{::options parse\_block\_html="true" /}。

### このページについて

{:.no\_toc .hidden-md .hidden-lg}<sub>o</sub>

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}.

# GitLabデータカタログインデックス

データカタログのページでは、アナリティクスソリューション、ダッシュボード、ワークフロー、キーワードなどのインデックスを掲載しています。リンクやリソースを追加したい場合は、ご自由にご投稿ください。

# Sisense ダッシュボード

Sisenseは、当社のエンタープライズスタンダードなBIツールです。

重要: 2021-09-08 SAFEダッシュボードの変更点

**2021-09-08**に、**GitLab**の成熟と**SAFE**データフレームワークの推進により、**Sisense**ダッシュボードのアクセスが変更されます。次のセクションでは、予定されているアクセスプロセスについて説明します。

Sisenseダッシュボードへのアクセスは、職務上の役割に基づいており、SAFEデータアクセスフレームワークによって管理されています。Sisenseでは、ダッシュボードは以下のスペースに分類されます。

- **SAFEダッシュボード**スペースには、**GitLab**の**SAFE基準**を満たすすべてのダッシュボードが置かれています。**SAFE**ダッシュボードは、役割や必要性に応じて**GitLab**チームメンバーが利用できます。
- **GitLab**スペースは、**SAFE**の取り扱いを必要としないすべてのダッシュボードを収容する一般アクセスエリアです。一般 アクセス用ダッシュボードは、**GitLab**チームメンバー全員がアクセスできます。

SAFEダッシュボードを含むダッシュボードの全リストは、GitLabDashboardIndexに掲載されています。

#### GitLab General Access Dashboardへのアクセス

一般アクセスダッシュボードは、すべてのGitLabチームメンバーが利用でき、アクセスリクエストの発行は必要ありません。アクセス方法は、Sisenseの概要と開発ページに記載されています。

#### SAFEダッシュボードへのアクセス

すべてのSAFEダッシュボードは、Sisense SAFE Dashboardスペース内に保存され、アクセスはSisenseスペースレベルで許可されます。1つのSAFEダッシュボードにアクセスすると、すべてのSAFEダッシュボードにアクセスできます。SAFEダッシュボード(およびスペース)へのアクセスには、以下が必要です。

- 1. 直属の上司の承認
- 2. 部署のVP (またはそれに相当するもの) の承認
- 3. GitLab Dashboard Indexで定義されたSAFE Dashboard Space Ownerの承認 SAFEダ

ッシュボードへのアクセス権の取得

- 1. アクセスリクエストを作成し、GitLabDashboardIndexからアクセスを必要とするSAFEダッシュボードを3つまでリストアップする。このSAFEダッシュボードのリストは、承認者があなたのニーズと意図を理解するのに役立ちます。
- 2. 直属の上司、所属部署の副社長(またはそれに相当する人物)、およびGitLabダッシュボードのインデックスヘッダーに定義されているSAFEスペースオーナーに承認を求める。SnowflakeでSAFEデータへのアクセスが60日以内に承認されていれば、承認は必要ない。このステップをスキップして、特定のARにリンクする。
- 3. リクエストが承認されたら、@gitlab-data/analystsをタグ付けすると、データチームがリクエストを処理します。
- 4. 処理が完了すると、Sisenseにログインして、ご希望のSAFEダッシュボードやSAFEダッシュボードスペース内の他のすべてのダッシュボードにアクセスできるようになります。

Sisense Spaces

SnowflakeでのSAFEデータへのアクセス方法については、「SnowflakeでのSAFEデータへのアクセス」をご覧ください。

SAFEデータのGoogle Sheetsファイルは、「SAFEダッシュボードへのアクセス」の手順でアクセスできます。Google SheetsのSAFEデータにアクセスするには、この手順に従ってください。

#### 2021-09-08終了 SAFEダッシュボード変更のお知らせ

# 分野別データ

マーケティング

#### ダッシュボード

- TD:マーケティングデータマート
- TD: SDRパフォーマンスダッシュボード

#### 分析プロジェクト

◆ 2021-10-01 SAO分析

売上 高

#### ダッシュボード

- TD:カスタマーセグメンテー
- ◆ ション TD: セールスファネル
- 手動使用データアップロードプロセス

#### 分析プロジェクト

◆ 買い方の傾向

#### ファイナンス

**◆ TD**:ファイナンス**ARR** 

# 製品

#### ダッシュボード

- TD: プロダクト・ジオロ
- ケーション TD: プライ シング・アナリシス

#### データモデルとプロセス

- ◆ 製品の使用状況デー
- タ SaaSサービスの
- Ping SaaS製品のイベント

#### 分析プロジェクト

• 2020-12 プロダクトアナリティクス オフサイト

#### 成長

#### ダッシュボード

- XMAU分析成長ダッ
- シュボード

分析プロジェクト

- 2021-08CustomerCentricProductInsights
- 2021-08ExperimentationWorkshop
- ◆ 2021-08Stage:Secureの採用と転換の分析 2021-08名前空
- ◆ 間の転換の分析
- FY22-Q1 Growth Team KPI Review
- 2021-08SSOログインディープダイブ
- 分析 Growth Insights
- ◆ 締めくくりの成長実験

#### カスタマーサクセス

◆ 建設中

#### エンジニアリング

• MRレート

#### 人々

- People Metrics Overview
- PTO By Roots (Slack)
- People Key Metrics
- People KPI Deck
- PromotionsReport
- タレント・アクイジション・メトリクス
- ◆ People Metrics Sisense Dashboardでのデータ発見

#### データチーム

- ◆ Sisenseの利用と導入 信頼でき
- るデータヘルス

# 指標と用語の索引

- ◆ セールス用語集
- ARR : Annual Recurring Revenue
- ATR : Available To Renewal
- CAC: Customer Aquisition Cost
- LTV: Customer LifeTime Value
- Namespace
- NDR: Net Dollar Retention
- PQL:Product Qualified Lead 製
- ◆ 品ステージ
- SM: Self-Managed aka Self-Hosted
- ◆ UPA:究極の親アカウント xMAU: x
- ◆ 月間アクティブユーザー数

#### 伝説

□は、そのソリューションが運用され、ハンドブックに組み込まれていることを示します。

**||** は、そのソリューションがWork In Progressであり、積極的に開発されていることを示します。この指標を使用する場合、課題はこのページからもリンクされている必要があります。

◆は、そのソリューションが近いうちに実用化される可能性が低いことを示しています。

layout: handbook-page-toc title:"カスタマーセグメンテーション"説明"顧客セグメンテーションとは、お客様が誰であるかを理解するために、共通の特徴に基づいてお客様をグループに分けるプロセスです。"

# このページについて

{:.no\_toc}

TOC {:toc}

# カスタマー・セグメンテーション分析

顧客セグメンテーションとは、お客様を理解し、素晴らしいカスタマー・エクスペリエンスを提供するために、共通の特徴に基づいてお客様をグループに分けるプロセスです。お客様を特定するためには、業界、製品カテゴリー、販売セグメント、配送、地域など、さまざまな特徴があります。Customer Segmentation Analysisページでは、GitLabチームメンバーが顧客データを調査し、顧客の洞察力を高めるために使用できる情報とツールを提供します。

このデータソリューションは、3つのSelf-Service Data機能を提供します。

- 1. ダッシュボードのユーザー業界、製品カテゴリー、販売セグメント、配送、アカウント・オーナー・チーム、テリトリーごとにARRと顧客数を視覚化するSisenseのダッシュボード
- 2. ダッシュボード開発者。新しいダッシュボードを構築したり、既存のダッシュボードを顧客セグメンテーション・データにリンクさせるための完全な次元モデル・コンポーネントを含む新しいSisenseデータ・モデルです。
- 3. SQL開発者。An Enterprise Dimensional Modelの対象領域 デー

