28,551	54,087	_			合計	_	25,536	28,551	54,087	_	
	日本人	17,570	14,415	17,093	31,508	13	14	赤坂地区総合支所管内	日本人	17,570	14,415	17,093	31,508	13	
赤坂地区総合支所管内	外国人		1,742	1,577	3,319		17		外国人	_	1,742	1,577	3,319	_	
MAX CERCO X / / I I I	合計	_	16,157	18,670			31		合計	_	16,157	18,670	34,827	_	
	日本人	29,000	23,751	29,574	53,325	△ 8	15	高輪地区総合支所管内	日本人	29,000	23,751	29,574	53,325	△ 8	
高輪地区総合支所管内	外国人	29,000	1,487	1,370	2,857		△ 19		外国人	_	1,487	1,370	2,857	_	1
同無地区心口又用目的	加国人		1,407	1,370	2,037		△ 19								

手順 **2** 

# データ構造を統一する

各種データの構造を統一します。

1. PowerQurery によるデータ変換 (参考 p. 7)

Excel を利用してデータ構造を統一します。

ケース[A] 保存形式 XLS/縦横多段クロス集計 (参考 p. 7)

ケース[B] 保存形式 XLS/単純クロス集計 (参考 p. 7)

ケース[C] 保存形式 XLS/単純表形式または複数ファイル <u>(参考 p. 7)</u>

ケース[D] 保存形式 XLS/1 シートに単純表形式が種類別に縦に3つ並んでいる (参考 p. 7)

2. AI-OCR によるデータ変換 (参考 p. 8)

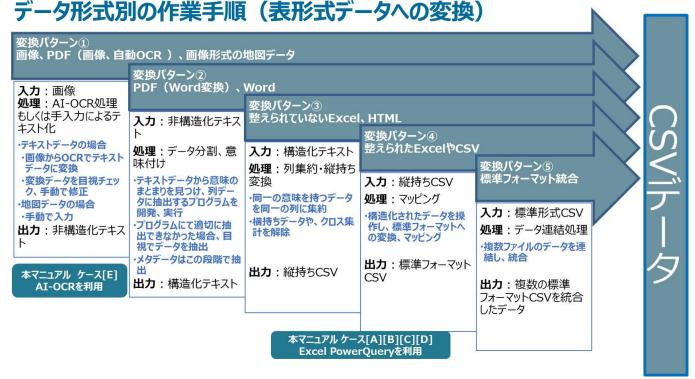
AI-OCR を利用してデータ構造を統一します。

ケース[E] 保存形式 PDF (参考 p. 7)

# データ形式別の作業フロー(表形式データ変換)

最終的に公開データに変換する為には、データの種類に応じて複数のパターンがあり、この章では行政内でよく使われる Excel データと」 PDF データを CSV 構造を統一する手順を説明しています。この手順を実施する事でその後のデータクレンジングなどのデータの精度向上手順を単純化する事が出来ます。

下記の図では「変換パターン①」の場合は変換パターン①から始まり、②、③、④、⑤までの作業がマッピング前までの CSV データ作成までに必要な手順となっており、「変換パターン③」の場合は、変換パターン③から始まり、④、⑤までの作業が必要な事を示しています。 (資料別目次(p. 3)では手順1、2として記載)



# データ形式別の作業フロー(表形式にならないデータ)

他にも既に標準化されたデータやオープンデータにできないデータも存在しますが、本資料では取り扱いません

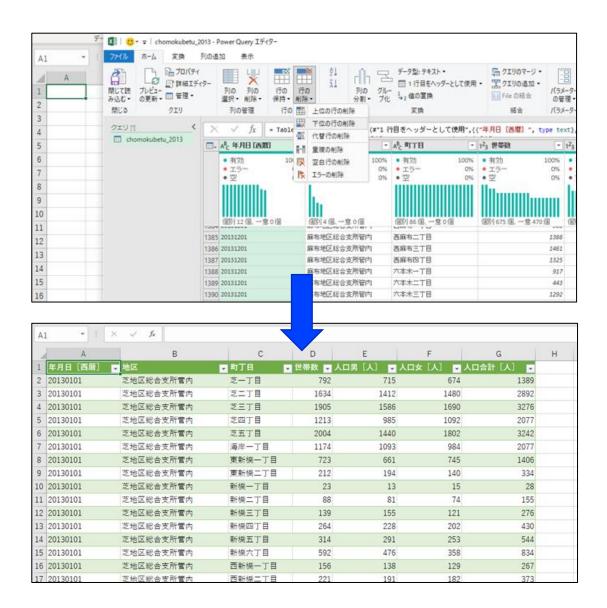
# Power Query によるデータ変換

Power Query (パワークエリ) とは、マイクロソフト社の エクセルの機能の一つで、外部データとの連携や、連携してからのデータの加工 (列の追加や抽出等) を自動化する機能です。

Excel の Power Query 機能を利用してデータを整形します。

その後、MySQL サーバーにデータをインポートし、検討した CSV データに変換します。

Power Query について詳細は、「Power Query (Microsoft Excel の標準またはアドイン)」(P. 5)をご覧ください。



取り扱うファイルのケースにあわせて、下記のケース別分類で Excel から CSV に変換します。

#### ケース[A] / 縦横多段クロス集計

適用データ例:

Excel 上でコピー、貼り付け等の操作や、結合セルを解除することで CSV データに変換できるデータのレイア タウト。

住民基本台帳のデータ

平成	23年1月	1日現在		ご注意:	印刷地	する場合	こしょく	て」を選択	択しますと	大部数	こなります	†のでご	主意くださ	U.											
「住民	基本台帳	·登鋒人口			品川区:																				
					Oi	複	16	ða.	21	aba.	31	裁	4章	b	51	歳	61	aba.	7	歳	81	ða.	91	故	1
「指定統計」 等コード (五十音順)	品川区 行政コード	丁目名	世帯数	人口総数	男	女	男	女	男	女	男	女	男	女	男	女	男	女	男	女	男	女	男	女	男
		総数	191,930	351,350	1,493	1,462	1,515	1,468	1,448	1,343	1,336	1,302	1,335	1,269	1,227	1,130	1,218	1,198	1,217	1,108	1,171	1,091	1,168	1,167	1,16
		品川地区	33,342	63,077	311	286	313	312	313	271	294	247	260	258	261	255	259	244	228	214	203	208	239	211	23:
060000100	1101	北品川1丁目	1,745	2,900	3	9	8	9	7	5	10	8	6	4	4	0	4	7	5	2	11	0	5	4	7
0060000200		北品川2丁目	1,681	3,006	7	9	7	11	10	15	22	11	10	11	10	8	12	7	7	6	7	13	13	8	8
060000300		北品川3丁目	1,674	3,166	27	21	23	20	17	13	16	15	8	19	16	17	19	12	14	9	9	17	10	10	12
060000400		北品川4丁目	653	1,349	8	8	10	8	10	12	9	5	9	7	11	9	4	5	5	5	8	7	6	5	1.0
0060000500	1105	北品川5丁目	2,266	4,490	25	18	19	20	24	28	20	14	15	23	16	17	19	14	14	13	22	9	16	17	12
060000600	1106	北品川6丁目	316	579	3	1	0	2	1	3	2	0	0	2	2	3	1	3	3	2	1	3	2	- 1	
180000100	1201	東品川1丁日	2.063	3.846	12	11	6	13	14	7	9	11	8	11	10	14	15	11	11	7	12	13	14	12	1.5
180000200	1202	東品川2丁目	926	1.464	3	7	4	4	3	5	3	1	4	3	1	1	3	2	4	3	2	1	4	3	
180000300	1203	東品川3丁目	4,761	10,168	66	62	57	70	85	53	65	51	73	64	75	53	53	69	56	54	43	46	56	45	47
180000400	1204	東品川4丁目	1,818	3,469	33	23	25	36	33	20	21	18	20	13	20	23	22	12	14	16	8	14	13	9	6
180000500	1205	東品川5丁目	14	14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
240000100	1301	南品川1丁目	727	1,381	4	2	7	5	2	6	4	3	4	- 1	3	1	4	4	3	3	0	6	5	8	
240000200	1302	南品川2丁目	1,052	1,877	5	3	4	4	6	5	7	7	2	3	1	4	6	5	7	4	6	4	3	5	1.4
240000300	1303	南品川3丁目	876	1,648	8	6	13	18	8	6	8	7	9	11	10	4	6	10	6	8	5	4	3	4	
240000400	1304	南品川4丁目	1,882	3,281	13	12	11	8	6	9	11	12	11	6	9	8	15	10	11	4	6	11	10	12	15
240000500	1305	南品川5丁目	2,828	5,269	25	19	27	24	23	17	13	27	14	21	14	15	17	20	13	16	16	17	14	19	12
240000600	1306	南品川6丁目	1,709	2,702	5	9	11	6	5	8	12	6	6	5	9	7	11	4	5	4	6	10	7	5	
130000100	1401	西品川1丁目	2,068	4,055	17	9	18	8	16	9	15	10	17	17	16	30	10	14	15	25	15	12	18	20	21
130000200	1402	西品川2丁目	2,082	3,749	8	11	19	15	15	13	16	13	12	13	11	16	12	12	15	14	11	12	20	11	13
130000300	1403	西品川3丁目	1,558	2,821	9	11	14	6	8	14	10	7	13	13	5	6	10	9	9	9	8	7	9	5	13
210000100	1501	広町1丁目	37	52	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
210000200	1502	広町2丁目	606	1,791	30	35	30	25	20	22	21	21	19	- 11	18	19	16	14	11	10	7	2	11	8	
		大崎地区	31,850	55,556	290	276	255	255	269	244	229	218	209	204	155	176	178	179	196	198	193	170	153	165	17:
0050000100		上大崎1丁目	1,484	2,769	8	14	11	10	10	12	8	9	6	18	9	4	14	10	9	6	9	13	8	14	10
050000200		上大崎2丁目	1,896	3,238	9	23	11	24	20	25	20	16	17	8	13	11	9	13	9	9	8	7	7	14	2
0050000300		上大崎3丁目	1,661	2,890	12	10	16	10	9	13	14	15	11	4	12	10	12	11	5	14	8	10	2	11	
0050000400		上大崎4丁目	684	1,202	6	6	6	5	6	6	2	2	3	4	2	3	3	4	2	0	1	2	0	3	2
170000100		東五反田1丁目	1,439	2,280	9	4	9	8	8	4	7	5	9	. 7	7	3	7	8	7	9	10	9	5	5	
170000200		東五灰田2丁目	2,110	3,889	32	29	26	34	24	22	18	23	20	14	12		9	10	13	17	13		14	14	- 11
170000300		東五反田3丁目	953	1,717	5	5	3	10	2	12	5	6	9		7	8	7	7	10	4	8	10	5	9	
170000400		東五反田4丁目	1,084	1,960	12	11	14		12	9	12	12	12	12					6	12	11	9	3		3
170000500		東五反田5丁目	1,458	2,666	14	10	7	10	11	8	11	8	13	7	7	11	8	4	9	11	8		8	6	12
120000100		西五反田1丁目	494	711	2	0	0	1	3	1	1	2	0	2	2	1	2	2	1	1	1	2	1	1	3
120000200		西五灰田2丁目	1,061	1,447	4	0	5	2	4	4	1	1	2	1	1	4	1	1	0	1	3	3	2	2	3
120000300		西五灰田3丁目	2,405	4,414	36	29	25	24	32	32	21	27	26	27	25	22	19	20	36	22	25	23	23	22	17
120000400		西五反田4丁目	2,282	4,071	26	21	15	15	18	23	16	13	11	15	9	9	11	13	16	12	12	11	7	9	18
120000500	2305	西五反田5丁目	2,565	4,457	21	25	21	20	22	15	21	13	16	15	9	18	18	15	26	18	32	11	13	11	1

#### ケース[B] / 単純クロス集計

1 行目が項目行、2 行目以降がデータ行となり、1 レコードが 1 行で保持されたデータ。

# 

**適用データ例:** 地域・年齢別人口の データ

#### ケース「C] / 単純表形式、または 複数ファイル

適用データ例:

医療機関データ

1 行目が項目行、2 行目以降がデータ行となった、表形式のデータ。

項目の粒度や目的が混在している場合は、変換したい CSV ファイルにあわせて表を分割する必要があります。



#### ケース[D] / 1シートに単純表形式が種類別に縦に3つ並んでいる、等

適用データ例:

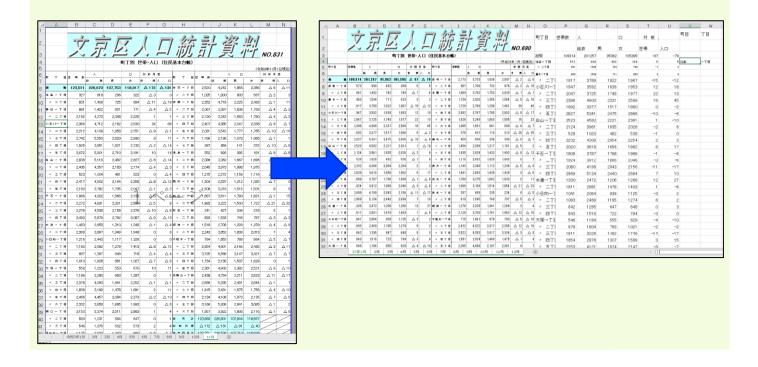
ハザードマップ

1シートに複数の単純表形式が混在しているレイアウト。

CSV データにするにはファイルを分割する必要がある。例えば同一シート内に単純表形式が縦に 3 つあれば、1 シートに 1 表、かつ 3 ファイルに分割する必要があります。



対象外のケース / シート毎に年月日が分かれている、1 枚のシート中にクロス集計表が横に並んでいる、等整えられていない Excel は、データをコピーと貼り付けすることにより、一括処理しやすい形式に変換します。



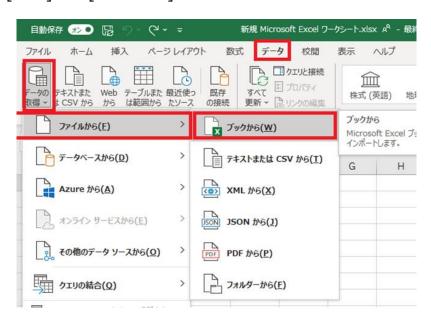
# 手順(ケース[A])

下記の手順で Excel から CSV にデータを変換します。

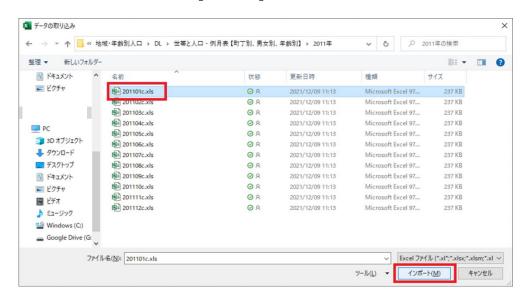
【お願い】: データ取得元(ソース)のファイルは閉じた状態で作業を行ってください。

- ① Excel で整形、加工したいファイルをインポートします。 (詳細: p.28)
- ② Excel の Power Query 機能で、データを整形します。 (詳細: p.29)
- ③ Excel の Power Query 機能で(年次、月次等)定期発生するデータを整形する (詳細: p.43)

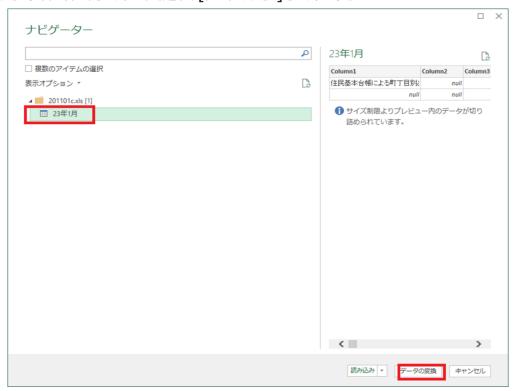
- ① Excel で整形、加工したいファイルをインポートする
- 1. データ取得元(ソース)のファイルを開いている場合は、閉じておく
- 2. [データ]タブ→[データの取得]→「ファイルから」→「ブックから」をクリックする



3. Excel ファイルをクリックして選び、[インポート]をクリックする



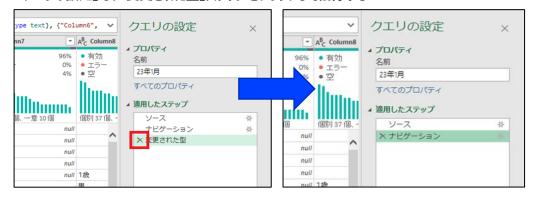
#### 4. 変換するアイテムをクリックして選び、「データの変換]をクリックする



ここでの「アイテム」は、 Excel のシートです。 シート単位で複数表 示されている場合は、 必要なシートを選びま す。

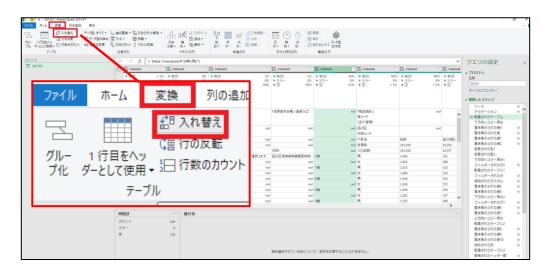
# ② Excel の Power Query 機能でデータを整形する

#### 1. 「クエリの設定」で、「変更された型」ステップをクリックして削除する



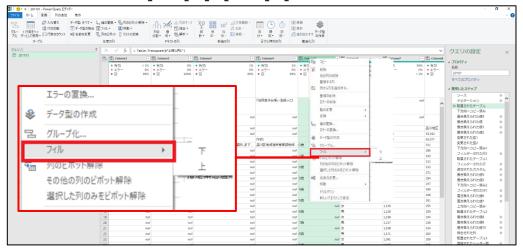
「変更された型」は操作を行うと自動的に追加されますが、手順は不要なので削除します。

## 2. [変換]タブ→[入れ替え]をクリックする



ここではヘッダー行を整 形するため、行列を入 れ替える処理を行いま す。

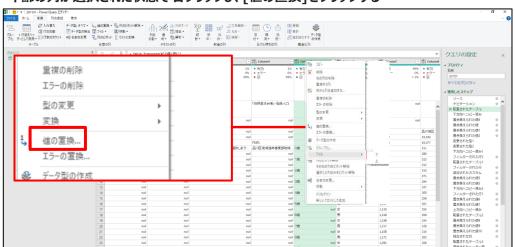
#### 3. 年齢の列を右クリックし、「フィル]→[下]または[上]をクリックする



[下]または[上]は、コ ピーしたい方向を示し ます。

[下]は上から下の方向にデータをコピーし、 [上]は下から上の方向にデータをコピーします。

## 4. 年齢の列が選択された状態で右クリックし、[値の置換]をクリックする



# 5. 置換する値に「null」、置換後に「""」(ダブルクォーテーション) と入力し、[OK]をクリックする



ここでの作業は、後の 手順で日付データをフィルするために行います。

「null」は何のデータも 含まれない、または長 さ 0 の空文字列という 意味です。

画!像

#### 6. 年月日の列が選択された状態で右クリックし、「値の置換」をクリックし、値を置換する

#### 年月日列で不要データの値を置換します。

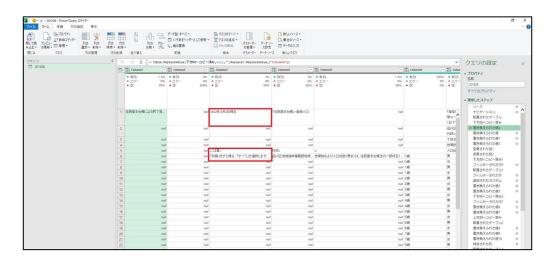
・対象列を選択した状態で右クリックし、「値の置換」を選ぶ。 置換する値「日付データ以外の文言」 → 置換後「null」

