

機械学習研修

Day 1

author = 'Toshifumi Tsutsumi' AND
presented_at = '2021-07-15'



堤 利史

[Twitter: @tosh2230](https://twitter.com/tosh2230)

データエンジニア

GMO Pepabo, inc.
技術部 データ基盤チーム
(2020/12~)

- 大規模データを ETL (Extract, Transform, Load)する方法を知る
- データウェアハウスの使い方やETLの基本的な考え方に慣れて、自分が利用するデータを抽出・加工できるようになる

NOTE:

演習では、一般社団法人 データサイエンティスト協会が提供する

“データサイエンス100本ノック（構造化データ加工編）” のデータを BigQuery にロードしています

<https://github.com/The-Japan-DataScientist-Society/100knocks-preprocess>

2021/7/15(木) 15:00 - 18:00

1. データエンジニアリングってなんですか？
2. Bigfoot 最速入門
3. Google BigQuery 入門
4. Google DataStudio 入門

2021/7/16(金) 15:00 - 17:00

5. ETLってなんですか？
6. ETL Ultraquick Tutorial
7. データエンジニアリングってなんですか？

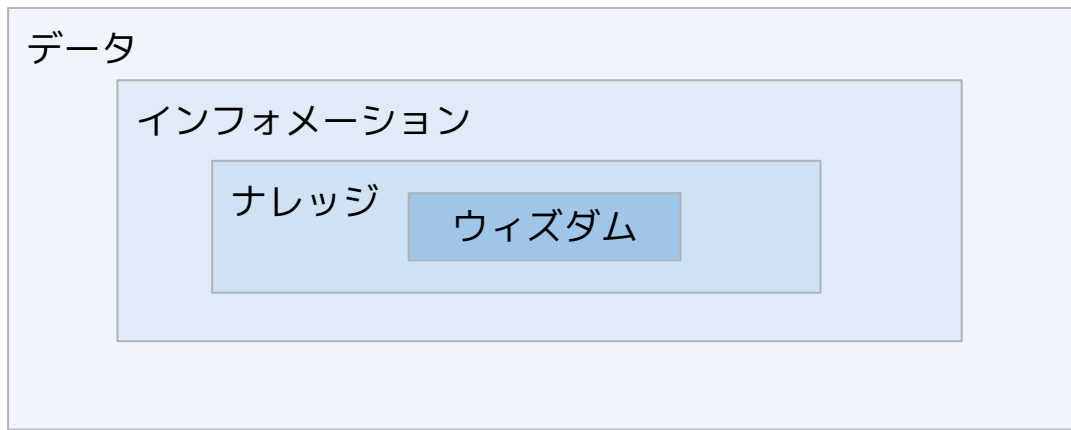


Section 1

データエンジニアリングって
なんですか？

データ × インフォメーション *

- データは「インフォメーションの原材料」
- インフォメーションは「コンテキストを持ったデータ」
 - コンテキストを持った = 意味付けされている



date-id-name-num0-num1-num2-num3

2021-07-15-0123456789-test0-12-34-567-890

2021-07-15-9876543210-test1-98-76-543-210

date-id-name-num0-num1-num2-num3

2021-07-15-0123456789-test0-12-34-567-890

2021-07-15-9876543210-test1-98-76-543-210

date,*id*,*name*,num0,num1,*num2*,num3

2021-07-15,*0123456789*,*test0*,12,34,*567*,890

2021-07-15,*9876543210*,*test1*,98,76,*543*,210

手段1: インフォメーションになりうるデータをつくる

- そもそも、データを集めるのはなぜか？
 - 何かを知りたいから
 - 適切な理解は適切なデータ作成から
- データをつくるにはナレッジ(=ドメイン知識)が必要
 - 目的に沿ったデータをつくることが大事
 - システムだけではなく人もデータをつくる

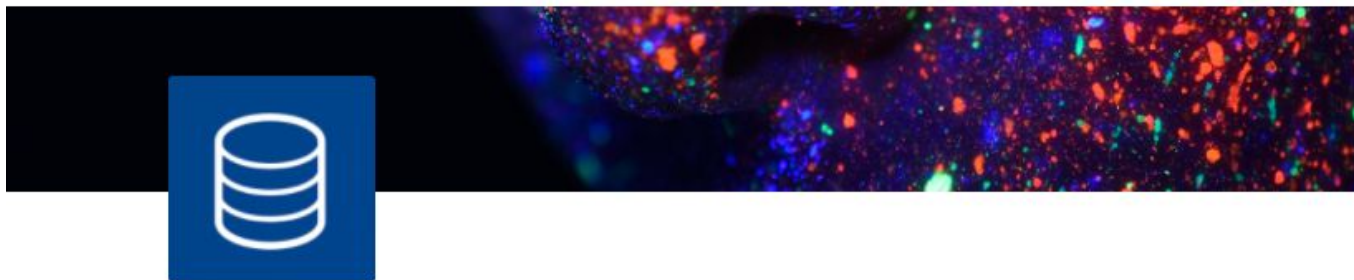
手段2: 既存のデータをインフォメーションに変換する

- コンテキストがわかりにくい場合
 - 誰もがわかる表現に加工する
 - enum(列挙型): ex, 0-6で曜日を表現
 - 信号処理
 - 別のデータとつなぎ合わせる
 - ユーザー情報 * 注文情報
 - メタデータ(データを説明するデータ)を追加する
 - 5W1H
 - 表現したいこと
 - 大元の発生源

手段2: 既存のデータをインフォメーションに変換する

- コンテキスト自体をまだ見いだせていない場合
 - データの中から、隠れたコンテキストを探す
 - 暗黙知を形式知に
 - 統計
 - パターン認識
 - 機械学習

DX Criteria (v202104) / データ駆動



データ駆動



「データの世紀」と呼ばれるように、企業の競争戦略にとってデータの利活用は必要不可欠なものです。しかし、そもそもデータの取得ができていなかったり、データのリテラシーが低くうまく経営に行かせないということも多くあります。