タプラットフォームの観点から、このソリューションは提供され

ます。

- 1. 顧客セグメンテーション分析のためのエンタープライズ次元モデルの拡張
- 2. Data Pipeline Healthダッシュボードのテストとデータ検証の拡張機能
- 3. ERD、dbtモデル、および関連するプラットフォームコンポーネント

クイックリンク

ガスタマーセグメンテーション ダッシュボード セグメンテーション

Dashboard Developer Certification - Customer Segmentation

SQL認定資格 - カスタマーセグメンテーション

Sisense Discoveryを使い始める

セルフサービス・ウォークスルー・ビデオ

#### データセキュリティの分類

Customer Segmentation Dashboardに含まれるデータの多くは、オレンジまたはイエローです。これには、オレンジ色のアカウントの顧客メタデータ、SalesforceやZuoraのコンタクトデータ、GitLabの非公開財務情報などが含まれますが、これらはすべて公開すべきではありません。このダッシュボードからデータを共有する際には、その詳細が GitLab の組織内に留まるようにし、外部に共有する際には適切な承認を得られるように注意しなければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保存しないよう、常に注意を払う必要があります。このデータはSnowflakeとSisenseに残すべきであり、特に承認されない限り、これらのアプリケーションを通じてのみ共有されることが理想的です。

#### ORANGE

- 説明行またはレコードレベルの顧客および個人データ。オブジェ
- クトです。
  - ⋄ dim\_billing\_accounts

- dim\_crm\_accounts
- dim crm persons

#### **YELLOW**

- 説明します。GitLab 財務データ。集約や合計を含む。オブジェクトです。
- dim subscriptions
  - fct charges
  - ♦ fct\_invoice\_items
  - fct mrr

#### ソリューション・オーナーシップ

- ソースシステムオーナー:
  - ◇ Salesforce@jbrennan1
  - Zuora:アンドリュー・マー
- ソースシステムのサブジェクト・マタ
  - o ー・エキスパート: Salesforce
  - @jbrennan1 Zuora:アンドリュー・マーレイ ( @andrew\_murray
- データチームのサブジェクト・マター・エキスパートPaul\_armstrong @jeanpeguero @jjstark @iweeks

#### 主な用語

- 1. 製品カテゴリー、製品階層、配送
- 2. 販売セグメント
- 3. アカウントオーナーチーム
- 4. テリトリー
- 5. お客様
- 6. 産業

#### キーメトリクス、KPI、PI

- 1. ARR
- 2. お客様の数

#### セルフサービス・データ・ソリューション

セルフサービス・ダッシュボード開発者

Sisenseでのチャート構築を始めるには、Sisenseの10分間のDataOnboardingVideoを見るのが良いでしょう。ダッシュボードを作成した後は、そのダッシュボードを簡単に見つけられるようにしたいものです。トピックは、ダッシュボードを1つの場所に整理し、簡単に見つけることができる素晴らしい方法です。ダッシュボードの右上にあるトピックへの追加アイコンをクリックすると、トピックを追加できます。ダッシュボードは、関連性のある複数のトピックに追加することができます。トピックには、Finance(財務)、Marketing(マーケティング)、Sales(販売)、Product(製品)、Engineering(エンジニアリング)、Growth(成長)などがあります。

#### セルフサービスのSQL開発者

#### キーフィールドとビジネスロジック

- ◆ データは、ZuoraとSalesforceから取得しています。
- 親顧客は1つ以上の製品を持つことができるため、製品カテゴリと配送ディメンションで複数回カウントすることができます。総顧客数のユニークなカウントを得るには、製品と配送をARRAYに集約して顧客のCOUNT DISTINCTを行うか、製品カテゴリや配送ディメンションを含まない顧客のCOUNT DISTINCTを行う必要があります。
- ◆ 料金が定期的に発生しているとみなされるには、データ上で有効な終了月が有効な開始月よりも大きい必要があります
- 。Zuoraでは、課金のeffective end dateおよびeffective end monthは、それぞれ更新の初日または月となります。

• 月間ARR計算では、有効終了月は解約が発生する時期を示しており、有効終了月はARR計算にカウントしません。例えば、有 効開始月=2020-07-01、有効終了月=2021-07-01の契約では、2020-07-01から2021-06-01までのARRを合計し、12ヶ月分のARR となります。

#### エンティティ・リレーションシップ・ダイアグラム

# Diagram/EntityGrainPurposeKeywords

ARR and Customer Count Analytics ERD	月間、サブ スクリプショ ン、製品カ テゴリー	ARRと顧客数について、様々 な角度からのインサイトの 提供	親顧客、製品カテゴリー、配送、業界、アカウントオーナー チーム、テリトリー、販売セグメント
Lead to Cash	以下の	リード・トゥ・キャッシュの全	親顧客、製品カテゴリー、配送、業界、アカウントオーナー
Overview ERD	すべて	プロセスの概要	チーム、テリトリー、販売セグメント、CRM、担当者、アカウント

#### リファレンスSQL

#### スニペットライブラリ説明

カスタマーセグメンテーション SQLスクリプト	ARRと顧客数を製品カテゴリー、配送、業界、アカウントオーナーチーム、地域、販売セグメント別にスライスするクエリ
顧客セグメンテーション TY四半 期とLY四半期の比較 SQLスクリ プト	TY対LYのARRと顧客数を四半期ごとに取り出し、製品カテゴリー、配送、業界、アカウントオーナーチーム、テリトリー、セールスセグメントごとに分類するクエリ。

# データプラットフォー ムソリューション

# データの系統

- データはSalesforce.comとZuoraから取得し、手動で管理されたZuoraの除外アカウントリストからアカウントを除外して
- ◆ います。データの系統図はdbtmart\_arr系統図に記載されています。

#### DBTソリューション

dbtソリューションは、RAWソースデータから次元モデルを生成します。ただし、以下のフィールドは例外で、特定のdbtモデルに 実装されたビジネスロジックに基づいて計算されます。

### フィールドビジネスロジック

product_	categoryZuoraproduct_rate_plan_nameに基づいて計算されています。	
delivery	製品カテゴリーに基づいて計算されます。	
サービスタイプ	製品レートに基づいて計算されます。	
ultimate_parent_account_	segmentSFDCのmultimate_parent_sales_segmentに基づいて、UnknownとNULLをグループ化して算出	
したもの。	セグメントをSMBに	

# トラステッドデータソリューション

#### TrustedDataFrameworkの概要を見る

dbtガイドの例では、さらなるテストを実施するための詳細と例を示しています。

#### ズーラ

◆ 信頼できるデータダッシュボード

。 tdfタグを含むすべてのZuoraデータテストを報告します。

・ データパイプラインの健全性を示すダッシュボード

• ZuoraアカウントとSubscriptionデータのデータの問題点を明らかにするための行数の報告。

EDM エンタープライズ ディメンション モデル バリデーション

- (WIP) エンタープライズ ディメンショナル モデル バリデーション ダッシュボード
  - 最新のEnterprise Dimensionalモデルのテストと実行に関するレポート

RAWソースデータパイプラインの検証

データパイプライン・ヘルスバリデーション

# layout: handbook-page-toc title: "WIP Email Marketing Data Mart"

このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}.