#### 年月日列を日付型でエラーとならない形式に整形します。

・対象列選択状態で 右クリック 「値の置換」

置換する値:不要な()や現在、末等の文言 → 置換後「""」(ダブルクォーテーション)

ここでの不要データ置換作業は、後の手順で日付データをフィルするために行います。



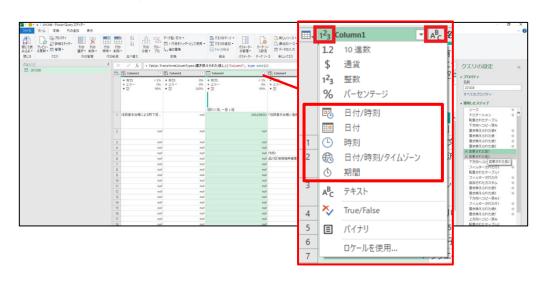
#### 日付型の列での注意点です。

和暦表記も可能ですが、「元号」表記、「全角数字」、「漢数字」は取り扱えません。

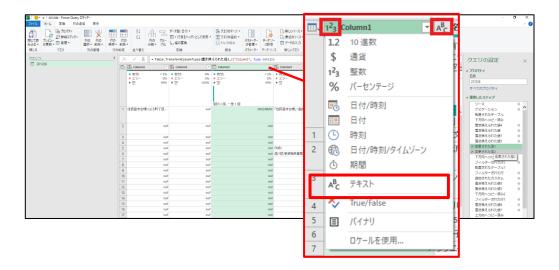
- × 令和元年1月1日 (元号)
- × 令和1年1月1日 (全角数字)
- × 令和一年一月一日 (漢数字)
- × 20190101
- $\times$  2019/1/1

- 〇 令和1年1月1日
- 〇 令和1年1月1日
- 〇 令和1年1月1日
- O 20190101
- O 2019/1/1

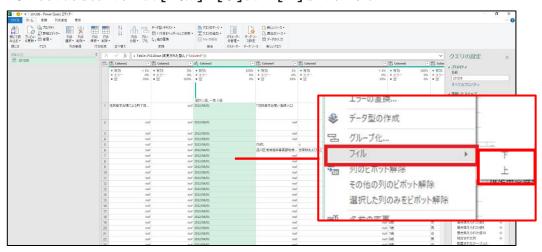
# 7. 年月日のヘッダー上部にある「ABC123」アイコンをクリックし、日付型を選択する



#### 8. 年月日のヘッダー上部にある「ABC123」アイコンをクリックし、テキスト型を選択する



9. 年月日の列を右クリックし、[フィル]→[下]または[上]をクリックする



[下]または[上]は、コ ピーしたい方向を示し ます。

[下]は上から下の方向にデータをコピーし、 [上]は下から上の方向にデータをコピーします。

#### ここから、不要な行を非表示にします。

不要な行を非表示にする方法は 2 種類あります。

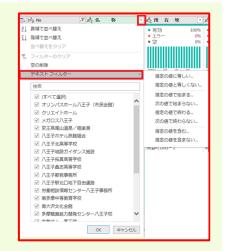
- 1. 対象列の右側「▼」をクリックし、数値フィルターまたはテキストフィルターを表示する
- 2. 「規定の値と等しくない」を選び、[OK]をクリックする

対象行が非表示になります。

数値フィルター

テキストフィルター

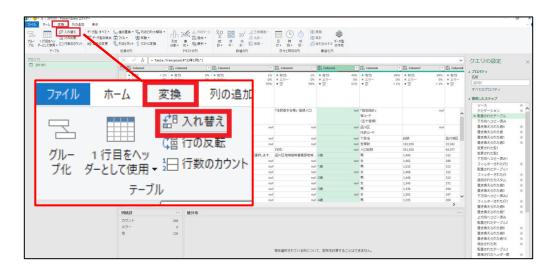




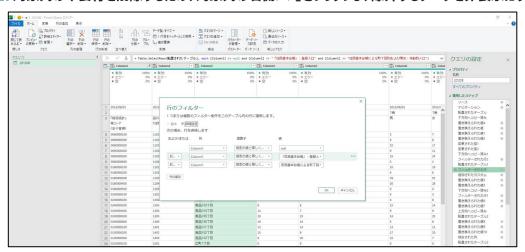
10. 年月日の列を右クリックし、除外するデータを非表示にする



## 11. [変換]タブ→[入れ替え]をクリックする

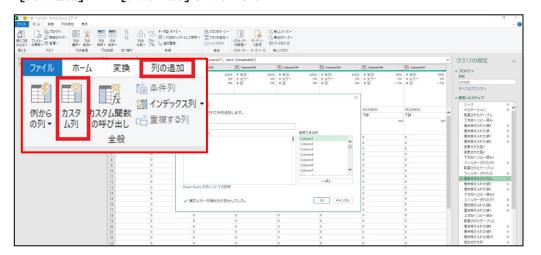


#### 12. 対象列の不要行を削除するため、対象列の右側「▼」をクリックし、除外するデータを非表示にする



#### データ整備マニュアル

#### 13. [列の追加]タブ→[カスタム列]をクリックする



今回は最終列をコピー し、年月日のカスタム 列を作成します。

#### 14. 追加したカスタム列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「不詳」 → 置換後「null」



ここでの不要データ置 換作業は、後の手順 で日付データをフィルす るために行います。

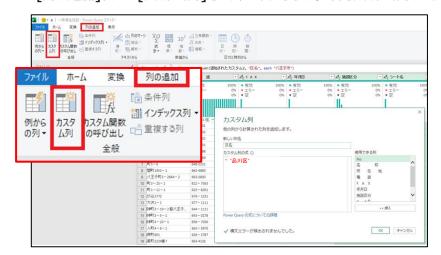
#### 15. 年月日の列を右クリックし、[フィル]→[下]または[上]をクリックする



[下]または[上]は、コ ピーしたい方向を示し ます。

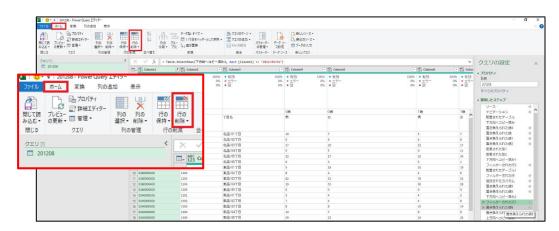
[下]は上から下の方向にデータをコピーし、 [上]は下から上の方向にデータをコピーします。

### 16. 「列の追加]タブ→「カスタム列]をクリックし、対象の区名をカスタム列の式に入力し、[OK]をクリックする

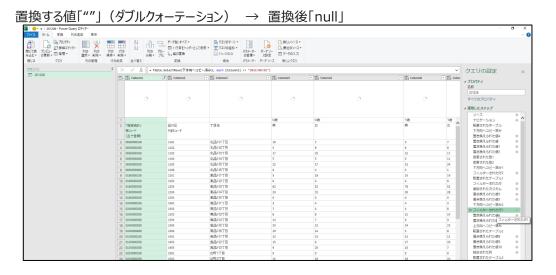


今回の例では"渋谷区"を入力します。

#### 17. [ホーム]タブ→[行の削除]→「上位の行の削除」をクリックし、削除する上位の行数を入力する



#### 18. 対象列が選択された状態で右クリックし、「値の置換」をクリックし、値を置換する



ここではヘッダー行ズレ を整形します。

## データ整備マニュアル

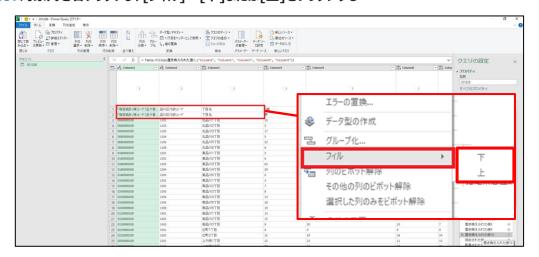
## 19. 対象列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「#(lf)」(改行コード) → 置換後「」(半角スペース) 「詳細設定オプション」の「特殊文字を使用した置換」にチェックを付けておきます。



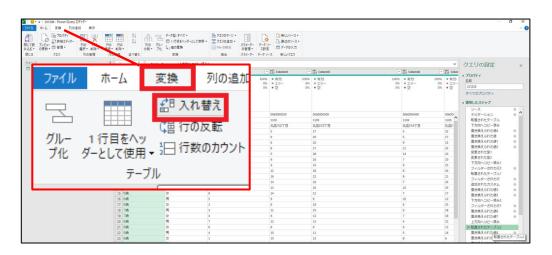
ここでは改行コードを 削除しています。

# 20. 対象列を右クリックし、「フィル]→[下]または[上]をクリックする



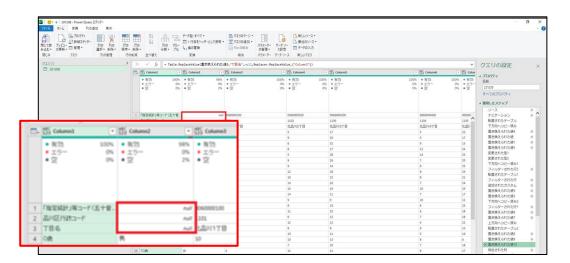
ここではヘッダー行デー タをコピーしています。

# 21. [変換]タブ→[入れ替え]をクリックする



# 22. 対象列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

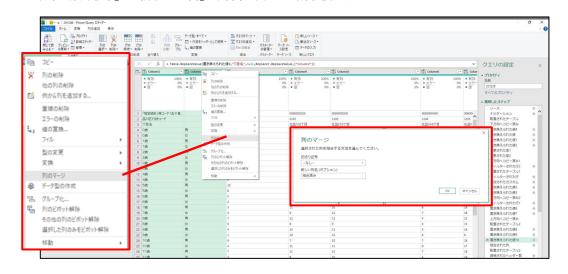
置換する値「Column1 と Column2 で同じ文言」 → 置換後「Column2 の値を null」



ここでは列をマージ (結合・統合) するため、不要な文言を削除しています。

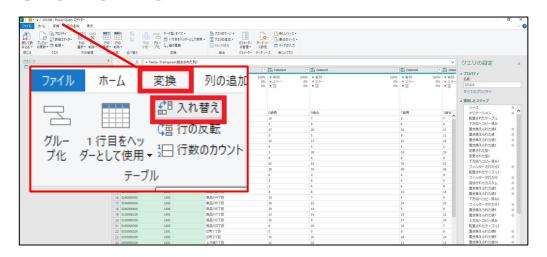
# 23. 対象列が選択された状態で右クリックし、「列のマージ]をクリックし、値を置換する

「段区切り記号」や「新しい列名」は内容に合わせて設定します。

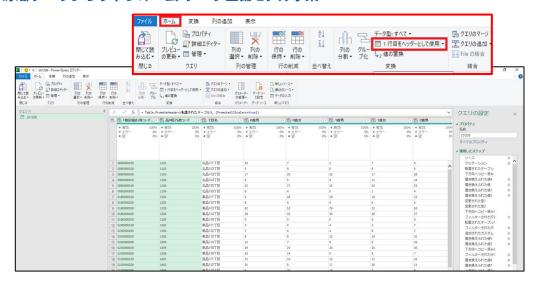


ここではヘッダー列を作 成するため、列をマー ジしています。

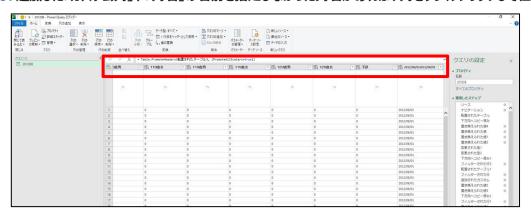
# 24. [変換]タブ→[入れ替え]をクリックする



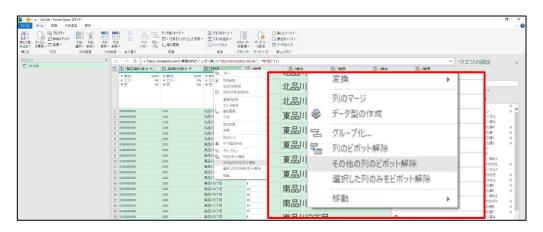
## 25. [ホーム]タブ→[1 行目をヘッダーとして使用]をクリックする



26. 追加した「カスタム列」や「列名」の名前を指定しなかった列名があれば、列をダブルクリックして任意の列名を入力する



27. 対象列が選択された状態で右クリックし、[そのほかの列のピボット解除]をクリックする



#### 28. ピボット解除により作成されたヘッダー列名に、任意の列名を入力する



# 29. 対象列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「全角数値0」 → 置換後「半角数値0」 <この間の数字分も、手順を繰り返します。>

置換する値「全角数値9」→ 置換後「半角数値9」

| \*\*\* | 201008 - Power Curry エチザラー | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | 201007 - \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* | \*\*\* |

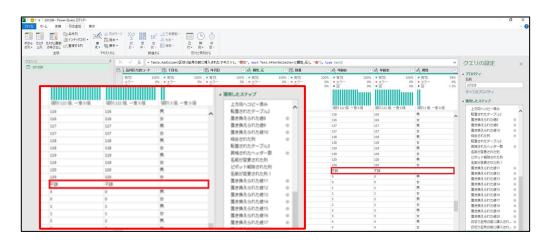
# 30. 対象列が選択された状態で右クリックし、[列の追加]タブ→[例からの列]→[選択範囲から]をクリックし、任意の値に変更して例を作成する

対象列:属性、年齢自、年齢列追加 属性\_元列 それぞれの列から必要範囲を切り出して例を作成します。 全データで想定通りの値が入力されているか確認してください。



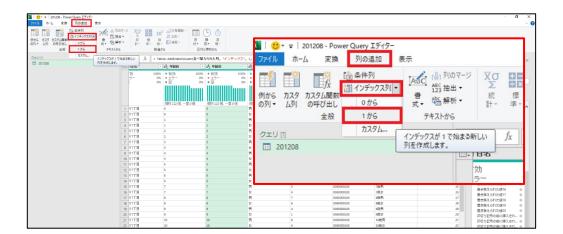
# 31. 対象列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「不詳」 → 置換後「999」 置換する値「""」(ダブルクォーテーション) → 置換後「その他」



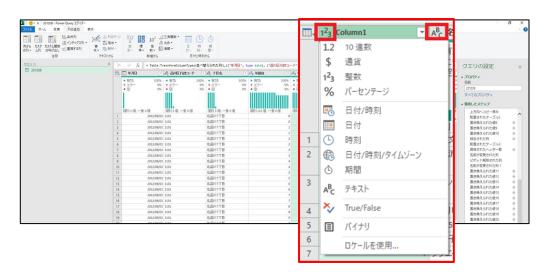
ここでは不詳データを 整形しています。

# 32. 追加する列が選択された状態で右クリックし、「列の追加]タブ→[インデックス列]→[1 から]をクリックする

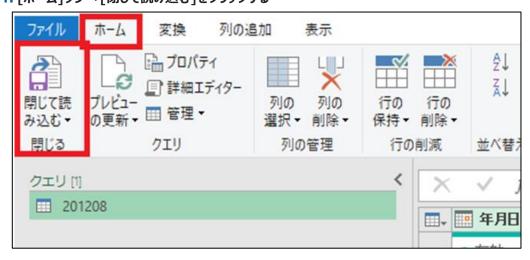


ここではインデックス列 を追加しています。

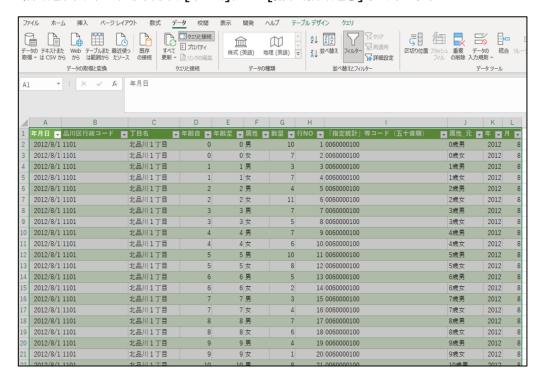
## 33. 各列のヘッダー上部にある「ABC123」アイコンをクリックし、任意の型を選択する



# 34. [ホーム]タブ→[閉じて読み込む]をクリックする



35. 読み込まれたデータを確認し、[ホーム]タブ→[閉じて読み込む]をクリックする



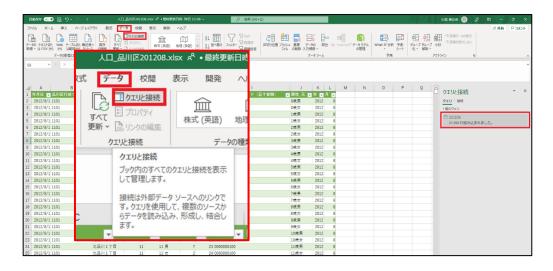
36. ファイルの種類欄で「テキスト(タブ区切り)」を指定して CSV を保存する

# ③ Excel の Power Ouerv 機能で(年次、月次等)定期発生するデータを整形する

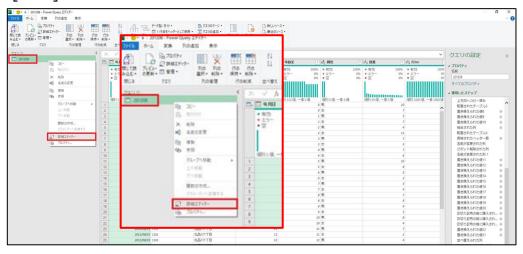
①②で処理した Power Query の詳細エディターを複製し、月毎等でエクセルブックが分かれている場合に発生する定期発生データも変換を行います。

【お願い】: データ取得元(ソース)のファイルは閉じた状態で作業を行ってください。

1. [データ]タブ→「クエリと接続」とクリックし、該当クエリをダブルクリックする



2. 「ホーム]タブ→「該当クエリ」を右クリックし、「詳細エディター」をクリックする



3. 詳細エディターのコードを全選択し、コピーする



## 4. 新規エクセルブックで詳細エディターの編集を行う

●B~H列: フォルダ名、シート名、ブック名、年月部分等、置換用の文字列を入力します。

·B列: YYYY/MM ※西暦(半角)·C列: YYYYMM ※西暦(半角)

·D列: YY 年 ※和暦(シート名用)(半角)

·E 列: YYYY 年 ※西暦(半角) ·F 列: MM 月 ※(半角)