また、機械学習やデータサイエンスの知見を利用したアプリケーションには、それを支えるビッグデータ処理の仕組みが合わせて必要になります。

<https://dxcriteria.cto-a.org/d44c9631ebfb4c8f8ff1ae69361dfa65>

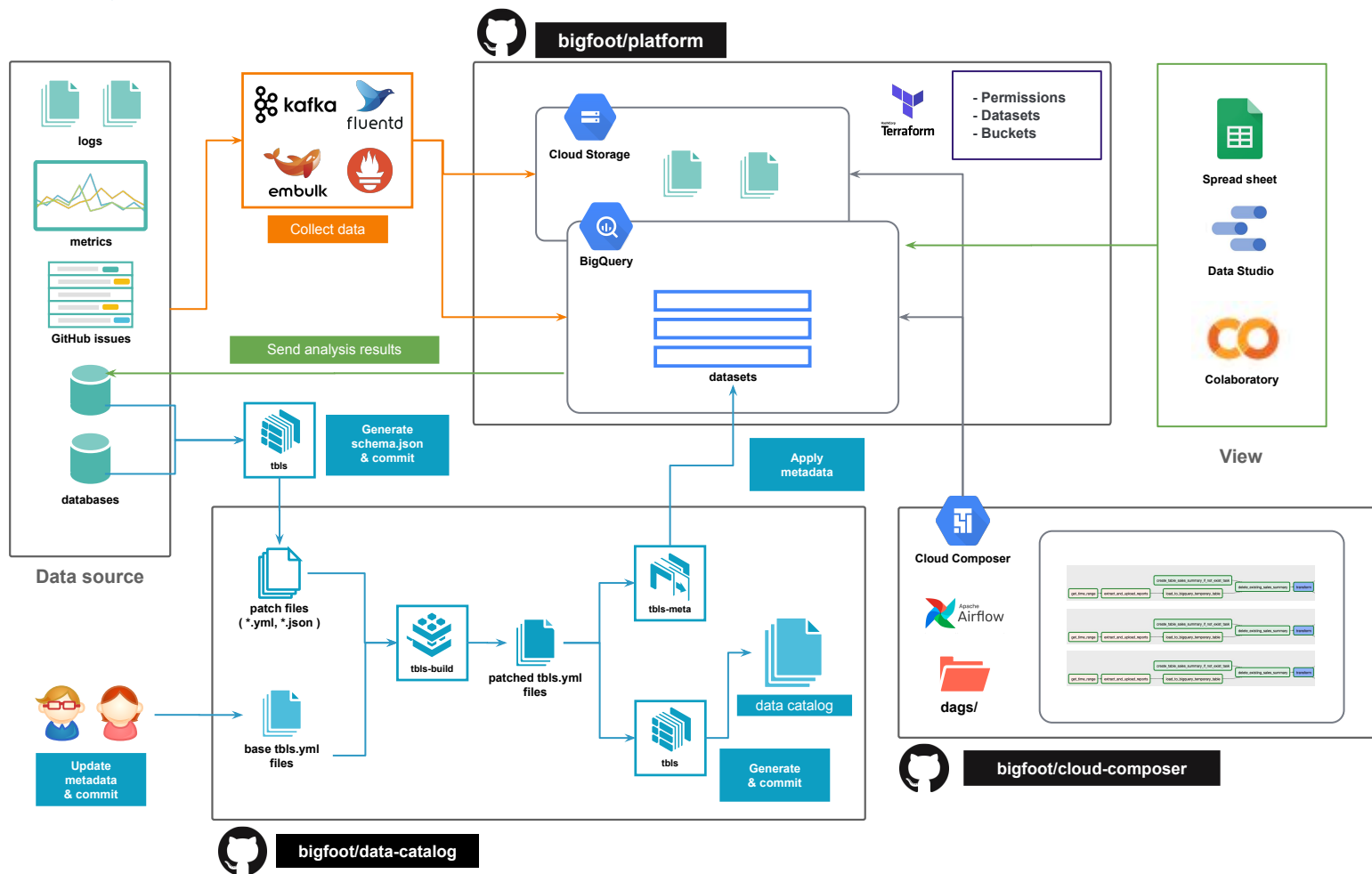


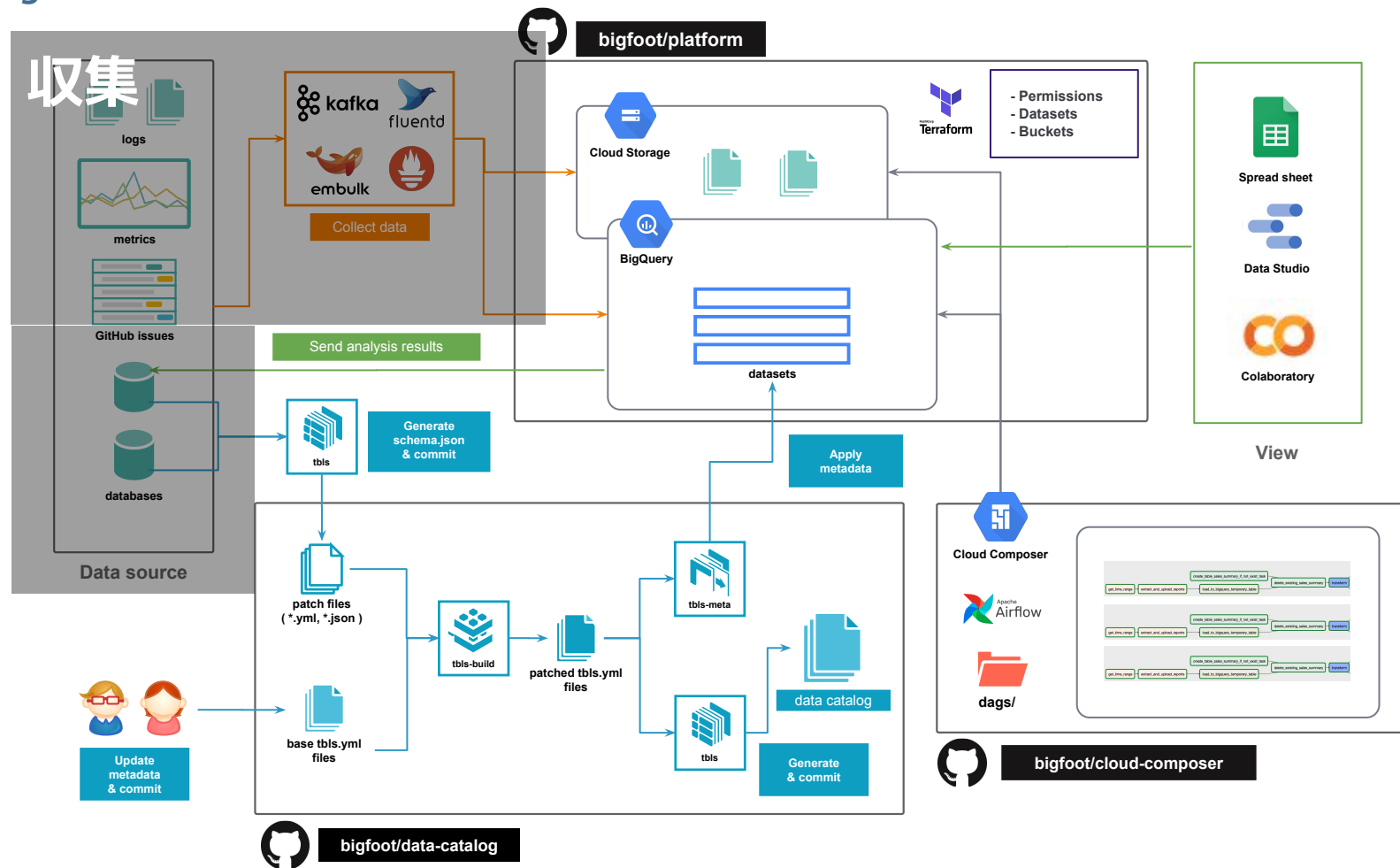
Section 2

Bigfoot 最速入門

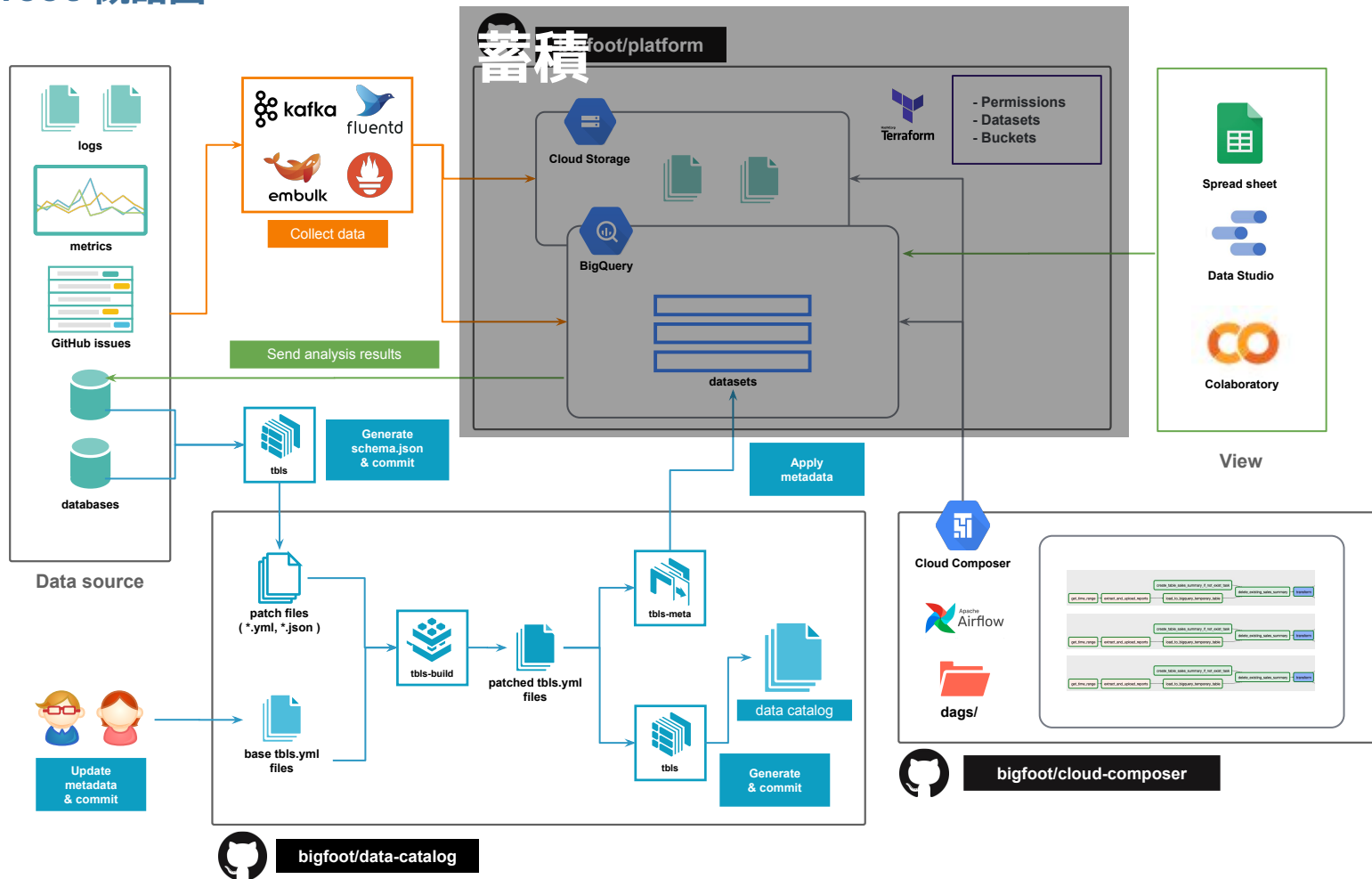
ペパボのデータ利活用基盤をご紹介します

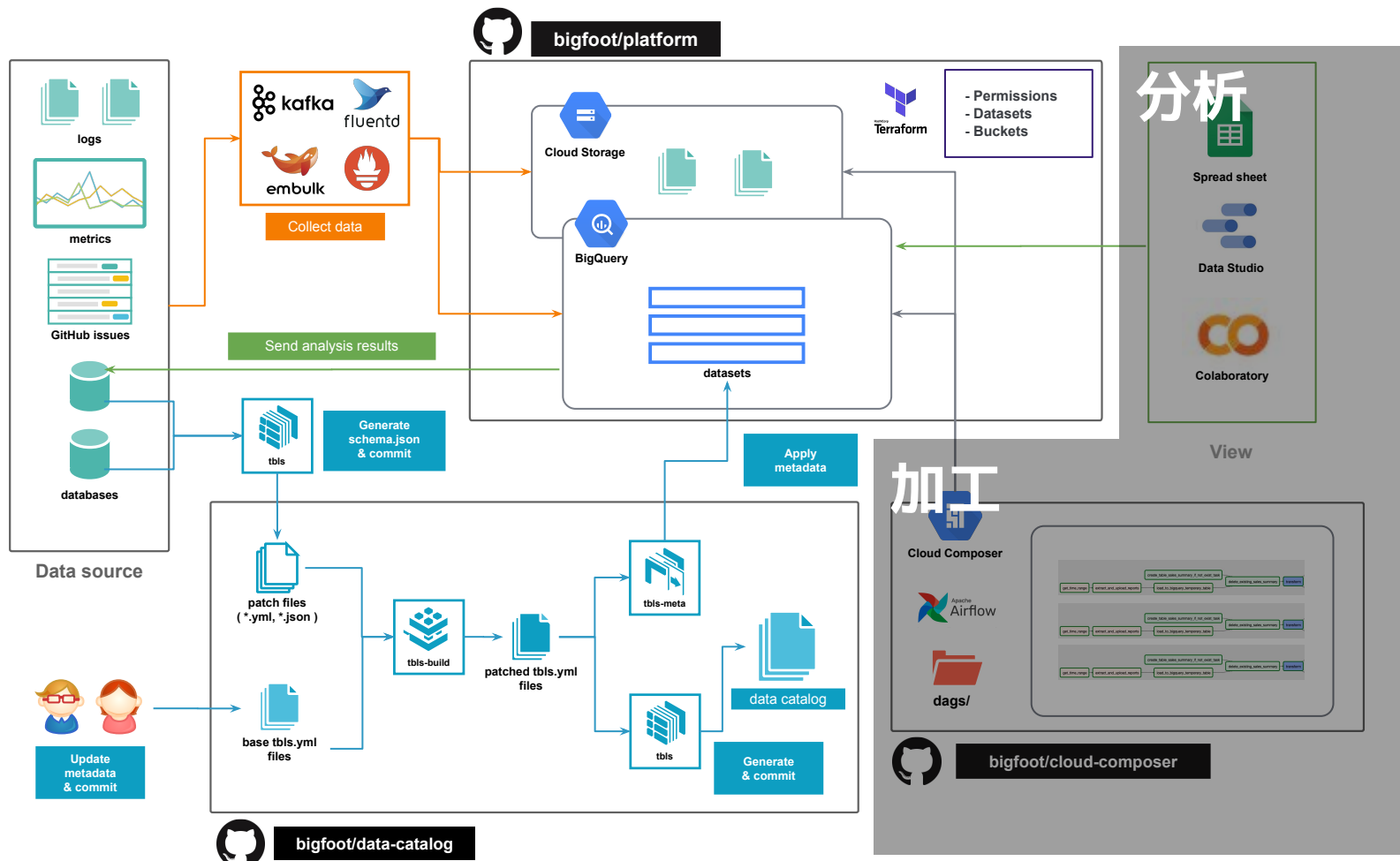
Bigfoot 概略図



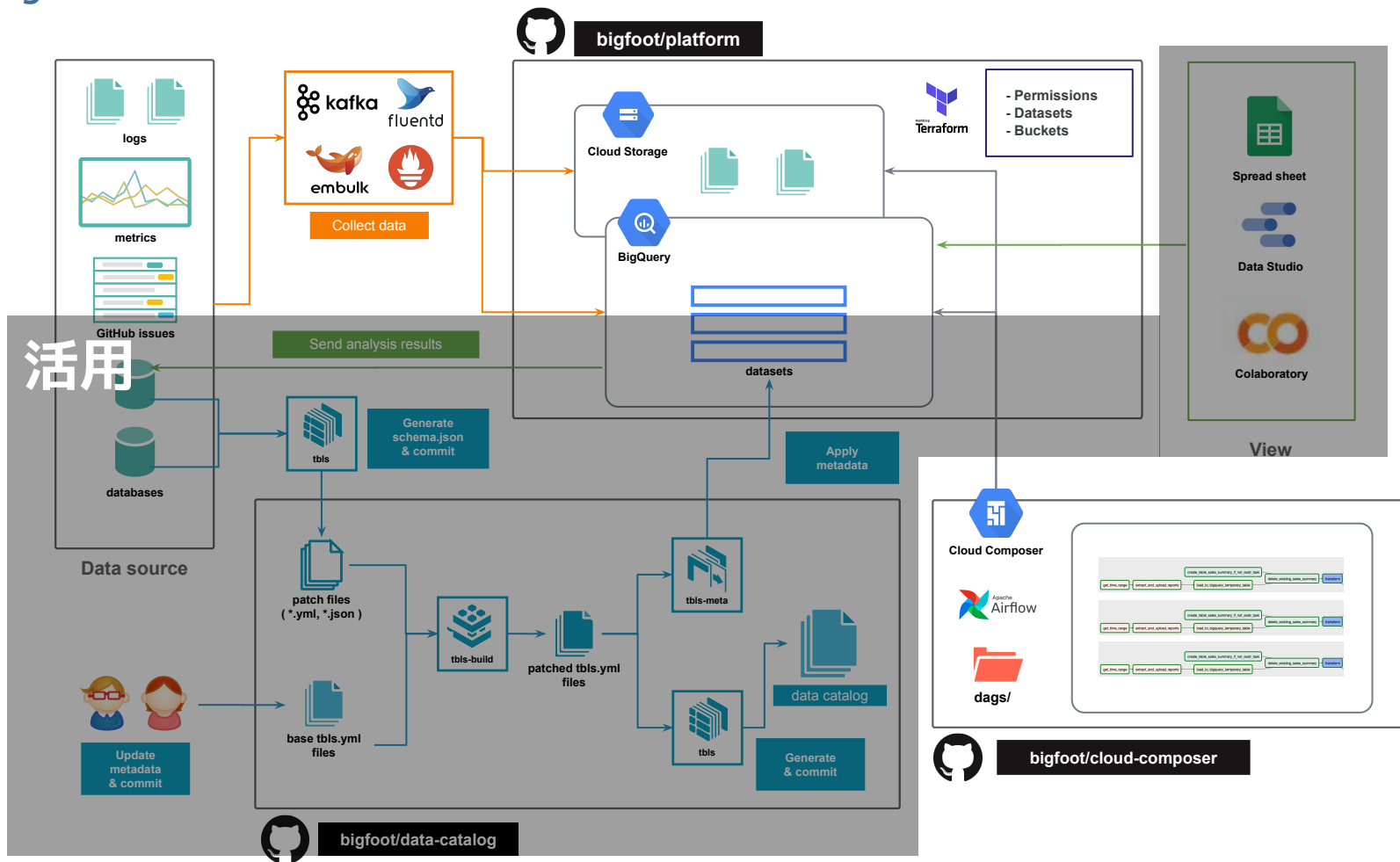


Bigfoot 概略図

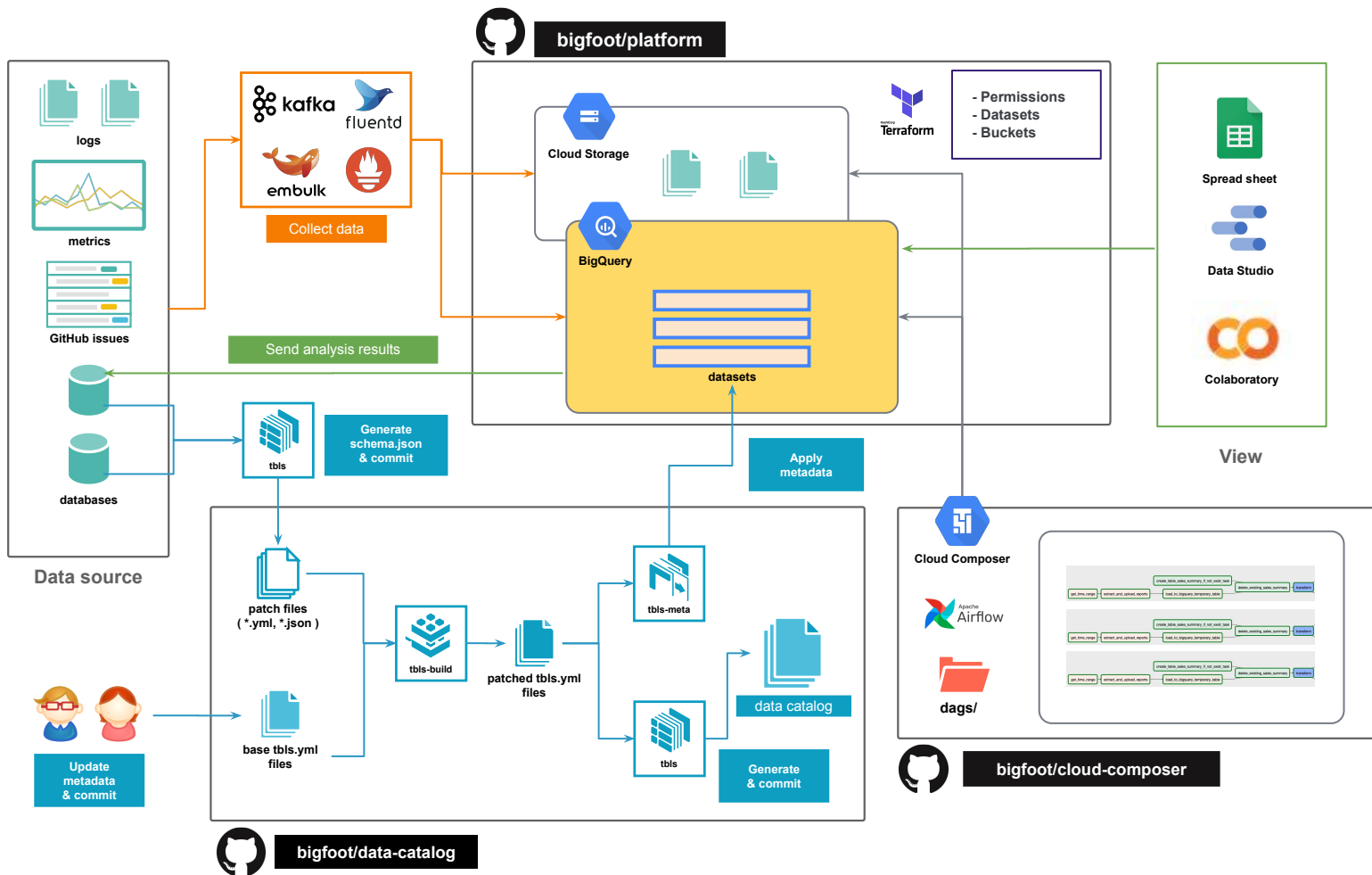




Bigfoot 概略図



Bigfoot 概略図



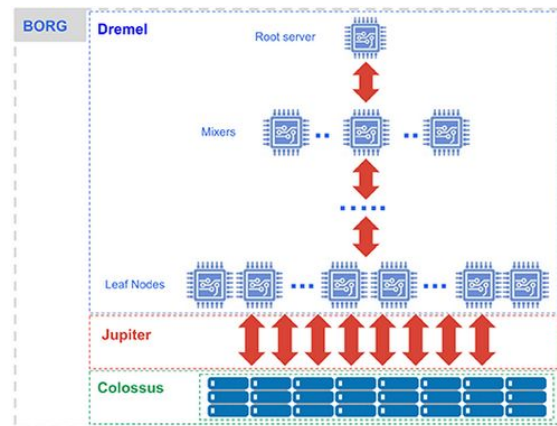


Section 3

Google BigQuery 入門

大規模データ分析対応のエンタープライズ向け フルマネージドデータウェアハウス