# メールマーケティングデータマート

昨年、私たちのチームは、データ要求とライフサイクル・マーケティング戦術を拡張するためのソリューションの必要性に気づきました。

メールマーケティングデータマート(メールマーケティングデータベースとも呼ばれる)は、複数の異なるデータソースを統合し、よりインテリジェントでスケーラブルな方法でユーザー、顧客、見込み客にアプローチできるようにします。

企業のアプリケーションチームやマーケティングチームと協力して、GitLab社が顧客やユーザーの重要なアップデートをより 迅速に伝え、見込み客とユーザーの重なりを理解し、より効率的なマーケティングを可能にするソリューションを作り上げま した。

#### このページの目的は

- TD:EmailMarketingDataMartREADME」を使ってメールキャンペーンを実施する際の参
- 考になります。TD: Email Marketing Data Mart」の作成に使用されたデータモデルの理解を助ける。
- ◆ Coming Soon GitLabでのあなたの役割に最も適した認定を受けることで、あなたの理解度を評価します。
- そして、みんなが貢献できるようにしましょう。

クイックリンク

ERD - TD: Email Marketing Database(メールマーケティングデータベー

# はじめに

◆ メトリクスをどのように説明するかを説明するキーワード

• データベースの背後にあるデータソース

# キーターム、キーフィールド、ビジネスロジック

- ▶ 主な用語
  - ◆ 販売セグメント
  - 販売地域 製品階
  - ◆ 層 製品納入 ユー
  - ◆ ザー役割
  - ◆ サブスクリプションの開始日と終了
  - ◆ 日トライアルの開始日と終了日
- ▶ キーメトリクス、KPI、PI
  - ◆ 使用方法 Ping
  - 試験終了までの日数
- ▶ キーフィールドとビジネスロジック

# データソースとデータモデルの理解

データソース

- Zuora
- Salesforce
- GitLab.com 顧
- ◆ 客DB

Entity Relationship Diagram(ERD)の略。

▶ クエリの例

# 追加リソース

▶ Trusted Data Solutionは近日

公開予定です。

- ► EDM Enterprise Dimensional Model Validations
- ► RAW Source Data Pipeline validations
- ▶ データセキュリティの分類

#### **ORANGE**

- 説明行またはレコード・レベルの個人データ。このため、このデータはPROD.COMMON\_MART\_MARKETINGではマスクされているが、PREP.SENSITIVEでは利用可能である。顧客メタデータ(顧客の会社名など)は
  PROD.COMMON MART MARKETINGで利用可能です。
- オブジェクトです。
  - prep.sensitive.dim\_marketing\_contact
  - o prep.sensitive.mart marketing contact
  - prod.common\_mart\_marketing.mart\_crm\_touchpoint
  - prod.common\_mart\_marketing.mart\_crm\_attribution\_touchpoint
- ▶ ソリューション・オーナーシップ

layout: handbook-page-toc title:"マージリクエスト (MR) レート"

#### このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

# マージリクエスト (MR) レート

MR率は、生産性と効率性の指標です。分子は、一連のプロジェクトに対するマージ要求の集まりです。分母は人の集まりです。どちらも時間をかけて(通常は月単位で)追跡される。

MR Rateの全体および全体の詳細については、こちらをご覧ください。

クイックリンク

MR率ダッシュボード

# はじめに

#### DBTドキュメント

MR率ダッシュボードは

rpt\_gitlab\_employees\_merge\_request\_metricsでは、2つの主要なフィルターを考慮しています。

- 1. ブレイクアウト。ダッシュボードを会社、部門、部署レベルで表示することができます。
- 2. Breakout\_division\_departmentの略。このフィルターは子フィルターで、ブレイクアウトの選択に応じて変化します。
- ▶ Breakout = Company --> Breakout\_division\_departmentでは、all\_companyが1つ選択されます。

- Breakout = Division --> Breakout\_division\_departmentはGitLabの全ての部門をリストアップします。
- Breakout = Department--> Breakout division departmentでは、すべての部門と部署の組み合わせが表示されます。

注:Breakout\_division\_departmentフィルターでは、データが事前に集計されているため、1つのオプションしか選択できません。 ダッシュボード自体は以下のように構成されています。

- 1. MRレートの推移
- 2. 少なくとも1つのMRトレンドを持つチームメンバーの%メンバーです。
- 3. 前月のチームメンバーによるMRレートの内訳
- 4. どのプロジェクトが計上されていて、どのプロジェクトが計上されていないかの内訳です。

#### 製品の一部であるプロジェクト

MR率とMR量の計算では、製品全体の取り組みに貢献するプロジェクトからのMRを考慮しています。現在のプロジェク

トのリストは、gitlab-data/analyticsプロジェクトで以下のシステム・データベースについて特定されています。

システムデータベース | ファイル | GitLab.com | projects\_part\_of\_product.csv | ops.gitlab.net |

projects\_part\_of\_product\_ops.csv | プロジェクトリストを更新するには、こちらの手順に従ってください。

# データモデル

基本となるデータモデルは、多くのデータモデルをあらかじめ考慮していますが、これをGitLabチームメンバーレベルの1つの集約されたデータモデルにまとめ、集約レベルのレポートを作成することで、分析しやすくしています。

gitlab\_employees\_merge\_requests\_xf -- GitLabチームメンバーのすべてのマージリクエストを表示し、そのマージリクエストがMR レートの計算に含まれているかどうかをis\_part\_of\_productフィールドを使って表示します。

さらに、gitlab\_bamboohr\_employee\_baseを使って、従業員を月別、組織属性(部門、部署など)別、有効期間別に分けたベースモデルを作成しています。

# layout: handbook-page-toc title: "Estimation Algorithm"

このページについて

 $\{:.no\_toc\}$ 

TOC {:toc}

#### 推定アルゴリズム

SaaS型のGitLab.com製品をホストするセルフマネージドインスタンスでは、Usage Pingsから実世界のデータを受け取ります。セルフマネージドインスタンスでは、お客様はUsage Pingを無効にするオプションを持っており、これらのインスタンスの使用量統計を受け取ることはありません。すべてのセルフマネージドカスタマーの使用量を計算するためには、Usage Ping を無効にしたインスタンスの使用量を予測するアルゴリズムを開発する必要があります。このアルゴリズムを「Estimation Algorithm」と呼んでいます。

このセクションでは、これらのxMAU推定に関する我々の最初の試みについて説明します。

#### 現在の方法論

#### 私たちが知っていることは?

- ◆ 月間アクティブサブスクリプション数
- Usage ping の特定の xMAU カウンタのリリース日(つまり月)。例えば、Create Source Code Groupの
- GMAUであるmerge\_requests\_usersがバージョン12.9でリリースされたことがわかります。

#### わからないことは?

アクティブなフリーインスタンスの数

#### 私たちにできることは?

実際の例を見てみましょう。Dev:CreateステージのSMAUとして使われているカウンタは、

action\_monthly\_active\_users\_project\_repoです。これはバージョン13.3でリリースされました。Estimated SMAUを算出するには、次のようなステップを踏みます。

- ◆ 記録されたGMAUを配信ごとに分割して算
- 出 1ヶ月ごとに算出
  - 10月の場合、5000のサブスクリプションのうち、Usage Ping Payloadを送信した3500のサブスクリプションがあることが わかります。
  - その中で、バージョン13.3以上のサブスクリプションと13.2以下のサブスクリプション、そしてペイロードを送信してこないサブスクリプションを分けています。
  - このSiSenseのチャートは、月ごとの分割を見たものです。
- 上のグラフは、全アクティブサブスクリプションのうち、ペイロードを送信し、13.4以上を使用しているサブスクリプションの割合を示しています。
- そして、Estimated SMAU = Recorded SMAU / % on 13.3 を導き出します。