・G列:cブック名 YYYYMM 以外文字列がc(小文字半角)以外の場合は個別対応必要

·H列:※上記年月で全角文字使用等例外時に設定

● ]列:生成後の文言出力用の数式を入力します。

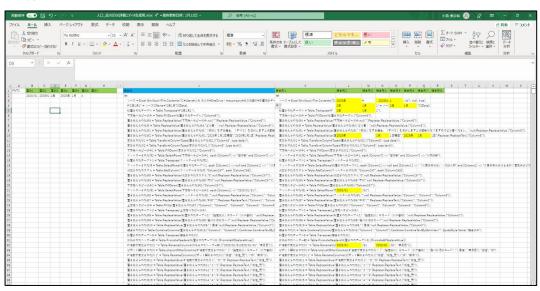
·J列:=K2&L2&M2&N2&O2&P2&Q2&R2

● K~R 列:コピーした詳細エディターを分割して入力

※B~H 列に設定した置換を行う必要のある文字列部分は列を分割します。

・置換文字列:=B2 ※該当の置換文字列が設定されているB~H列を設定すること。

上記生成用の範囲をコピーアンドペーストし、B~H 列の年月を該当年月に変更し、詳細エディターを 生成します。

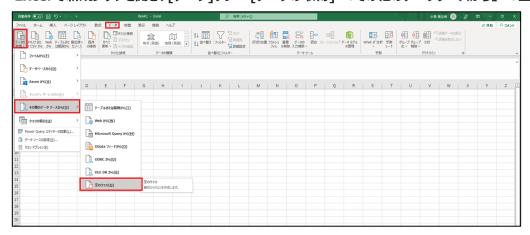


## 行については下記のように表記されます。

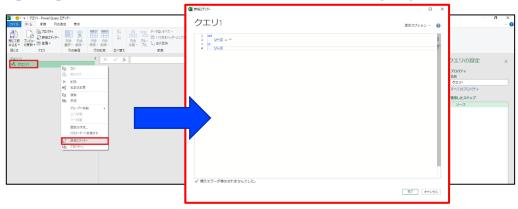
2 行目 ソース=...:フォルダ名、ブック名記載箇所

3 行目 #"2 年 1 月 1" = ソース{[Name="2 年 1 月"]}[Data],:シート名記載箇所

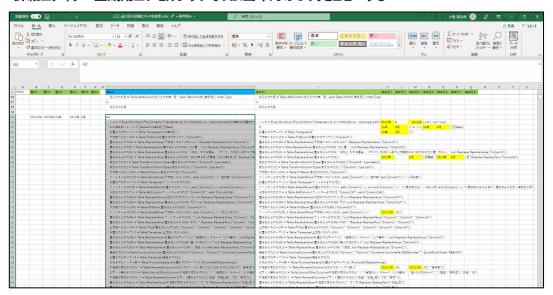
ブック名、シート名、年 月日に関わる部分で、 記載内容と異なる状 況がある場合(品川 区の月ごとデータ 等)、FAQ「Power Query」(P. 117) もご参考ください。 5. Excel で新規ブックを開き、「データ」タブ→「データの取得]→「その他のデータソースから」→「空のクエリ」をクリックする



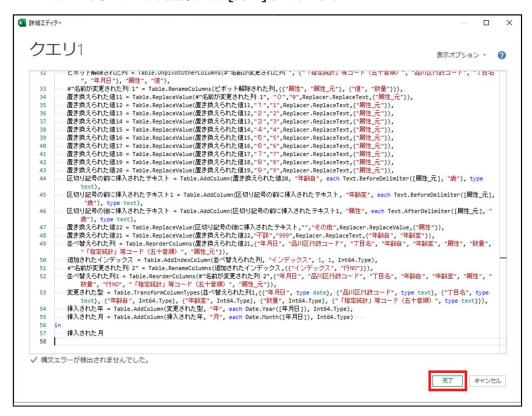
6. [ホーム]タブ→「該当クエリ」を右クリックし、新規ブックの Power Query 詳細エディターを表示する



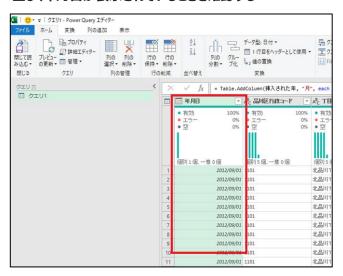
7. 詳細エディター生成用エクセルブックで、該当年月の J 列をコピーする



## 8. コピーした該当年月のJ列を貼り付け、「完了」をクリックする



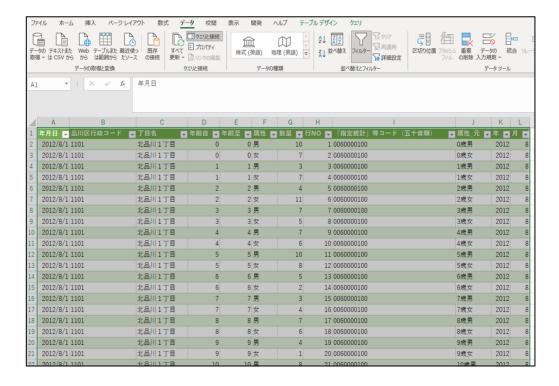
# 9. 正しく年月日が表示されていることを確認する



## 10. [ホーム]タブ→[閉じて読み込む]をクリックする



#### 11. 読み込まれたデータを確認、データを全選択し、コピーする



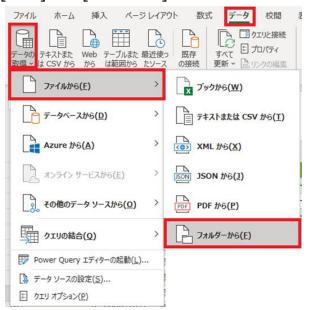
- 12. ②で作成した CSV ファイルを開き、最終行の次行に貼り付ける
- 13. ファイルの種類欄で「テキスト(タブ区切り)」を指定して CSV を保存する
- 14. 全データの CSV 更新が完了したら、途中行に追加されたヘッダー行を検索し行を削除する

# 手順(ケース[B])

下記の手順で Excel から CSV にデータを変換します。

## 【お願い】: データ取得元(ソース)のファイルは閉じた状態で作業を行ってください。

- ① Excel で整形、加工したいファイルをインポートします。 (詳細: p. 48)
- ② Excel の Power Query 機能で、データを整形します。 <u>(詳細: p. 50)</u>
- ① Excel で整形、加工したいファイルをインポートする
- 1. データ取得元 (ソース) のファイルを開いている場合は、閉じておく
- 2. [データ]タブ→[データの取得]→「ファイルから」→「フォルダから」をクリックする



# 3. Excel ファイルをクリックして選び、[開く]をクリックする



# 4. [データの結合と変換]をクリックする

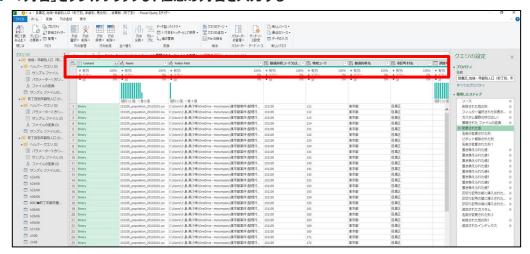


## 5. [OK]をクリックする

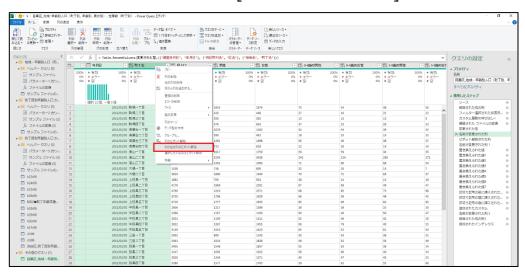


# ② Excel の Power Query 機能でデータを整形する

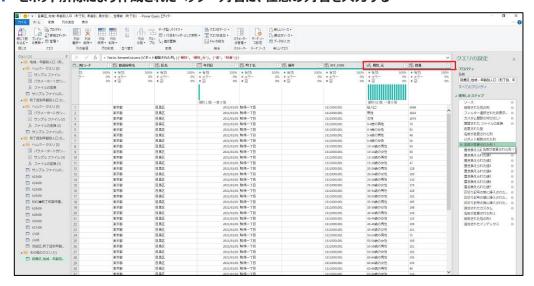
1. 「列名」をダブルクリックし、任意の列名を入力する



2. 対象列が選択された状態で右クリックし、[そのほかの列のピボット解除]をクリックする



3. ピボット解除により作成されたヘッダー列名に、任意の列名を入力する

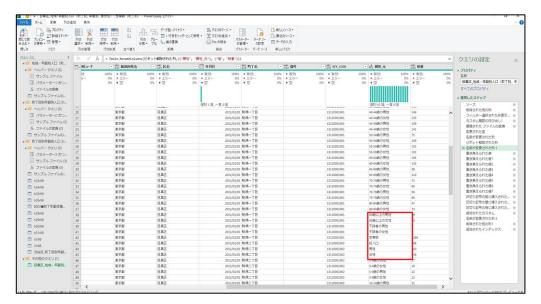


## データ整備マニュアル

## 4. 対象列が選択された状態で右クリックし、「値の置換」をクリックし、値を置換する

置換する値「000 歳-000 歳の〇〇〇形式となっていないデータ」

→ 置換後「000歳-000歳の○○○」



ここでは、属性\_元列を年齢(開始)、年齢(終了)、属性に分割するため、入力フォーマットを揃えます。

属性\_元列の値の置 換を行います。

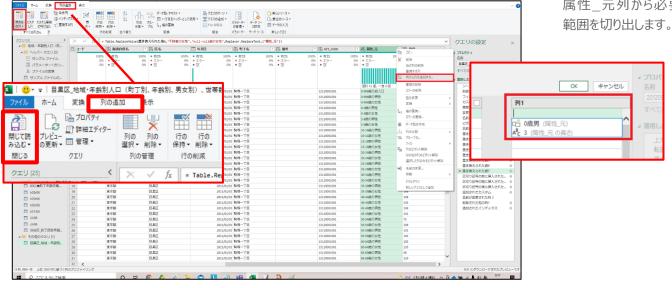


値の置換の際、"男性"、"女性"等、置換を行いたくない他のセル内容の一部と合致してしまう場合は、「セルの内容全体の照合」のチェックを ON にします。

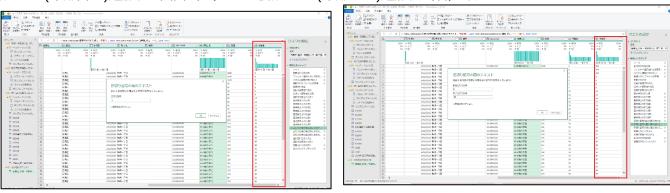
# 5. 対象列が選択された状態で右クリックし、「列の追加]→「例からの列]→「選択範囲から]をクリックし、例を作成する

追加列のセルをダブルクリックし、選択セルのデータ内容をクリック、任意の値に変更して例を作成します。 作成後は、全データ想定通りの値が入力されている事を確認します。

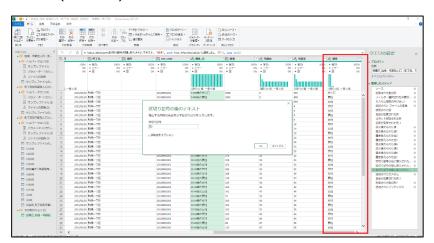
ここでは、年齢自、年齢至、属性列追加属性\_元列から必要節囲を切り出します。



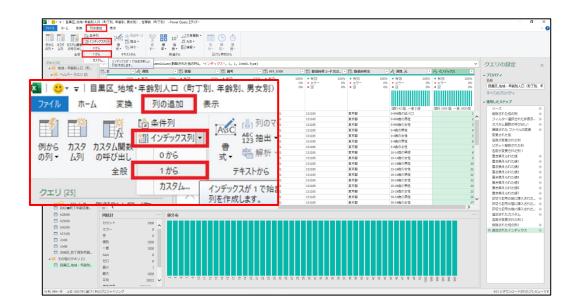
例からの列(年齢自列)追加後の画面(左)および例からの列(年齢至列)追加後(右)です。



例からの列(年齢自列)追加後の画面です。

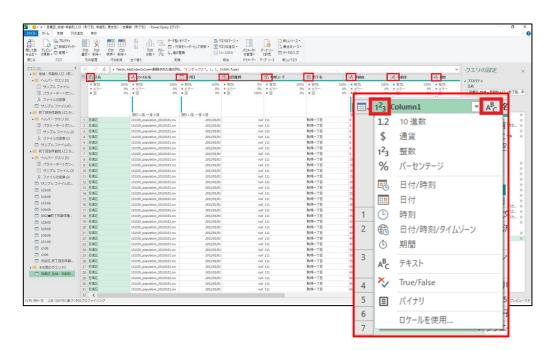


6. 追加する列が選択された状態で右クリックし、「列の追加]タブ→[インデックス列]→[1 から]をクリックする



ここではインデックス列 を追加しています。

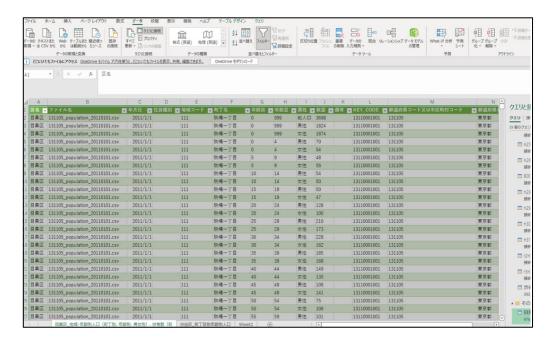
7. 各列のヘッダー上部にある「ABC123」アイコンをクリックし、任意の型を選択する



8. [ホーム]タブ→[閉じて読み込む]をクリックする



9. 読み込まれたデータを確認する



10. ファイルの種類欄で「テキスト(タブ区切り)」を指定して CSV を保存する

# 手順(ケース[C])

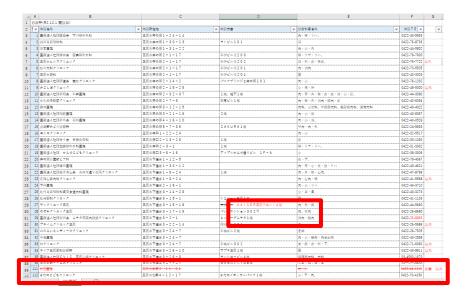
下記の手順で Excel から CSV にデータを変換します。

## 【お願い】: データ取得元(ソース)のファイルは閉じた状態で作業を行ってください。

- ① Excel マクロにより、文字装飾等を削除する前処理を実施します。 (詳細: p. 55)
- ② Excel で整形、加工したいファイルをインポートします。 (詳細: p. 58)
- ③ Excel の Power Query 機能で、データを整形します。 (詳細: p. 60)クエリパターン 1 (詳細: p. 92)クエリパターン 2 (詳細: p. 64)
- ④ Excel の Power Query 機能で、複数パターンが発生するクエリを結合します。 (詳細: p. 71)

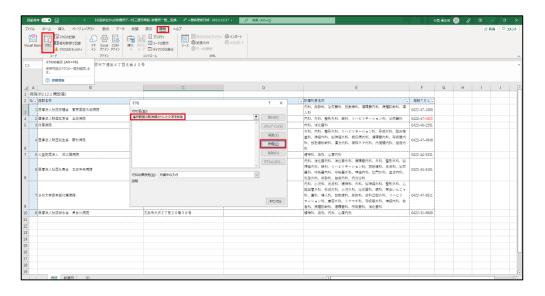
# ① Excel マクロにより、文字装飾等を削除する前処理の実施

#### 1. 取消線や装飾文字のある箇所を確認する



Excel に直接打ち込んだ、取り消し線等をマクロで削除します。

2. [開発]→[マクロ]をクリックし、マクロ名を入力した後、[作成]をクリックする



3. 選択範囲を対象とした、取消線のついた文字の装飾を削除するマクロを入力する

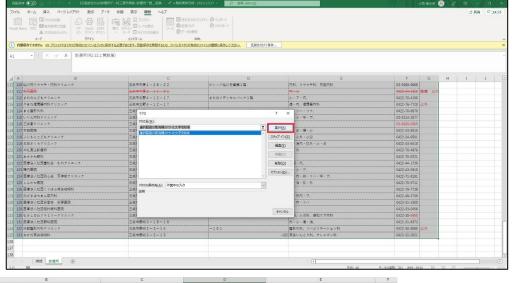


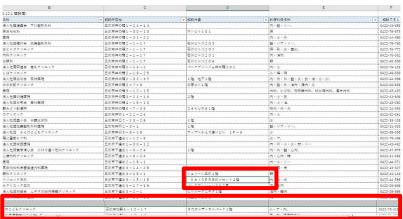
マクロの内容は下記です。

4. [ファイル]タブ→[終了して Microsoft Excel へ戻る]をクリックする



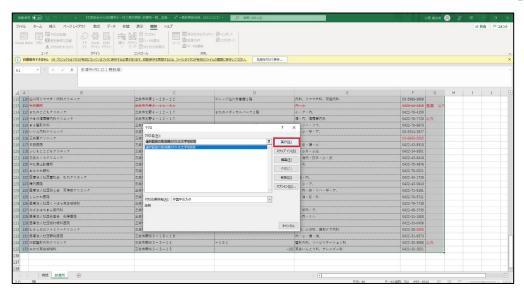
5. マクロを適用する範囲を選択し、「選択範囲の取り消し線の付いた文字を削除」マクロを選択し、「実行」をクリックする



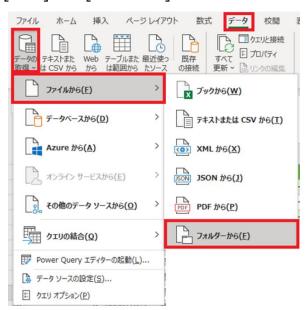


マクロが実行され、選択範囲の文字装飾がなくなっていることを確認します。

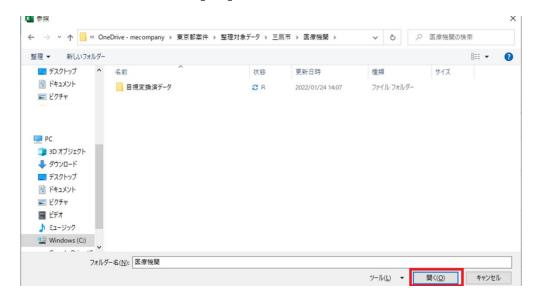
6. マクロを適用する範囲を選択し、「選択範囲の取り消し線の付いた文字を削除」マクロを選択し、[実行]をクリックする



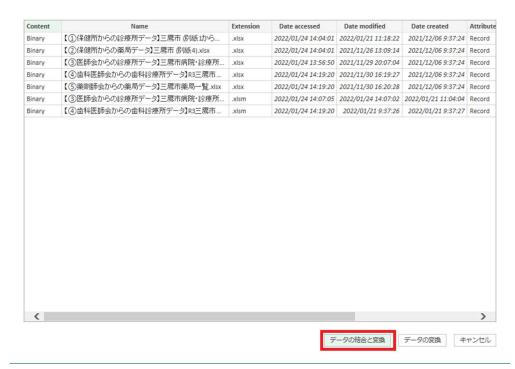
- ② Excel で整形、加工したいファイルをインポートする
- 1. データ取得元 (ソース) のファイルを開いている場合は、閉じておく
- 2. [データ]タブ→[データの取得]→「ファイルから」→「フォルダーから」をクリックする