- コンピュート、ストレージ、メモリがそれぞれ分離している構造
- それぞれ自動でスケールアウトしてくれるので、
 - データの保存量は事実上無制限
 - クエリ実行速度が超速(特に集計処理)



<https://cloud.google.com/blog/products/data-analytics/new-blog-series-bigquery-explained-overview>

- [Google Data Studio](#)
- [Google Sheets](#)
- [BigQuery Web Console](#)
- [Google Apps Script](#)
- [Google Colaboratory](#)
- [bq command-line tool](#)
- [API Client Libraries](#)
 - C#, Go, Java, Node.js, PHP,
 - Python, Ruby
- [REST API](#)



```
require "google/cloud/bigquery"

def query
  bigquery = Google::Cloud::Bigquery.new
  sql = "SELECT name FROM `bigquery-public-data.usa_names.usa_1910_2013` " \
        "WHERE state = 'TX' " \
        "LIMIT 100"

  # Location must match that of the dataset(s) referenced in the query.
  results = bigquery.query sql do |config|
    config.location = "US"
  end

  results.each do |row|
    puts row.inspect
  end
end
```

- ストレージ料金
 - アクティブ ストレージ: \$0.020 per GB
 - 長期保存: \$0.010 per GB
 - 90日間連続で使用していない場合は、自動的に“長期保存”と判定される
 - (やさしさ)
- スキャン料金
 - 通常は 1TB スキャンで \$5
 - 2021年6月より、Bigfoot はクエリ実行定額プランに移行しました
 - いまならなんと！クエリ投げ放題！お得！！1

RDB: 行指向

1. レコードに対する細かい操作が得意
 - トランザクション処理
 - インデックスによる行の特定

2. 列方向の集計処理効率はいくつか悪い
 - もちろん集計はできますが、リソース効率的に向いていない
 - 行データは、一定の大きさのブロックとしてファイルストレージに保存される