カウンターはGitLabの異なるバージョンでリリースされているので、推定値はカウンターごとにカスタマイズされます。

この数値は、実際には有料のXMAUを推定するのに非常に正確な数値ですが、最初は推定XMAUを計算するために使用することができます。しかし、この推定値をより強固なものにするためには、いくつかの改善点が考えられます。

- 最初の改善点は、サブスクリプションの数ではなく、注文したシートの数量の合計を使用することです。また、プラ
- ◆ ンごとに見積もりを分割することもできます(現在の問題:1つのサブスクリプションに2つの異なる製品層の2つのインスタンスが含まれる可能性があります)。
- 主な改善点は、Coreの推定に関するものです。現在のところ、バージョンアップや使用量pingのオプトインに関して、有料と無料で同じパターンがあると仮定しています。使用率の推定は非常に複雑ですが、コア・インスタンスのバージョンアップと使用率のデータ傾向については、重要な改善点となるでしょう。

# layout: handbook-page-toc title:"TD : スナップショット 年間経常収益(ARR)"こ

のページについて

TOC {:toc}

{::options parse\_block\_html="true" /}<sub>o</sub>

#### このページの目的は

- Snapshot ARR Dashboardsの操作方法を理解するのに役立ちます。
- Snapshot ARRダッシュボードの作成に使用されるデータモデルの理解を助ける。
- そして、全体的にみんなが貢献できるようにします。

#### スナップショット ARR

ZuoraデータのARRは、契約の更新や修正に伴って日々変化しています。そのため、毎日データのスナップショットを作成し、データウェアハウスにARRの過去の記録を残して、レポートや分析に利用する必要があります。ARRのスナップショットとレポート作成に使用する3つの方法を以下に示します。

#### 方法は1つ。

dbtスナップショットは、変更可能なソーステーブルにタイプ2のSlowly Changing Dimensionsを実装しています。Slowly Changing Dimensions (SCD)は、テーブルの行が時間の経過とともにどのように変化するかを示すものです。方法1では、ライブARRメトリクスを生成するために使用されるマートであるmart\_arrテーブルをスナップショットします。この方法では、毎日 mart\_arr の正確なレプリカを作成し、日付スパイン技術を使用してモデル化します。この方法では、スナップショットされた生データを使用せず、

martテーブルをスナップショットしています。これにより、スナップショットされた mart\_arr テーブルのシンプルなモデルが作成され、常に mart\_arr テーブルのカラムが最新の状態に保たれるという利点があることがわかりました。

mart\_arrを完全にリフレッシュする必要はありません。これは、カラムに追従するmart\_arrの正確な3リングバインダーコピーが必要で、完全にリフレッシュする必要がないというユースケースにぴったりです。モデルが完全にリフレッシュされたとしても、常に同じ結果を返すことができます。

#### 方法は2つ。

SnowflakeClone機能を使用して、mart\_arrのゼロコピークローンを作成します。方法2は、毎日mart\_arrの正確な3リングバインダーコピーが必要な方法1と同じユースケースを満たしています。この方法は、dbtから独立した方法を使用することで、方法1のバックアップとして機能します。これらのクローン化されたテーブルは、クエリやレポーティングに使用される主要なテーブルではありませんが、特定の日にdbtスナップショットのパフォーマンスに問題があった場合に、冗長性とフェイルセーフを提供します。

#### 方法は3つ。

方法3では、より伝統的なボトムズアップ式のスナップショットを使用します。のスナップショットを使用して、変更可能なソーステーブルにタイプ2のSlowly Changing Dimensionsを実装し、生のソースデータをスナップショットします。その後、Date Spinningモデリング技術を使用して、スナップショットされた $\texttt{mart\_arr}$ テーブルをボトムアップで構築します。このモデルでは、スナップショットされたARRファクトと、ライブのCRM AccountおよびProduct Detailsディメンションを使用します。これにより、Sales CRM の Account Hierarchy と Product Tier や Product Delivery などの Product Details 属性の現在の状態で ARR がどのように見えるかという質問に答えることができます。この方法ではフルリフレッシュが可能で、ある日のARRの合計金額はフルリフレッシュ中に変更されませんが、そのARRのスライスは、ライブのSales CRM Account HierarchyとProduct Detailsディメンションが報告する内容に応じて変更され、更新されます。

#### Release Train Cadence:

1.近日公開

メンテナンスの予定です。

1.必要に応じて、2週間に1度、金曜日の午前9時から午前11時(米国東部時間)に定期メンテナンスを実施します。

クイックリンク

マーケット分析ハブへ

Sisense Discoveryを使い始めるには

# はじめに

そのためには、お客様に理解していただく必要があります。

- ◆ このダッシュボードでサポートされるKPI/PIは?
- ◆ メトリクスをどのように説明するかを説明するキーワー
- ドダッシュボードを支えるデータソース
- さらに探求するために、Sisenseでビジュアルや分析を自分で作成することができます。Sisenseディスカバリーツールから 始めるのが良いでしょう。Sisenseを始めたい方はこちらをご覧ください。
- さらに深く掘り下げるために、snowflakeでデータを探索することができます。Snowflakeでの探索の利点は、追加の情報 (つまり他のデータソース) に結合できることです。Snowflakeでの探索に関する追加情報はこちらをご覧ください。

# キーターム、メトリクス、KPI/PI、そしてキーフィールドとビジネスロジック

- ▶ 主な用語
  - ◆ 近日公開
- ▶ キーメトリクス、KPI、
  - PIは近日公開予定で

す。

▶ キーフィールドとビジネスロジック

# データソースとデータモデルの理解

ARRのダッシュボードとデータモデルは、ARRのERDにあるデータモデルを使用しています。

- ▶ データリネージの
- ▶ クエリ例

近日公開

# 追加リソース

▶ トラステッドデータソリューション

ARR モデルは、arr、arr\_snapshots、MRR、zuora、billing\_account、crm\_account タグを Trusted Data のテストとその結果に使用します。これは、Trusted Data Dashboard を使用して簡単に確認できます。

TrustedDataFrameworkの概要を見る

dbtガイド例 さらなるテストの実施に関する詳細と例 EDM Enterprise

- ► Dimensional Model Validations
- ▶ RAWソースデータ パイプラインバリデ
- ▶ ーション データセキュリティ分類

近日公開

#### **ORANGE**

- 説明行またはレコードレベルの顧客および個人データ。
- オブジェクトです。

0

#### **YELLOW**

- 説明します。GitLabの財務データで、集約や合計を含みます。
- オブジェクトです。

0

**ソリューション・オーナーシップ** ▶

merged.md 2021/11/8

# layout: handbook-page-toc title: "License Utilization Analysis"

このページでは

{:.no\_toc}

TOC {:toc}

# ライセンス使用状況の分析

お客様が注文したライセンスをどのように消費しているかを理解することは、お客様の全体像を把握する上で非常に重要なステップです。これは、製品が潜在的な導入問題を特定するのに役立ち、営業チームやTAMチームが顧客をよりよく理解し、解約やダウングレードの潜在的なリスクを特定するのに役立ちます。

このページはMVCであり、次の四半期に向けての意思表示と捉えていただく必要があります。

#### キーメトリクス、KPI、PI

- ライセンス稼働率は、アクティブユーザー数/注文したライセンス数として算出
- ◆ されます。ライセンス使用率は、使用されたライセンス数 / 注文されたライセンス数 / 注文されたライセンス数として計算されます。

#### 将来の分析が必要

これは、SaaSとセルフマネージドの両方のライセンス・アクティビティ・レートを示す予備的な分析でした。データ・チームは、FY21 Q1ではこの分析にさらに焦点を当て、より深い分析を行いたいと考えています。焦点を当てるべき重要な分野は以下の通りです。