3. Excel ファイルをクリックして選び、[開く]をクリックする

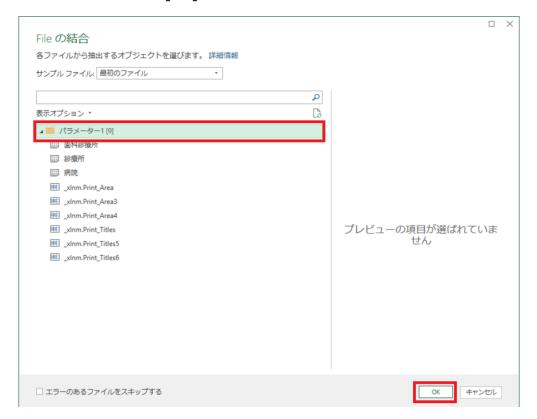


# 4. [データの結合と変換]をクリックする



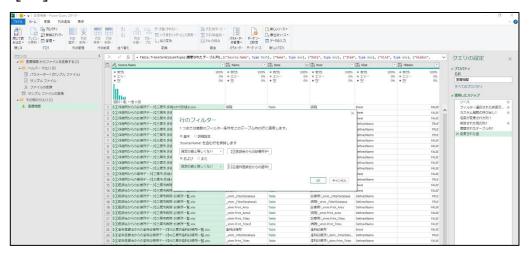
ここでの「アイテム」は、 Excel のシートです。

# 5. フォルダをクリックして選び、[OK]をクリックする



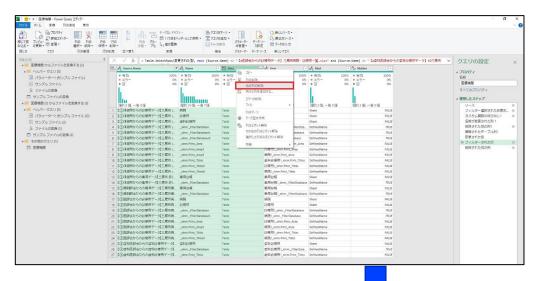
# ③ Excel の Power Query 機能でデータを整形する

# 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする

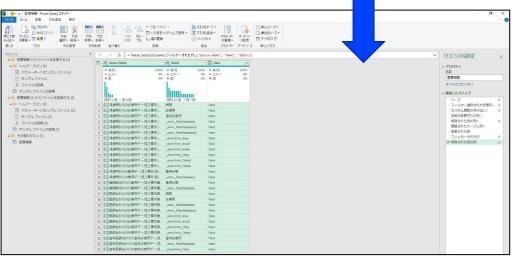


手順①でマクロ有(拡張子.xlsm)で保存したブックが存在する場合はマクロ無(拡張子.xlsx)を除外します。

2. 対象列を選択した後、右クリックし、「他の列の削除」をクリックする

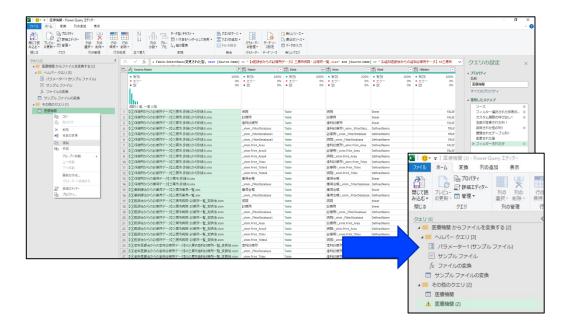


ここでは、対象列は、 Source.Name 、 Source、Data 例を 指します。



対象列が残り、不要列を削除した状態です。

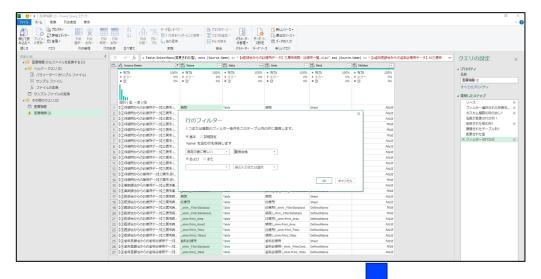
3. 該当するクエリを選択した後、右クリックし、「複製」をクリックする



ヘッダーパターン分のク エリを複製します。

# くクエリパターン 1>

4. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しい」を選び、[OK] をクリックする



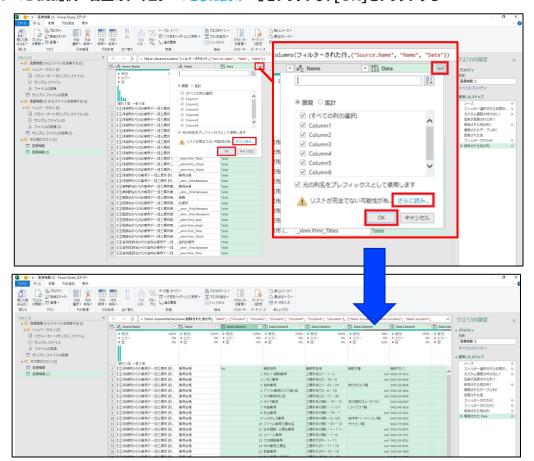
ここでは必要行を残す フィルターを設定してい ます。



必要行が残った状態です。

## データ整備マニュアル

5. 「Data」例 右上のアイコン→「さらに読み…」をクリックし、「OK]をクリックする



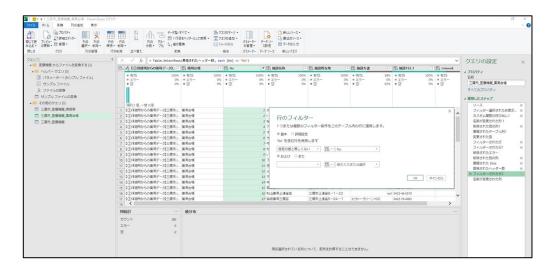
ここでは、Data 例を展開します。

6. [ホーム]タブ→「1 行目をヘッダーとして使用」をクリックする



ここでは1行目をヘッダーとする整形を行っています。

7. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする



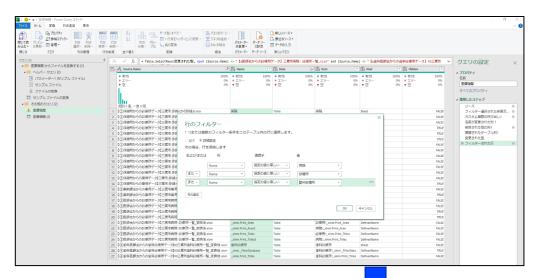
ここでは不要行を非表 示する処理をしていま す。

8. 追加した「カスタム列」や「列名」の名前を指定しなかった列名があれば、列をダブルクリックして任意の列名を入力する



# **<クエリパターン 2>**

9. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しい」を選び、[OK] をクリックする

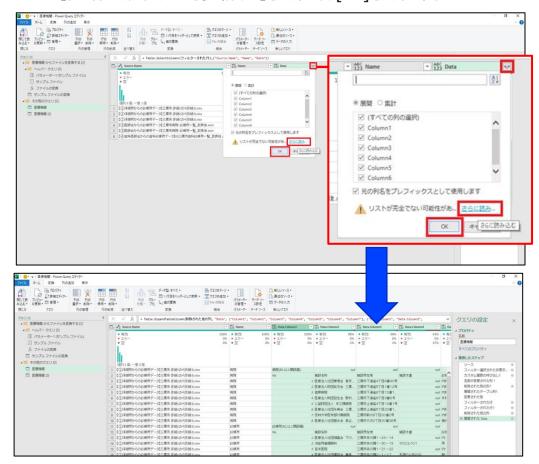


ここでは必要行を残す フィルターを設定してい ます。



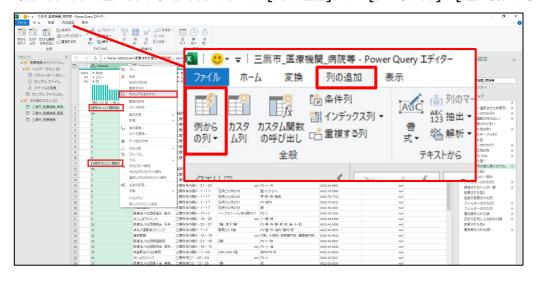
必要行が残った状態です。

## 10. 「Data」例 右上のアイコン→「さらに読み…」をクリックし、「OK]をクリックする



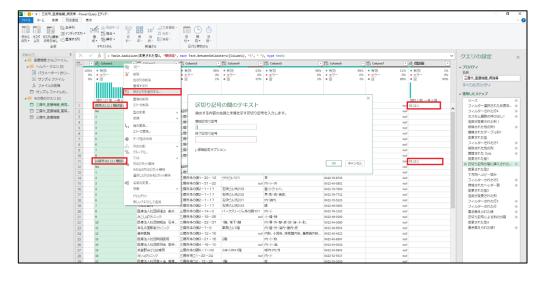
ここでは、Data 例を展開します。

# 11. 開設届の列が選択された状態で右クリックし、「列の追加]タブ→「例からの列]→「選択範囲から]をクリックする

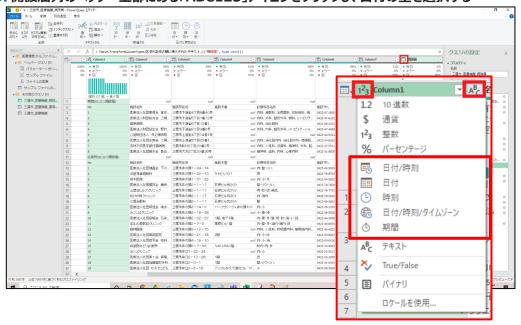


ここでは、No、開設届年月日が格納されている Column1 列から必要範囲を切り出し、開設届列を追加します。

#### 12. 追加列のセルをダブルクリック、選択セルのデータ内容をクリック、任意の値に変更して例を作成する



## 13. 開設届列のヘッダー上部にある「ABC123」アイコンをクリックし、日付の型を選択する



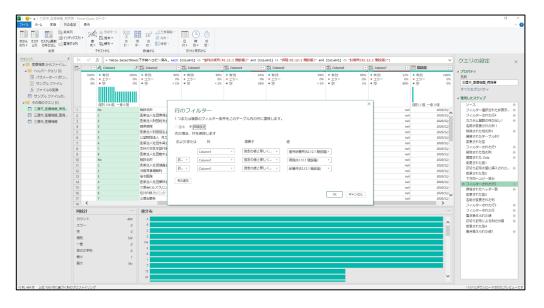
14. 開設届の列を右クリックし、[フィル]→[下]または[上]をクリックする



15. 開設届の列で、方を追加列のセルをダブルクリック、選択セルのデータ内容をクリック、任意の値に変更して例を作成する



16. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする



ここでは不要行を削除 する処理をしていま

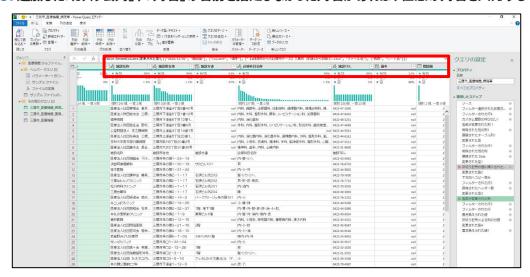
前の手順で別列へ分離した開設届年月日情報を含むデータ行は除外します。

17. [ホーム]タブ→「1 行目をヘッダーとして使用」をクリックする

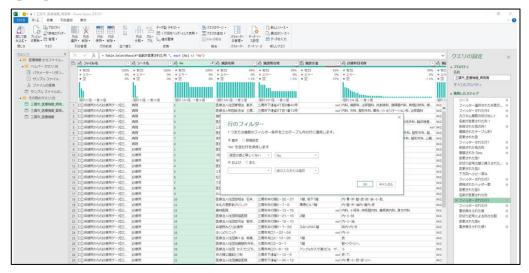


ここでは1行目をヘッダーとする整形を行っています。

18. 追加した「カスタム列」や「列名」の名前を指定しなかった列名があれば、任意の列名を入力する



19. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする



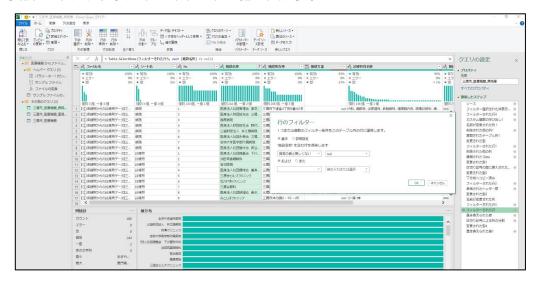
ここでは不要行を削除 する処理をしていま す。

ヘッダー行列名情報の データ行は除外しま す。

20. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする

## 東京都データプラットフォーム データ整備モデル事業

# データ整備マニュアル

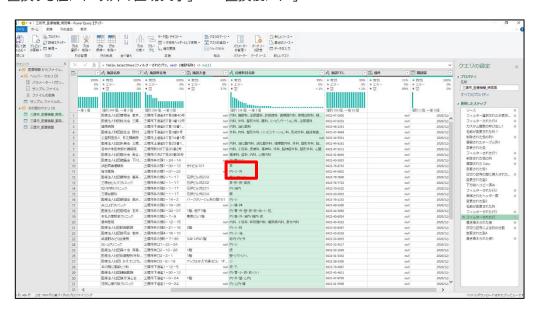


ここでは不要行を削除 する処理をしていま

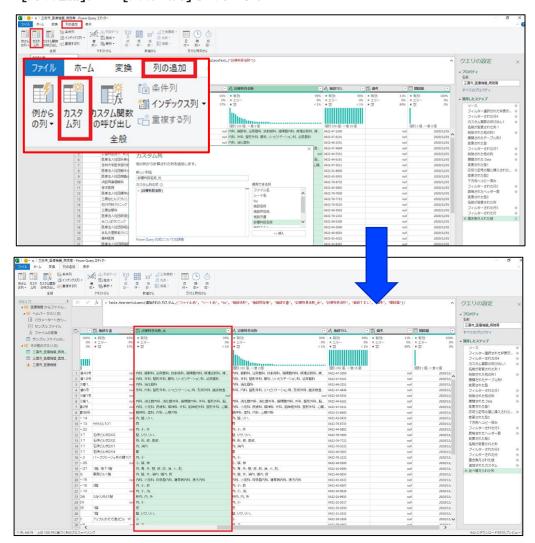
前の手順で、全て取 消線で「null データ」と なったデータ行は除外 します。

# 21. 診療科目名称の列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「"、"以外の区切文字」 → 置換後「"、"」



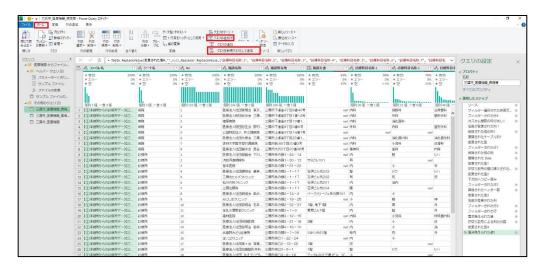
# 22. [列の追加]タブ→[カスタム列]をクリックする



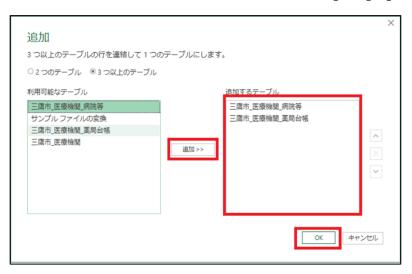
今回は元データ保持 のため、診療科目名 称をコピーした診療科 目名称\_元カスタム列 を作成します。

カスタム列を追加後の状態です。

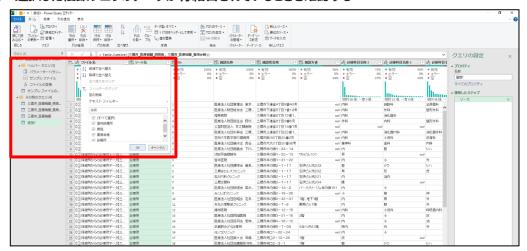
- ④ Excel の Power Query 機能で、複数パターンが発生するクエリを結合する
- 1. [ホーム]タブ→「クエリの追加」→「クエリを新規クエリとして追加」とクリックする



1. ①②③で作成したクエリを、利用可能なテーブルより選択し、[追加]→[OK]をクリックする



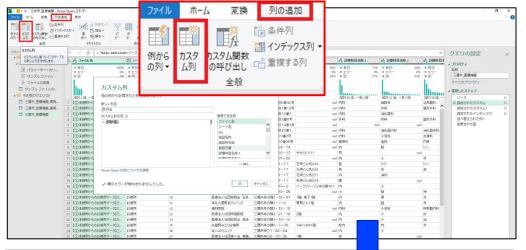
2. 選択した複数クエリのデータが行結合されていることを確認する



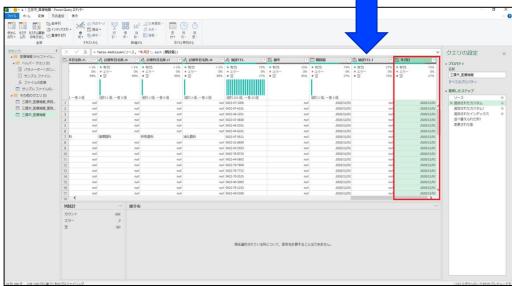
## 東京都データプラットフォーム データ整備モデル事業

## データ整備マニュアル

## 3. [列の追加]タブ→[カスタム列]をクリックする

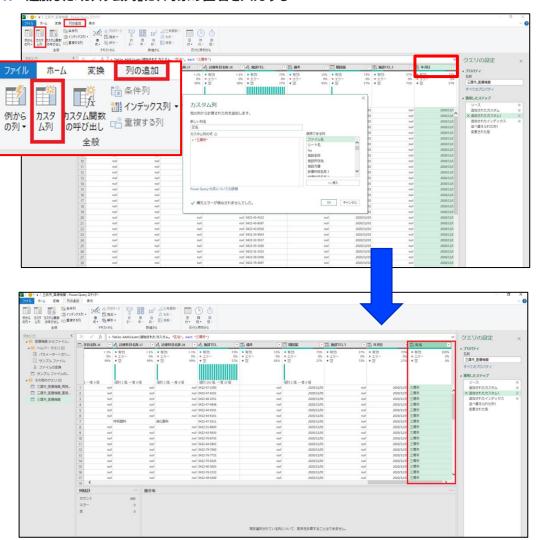


今回はカスタム列の式 に開設届列を作成し ます。



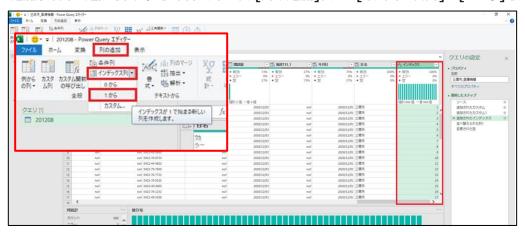
年月日カスタム列を追 加後の状態です。

## 4. 追加した「カスタム列」に、対象の区名を入力する



区名カスタム列を追加 後の状態です。

5. 追加する列が選択された状態で右クリックし、「列の追加]タブ→[インデックス列]→[1 から]をクリックする



ここではインデックス列 を追加しています。

## 6. 追加した列をダブルクリックし、任意の列名を入力する



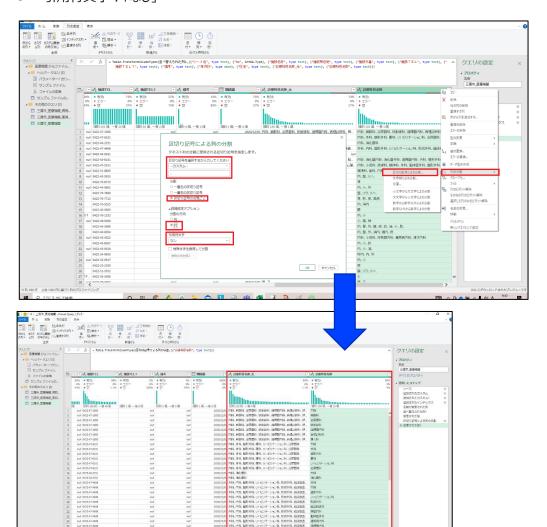
## 7. 対象列が選択された状態で右クリックし、「列の分割」→「区切記号による分割」をクリックする