 - 抽出対象の行が格納されているブロックをすべて参照して、目的の行を探す
 - たとえ1列だけ欲しいとしても、構造的に全列を走査する

DWH: 列指向(が多い)

1. 特定列に対する集計処理が得意
 - 例: 1億行・10列のテーブルから、列Aの平均を算出
 - 列ごとにデータを保存しているので、列Aのみ走査する
2. Primary Key, Foreign Key の概念がない
 - 行の一意性はテーブル設計者が担保
 - テーブル結合はできる

DWH: 列指向(が多い)

3. “あえて” 非正規化して保存することが多い

- そもそも、なぜ RDB では正規化するのか？
 - Insert, Update, Deleteで発生する更新時異状(Update anomaly)を防ぐため
- 可視化用/分析用に加工したデータを、行レベルで更新することは基本的でない
 - データ加工については、Day2でお話しします
- ひとつのテーブルとして保存しているほうが、使う人にとってわかりやすい
- そしてテーブル結合が不要なので、(一般的に)クエリ速度は速くなる

ではそろそろ実物を...

<https://console.cloud.google.com/bigquery>

エクスプローラ

プロジェクト
└ データセット
└ テーブル
ビュー

Google Cloud Platform

エクスプローラ

クエリエディタ

クエリ結果

クエリ履歴

行	ID
1	1
2	2
3	3

タブ領域

- クエリエディタ
- クエリスキャン見積
- データセット情報参照
- テーブル情報参照

クエリ結果

- クエリ実行結果の表示
- クエリ実行結果の保存
- ジョブ実行情報

履歴関連

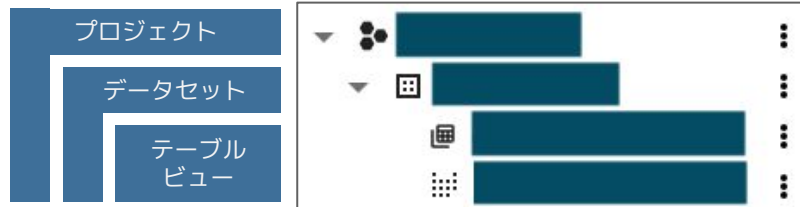
- ジョブ実行履歴
- クエリ実行履歴
- 過去に保存したクエリ

プロジェクト

- GCPのリソースやコストの管理をするためのグループ

データセット

- テーブルやビューをグループ化する概念
- Bigfootでは、サービス単位で作成している
- データセットでアクセス権限を管理



列を指定するとスキャン量が減る → 速

- 全ての列を指定すると
表全体をスキャン
- 特定の列を指定すると
その列のみをスキャン

receipt_join		クエリ
スキーマ	詳細	プレビュー
テーブル情報		
テーブル ID	[REDACTED] training.receipt_join	
表のサイズ	41.88 MB	
長期保存のサイズ	0 B	
行数	104,681	
作成	2021. 7月 9., 18:28:45 UTC+9	
最終更新	2021. 7月 13., 15:19:27 UTC+9	
テーブルの有効期限	常にオフ	
データのロケーション	US	
説明		