- ◆ ライセンス利用率にも同じグラフを表示してみる
- ◆ ライセンス使用率/アクティベーション率の低さと解約/ダウングレードの相関関係を探る 製品
- ◆ 層、案件規模、業界別に分析を深める

#### さらなる努力が必要

このページを当社の品質基準に合わせるためには、以下のことが必

- ◆ 要になります。EDMモデルの作成とERDの作成
- ◆ L2ソリューションへの取り組み ページ

# layout: handbook-page-toc title:"Manual Upload of Usage Payload" $\subset$

のページについて

{:.no\_toc}

TOC {:toc}

### 利用データを用いたユーザーへの取り組み

- 使用状況のPingデータを初めて手動でアップロードする場合は、snowflake\_imports GCSバケットへのアクセスが必要です。アクセスリクエストを作成してください。アクセスの準備ができたら、アクセスリクエストに @gitlab-data/engineers グループをタグ付けします。
- ◆ snowflake\_importsGCSバケットに顧客用のフォルダがまだ存在しない場合は、作成します。
- snowflake\_importsGCSバケットの、そのお客様のフォルダ内に、使用方法のpingペイロードファイルを配置します。まず、そのフォルダをクリックします。次に、「UPLOAD FILES」ボタンを押し、ファイルを選択して「Upload」を押します。
- ファイルを整理する。理想的には、ファイルはjson形式で、例えばcompanyname\_usage\_payload\_00\_00\_0000.jsonのようになる。
- ◆ 使用状況のpingデータは、GitLabインスタンス内でどの機能やサービスが利用されているかを示す情報です。詳しくは

こちらのリンクをご覧ください。

# layout: handbook-page-toc title:"グループ名空間の変換メトリクス" この

ページの目次

{:.no\_toc}

TOC {:toc}

# グループ名空間の変換メトリクス。GitLab.com

新しいグループのネームスペースがGitLab.comに作成されると、ネームスペースの作成者とユーザーは、自分のネームスペースにメンバーを追加したり、特定のステージイベントを完了したり、最終的には無料トライアルから有料プランに変更したりして、GitLab.comとの関わりをさらに深めていくことが非常に重要になります。また、GitLab.comでネームスペースが作成された後、これらのイベントやアクションがどれだけ早く行われているかを知ることも重要です。

ダッシュボードで答えた質問

- ◆ ある期間内に、新規グループの名前空間のうち、無料トライアルから有料プランにアップグレードされ
- た割合はどのくらいですか?一定期間内に新規のグループネームスペースのうち、2名以上のメンバーで構成されている割合は?
- 新規グループのネームスペースのうち、一定期間内に「作成」ステージに参加した割合は?新規
- ◆ グループのネームスペースのうち、一定期間内に「Verify」ステージに参加した割合は?
- ステージの採用数、転換数、ネームスペースの作成数に大きな変化はありましたか?

#### このページの目的は

- グループネームスペースの操作方法の理解を助ける変換指標の理解を助けるデータモデルの
- ◆ 理解を助ける (WIP)
- ◆ そして、全体的にみんなが貢献できるように。

クイックリンク

新しいグループ名空間の変換メトリクス・ダッシュボード

グループ名空間コホート分析ダッシュボードビュー(WIP

データの注意点と制約事項

• 無料」または「有料」でフィルタリングすると、グループネームスペースの現在のステータスが反映され、以下のようになります。

は、採用イベントの時点でのネームスペースのアカウント状況を把握していません。

- これらのレポートには、新規ユーザーおよび既存の有料アカウントのユーザーが作成したグループ
- ◆ 名前空間が含まれます。データには、スパムユーザーによって作成された名前空間が含まれる場合があります。

このデータはGitLab.comの製品のみを反映しており、自己管理の利用状況は報告されていません。

#### データセキュリティの分類

グループ名空間変換ダッシュボードのデータの多くはオレンジ色であり、一般に公開すべきものではありません。このダッシュボードのデータを共有する際には、詳細をGitLabの組織内にとどめ、外部に共有する際には適切な承認を得るように注意してください。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保存しないよう、常に注意を払う必要があります。このデータはSnowflakeとSisenseに残すべきであり、特に承認されない限り、これらのアプリケーションを通じてのみ共有されることが理想的です。

ソリューション・オーナーシップ

- ◆ GitLab.comクレイグ・ゴメス (@craig-gomes
- データチームのサブジェクト・マター・エキスパート@dpeterson1 @kathleentam @mpeychet\\_ @iweeks @jeanpeguero @m walker @pluthra

新しいグループ名空間 作成段階 採用率@mkarampalas

新しいグループ名空間の検証段階の採用率。@jstava 少なくとも2人のユーザーが追加された新しいグループ・ネームスペース。 @s\_awezec

٠

•

٠

٠

◆ 新しいグループ名空間の試みから有料コンバージョン率へ。@s awezec

#### 主な用語

- グループネームスペース。トップレベルのネームスペースで、サブグループやプロジェクトを保持できます。これらのKPI および関連ダッシュボードのデータには、個人のネームスペースは含まれていません。
- 名前空間の作成日。新しいグループネームスペースが作成された日付(週
- ステージを検証する。グループネームスペースが最初のX日間でci\_pipelinesイベントを登録すれば成功
- 作成段階。グループネームスペースが最初のX日間でaction\_monthly\_active\_users\_project\_repoイベントを完了した場合に成功
- 有料プラン。名前空間が現在有料か無料かを示します。

#### ダッシュボードのフィルター

- 有料プランです。無料プラン、有料プランでの絞り込み
- 最初の X 日間フィルター。グループネームスペースが作成されてから X 日以内の、これらの KPI のコンバージョン率と 採用率を示します。現在のオプションには、最初の3日、7日、14日、30日、60日、90日、180日があります。
- 集約。ネームスペースの作成日を、選択した日付間隔(日、週、月、四半期、年)で集計します。 日付範囲。
- ◆ 過去X日以内に作成されたグループネームスペースのみを表示するようにチャートをフィルタリングします。

#### キーメトリクス、KPI、PI

新規グループ名空間試用から有料転換率 新規グルー

プ名空間 作成段階での採用率 新規グループ名空間

検証段階での採用率 2人以上のユーザーが追加され

た新規グループ名空間

#### キーフィールドとビジネスロジック

◆ データは主にSnowflakeから取得しています。

#### スニペット

スニペットは、異なるチャートで再利用可能なSQLコードの文字列を作成するために使用されます。詳細については、Sisense SQL Snippetsのページをご覧ください。

#### リファレンスSQL

	スニペットライブラリ説明	
分母_グループ_名前空間_on_trial	このスニペットは、「New Group Namespace Trial to Paid」の分母として使用されます。 コンバージョン率のKPIで、作成されたグループネームスペースのうち、最初の X日以内に試用を開始したグループネームスペースの総数を算出します。	
分水嶺_全グループ名空間	クエリは、上記の残りの3つの KPI の分母として使用されます。分母は、作成されたグループネームスペースの数を計算します。 denominator_group_namespaces_on_trialのスニペットとは異なり、このスニペットは	
	トライアルを開始していない新規グループネームスペースを除外しません。	

#### エンティティ・リレーションシップ・ダイアグラム

# Diagram/EntityGrainPurposeKeywords

グループ名空間 変 換指標 ERD

日付、場 所、プラ ンの種類 キーとなるステージのイベントや アクションの新ネームスペース採用 に関する洞察を提供する

Stage Adoption, Invite Members, New Groups, Paid Conversion, New Namespaces, GitLab.com

データプラットフォームソリューション

データの系統

近日公開

DBTソリューション

近日公開

トラステッドデータソリューション

▶ トラステッドデータフレームワーク

近日公開

EDM エンタープライズ ディメンション モデル バリデーション

◆ エンタープライズ次元モデル検証ダッシュボード

近日公開

RAWソースデータパイプラインの検証

近日公開

カスタマーDB

- ◆ 信頼できるデータダッシュボード
  - tdfタグを含むすべてのCustomerDBデータテストのレポート。

近日公開

# layout: handbook-page-toc title: "People Analytics Overview"

このページについて

{:.no\_toc .hidden-md .hidden-lg} Ø

ように。TOC {:toc .hidden-md

.hidden-lg}.