「区切り記号による列の分割」では、下記の変更を行った後、「OK]をクリックしてください。

区切り記号:「カスタム」→「"、"」分割:「区切り記号の出現ごと」

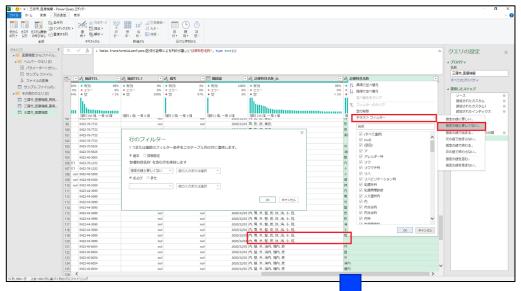
分割の方向:「行」引用符文字:「なし」

ここでは、同一列に複数列記されているデータで列を分割します。

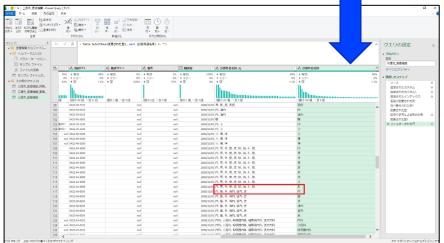


列が分割されます。

8. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする



ここでは、列の分割前 データの末尾に付与されていた区切り文字による空白データ(「皮」 「ア」等も含む)を削除します。



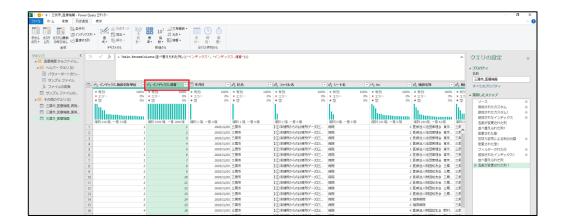
不要行(空白行)が削除されます。

9. 追加する列が選択された状態で右クリックし、「列の追加]タブ→[インデックス列]→[1 から]をクリックする

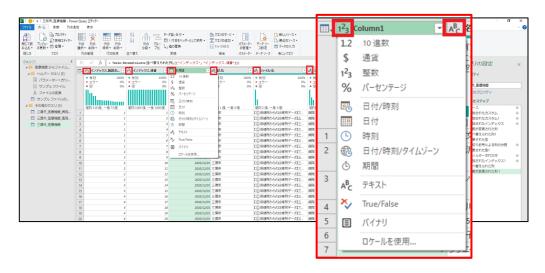


ここではインデックス列 を追加しています。

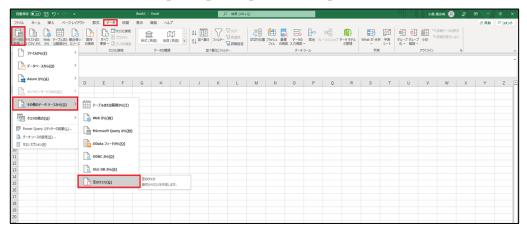
#### 10. 追加した列をダブルクリックし、任意の列名を入力する



## 11. 各列のヘッダー上部にある「ABC123」アイコンをクリックし、任意の型を選択する



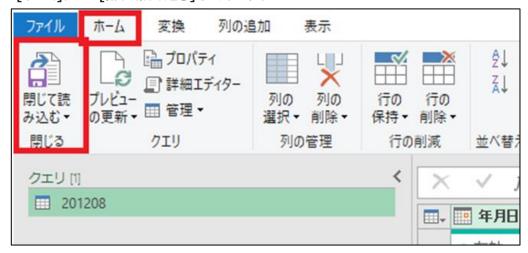
12. Excel で新規ブックを開き、[データ]タブ→[データの取得]→「その他のデータソースから」→「空のクエリ」をクリックする



## 13. 「ホーム]タブ→「該当クエリ」を右クリックし、新規ブックの Power Query 詳細エディターを表示する

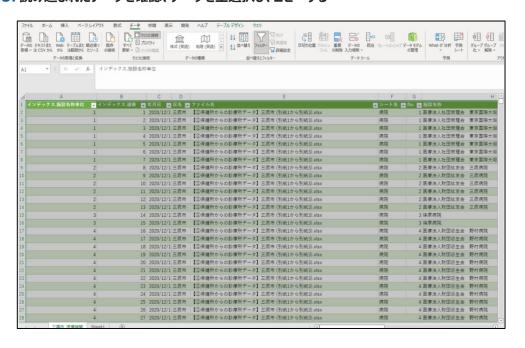


## 14. [ホーム]タブ→[閉じて読み込む]をクリックする



【お願い】 データ取得元(ソース)のファイルは閉じ た状態で作業を行っ てください。

## 15. 読み込まれたデータを確認、データを全選択し、コピーする



- 16. CSV ファイルを開き、貼り付ける
- 17. ファイルの種類欄で「テキスト(タブ区切り)」を指定して CSV を保存する

## 手順(ケース[D])

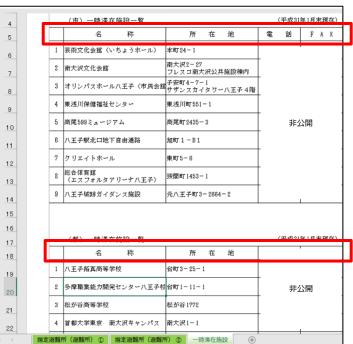
下記の手順で Excel から CSV にデータを変換します。

## 【お願い】: データ取得元 (ソース) のファイルは閉じた状態で作業を行ってください。

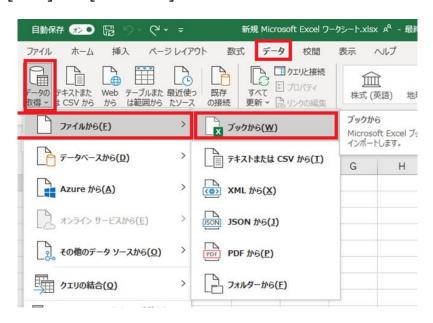
- ① Excel の表で、列名の状況を確認します。 (詳細: p. 78)
- ② Excel で整形、加工したいファイルをインポートします。 (詳細: p. 79)
- ③ Excel の Power Query 機能で、ブックごとにデータを整形します。 (詳細: p. 80)
- ④ Excel の Power Query 機能で、Excel ブック用の処理を行います。※④はケース[C]の「Excel の Power Query 機能で、複数パターンが発生するクエリを結合する」 (詳細: p. 71)をご参照ください。

## ① Excel の表で、列名の状況を確認する

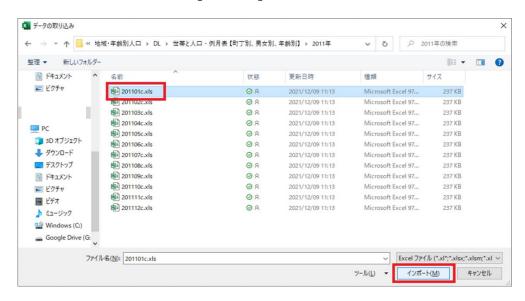
1. 複数ある表の、それぞれの列名が一致しているか目視で確認する



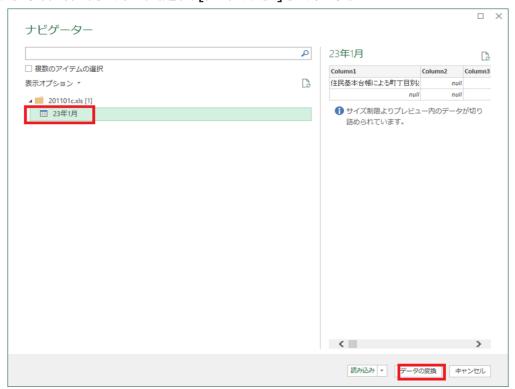
- ② Excel で整形、加工したいファイルをインポートする
- 1. データ取得元(ソース)のファイルを開いている場合は、閉じておく
- 2. [データ]タブ→[データの取得]→「ファイルから」→「ブックから」をクリックする



3. Excel ファイルをクリックして選び、[インポート]をクリックする



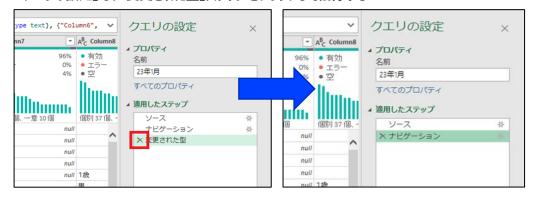
## 4. 変換するアイテムをクリックして選び、「データの変換]をクリックする



ここでの「アイテム」は、 Excel のシートです。 シート単位で複数表 示されている場合は、 必要なシートを選びま す。

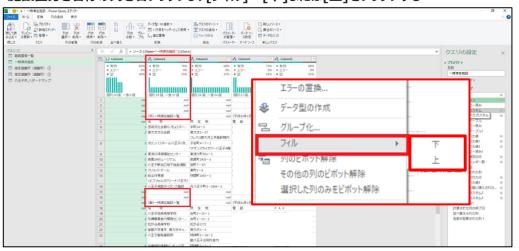
## ③ Excel の Power Query 機能でブックごとにデータを整形する

## 1. 「クエリの設定」で、「変更された型」ステップをクリックして削除する



「変更された型」は操作を行うと自動的に追加されますが、手順は不要なので削除します。

## 2. 施設区分と名称の列を右クリックし、[フィル]→[下]または[上]をクリックする



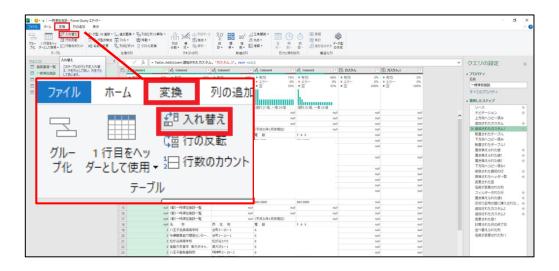
施設区分と名称が同列に混在するため分離します。

## 3. [列の追加]タブ→[カスタム列]をクリックする



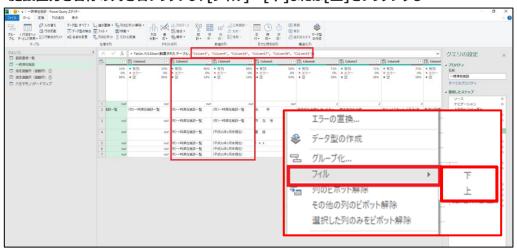
今回は最終列をコピー し、年月日のカスタム 列を作成します。

## 4. [変換]タブ→[入れ替え]をクリックする

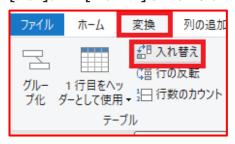


ここでは追加したカスタム列に値をコピーするため、行列を入れ替える処理を行います。

## 5. 施設区分と名称の列を右クリックし、[フィル]→[下]または[上]をクリックする



## 6. [変換]タブ→[入れ替え]をクリックする



## 7. 年月日の列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

## 施設区分列にフィルされている年月日データの値を置換します。

・対象列を選択した状態で右クリックし、「値の置換」を選ぶ。

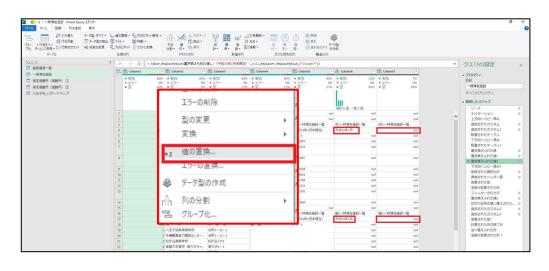
置換する値:年月日の文言 → 置換後「null」

## 年月日列を日付型でエラーとならない形式に整形します。

・対象列選択状態で 右クリック 「値の置換」

置換する値:不要な()や現在、末等の文言 → 置換後「""」(ダブルクォーテーション)

ここでの不要データ置換作業は、後の手順で日付データをフィルするために行います。



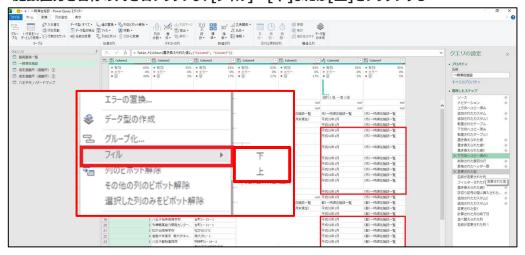
#### 日付型の列での注意点です。

和暦表記も可能ですが、「元号」表記、「全角数字」、「漢数字」は取り扱えません。

- × 令和元年1月1日
- × 令和1年1月1日
- × 令和一年一月一日
- × 20190101
- $\times$  2019/1/1

- 令和1年1月1日
- 〇 令和1年1月1日
- 〇 令和1年1月1日
- O 20190101
- O 2019/1/1

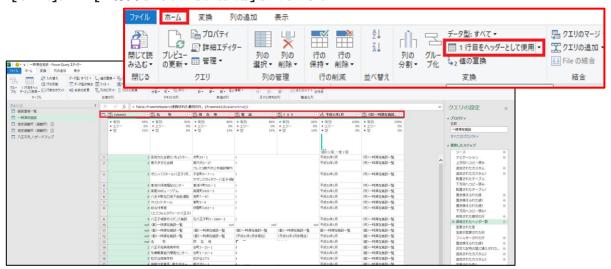
8. 施設区分と名称の列を右クリックし、「フィル]→[下]または[上]をクリックする



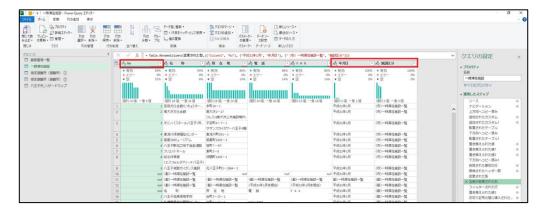
9. [ホーム]タブ→[行の削除]→「上位の行の削除」をクリックし、削除する上位の行数を入力する



10. [ホーム]タブ→[1 行目をヘッダーとして使用]をクリックする



11. 追加した「カスタム列 |や「列名 |の名前を指定しなかった列名があれば、列をダブルクリックして任意の列名を入力する



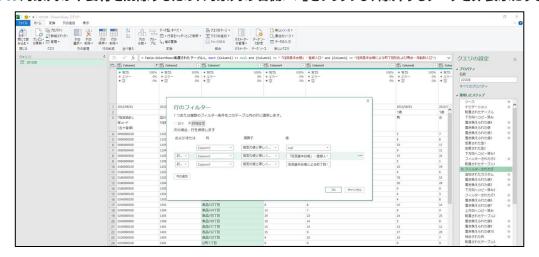
12. 対象列の右側「▼」をクリックして、数値フィルターまたはテキストフィルターを表示した後、「規定の値と等しくない」を選び、 [OK]をクリックする



13. 年月日の列を右クリックし、除外するデータを非表示にする

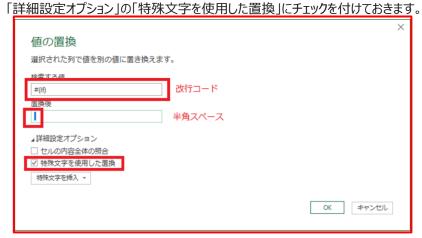


#### 14. 対象列の不要行を削除するため、対象列の右側「▼ |をクリックし、除外するデータを非表示にする



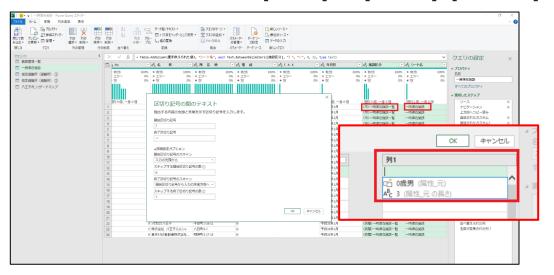
## 15. 対象列が選択された状態で右クリックし、[値の置換]をクリックし、値を置換する

置換する値「#(If)」(改行コード) → 置換後「 」(半角スペース)

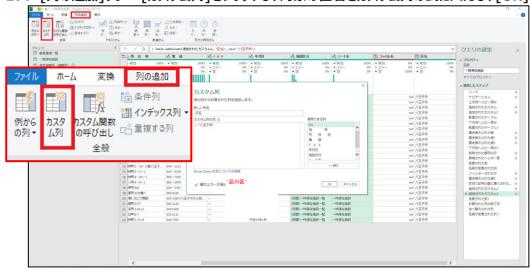


ここでは改行コードを 削除しています。

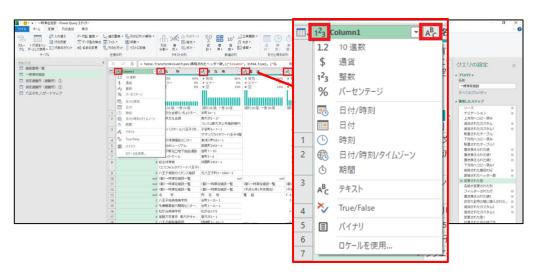
## 16. 開設届の列が選択された状態で右クリックし、「列の追加]タブ→「例からの列]→「選択範囲から]をクリックする



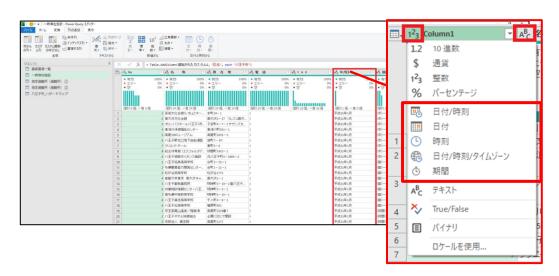
ここでは、施設区分列 から必要範囲を切り出 します。 17. [列の追加]タブ→[カスタム列]をクリックし、対象の区名をカスタム列の式に入力し、「OK]をクリックする



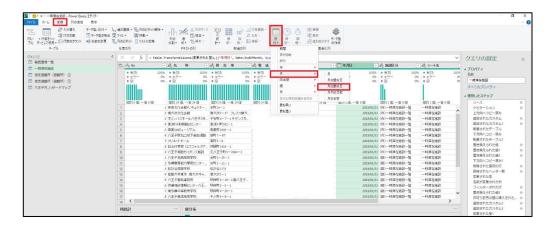
18. 列のヘッダー上部にある「ABC123」アイコンをクリックし、任意の型を選択する



19. 年月日のヘッダー上部にある「ABC123」アイコンをクリックし、日付型を選択する

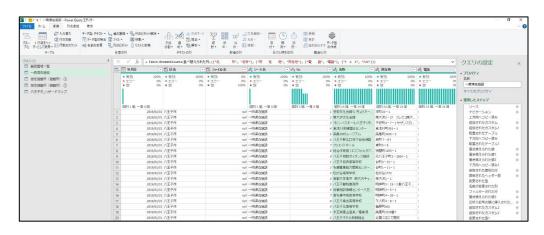


## 20. [変換]タブ $\rightarrow$ [日付] $\rightarrow$ [月] $\rightarrow$ [該当の日付を選択]をクリックし、年月日列に日付を付与する



例に使用したデータは 日付が末となっていた ため、「月の最終日」を 選択しています。

## 21. ヘッダー列名に不要な空白が存在する場合は、ヘッダー列名をダブルクリックし、任意の列名を入力する



## AI-OCR によるデータ変換

下記の手順で PDF から CSV データを生成する変換を進めます。

【準備】: Adobe Acrobat DC および CLOVA OCR READER を利用できる環境を確保します。(初回のみ)

- ① Adobe Acrobat DC で、AI-OCR に読み込むアクションを登録します。 (詳細: p. 88)
- ② アクションを実行し、データを作成します。 (詳細: p. 92)
- ③ CLOVA OCR READER で、AI-OCR 処理するテンプレートを作成します。 (詳細: p. 94)
- ④ CLOVA OCR READER で、AI-OCR 処理を行い、CSV ファイルを作成します。 (詳細: p. 97)

## ご参考

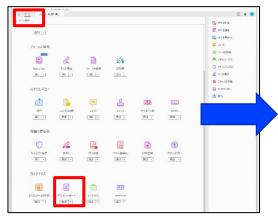
本書記載以外の、CLOVAの詳しい操作については、下記をご参照ください。

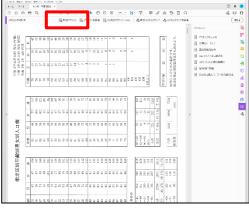
https://clova.line.me/clova-ocr/

## 手順(ケース[E])

## ① Adobe Acrobat DCで、AI-OCR に読み込むアクションを登録する

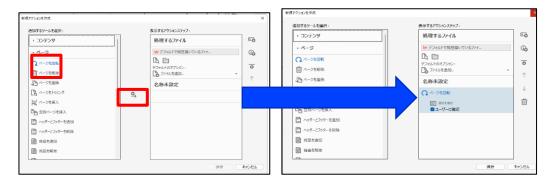
- 1. Adobe Acrobat DC を開く
- 2. ツールからアクションウイザードを開き、新規アクションをクリックする