実行 展開 ⋮

このクエリを実行すると、41.9 MiB が処理されます。

```
1 SELECT
2   *
3 FROM
4   training.receipt_join
5
```

実行 展開 ⋮

このクエリを実行すると、817.8 KiB が処理されます。

```
1 SELECT
2   SUM(amount) AS total_amount
3 FROM
4   training.receipt_join
5
```

以下のクエリを実行してみてください(画面右上のスキャン量に注目)

1. `SELECT * FROM training.receipt_join`
2. `SELECT receipt_no FROM training.receipt_join`
3. `SELECT receipt_no FROM training.receipt_join LIMIT 100`

- 公式の説明:

“ARRAY 型ではないゼロ以上の要素の順序付きリスト。”

トートロジー...

- 配列をそのままフィールドへ挿入できる
- 配列内の値は同一の型でなければならない

```
SELECT 'a' AS char, [1, 2, 3] AS nums
```

ジョブ情報 結果 JSON 実行の詳細

行	char	nums
1	a	1
		2
		3

クエリ結果は1行

```
WITH base AS (
  SELECT 'a' AS char, [1, 2, 3] AS nums
)
```

```
SELECT char, u_nums FROM base, UNNEST(nums) AS u_nums
```

ジョブ情報 結果 JSON 実行の詳細

行	char	u_nums
1	a	1
2	a	2
3	a	3

UNNEST関数で配列を展開してCROSS JOIN
クエリ結果は3行

- 公式の説明:

“ 順序付きフィールドのコンテナ。

各フィールドはデータ型（必須）とフィールド名（オプション）を持ちます。”

- 子テーブルが列の中に存在するイメージ

行	_sdc_sequence	unique_actions._7d_view	unique_actions._7d_click
1	1614217069695550136	<i>null</i>	1.0
		<i>null</i>	1.0
		<i>null</i>	1.0
		1.0	<i>null</i>
		1.0	<i>null</i>
		1.0	<i>null</i>
		3.0	<i>null</i>
		1.0	<i>null</i>
		1.0	<i>null</i>
		1.0	<i>null</i>
		3.0	<i>null</i>
		4.0	<i>null</i>

- 1:N の関係にあるテーブル同士が、
他の列をキーとしてあらかじめ
結合された状態

この例では、
unique_actions が STRUCT型

unique_actions のフィールドは2つ

- unique_actions._7d_view (ARRAY型)
- unique_actions._7d_click (ARRAY型)

training.receipt_join テーブル (receipt テーブルをベースに非正規化したもの) を使って計算してみましょう

列情報は テーブルの "スキーマ" タブを参照ください

1. `customer_id = 'CS029512000063'` の方がこれまで購入した金額は合計でいくらでしょうか？
2. `category_small_name = 'その他駄菓子'` の商品を販売した実績のある店舗はいくつあるでしょうか？

training.receipt_join テーブル (receipt テーブルをベースに非正規化したもの) を使って計算してみましょう

列情報は テーブルの "スキーマ" タブを参照ください

1. `customer_id = 'CS029512000063'` の方がこれまで購入した金額は合計でいくらでしょうか？

行	total_amount
1	776

```
SELECT
    SUM(amount) AS total_amount
FROM
    training.receipt_join
WHERE
    customer.customer_id = 'CS029512000063'
```

training.receipt_join テーブル を使って計算してみましょう

列情報は テーブルの “スキーマ” タブを参照ください

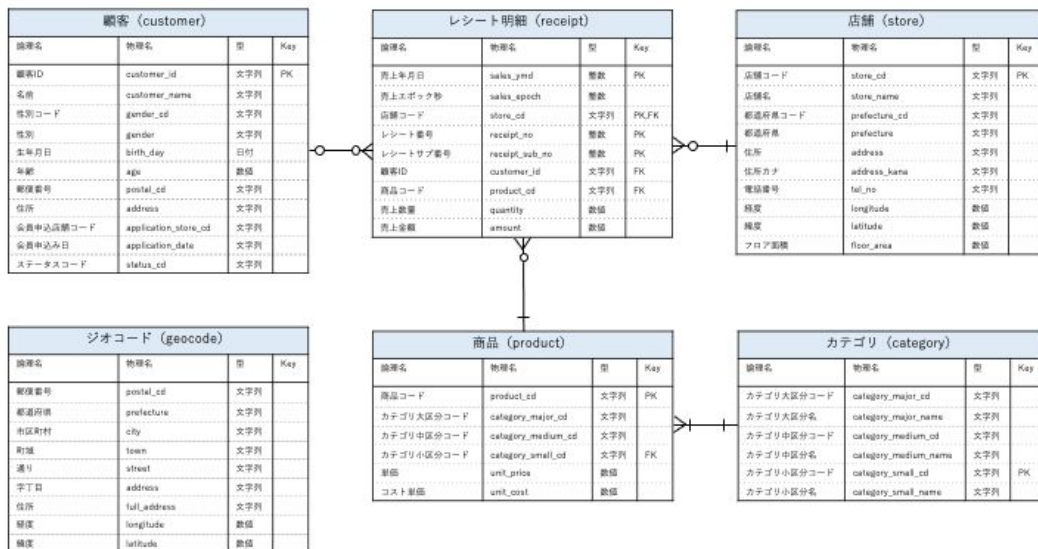
2. `category_small_name = 'その他駄菓子'` の商品を販売した実績のある店舗は
いくつあるでしょうか？