{::options parse\_block\_html="true" /}。

# ピープルアナリティクス

仕事をしている人の研究です。人事部は組織行動へのアプローチを反転させ、判断や意見を用いて意思決定を行う代わりに、データの力で事実に基づいた意思決定を行っています!

#### ピープルスペースにおけるデータのメリット

- より合理的な人材獲得プロセスを構築することで、GitLabが強固で多様性のあるチームを構築することを支援し、プロセスを経る候補者にも充実した経験を提供します。
- GitLabチームメンバーの体験を常に向上させるために、チームを推進しています。 "Everyone can contribute "という言葉は、私たちが本気で言っていることです。エンゲージメント調査やKPIを用いてチームの目標を達成するために、私たちは人を最優先に考えています。
- 他のデータと結びつけようセールスデータであれ、エンジニアリングデータであれ、私たちにとってはすべてが重要です。私 たちは、チームメンバーの成長をどのように支援できるか、またリーダーが組織の影響力をどのように理解するかを理解 したいと考えています。

### 一般的なルールとガイダンス

ピープルグループは、運用ソフトウェア(BambooHR、Greenhouseなど)の中で日々の業務を行うことができなければ

- なりません。ウェアハウス内のピープルデータは、組織内の一般的なピープル情報を報告するためのものでなけれ
- ばなりません。
- リスクを把握するために、「ゆりかごから墓場まで」データにアクセスする役割とユーザーを常に把握する。

◆ レポート作成に必要なものだけをデータウェアハウスに取り込み、その他の機密性の高いデータは運用ツールに残しておくべきです。

• 可能な限りリスクを軽減するために、メタリックな計算や報告に使用されるセンシティブなデータを匿名化する。

# 人々のデータソース

#### BambooHR

人事管理システムです。

温室

採用・応募者管理システム

#### PTO By Roots

チームメンバーの休暇を記録するslackアプリケーション

layout: handbook-page-toc title: "People KPI Deck Report"

# このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

# People KPI Deck Report

このレポートは、People Groupの月例KPIミーティングで説明された指標を考慮しています。こ

のレポートは、このGoogleレポートに掲載されています。

#### 基礎となるデータ

rpt\_people\_kpisデータモデルのドキュメントは、こちらをご覧ください

。これは、過去12ヶ月間のデータをキャプチャします。

layout: handbook-page-toc title:"プロモーションレポート"

# このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

### プロモーションレポート

このレポートは、トータルリワードチームのため

に作成されます。レポートは、このgoogleフォ

ルダにあります。

#### プロモーションレポートの更新

1. Promotions Last 12 Months」タブを更新します。A. snowflakeに向かう(Oktaでログイン) B. 以下のコードを実行する。

```
SELECT *
from "prep". "sensitive". "bamboohr_promotions_xf"
WHERE promotion_month BETWEEN DATEADD(year,-1, DATE_TRUNC(month, CURRENT_DATE()))AND
DATEADD(moon,-1, DATE_TRUNC(moon, CURRENT_DATE()))
ORDER by promotion_date
```

C.結果をダウンロードして、「過去12ヶ月間のプロモーション」タブに配置します。

2. Budget Spend」タブを更新する A. snowflakeで、以下のコードを実行する。

```
SELECT *
"prep""sensitive""bamboohr_budget_vs_actual "より
ORDER BY 会計年度,会計四半期,部門
```

B.結果をダウンロードして、「予算支出」タブに入れます。

- 3. 予算と実績のピボットを更新する 注:予算の合計を更新するには、ここをクリックします。People Budget Sheetload。このデータは全体のレポートを更新するために使用されます。
- 4. リフレッシュ時にプロモの予算を更新するには、以下の方法があります: sheetload.budget\_people

layout: handbook-page-toc title:"スラックダッシュボード"

# Slackダッシュボード

このダッシュボードは、Chief of Staff Team to the CEOのKPIに関連して、DMSaではないメッセージの割合を監視するために使用されます。このダッシュボードは、パブリックチャネルの利用を増やし、ハンドブックファーストにする方法を理解するのに役立ちます。

ダッシュボードはこちらからご

覧いただけます###Updating

Data

- 1. Head to https://gitlab.slack.com/stats
- 2. 日付範囲を全時間帯に変更してデータをエクスポート
- 3. 以下のコードを使用してデータをクリーンアップします (これはR Studioから行うことができます)
  - https://gitlab.com/gitlab- data/analytics/-blob/master/transform/general/clean\_slack\_data.R
- 4. クリーニングされたデータを入手したら、sheetload.slack\_stats(https://docs.google.com/spreadsheets/d/15a2PVvSs7K\_C-)に向かいます。

EsGq2hKLUYNywxPmE74Ldx-YbXqN-w/edit#gid=673732546) を作成し、現在のタブのデータを新しいデータで置き換える。

layout: handbook-page-toc title:"Talent Acquisition Metrics"

#### このページについて

TOC {:toc}

 ${::options\ parse\_block\_html="true"\ /}_{\circ}$ 

# タレント・アクイジション・キー・メトリクス

GitLabはチームメンバーをサポートすることに重点を置いています。他のビジネスステークホルダーと同様に、People Group内の Talent Acquisitionチームは、チームが社内で適切なチームメンバーを採用できるように採用プロセスを常に改善し、またプロセ

スを経た候補者が良い経験をできるようにすることに重点を置いています。

タレント・アクイジションのKPIの多くは、KPIインデックスに掲載されています。

2021/11/8 merged.md

TalentAcquisitionMetricsDashboardは、他の指標に加えて、当社のKPIを考慮しています。

クイックリンク

タレント・アクイジション・メトリクス・ダッシュボード

モデル・ドキュメンテ

# ダッシュボードの使い方

# **People Metrics**

# Dashboard

人材獲得指標ダッシュボードでは、先月、前四半期末のGitLab社の状況を簡単に把握することができ、選択した日付範囲(デフ オルトは過去12ヶ月間)でのトレンドを確認することができます。

ダッシュボードでは、ユーザーは以下の項目でフィルターをかけることができます。

- ◆ タレント・アクイジション部門
- Talent Acquisition Department (部門レベルに移行するには、部門が選択されている必要があり
- ます) Talent Source
- ◆ タレントソース名 (タレントソースが選択されている必要があります
- ◆ タレント・アクイジション・チーム (エンジニアリング・タレント・アクイジション・チーム、セールス
- ◆・タレント・アクイジション・チームなど) タレント・アクイジション・チームのリクルーター (タレ ント・アクイジション・チームのフィルターに依存する
- ◆ 候補者コーディネーター
- ◆ Is Prospect(候補者が、当社がアプローチした潜在的な見込み客であるのか、それとも他のソースから応募して採用プロ セスを開始した候補者であるのかを確認するため
- ソーシング・チーム・フラッグ
- Sourcing Team Sourcer (sourcing team flag set = sourcing teamに依存) このダ