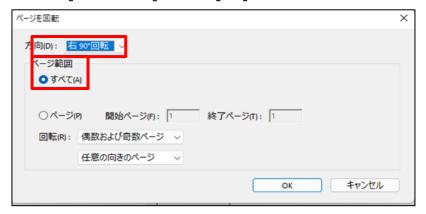
**アクション**は、一定の 操作を記録し、繰り返 し使うためのメニューで す。

3. [ページ]→[ページを回転]をクリック、[→]をクリックして設定を右側パネルに移動させる



右側パネルに移動した操作が、上から順に実行されます。

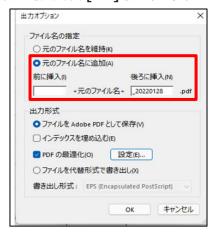
4. 「すべて」と [右 90 度回転]を選択し、[OK]をクリックする



5. [保存と書き出し]→[保存]をクリック、[→]をクリックして設定を右側パネルに移動させる



6. 設定を変更し、[OK]をクリックする



ファイル名の指定:元のファイルに追加

後ろに挿入:\_日付

## 7. PDF の最適化設定を変更し、[OK]をクリックする



## 互換性を確保: Acrobat10.0 及びそれ以降

ダウンサンプルの解像度:300ppi

カラーの圧縮: JPEG2000



オブジェクトを破棄: すべてにチェックを入れる



ユーザーデータを破棄: すべてにチェックを入れる

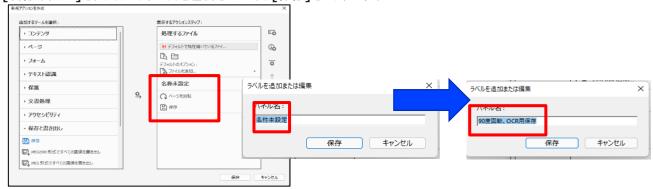


最適化: すべてにチェックを入れる

8. 「保存」をクリック後、保存名を「OCR用」と入力し、「OK」をクリックする

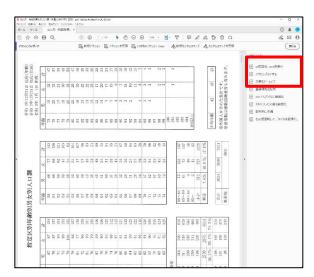


9. [名称未設定]をダブルクリック、処理名を入力し、[保存]をクリックする



## 10. アクション名を入力し、[保存]をクリックする

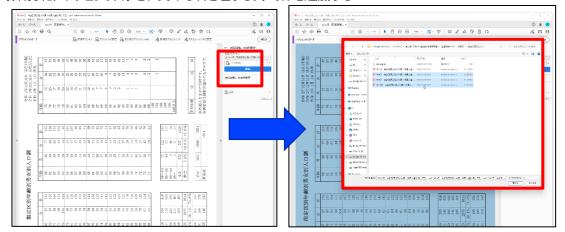




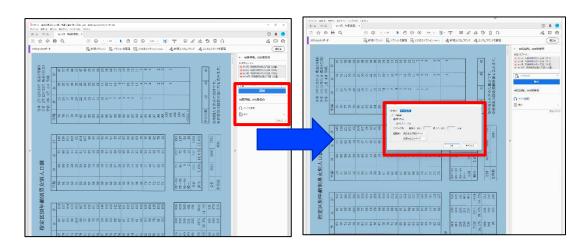
アクション名に追加されます。

## ② アクションを実行し、データを作成する

## 1. 作成したアクションリストをクリックし、処理するファイルを追加する

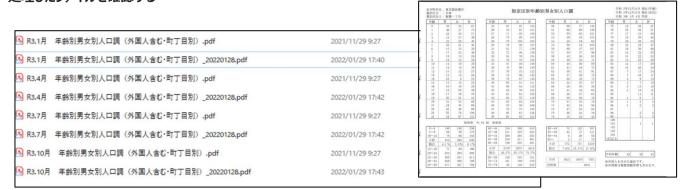


## 2. [開始]をクリックし、ページの回転で[OK]をクリックする



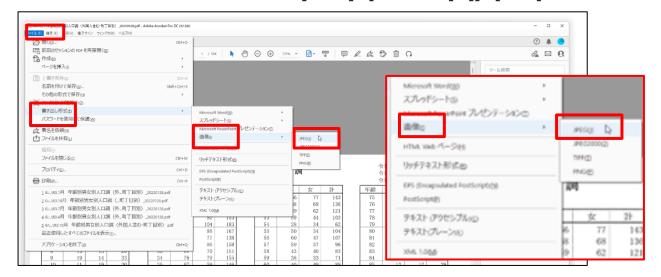
処理が実行され ます。処理に時 間がかかること があります。

## 3. 処理したファイルを確認する

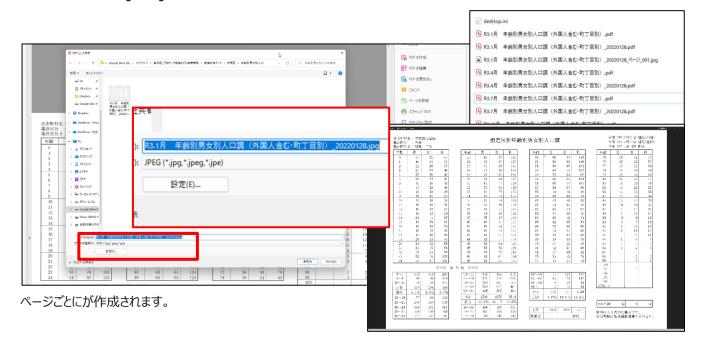


ファイル名末尾に指定した日付が入ったファイルが今回作成したデータです。

4. OCR 処理用の JPEG ファイルを 1 ページ出力する (「ファイル]  $\rightarrow$  [書き出し形式]  $\rightarrow$  []  $\rightarrow$  [JPEG])

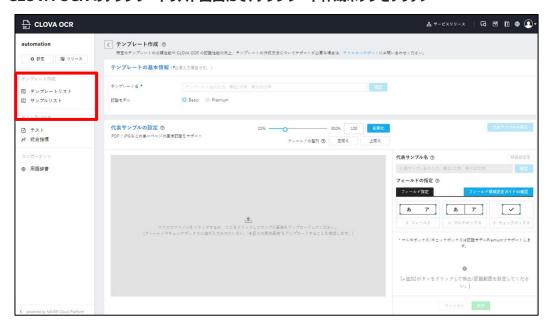


## 5. ファイル名を入力し、[保存]をクリックする



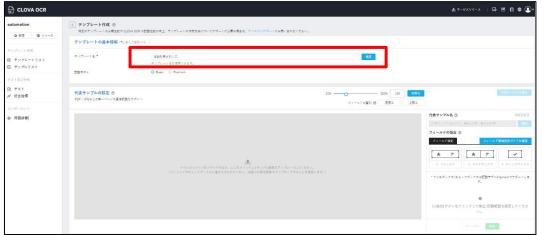
## ③ CLOVA OCR READER で、AI-OCR 処理するテンプレートを作成する

1. CLOVA OCR のテンプレートリスト画面にて、テンプレート作成ボタンをクリック



**※CLOVA OCR の** 初期設定がすべて完 了していることを前提 に手順を記載してお います。

2. テンプレート名を入力し、[確認]をクリック

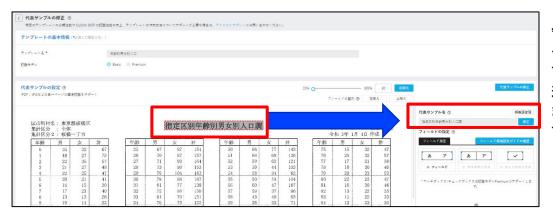


テンプレート名が確定 されたことを確認してく ださい。 3. ②で作成したファイル (JPG 形式) をドラッグ&ドロップで追加する



領域色が反転すると、ファイルが認識されています。そのまま保持状態を解除します。

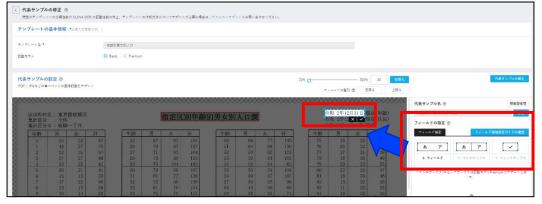
4. PDF のタイトルにする領域を選択した後、「代表サンプル名」に『指定区別年齢男女別人口調』と入力、「確認]をクリックする



領域はマウスのドラッ グ操作で選択できま す。

選択部分は破線で囲 まれ、色が反転しま す。

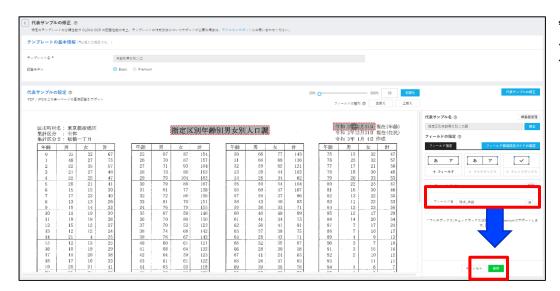
5. [+フィールド]をクリックし、読み取りたい領域を選択する



領域はマウスのドラッグ操作で選択できます。

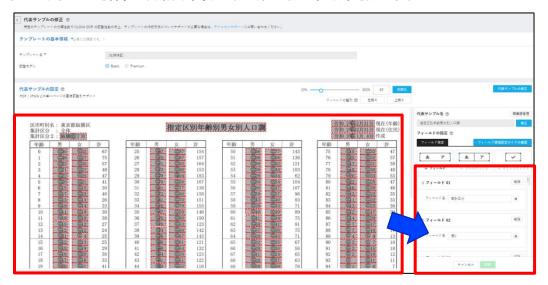
選択部分は破線で囲 まれ、色が反転しま す。

## 6. フィールド名を設定し、[保存]をクリックする



領域はマウスのドラッ グ操作で選択できま す。

#### 7. 手順 5、手順 6 を繰り返し、読み取りたい領域すべてにフィールドを設定する



最後に[保存]をクリッ クするのを忘れないよ うにしてください。

## **④ CLOVA OCR READER で、AI-OCR 処理を行い、CSV ファイルを作成する**

## 1. CLOVA OCR READER を開き、ファイルをアップロードする



※ CLOVA OCR READER の初期設 定がすべて完了してい ることを前提に手順を 記載しておいます。

画面右側が、ファイル をアップロードする領 域です。ここにドラッグ 操作か、領域をクリッ クしてファイルをアップ ロードします。



ファイル選択画面では、[+]をクリックしてファイルの追加を繰り返します。

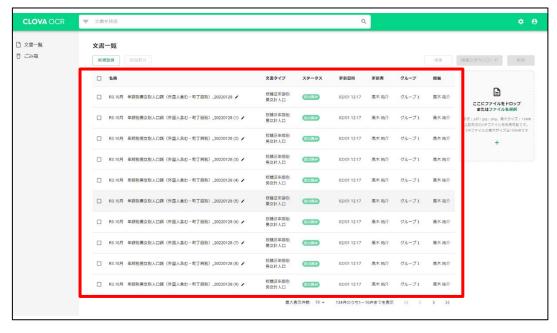
## 2. CLOVA OCR READER を開き、読み込む PDF ファイルをアップロードする



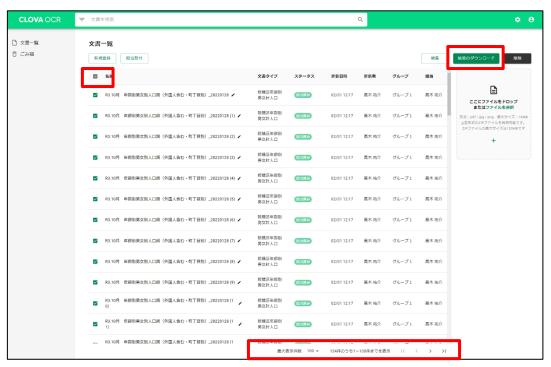
ファイルの追加が終わったら、[登録]をクリックします、



読み込みがはじまりま す。



読み込み後は、1 ファ イルにつき、1 行の文 書情報が表示されま す。 3. CSV データを入手したい文書のチェックボックスにチェックを付け、「結果のダウンロード」をクリックする



数ページにわたる数の 文書を一度に表示す るには、画面下部の 表示件数を変更しま す。

文書先頭のチェックボックスは、画面上部の「名前」左側のチェックボックスでまとめて 選択できます。

Α	В	С	D	E	F	G	Н	1	J	
title	集計区分	男0	男1	男2	男3	男4	男5	男6	男7	男8
指定区別年	板橋一丁目	37	33	29	19	24	17	22	11	
指定区別年	板橋二丁目	18	14	22	18	19	16	9	12	
指定区別年	板橋三丁目	16	26	23	19	22	18	22	31	
指定区別年	板橋四丁目	22	6	14	17	13	13	14	13	
指定区別年	加賀一丁目	33	48	59	42	58	67	58	61	
指定区別年	加賀二丁目	12	10	14	7	16	21	25	17	
指定区別年	大山東町	19	19	18	16	18	15	18	9	
指定区別年	大山金井町	25	12	16	14	10	20	8	14	
指定区別年	熊野町	7	17	10	9	15	22	13	12	
指定区別年	中丸町	19	28	30	26	16	18	31	19	
指定区別年	南町	16	9	7	10	10	22	5	14	
	title 指定区别与指定区别与指定区别与有指定区别与有指定区别与有指定区区别与有指定区区别与有指定区区别与有指定区区别与有指定区别与有限。	性itle 集計区分 指定区別4板橋一丁目 指定区別4板橋三丁目 指定区別4板橋三丁目 指定区別4板橋四丁目 指定区別4板橋四丁目 指定区別4加賀一丁目 指定区別4加賀二丁目 指定区別4大山東町	title         集計区分         男0           指定区別4板橋一丁目         37           指定区別4板橋二丁目         18           指定区別4板橋三丁目         16           指定区別4板橋四丁目         22           指定区別4加賀一丁目         33           指定区別4加賀二丁目         12           指定区別4大山東町         19           指定区別4株野町         7           指定区別4中丸町         19	title     集計区分     男0     男1       指定区別4板橋一丁目     37     33       指定区別4板橋二丁目     18     14       指定区別4板橋三丁目     16     26       指定区別4板橋四丁目     22     6       指定区別4板橋四丁目     33     48       指定区別4加賀二丁目     12     10       指定区別4大山東町     19     19       指定区別4大山金井町     25     12       指定区別4帐野町     7     17       指定区別4中丸町     19     28	title         集計区分         男0         男1         男2           指定区別4板橋一丁目         37         33         29           指定区別4板橋二丁目         18         14         22           指定区別4板橋四丁目         16         26         23           指定区別4板橋四丁目         22         6         14           指定区別4加賀一丁目         33         48         59           指定区別4加賀二丁目         12         10         14           指定区別4大山東町         19         19         18           指定区別4大山金井町         25         12         16           指定区別4         中丸町         7         17         10           指定区別4         中丸町         19         28         30	title         集計区分         男0         男1         男2         男3           指定区別4板橋一丁目         37         33         29         19           指定区別4板橋二丁目         18         14         22         18           指定区別4板橋三丁目         16         26         23         19           指定区別4板橋四丁目         22         6         14         17           指定区別4板橋四丁目         33         48         59         42           指定区別4加賀二丁目         12         10         14         7           指定区別4大山東町         19         19         18         16           指定区別4株野町         7         17         10         9           指定区別4中丸町         19         28         30         26	title     集計区分     男0     男1     男2     男3     男4       指定区別4板橋一丁目     37     33     29     19     24       指定区別4板橋二丁目     18     14     22     18     19       指定区別4板橋三丁目     16     26     23     19     22       指定区別4板橋四丁目     22     6     14     17     13       指定区別4加賀一丁目     33     48     59     42     58       指定区別4加賀二丁目     12     10     14     7     16       指定区別4大山東町     19     19     18     16     18       指定区別4大山金井町     25     12     16     14     10       指定区別4中丸町     7     17     10     9     15       指定区別4中丸町     19     28     30     26     16	title         集計区分         男0         男1         男2         男3         男4         男5           指定区別4板橋一丁目         37         33         29         19         24         17           指定区別4板橋二丁目         18         14         22         18         19         16           指定区別4板橋三丁目         16         26         23         19         22         18           指定区別4板橋四丁目         22         6         14         17         13         13           指定区別4板橋四丁目         33         48         59         42         58         67           指定区別4加賀二丁目         12         10         14         7         16         21           指定区別4大山東町         19         19         18         16         18         15           指定区別4株野町         7         17         10         9         15         22           指定区別4中丸町         19         28         30         26         16         18	title     集計区分     男0     男1     男2     男3     男4     男5     男6       指定区別4板橋一丁目     37     33     29     19     24     17     22       指定区別4板橋二丁目     18     14     22     18     19     16     9       指定区別4板橋四丁目     16     26     23     19     22     18     22       指定区別4板橋四丁目     22     6     14     17     13     13     14       指定区別4加賀一丁目     33     48     59     42     58     67     58       指定区別4加賀二丁目     12     10     14     7     16     21     25       指定区別4大山東町     19     19     18     16     18     15     18       指定区別4株野町     7     17     10     9     15     22     13       指定区別4中丸町     19     28     30     26     16     18     31	title         集計区分         男0         男1         男2         男3         男4         男5         男6         男7           指定区別4板橋一丁目         37         33         29         19         24         17         22         11           指定区別4板橋二丁目         18         14         22         18         19         16         9         12           指定区別4板橋三丁目         16         26         23         19         22         18         22         31           指定区別4板橋四丁目         22         6         14         17         13         13         14         13           指定区別4板橋四丁目         33         48         59         42         58         67         58         61            指定区別4加賀二丁目         12         10         14         7         16         21         25         17           指定区別4大山東町         19         19         18         16         18         15         18         9           指定区別4株野町         7         17         10         9         15         22         13         12           指定区別4中東町         19         28         30         26         16         18         31 </td

[結果のダウンロード] クリック後に、CSV ファイルがダウンロード されます。 戶順

# データを適切に配置し、データの精度を高める

手順3はデータのマッピングとクレンジングの作業を行います。

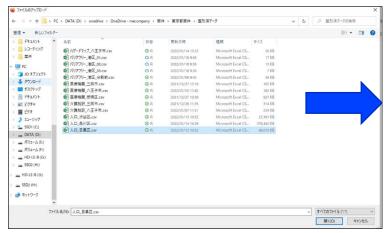
この作業は、目視で実施すると膨大な作業が発生し、手入力で実施すると誤りが発生する原因になりますので、基本的には API や外注作業にて対応を行います。この章に記載している作業を参考に、調達を実施して適切なデータに変換を行ってください。 この章ではその際に自動化が難しいケースについても記載をしています。

- ① これまでの手順で整備した CSV データをデータベースが読み込める場所に格納する (参考 p. 101)
- ② データをマッピングする (参考 p. 103)
- ③ マッピングデータを目視で確認する
- 4 ドメインデータを判別する
- ⑤ ドメインごとにクレンジングを行う (参考「クレンジング難度が高い例」p. 105)
- ⑥ データを正規化する (参考 p. 108)
- アメインデータをチェックする