行	unique_store_count
1	18

```
SELECT
    COUNT(DISTINCT store.store_cd) AS unique_store_count
FROM
    training.receipt_join
WHERE
    product.category.category_small_name = 'その他駄菓子'
```

training.receipt_join テーブルと同じデータを取得できる SELECT 文を書いてみてください。

training.receipt_join テーブルは、下のER図のリレーションをもとに結合しています。



https://github.com/The-Japan-DataScientist-Society/100knocks-preprocess/blob/master/docker/doc/100knocks_ER.pdf

training.receipt_join テーブル と同じデータを取得できる SELECT 文を書いてみてください。

```
WITH product_with_category AS (  
  SELECT  
    product.* EXCEPT(category_major_cd, category_medium_cd, category_small_cd),  
    category,  
  FROM  
    training.product AS product  
    LEFT OUTER JOIN training.category AS category  
      USING (category_major_cd, category_medium_cd, category_small_cd)  
)  
  
SELECT  
  receipt.* EXCEPT (customer_id, product_cd, store_cd),  
  customer,  
  store,  
  product  
FROM  
  training.receipt AS receipt  
  LEFT OUTER JOIN training.customer AS customer USING (customer_id)  
  LEFT OUTER JOIN training.store AS store USING (store_cd)  
  LEFT OUTER JOIN product_with_category AS product USING (product_cd)
```



Section 4

Google DataStudio 入門

- Google マーケティングプラットフォーム で提供されているBIサービス
- Google アナリティクスも、Google マーケティングプラットフォームのサービスのひとつ
- 日本では「データポータル」と呼ばれていますが、ここでの表記はData Studio に統一します。



どうやってアクセスするの？

こちらへ

<https://datastudio.google.com/>



- 自分が作成したものや共有されているものが表示される
- Google Drive に似ているが、別管理になっている
- 個人の所有物扱いなので、他の人へ見せるには共有の設定をする



新規作成で表示される編集画面

共有

閲覧画面を表示

ページ

データ

コントロール

他のユーザーと共有

 利史 堤 として共有

[ユーザーを追加](#)

[アクセスを管理する](#)

名前またはメールアドレスを入力...

閲覧者

☒ ユーザー通知の送信元メール

キャンセル

送信

- ページ
 - グラフや画像、コントロールを載せる場所。
 - スプレッドシートでいうシートにあたる。
- レポート
 - ページをまとめたもの。ダッシュボード=レポート。
 - スプレッドシートでいうファイル(ブック)にあたる。
- データ
 - グラフが参照するデータソース。
 - グラフにつき一つだけ指定できる。
- コントロール
 - プルダウンやチェックボックス、期間指定、などフィルタリングを行うためにレポートに設置する部品。

実際につくってみましょう

[公式のヘルプページ](#) の解説がわかりやすいです。

Data Studio のホーム画面に表示されているチュートリアルもおすすめです。



1. 全店舗の売上を集計した総売上推移を、日付ごとにわかるようにしたいです。
期間指定とか、月次・年次に切り替えられたりできるといいなあ
2. 1のグラフとは別に、もうひとつお願いします！！
地域差があるかを知りたいので、売上の累計を都道府県別・店舗別にみたいです。
3. 全店舗行ったことないのですが、女性のお客が多い気がする...
2のグラフの内訳で、性別がみれるようにしたいです。

1. 全店舗の売上(receipt.amount)を集計した総売上推移を、日付(receipt.sales_ymd)ごとにわかるようにしたいです。
期間指定とか、月次・年次に切り替えられたりできるといいなあ
2. 1のグラフとは別に、もうひとつお願いします！！
地域差があるかを知りたいので、売上(receipt.amount)の累計を都道府県(store.prefecture)別・店舗(store.store_name)別にみたいです。
3. 全店舗に行ったことはないのですが、女性のお客が多い気がする...
2のグラフの内訳で、性別(customer.gender)がみれるようにしたいです。

つくったグラフを
みんなで見てみましょう



閲覧権限

編集権限

ダッシュボードを Notion に埋め込む

The image shows a Notion dashboard titled "エンジニア研修2021_機械学習_回答例" with a menu open. The menu includes options like "共有...", "レポート設定", "変更履歴", "公開設定", "新しいレポート", "コピーを作成...", "ダウンロード形式", and "レポートを埋め込む". An orange box highlights the "レポートを埋め込む" option. An arrow points from this option to a "レポートの埋め込み" (Embed Report) dialog box. This dialog has a checked box for "埋め込みを有効にする" (Enable embedding) and radio buttons for "コードを埋め込む" (Embed code) and "URL を埋め込む" (Embed URL). Below these is a text input field with the placeholder "ここに URL が表示されます". An orange box highlights the "クリップボードにコピー" (Copy to clipboard) button. An arrow points from this button to the text "Notion に貼り付け" (Paste into Notion). Below this text is a Notion modal menu with options "Dismiss", "Create bookmark", and "Create embed", where "Create embed" is highlighted with an orange box.

エンジニア研修2021_機械学習_回答例

ファイル 編集 表示 挿入 ページ 配置 リソース

3/3 ページ

共有...

レポート設定

変更履歴

公開設定

新しいレポート

コピーを作成...

ダウンロード形式

レポートを埋め込む

データを追加

店舗別売上

レポートの埋め込み

☒ 埋め込みを有効にする

☐ コードを埋め込む ☒ URL を埋め込む

次の情報をサイトに貼り付けます。

ここに URL が表示されます

クリップボードにコピー

完了

Notion に貼り付け

Dismiss

Create bookmark

Create embed