ッシュボードを使用する際の重要なヒントがあります。

1.メトリクス名が青で表示されている場合は、データをより深く掘り下げることができる付属のダッシュボードがあります。 リンクをクリックすると、ダッシュボードに移動します。

# タレント・アクイジション・メトリクス・デー ・モデル

このダッシュボードは greenhouse stage analysis をベースに作られています。ソースデータはGreenhouseからのものです。

このデータモデルは、採用プロセスを経た各候補者の概要と、主要な属性(ソース、リクルーター、求人部門など)に加えて 、どのステージに到達したかを示しています。

このモデルを作成するための基礎データは、温室ロールにアクセスできるチームメンバーがアクセスでき

ます。Sisenseでは、基礎となるクエリが次のように表示されます。

```
WITH talent acquisition_data AS

( SELECT * )
FROM LEGACY.GHGENUSIVEAN STAGE_ANalysis
WHERE [division_modified=Talent Acquisition Division]
AND [department_name=talent acquisition_department]
AND [source_type=Talent_Source] .
AND [sourcing_team=Sourcing_Team_Flag]
AND [sourcer_name=Sourcing_Team_sourcer]

AND [talent acquisition_team=talent acquisition_team]
```

```
AND [candidate_recruiter=Talent Acquisition_Team_Recruiter]
AND [source_name=Talent Acquisition_Source_Name] 。
AND [candidate_coordinator=candidate_coordinator]
AND [is_prospect=Is_Prospect]となります。

[talent acquisition_metrics] --- これは、すべての主要なメトリクスを識別するsisenseスニペットを使用しています。
例えば、どのような候補者が応募者とみなされるのか、応募から審査までの割合はどのくらいか、応募から採用までの割合はどのくらいか、などを計算します。

スニペット全体は[こちら](https://app.periscopedata.com/app/gitlab/snippet/recruiting_metrics/55081acd27f44d7fb1706d47f44b5ae 8/edit)でご覧いただけます。

SELECT prospected
FROM metrics
WHERE-month_stage_entered_on = DATE_TRUNC(month, CURRENT_DATE())
```

layout: handbook-page-toc title:"チームメンバーの分離"

### このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

# チームメンバーの分離レポート

このレポートは、分離されたすべてのチームメンバーを特定します。注:このデータは、roles.ymlで bamboohr\_sensitive\_separationロールを持つユーザーがsnowflakeでのみアクセスできます。このデータはSisenseではアクセスできません。

```
ランニングでデータにアクセスできます。
、、
SELECT
*.
from "prep". "sensitive". "bamboohr_separations"
、、、
```

# アクセス方法

このレポートにアクセスするためには、アクセスリクエストを作成する必要があります。

- 1. スノーフレーク
- 2. bamboohr\_sensitive\_separationロールへのアクセス。

このアクセスリクエストは、ピープルチーム(トータルリワードマネージャーまたはピープルサクセス担当シニアディレクター)が確認する必要があります。

# layout: handbook-page-toc title: "PTO by Roots"

# このページについて

{:.no\_toc .hidden-md .hidden-lg}.

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}.

# PTO by Roots (

# Slack) ゴール

GitLabチームメンバーはslackを使って休暇を記録しています。このデータを他のデータソースと結びつけることで、 理解を深めるのにとても役立ちます。

- ◆ 休みを取っていないメンバーの割合は?家族と友人
- ◆ の日の計画を手伝う
- ◆ 今後のキャパシティを把握し、マイルストーンやOKRを計画するパ
- ◆ フォーマンス指標への影響を測定する

具体的な例としては、開発部の狭義のMR率(チームメンバーの生産性を示す指標)の変動を把握するためにPTOデータを使用しています。当社では、メンバーに充電のための休暇を取ることを奨励しており、休暇を考慮することで、時系列で見たMR率の低下や上昇を説明することができます。例えば、前月に比べて休みが増えたことで、今月の狭義のMR率が下がったことを説明できるかもしれません。

このデータを使って、個々のチームメンバーの欠席時間や理由を評価しないでください。個人の出席状況に懸念がある場合は、担当のピープルビジネスパートナーに連絡してください。

#### データプロセス

私たちが取得するデータ

- データレベルでは、理由ではなく、休んだ日が気になります。したがって、我々のデータでは、extract層にある休んだ理由やsnowflakeで利用可能なものはすべて削除しています。
- このデータは、様々なデータソリューションに使用するために集計されます。例えば、MR率については、月、部門 、部署のレベルでKPIを表示します。休み時間やチームメンバーの情報は、Sisenseではアクセスできません。

このタイムオフデータへのアクセスは、データチーム(Parul Luthra, Data Engineer-未定)に限定されています。アクセスを希望する場合は、アクセスリクエストフォームを提出し、SlackのPTO管理者、Senior Director, People Successの承認を得る必要があります。

layout: handbook-page-toc title: "People Metrics"

# People Key Metrics

ギットラボはチームメンバーのサポートに力を入れています。他のビジネスステークホルダーと同様に、ピープルグループは、チームメンバーがギットラボでの貢献に集中できる環境を作るために、データを重視しています。

People KPIの多くは、KPIIndexに掲載されています。

PeopleMetricsDashboardは、当社のKPIに加え、ピープルデータの意思決定に役立つその他の指標も考慮しています。

クイックリンク

ピープルメトリクス ダッシュボード

ピープルメトリクスの概要

モデルドキュメント

人物データ発見ダッシ ュボード

データモデルのデモ

データモデルのビデオ・デモと ダッシュボード

# ダッシュボードの使い方 People Metrics

# Dashboard

PeopleMetricsのダッシュボードでは、先月、前四半期末のGitLabの状況を簡単に把握することができ、選択した日付範囲(デフォルトでは過去12ヶ月)での会社の傾向を確認することができます。

ダッシュボードは、3つのレベルにフィルタリングすることができます。

**ダッシュボ** ードビュー

ードビューフィルター	1フィルター2
GitLab OverallBreakout =	CompanyBreakout_Division_Department = All Company
GitLab OverallBreakout =	DivisionBreakout_Division_Department = 該当する部門のフィルタリング
OverallBreakout = GitLab	Breakout_Division_Department = 該当する部門でのフィルタリング

注:この画面では、一度に1つの部門または1つの部署しか選択でき

ません。このダッシュボードを使用する際の重要なヒントがありま

す。

ダッシュボ

- 1. メトリクス名が青色で表示されている場合は、データをさらに詳しく調べることができるダッシュボードが表示されます。リンクをクリックすると、ダッシュボードに移動します。
- 2. すべてをレポートにまとめて見る必要がある場合は、「PeopleMetricsOverview-CurrentMonth」をご覧ください。 これはダッシュボードからアクセスできます。
  - これは、ドリルダウンレポートと同じように機能します。たとえば、People\_Metrics\_View = 'Division`でフィルタリングすると、すべての kpis の内訳がレポートに表示され、上部には部門が、最後には GitLab Overall が表示され、部門レベルでの比較が簡単にできます。
    - これは、どの部門の離職率が高いのか、どの部門の採用が遅れているのかといった質問に最適です。
  - People\_Metrics\_View = 'Department'でフィルタリングすることにより、ビューアーは上記の質問をより深く掘り下げ、部門レベルで問題があるかどうかを判断することができます。