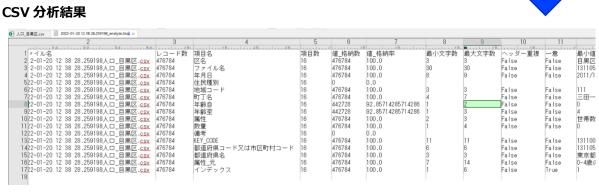
## 機能紹介/CSV 外形一括チェック

項目値の詳細チェック等が不要な場合等、簡易的に CSV ファイルの構造チェックが可能です。

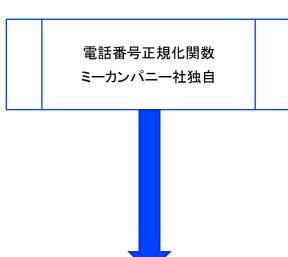
#### 対象ファイルを選択



※複数ファイル一括チェックが可能です。 出力最終段にて全一括チェック等に利用します。



※格納率、一意、最大値、最小値、最大桁、最少桁等から、外れ値や作業ミスの簡易検出が可能です。



# データの格納

整備対象データをデータベースに格納します。

## 参考ロジック

SELECT 項目 1、項目 2、項目 3、項目 4 FROM テーブル

名前	自動増加値	更新時間	データの長さ	エンジン	行	注积行
ハザードマップ_八王子市	0	2022-01-14 1	96 KB	InnoDB	274	
… バリアフリー_港区_01	0		16 KB	InnoDB	12	
	0		16 KB	InnoDB	12	
	0		16 KB	InnoDB	12	
■ 医療機関_三鷹市	0	2022-01-14 1	224 KB	InnoDB	668	
医療機関_八王子市	0	2022-01-14 1	320 KB	InnoDB	1175	
■ 医療機関_板橋区	0	2022-01-14 1	1552 KB	InnoDB	1927	
■ 介護施設_三應市	0	2022-01-20 0	352 KB	InnoDB	345	
₩ 介護施設_八王子市	0	2022-01-14 1	304 KB	InnoDB	835	
■ 人口_渋谷区	0	2022-01-14 1	36416 KB	InnoDB	285845	
■ 人口_品川区	0	2022-01-14 1	483328 KB	InnoDB	3724788	
₩人口_目黒区	0	2022-01-14 1	85632 KB	InnoDB	471877	

#### データの格納例(人口:目黒区)

区名	ファイル名 年月日	住民種別	地域コード	町丁名	年齢自	年齢至	屋性	数量	備者	KEY CODE	都道府県コード又は市区町村コー	都道府県名	属性 元	インデックス
目里区	131105 population 20110101, 2011/1/1	(Null)	111	駒場一丁目	0	999	総人口	3698	(Null)	13110001001	131105	東京都	0-999歳の総人口	1
日里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	0	999	男性	1824		13110001001	131105	東京都	0-999歳の男性	2
日里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	0	999	女性	1874		13110001001	131105	東京都	0-999歳の女性	3
目里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	0	4	男性	70		13110001001	131105	東京都	0-4歳の男性	4
目里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	0	4	女性	54		13110001001	131105	東京都	0-4歳の女性	5
日里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	5	9	男性	48		13110001001	131105	東京都	5-9歳の男性	6
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	5	9	女性	56		13110001001	131105	東京都	5-9歳の女性	7
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	10	14	男性	54		13110001001	131105	東京都	10-14歳の男性	8
里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	10	14	女性	50		13110001001	131105	東京都	10-14歳の女性	9
里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	15	19	男性	50		13110001001	131105	東京都	15-19歳の男性	10
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	15	19	女性	47	(Null)	13110001001	131105	東京都	15-19歳の女性	11
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	20	24	男性	128	(Null)	13110001001	131105	東京都	20-24歳の男性	12
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	20	24	女性	100	(Null)	13110001001	131105	東京都	20-24歳の女性	13
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	25	29	男性	210	(Null)	13110001001	131105	東京都	25-29歳の男性	14
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	25	29	女性	173	(Null)	13110001001	131105	東京都	25-29歳の女性	15
黒区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	30	34	男性	228	(Null)	13110001001	131105	東京都	30-34歳の男性	16
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	30	34	女性	162	(Null)	13110001001	131105	東京都	30-34歳の女性	17
黒区	131105_population_20110101. 2011/1/1		111	駒場一丁目	35	39	男性	185		13110001001	131105	東京都	35-39歳の男性	18
里区	131105_population_20110101, 2011/1/1	(Null)	111	駒場一丁目	35	39	女性	168		13110001001	131105	東京都	35-39歳の女性	19
黒区	131105_population_20110101, 2011/1/1	(Null)	111	駒場一丁目	40	44	男性	149	(Null)	13110001001	131105	東京都	40-44歳の男性	20
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	40	44	女性	135		13110001001	131105	東京都	40-44歳の女性	21
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	45	49	男性	109	(Null)	13110001001	131105	東京都	45-49歳の男性	22
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	45	49	女性	141	(Null)	13110001001	131105	東京都	45-49歳の女性	23
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	50	54	男性	75		13110001001	131105	東京都	50-54歳の男性	24
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	50	54	女性	109	(Null)	13110001001	131105	東京都	50-54歳の女性	25
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	55	59	男性	101	(Null)	13110001001	131105	東京都	55-59歳の男性	26
黒区	131105_population_20110101. 2011/1/1		111	駒場一丁目	55	59	女性	102	(Null)	13110001001	131105	東京都	55-59歳の女性	27
黒区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	60	64	男性	105	(Null)	13110001001	131105	東京都	60-64歳の男性	28
黒区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	60	64	女性	133	(Null)	13110001001	131105	東京都	60-64歳の女性	29
黒区	131105_population_20110101. 2011/1/1		111	駒場一丁目	65	69	男性	98		13110001001	131105	東京都	65-69歳の男性	30
黒区	131105_population_20110101, 2011/1/1	(Null)	111	駒場一丁目	65	69	女性	110	(Null)	13110001001	131105	東京都	65-69歳の女性	31
黒区	131105_population_20110101, 2011/1/1	(Null)	111	駒場一丁目	70	74	男性	72	(Null)	13110001001	131105	東京都	70-74歳の男性	32
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	70	74	女性	85		13110001001	131105	東京都	70-74歳の女性	33
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	75	79	男性	56	(Null)	13110001001	131105	東京都	75-79歳の男性	34
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	75	79	女性	86	(Null)	13110001001	131105	東京都	75-79歳の女性	35
里区	131105_population_20110101. 2011/1/1		111	駒場一丁目	80	84	男性	47		13110001001	131105	東京都	80-84歳の男性	36
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	80	84	女性	73	(Null)	13110001001	131105	東京都	80-84歳の女性	37
里区	131105_population_20110101. 2011/1/1	(Null)	111	駒場一丁目	85	999	男性	39	(Null)	13110001001	131105	東京都	85-999歳の男性	38
里区	131105 population 20110101, 2011/1/1		111	駒場一丁目	85	999	女性	90		13110001001	131105	東京都	85-999歳の女性	39

## データの格納例 (医療機関:八王子)

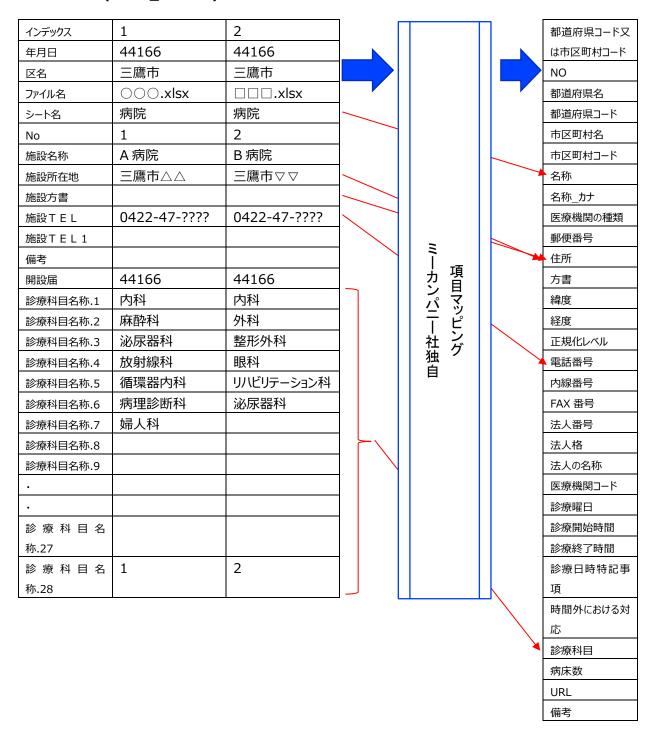


## データのマッピング

## データのフォーマット統一化

出典各様式を、共通フォーマットへ統一する作業です。

## 例: 医療機関(八王子\_医療機関)



#### 参考ロジック

```
select
 市区町村コード デジット as 都道府県コード又は市区町村コード
 , `インデックス` as NO
,' 東京都'as 都道府県名
, 13 as 都道府県コード
, `区名` as 市区町村名
, 市区町村コード as 市区町村コード
 , `施設名称` as 名称
, null as 名称_カナ
, null as 医療機関の種類
, 郵便番号 as 郵便番号
, `施設所在地` as 住所
, null as 方書
, null as 緯度
, null as 経度
, null as 正規化レベル
, case when `施設TEL` IS not null then `施設TEL` else `施設TEL1` end as 電話番号
, null as 内線番号
, null as FAX 番号
, null as 法人番号
, null as 法人格
, null as 法人の名称
, null as 医療機関コード
, null as 診療曜日
, null as 診療開始時間
, null as 診療終了時間
, null as 診療日時特記事項
, null as 時間外における対応
, concat_ws('|', `診療科目名称.1`, `診療科目名称.2`, `診療科目名称.3`, `診療科目名称.4`, `診療科
目名称.5`, `診療科目名称.6`, `診療科目名称.7`, `診療科目名称.8`, `診療科目名称.9`, `診療科目名
称.10`,
              `診療科目名称.11`, `診療科目名称.12`, `診療科目名称.13`, `診療科目名称.14`,
`診療科目名称.15`, `診療科目名称.16`, `診療科目名称.17`, `診療科目名称.18`, `診療科目名称.19`,
`診療科目名称.20`,
              `診療科目名称.21`, `診療科目名称.22`, `診療科目名称.23`, `診療科目名称.24`,
`診療科目名称.25`, `診療科目名称.26`, `診療科目名称.27`, `診療科目名称.28`
   ) as 診療科目
, null as 病床数
, null as URL
, `備考` as 備考
 from `医療機関_三鷹市`
```

## データのクレンジング

## クレンジング難度が高い例(データ構造の課題)

## 「1 セル 1 データ」となっていない。



## 営業日・営業時間の表記ゆれ

営業時間
窓口受付時間
受付営業時間
サービス提供時間
通いサービス提供時間
訪問可能時間帯
24 時間対応
休業日
窓口定休日
受付休業日
通いサービスの定休日



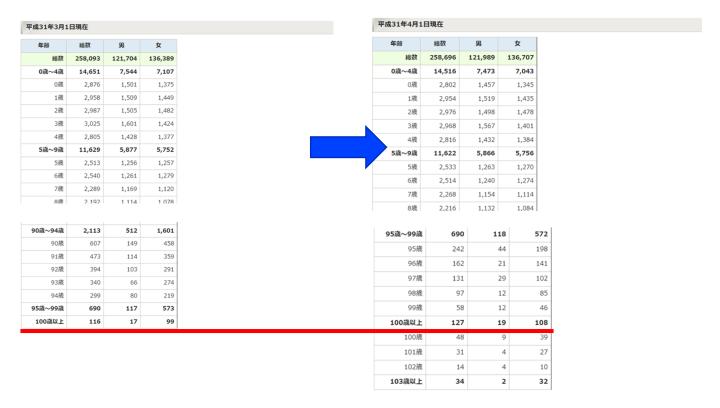
推奨データセット

介護事務所営業時間項目

利用可能曜日

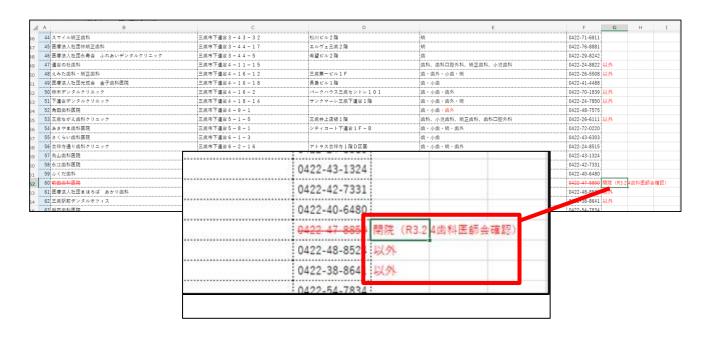
利用可能曜日特記事項

## 経年でデータが変わる



2019 年以降は 100 歳以上も 1 歳刻みで表示されている。

## 枠外にデータが表記されている



## クレンジング難度が高い例(データの課題)

ここでは単純なデータクレンジング作業では対応できず、目視を伴う作業が発生するケースを記載しています。

以下のようなケースでは、100%自動化は難しく、目視作業が必要になります。作業を依頼する場合には、こういった作業が発生する 事を念頭に依頼を行ってください。

## 枠外に閉院データが表記されている



#### 打ち消し線でデータが修正されている



## 異表記がデータに含まれる



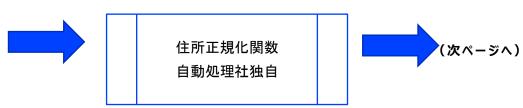
## データの正規化

### 住所

データ最終化前の、文字欠け、文字化け、表記揺れ対応、対応方針の正確性確認、不正値検知等のために必要な作業です。

### 例(住所:八王子)





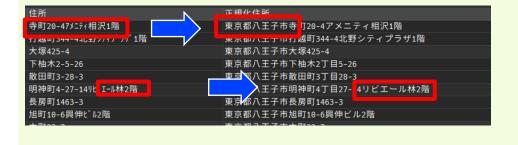
#### 住所正規化関数の役割

SOL 関数にて郵便番号を追加します。

住所を都道府県、市区町村を補完し、市区町村合併を反映、緯度経度を付与します。 ※住所については日本に標準となるものがない為、100%精度の適切な住所が存在しませ ん。その為、住所の正規化は非常に難度の高い作業になります。また日本全国の建物や番地 に緯度経度がついていない為、緯度経度を100%つける事はほぼ不可能です。緯度経度を 公開する際には、正規化精度を必ず一緒につけて公開するようにしてください。

※Google Place API など緯度経度の情報を安く取得できる方法はありますが、ほとんどの場 合は、取得したデータをそのまま公開してはいけない利用規約になっており、公開すると利用規 約違反になってしまう事が多く、莫大な損害賠償を請求される可能性があります。緯度経度を オープンデータとして公開する際には、許諾を適切に取ったデータのみ公開するようにしてくださ (1<sub>o</sub>

例:住所正規化



※前ページからの続き

### 元情報住所

### 表記揺れ統一後の正規化データ

### 元住所

住所	i 規化住所結果JSON	@json_normalize_string_wk	正規化住所	@json_normalize_string_wk:
京都寺町20-47メニティ相沢1階	{ esults": {"id": 12428729	, " {"results": {"normalize_string"	東京都八王子市寺町20-4アメニティ相沢1階	{"results": {"normalize_string"
夏京都打越町344-4北野シティブラザ1階	{ results": {"id": 12375697	, " {"results": {"normalize_string"	東京都八王子市打越町344-4北野シティプラザ1階	{"results": {"normalize_string"
<b>[京都大塚425-4</b>	{ results": {"id": 12380392	, " {"results": {"normalize_string"	東京都八王子市大塚425-4	{"results": {"normalize_string"
京都下柚木2-5-26	{ results": {"id": 12416142	, " {"results": {"normalize_string"	*東京都八王子市2-5-26	("results": ("normalize_string"
京都散田町3-28-3	{ results": {"id": 12411251	, " {"results": {"normalize_string"	東京都八王子市3-28-3	("results": ("normalize_string"
夏京都明神町4-27-14リピエール林2階	{ results": {"id": 12462204	, " {"results": {"normalize_string"	*東京都八王子市4-27-14リビエール林2階	{"results": {"normalize_string"
原都長房町1463-3	{ results": {"id": 12437918	, " {"results": {"normalize_string"	'東京都八王子市長房町1463-3	("results": ("normalize_string"
原都旭町10-6興伸ビル2階	{ results": {"id": 12367358	, " {"results": {"normalize_string"	東京都八王子市旭町10-6興伸ビル2階	("results": ("normalize_string"
京都本町33-7	{ results": {"id": 0, "area":	", {"results": {"normalize_string"	*東京都本町33-7	0
京都大塚622-9ニューハイム井上201	{results": {"id": 12379784	, " {"results": {"normalize_string"	・東京都八王子市大塚622-9ニューハイム井上201	{"results": {"normalize_string"
京都高尾町1602	{ results": {"id": 12419282	, " {"results": {"normalize_string"	東京都八王子市高尾町1602	("results": ("normalize_string"
京都元横山町2-1-20	{ results": {"id": 12470417	, " {"results": {"normalize_string"	東京都八王子市町2-1-20	("results": ("normalize_string"
原京都狭間町1450-1めじろ台コーオ・ラス	{ results": {"id": 12446314	, " {"results": {"normalize_string"	東京都八王子市狭間町1450-1めじろ台コーポラス	("results": ("normalize_string"
京都みなみ野3-9-7	{ results": {"id": 12489190	, " {"results": {"normalize_string"	東京都八王子市野3-9-7	{"results": {"normalize_string"
京都八日町4-9幸ビル3階	{ results": {"id": 12473681	, " {"results": {"normalize_string"	東京都八王子市八日町4-9幸ビル3階	{"results": {"normalize_string"
京都長沼町200-3旗野コーポーラス 1階	{results": {"id": 12434268	, " {"results": {"normalize_string"	・東京都八王子市長沼町200-3旗野コーポラス1階	{"results": {"normalize_string"
京都片倉町443-4	{ results": {"id": 12388043	, " {"results": {"normalize_string"	東京都八王子市片倉町443-4	{"results": {"normalize_string"
京都子安町4-6-1Phit 1/3階	{results": {"id": 12408873	, " {"results": {"normalize_string"	・東京都八王子市4-6-1Phiビル3階	{"results": {"normalize_string"
京都子安町4-6-1Phit 1/2階	{ results": {"id": 12408873	, " {"results": {"normalize_string"	東京都八王子市4-6-1Phiビル2階	{"results": {"normalize_string"
『京都明神町3-26-10土屋ピル1階	esults": ("id": 12462013	" ("result ("normalize_string"	東京都八王子市3-26-10土屋ビル1階	{"results": {"normalize_string"
京都東浅川町879-6	{ results": {"id":	rmalize_string	東京都八王子市東浅川町879-6	{"results": {"normalize_string"
原京都子安町1-2-8守屋ピル3階	{ esults": {"id":	normalize_string	東京都八王子市1-2-8守屋ビル3階	{"results": {"normalize_string
京都元本郷町3-16-7	esults": {"id": 12469867	, " {"results": {"normalize_string"	東京都八王子市町3-16-7	{"results": {"normalize_string"
京都中野上町5-5-3	{results": {"id": 0, "area":	", {"results": {"normalize string"	東京都中野区上町5-5-3	0
夏京都松木14-3ザ・テラス松木103	{results": {"id": 12456134	, " {"results": {"normalize_string"	東京都八王子市松木14-3ザ・テラス松木103	{"results": {"normalize_string"
京都高月町365-1 2階	{results": {"id": 12420849	, " {"results": {"normalize string"	東京都八王子市高月町365-1-2階	{"results": {"normalize string"
京都子安町4-21-7	{results": {"id": 12409000	, " {"results": {"normalize string"	東京都八王子市4-21-7	{"results": {"normalize string"
京都横山町8-19牛久保ピル3階	{results": {"id": 12476628	, "{"results": {"normalize string"	東京都八王子市横山町8-19牛久保ビル3階	{"results": {"normalize string
京都久保山町2-43-4	{results": {"id": 12483042	, " {"results": {"normalize string"	東京都八王子市町2-43-4	{"results": {"normalize string
京都大和田町6-5-17	{results": {"id": 12384179	, " {"results": {"normalize_string"	'東京都八王子市町6-5-17	{"results": {"normalize_string"
原京都南大沢2-277レスコ南大沢4階	esults": {"id": 12458559	, " {"results": {"normalize_string"	東京都八王子市2-27フレスコ南大沢4階	{"results": {"normalize_string"
原京都横山町11-4井藤ビル4階FC	esults": {"id": 12476589	, " {"results": {"normalize string"	東京都八王子市横山町11-4井藤ビル4階FC	{"results": {"normalize string"
原京都別所1-17-1クレセントヒルス 2階-B	esults": {"id": 12452677	, " {"results": {"normalize string"	東京都八王子市-17-1クレセントヒルズ2階-B	{"results": {"normalize string"
京都日吉町13-31		. " {"results": {"normalize string"	東京都八王子市日吉町13-31	{"results": {"normalize string"
京都北野町546-6小山ピル2階	{ esults": {"id": 12399888	, " {"results": {"normalize string"	東京都八王子市北野町546-6小山ビル2階	{"results": {"normalize string"
京都長房町375	{ results": {"id": 12436119	, " ("results": ("normalize_string"	市京都八千子市長房町375	{"results": {"normalize string"
京都鹿島6		, " {"results": {"normalize string"	東京都八王子市鹿島6	{"results": {"normalize string
京都東浅川町1028		, " {"results": {"normalize string"	東京都八王子市東浅川町1028	{"results": {"normalize string
京都本町3-4		, " {"results": {"normalize_string"	東京都八王子市本町3-4	{"results": {"normalize_string"
京都めじろ台2-49-19		, " {"results": {"normalize string"	東京都八王子市台2-49-19	{"results": {"normalize string"
京都子安町1-4-5	{ esults": {"id": 12407438	, " {"results": {"normalize string"	東京都八千子市1-4-5	{"results": {"normalize string
京都大和田町4-11-19		. " {"results": {"normalize string"	東京都八王子市町4-11-19	{"results": {"normalize_string"
「京都片倉町492-1綱木ピル1階	{ results": {"id": 12386681	. " {"results": {"normalize string"	東京都八王子市片倉町492-1綱木ビル1階	{"results": {"normalize string"
京都明神町4-12-2	,	, " {"results": {"normalize_string"	東京都八王子市4-12-2	{"results": {"normalize_string
京都川口町1088-4		, " ("results": ("normalize string"		{"results": {"normalize string"
i 京都北野台4-27-2		, " {"results": {"normalize string"	東京都八王子市4-27-2	{"results": {"normalize string
京都高倉町48-2高橋ビル2階		, "{"results": {"normalize_string"	東京都八王子市高倉町48-2高橋ビル2階	{"results": {"normalize string
京都北野町555-17		, " ("results": ("normalize string"	東京都八王子市北野町555-17	{"results": {"normalize string"
京都散田町5-1-6	,	, "{"results": {"normalize_string"	* 車京都八王子市5-1-6	{"results": {"normalize_string"
E京都横川町536-3	{ results": {"id": 12474858			("results": ("normalize_string"

「表記揺れ統一後の正規化データ」では、先頭の「都道府県」は、前処理にて付与されます。 元住所に東京都が記載されており、「東京都東京都板橋区」となっても問題は発生しません。 ※余計な空白は削除(全半角統一等)されます。

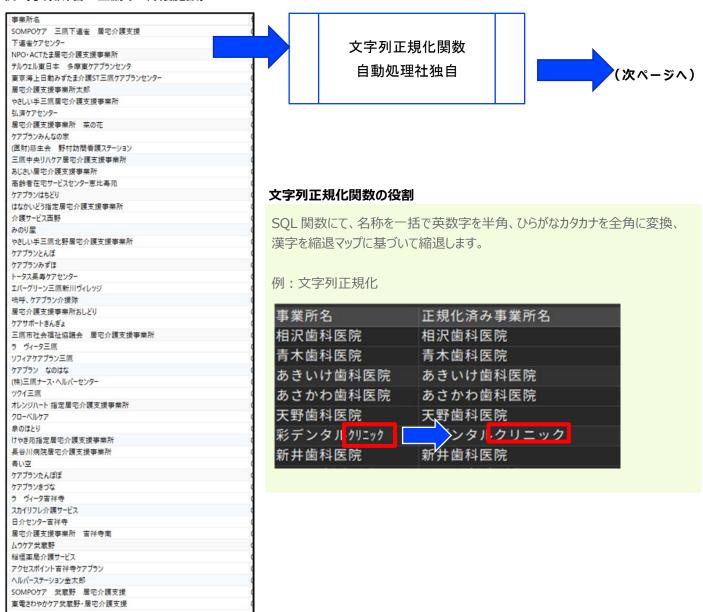
### 参考

「元住所」では、郵便番号マスタ比較にて、有効であるか判別ができます。

### 事業所名(文字列)

データ最終化前の、表記揺れ対応、対応方針の正確性確認、不正値検知等のために必要な作業です。

### 例(事業所名:三鷹市 介護施設)



出所:三鷹市・三鷹市介護保険事業者連絡協議会「介護サービス事業者ガイドブック」

※前ページからの続き

### 元事業所名

### 表記揺れ統一後の正規化データ

<b>拉墨尼克</b>	To the second se
事業所名	results
SOMPOケア 三鷹下連雀 居宅介護支援	SOMPOケア 三鷹下連雀 居宅介護支援
下連省ケアセンター	下連省ケアセンター
NPO・ACTたま居宅介護支援事業所	特定非営利活動法人・ACTたま居宅介護支援事業所
テルウェル東日本 多摩東ケアプランセンタ	テルウェル東日本 多摩東ケアプランセンタ
東京海上日動みずたま介護ST三鷹ケアプランセンター	東京海上日動みずたま介護ST三鷹ケアプランセンター
居宅介護支援事業所太郎	居宅介護支援事業所太郎
やさしい手三鷹居宅介護支援事業所	やさしい手三鷹居宅介護支援事業所
弘済ケアセンター	弘済ケアセンター
居宅介護支援事業所 菜の花	居宅介護支援事業所(菜の花)
ケアプランみんなの家	ケアブランみんなの家
(医財)慈生会 野村訪問看護ステーション	医療法人財団慈生会 野村訪問看護ステーション
三鷹中央リハケア居宅介護支援事業所	三鷹中央リハケア居宅介護支援事業所
あじさい居宅介護支援事業所	あじさい居宅介護支援事業所
高齢者在宅サービスセンター恵比寿苑	高齢者在宅サービスセンター恵比毒苑
ケアプランはちどり	ケアプランはちどり
はなかいどう指定居宅介護支援事業所	はなかいどう指定居宅介護支援事業所
介護サービス西野	介護サービス西野
みのり屋	みのり屋
やさしい手三鷹北野居宅介護支援事業所	やさしい手三鷹北野居宅介護支援事業所
ケアプランとんぼ	ケアブランとんぼ
ケアプランみずほ	ケアブランみずほ
トータス長寿ケアセンター	トータス長寿ケアセンター
エバーグリーン三鷹新川ヴィレッジ	エバーグリーン三應新川ヴィレッジ
嗚呼、ケアプラン介援隊	嗚呼、ケアプラン介援隊
居宅介護支援事業所おしどり	居宅介護支援事業所おしどり
ケアサポートきんぎょ	ケアサポートきんぎょ
三鷹市社会福祉協議会 居宅介護支援事業所	三應市社会福祉協議会 居宅介護支援事業所
ラ ヴィータ三鷹	ラ ヴィータ三鷹
ソフィアケアプラン三鷹	ソフィアケアプラン三鷹
ケアブラン なのはな	ケアブラン なのはな
(株)三鷹ナース・ヘルパーセンター	株式会社三鷹ナース・ヘルパーセンター
ツクイ三鷹	ツクイ三鷹
オレンジハート 指定居宅介護支援事業所	オレンジハート 指定居宅介護支援事業所
クローベルケア	クローベルケア
泉のほとり	泉のほとり
けやき苑指定居宅介護支援事業所	けやき苑指定居宅介護支援事業所
<b></b> 長谷川病院居宅介護支援事業所	長谷川病院居宅介護支援事業所
青い空	青い空
ケアプランたんぽぽ	ケアブランたんぽぽ
ケアプランきづな	ケアブランきづな
ラ ヴィータ吉祥寺	ラ ヴィータ吉祥寺
スカイリフレ介護サービス	スカイリフレ介護サービス
THE PERSON NAMED IN COLUMN NAM	CONTRACTOR WAY

出所:三鷹市・三鷹市介護保険事業者連絡協議会「介護サービス事業者ガイドブック」

「表記揺れ統一後の正規化データ」では、法人格表記統一、スペース統一、記号全半角統一等、正規化内容の詳細な制御が可能です。

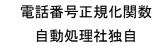
例: 法人格は削除等

### 電話番号

データ最終化前の、表記揺れ対応、対応方針の正確性確認、不正値検知等のために必要な作業です。

### 例(電話番号:八王子 医療機関)







### 電話番号正規化関数の役割

SQL 関数にて、電話番号は住所を元に市外局番を補完し、行政基本情報データ連携モデル準拠に変換します。

例:電話番号正規化



一部地域で 03 以外の番号を付与しています。

- 三鷹市中原一丁は03、それ以外は0422
- 八王子市は 042※2006/3/5 までは 0426

※前ページからの続き

電話番号 表記揺れ統一後の正規化データ 電話番号 convert\_domain\_data('001', 電 convert\_domain\_data('001', co convert\_domain\_data('001', co 626-4188 042-626-4188 648-3566 042-648-3566 677-0480 042-677-0480 679-1188 042-679-1188 661-0249 042-661-0249 645-3116 042-645-3116 661-6500 042-661-6500 644-3662 042-644-3662 623-2052 042-623-2052 676-8288 042-676-8288 661-0128 042-661-0128 642-1827 042-642-1827 664-5862 042-664-5862 632-7676 042-632-7676 625-5484 042-625-5484 636-4618 042-636-4618 637-2622 042-637-2622 625-7520 042-625-7520 625-7142 042-625-7142 646-9939 042-646-9939 664-5901 042-664-5901 656-1631 042-656-1631 625-6491 042-625-6491 626-3944 042-626-3944 675-4182 042-675-4182 696-5373 042-696-5373 622-1427 042-622-1427 646-6547 042-646-6547 691-7100 042-691-7100 644-8852 042-644-8852 677-5383 042-677-5383 646-5497 042-646-5497 679-0980 042-679-0980 622-4019 042-622-4019 646-3918 042-646-3918 661-4552 042-661-4552 676-2642 042-676-2642 669-7666 042-669-7666

「表記揺れ統一後の正規化データ」で、値が空欄な箇所は、元データが電話番号として無効な場合です。(桁不足等)

赤枠内1列目: 元データを電話番号正規化した場合

赤枠内 2 列目 : 元データ先頭に 03 を付記して、電話番号正規化した場合 赤枠内 3 列目 : 元データ先頭に 042 を付記して、電話番号正規化した場合 手順 **4** 

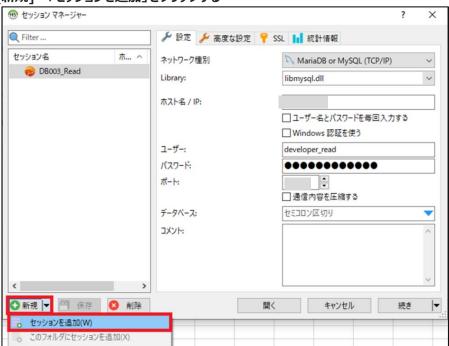
# データの保存(エクスポート)

## データベース接続の設定

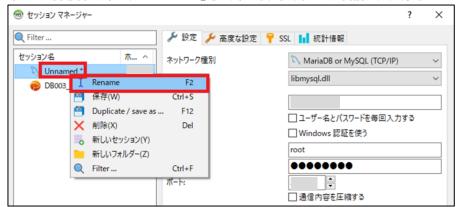
ここでは、HeidiSQLを用いて、データベースに接続する設定を行います。

任意のデータベース接続アプリケーションを使用される場合は、説明文をお使いのアプリケーションに適宜読み替えてご参照ください。

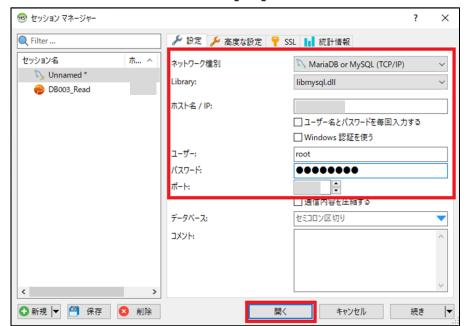
- 1. HeidiSQL等、データベースに接続できるアプリケーションを開く
- 2. [新規]→「セッションを追加」をクリックする



3. セッション名を右クリック、「Rename」をクリックした後、任意の名前に変更する



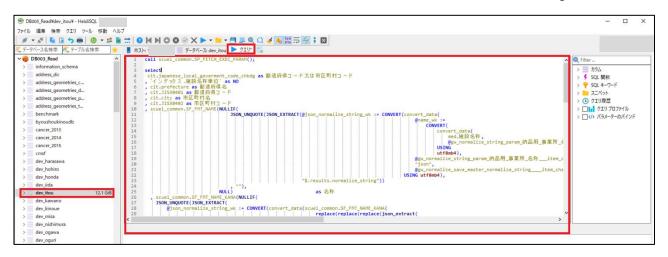
4. 「設定」タブをクリック、設定を変更した後、「開く」をクリックする



ネットワーク種別: MariaDB or MySQL(TCP/IP)を選択
Library: libmysql.dllを選択
ホスト名/IP:
接続先 DB のホスト名/IPを入力
ユーザー:
接続先 DB のユーザーを入力
パスワード:
接続先 DB のパスワードを入力
ポート:
接続先 DB のポートを入力

# データのエクスポート

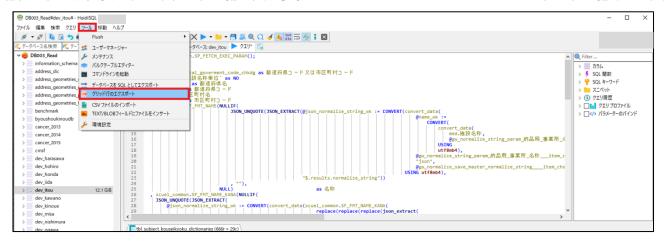
1. エクスポート対象のデータベースをクリックして選択、「▶クエリ」タブをクリックし、整形・正規化を行う SQL クエリを入力する



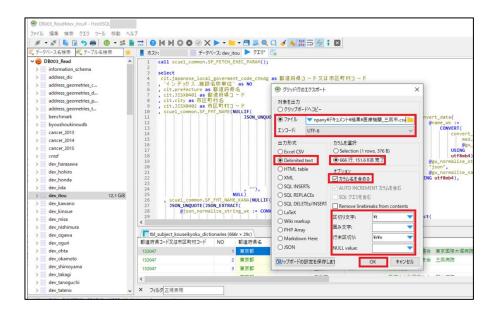
2. 「▶」→「実行」をクリックして、クエリを実行する



3. 「ツール」タブ→「グリッド行のエクスポート」をクリックする



### 4. 「対象を出力」で任意の出力先を指定した後、設定を変更し、[OK]をクリックする



**エンコード:** 「UTF-8」を選択

出力形式:「Delimitied text」を

選択

カラムを選択:下側(複数行)を

選択

オプション: 「カラム名を含める」に

チェック

区切文字:「¥t」を選択

**囲み文字:** (なし)

行末区切り:「¥r¥n」

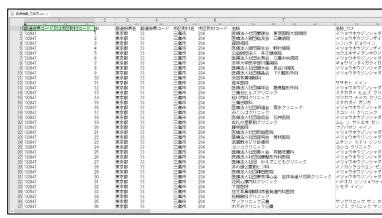
(Windows 形式の改行)を選

択

**NULL value:** (なし)

### 5. 指定したパスにファイルがエクスポートされるので、ファイルの内容を確認する





エクスポートされた内容(例)

# **FAQ**

# アプリケーション

項目名	· · · · · · · · · · · · · · · · · · ·
CSV ファイルの保存形式について	CSV は、「テキスト(タブ区切り)」で保存しています。
教えてください。	データにカンマが含まれる可能性を考慮し、タブ区切りを指定しています。

# データ

項目名	·····································
重複データが存在する。	重複データが目視で見つかった場合は、1 行分を手作業で削除してください。(存在が確認できている資料:板橋区データ)
名寄せや ID 付与がわからない。	東京都では来年度以降にもデータ整備モデル事業の実施を予定しています。名寄せや ID
	付与については来年度のマニュアルに実装される予定です。

# **Power Query**

項目名	·····································
	個別対応が必要となる場合、下記の情報をご参考ください。
	●ブック名 YYYYMM 以外文字列が c ではない:2012/05、2012/12、2019/05、
	2021/08、2021/11、2021/12
	●シート名の和暦年が全角: 2016/01、2021/12
	●シート名の月が全角: 2014/07、2020/10、2021/09
	●シート名の和暦年月が全角:2020/05、2020/06、2020/08、2020/09、
	2020/11、2020/12、2021/01、2021/02、2021/03,2021/06、2021/08
	● シート名和暦年月以外: 2018/08(201808c)、2018/09(201809c)、
	2019/05(2019 年 5 月)、2019/06(C,町丁別年齢別《保存用》)、2019/07(2019
ブック名、シート名、年月日に関わ	年7月)、2019/08(2019年8月)、2019/09(2019年9月)、2019/10(2019年
る部分で、記載内容と異なる状況	10月)、2019/11(2019年11月)、2019/12(2019年12月)、2020/02(2020
がある場合(品川区の月ごとデー	年 2 月)
夕等)	● シート名に元号が追加: 2015/01~2015/04
, ,	●令和元年のため、手順3年月日列の整形変更有:2019/05~2019/12
	●2019/05、2019/06 :値の置換「令和元年 5 月 1 日現在」 →「2019 年 5
	月1日」
	●2019/07、2019/08、2019/10、2019/12:値の置換「令和元年 7 月 1 日
	現在」 →「2019年7月1日」
	● <b>2019/09 :</b> 値の置換「令和元年 9月1日 現在」→「2019年9月1日」
	● <b>2019/11 :</b> 値の置換「令和元年11月1日 現在」→「2019年11月1日」
	●元号()で記載のため、手順3年月日列の整形変更有:2021/11
	●2021/11 :値の置換「2021 年(令和3年)11 月 1 日現在」→「2021 年 11
	月1日」

項目名	詳細
マッピング操作でエラーが出る	データ群が多い場合(例:八王子のハザードマップでは約 270 件)は、マッピング時にデータ群を分割する必要があります。 データを分割の後、再度マッピングを実施してください。
市外局番付与が難しい	八王子市 市外局番は、「0426」であるが、 市外局番末尾の「6」は元データに付与されている。 混在しており、個別判断は困難な状態。 三鷹市 23 区外であるが、「03」「0422」

# 東京都データプラットフォーム

## 行政データ整備モデル事業

## データ整備マニュアル

### マニュアルの改善に向けて

本マニュアルの改善に向けて、気づいた点などございましたら、データ利活用の推進のページにある「ご意見・ご感想」フォームからお願いいたします。

https://www.digitalservice.metro.tokyo.lg.jp/society5.0/index.html

2022 年 3 月 第 1 版 発行元·著作 東京都 委託先 日本総合研究所、株式会社自動処理

本書掲載の内容は 2022 年 3 月現在の情報です。

内容には万全を期しておりますが、システム改変等により実際の内容と異なる可能性があることをご了承ください。 操作説明に使用している画面は代表的なデータを使用するため、実際の利用画面と異なることがあります。

本書に記載されている会社名、システム名、製品名は一般に各社の登録商標または商標です。 なお、本文および図表中では、「™」、「®」は明記しておりません。