# People Metricsデータモデル

ダッシュボードの大部分は、bamboohr\_rpt\_headcount\_aggregationをベースに

構築されています。このデータモデルは、4つのレベルですべての関連する人数メ

トリクスを月ごとに集計します。

- 1. kpi breakout -- 会社レベルまで集計されたすべての指標を表示する。
- 2. ディビジョン\_ブレイクアウト
- 3. デパートメント・ブレークアウト
- 4. eeoc\_breakout -- ダイバーシティ指標を理解するために集計された指標を示します(性別、民族、地域など)。

このモデルを作成するための基礎データは、チームメンバーがbamboohr\_analyticsロールにアクセスし、よりセンシティブなデータ (例:eocデータや報酬) についてはbamboohrロールにアクセスします。

独自のレポートを作成します。

自分でレポートを作成したい方は、Sisenseのディスカバリー機能の使い方を説明したデモビデオをご覧ください。これにより、ユーザーはフィールドをドラッグ&ドロップすることができます。

この機能のプレビューは、PeopleDataDiscoveryFeatureDashboardでもご覧いただけます。

layout: handbook-page-toc title:"価格分析"

### このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}.

価格分析

価格設定とは、特定の価格帯でお客様がGitLabから受け取る価値と、その価格帯の収益性を分析するプロセスです。また、これらの価格がビジネス全体にどのような影響を与えるかを理解し、最適な価格を決定することも含まれます。

の価格設定は、顧客とGitLabの収益性を左右するものです。価格分析ページでは、GitLabチームメンバーが現在の価格戦略を検討し、最適化のための洞察を深めるために使用できる情報とツールを提供します。

このページの目的は

- ◆ 価格設定-顧客割引ダッシュボードの操作方法の理解を助けます。
- Pricing Customer Discounts Dashboard」の作成に使用されたデータモデルの理解を助けます
- 。GitLabでのあなたの役割に最も適した認定を受けることで、あなたの理解度を確認します。
  - o ダッシュボードの使い方を詳しく知りたい方は、ダッシュボードのユーザー認定を受けてください。
  - o Sisenseダッシュボードの開発について詳しく知りたい方は、DashboardDeveloper認定資格を取得し
- てください。 そして、全体的にみんなが貢献できるようにしてください。

クイックリンク

価格 - 顧客割引ダッシュボード

PnPテスト&リサーチ ダッシュボード

はじめに よな様の概要がNo. T ボード Signes Discovenyの使用関始について

> 始めるにあたり、以下のことを確認したいと思 います。

を理解しています。

- ◆ このダッシュボードでサポートされるKPI/PIは?
- ◆ メトリクスをどのように説明するかを説明するキーワー
- ドダッシュボードを支えるデータソース
- さらに探求するために、Sisenseでビジュアルや分析を自分で作成することができます。Sisenseディスカバリー・ツールを 使用するのが良い方法です。Sisenseを始めたい方はこちらをご覧ください。
- さらに深く掘り下げるために、snowflakeでデータを探索することができます。Snowflakeでの探索の利点は、追加の情報 (つまり他のデータソース) に結合できることです。Snowflakeでの探索に関する追加情報はこちらをご覧ください。

目で見るのがお好きな方は価格分析ソリューションのトレーニングをご覧ください。

# キーターム、メトリクス、KPI/PI、そしてキーフィールドとビジネスロジック

- ▶ 主な用語
- ▶ キーメトリクス、KPI、PI
- ▶ キーフィールドとビジネスロジック

## データソースとデータモデルの理解

このダッシュボードは、Sisenseの`base\_pricing\_discount`スニペットから作成されています。Sisenseのスニペットやビューを使うと、シンプルなSQL文字列を好きなだけチャートで参照することができます。これにより、エンドユーザーがARRやシート数などの計算の裏にあるロジックを理解したり、どのような結合が可能かを推測で理解したりするのに役立ちます。

このモデルから独自のダッシュボードを作成するには、Sisenseのクエリに次のように入力するだけです。

### [ベース\_プライシング\_ディスカウント]

Base\_pricing\_discount は mart\_discount\_arr を基に構築されており、様々な顧客ディメンション(親顧客、製品カテゴリー、配送、アカウントオーナーチーム、再販業者/そうでないもの、サブスクリプション、アカウントなど)による割引や再販取引についての洞察を提供します。

merged.md

2021/11/8

このマートは、ERD(Entity Relationship Diagram)に見られるようなデータモデルを考慮しています。

```
▶ データリネージの▶ 例 クエリ[ベース_プライシング_ディスカウント]
```

)

```
</details> </details
<br>
<style> #headerformat {
background-color:6666c4; color: black; padding:5px; text-align: center;
<h1 id="headerformat">追加のリソース</h1>。
<詳細
<summary markdown='span'> トラ
  ステッドデータソリューション
</summary>
概要は[Trusted Data Framework](https://about.gitlab.com/handbook/business-technology/data-
team/platform/#tdf)をご覧ください。
[dbtガイドの例](https://about.gitlab.com/handbook/business-technology/data-team/platform/dbt-
guide/#trusted-data-framework)を参照してください。
さらなるテストの実施に関する詳細と例
</details> </details
<詳細
<summary markdown='span'>
 EDM エンタープライズ ディメンション モデル バリデーション
</summary>
```

The [(WIP) Enterprise Dimensional Model Valida独向p24Bashboard] (エンタープライズ ディメンショナル モデル バリデーション ダッシュボード)

```
(https://app.periscopedata.com/app/gitlab/760445/WIP:-Enterprise-Dimensional-Model-Validation-
Dashboard) では、最新のEnterprise Dimensionalモデルのテストと実行について報告しています。
</details> </details
<詳細
<summary markdown='span'>
  RAWソースデータパイプラインの検証
</summary>
[Data Pipeline Health Validations](https://app.periscopedata.com/app/gitlab/715938/Data-Pipeline-
Health-Dashboard)
</details> </details
<詳細
<summary markdown='span'>
   データセキュリティの分類
</summary>
価格分析のデータの多くは、[オレンジ] (/handbook/engineering/security/data-classification-
standard.html#orange) stan
standard.html#yellow) となっています。これには、ORANGEのアカウントの顧客メタデータや、GitLabのNon
publicの財務情報が含まれており、これらはすべて一般に公開されるべきではありません。このダッシュボードからデ
ータを共有する際には、その詳細が GitLab の組織内に留まるようにし、外部に共有する場合には適切な承認を得られ
るように注意しなければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラッ
プトップにデータを保存しないよう、常に注意を払う必要があります。これらのデータは、
[Snowflake](/handbook/business-technology/data-team/platform/#data-warehouse) &
[Sisense](/handbook/business-technology/data-team/platform/periscope/)に残すべきであり、特別な承認がない
限り、これらのアプリケーションを通してのみ共有されるのが理想的です。
**ORANGE**
- 説明行またはレコードレベルの顧客および個人データ。
- オブジェクトです。
    - dim_billing_accounts`.
   - dim_crm_accounts`.
YELLOW** (イエロー
- 説明します。GitLabの財務データで、集約や合計を含みます。
- オブジェクトです。
   - dim subscriptions
   - `fct_charges`
   - `fct_invoice_items`
    - `fct_mrr`
</details> </details
く詳細
<summary markdown='span'>
   ソリューションの所有権
</summary>
* ソース・システム・オーナー
   * セールスフォース: `@jbrennan1`.
   * Zuora: `@andrew_murray` (英語)
* ソースシステムのサブジェクト・マター・エキスパート。
   * セールスフォース: `@jbrennan1`.
    * Zuora: `@andrew_murray` (英語)
* データチーム・サブジェクト・マター・エキスパート:`@paul_armstrong` `@jeanpeguero` `@jjstark` `@iweeks`
</details> </details
レイアウト:ハンドブック-ページ-トック
のタイトルを付けました。"Product Geolocation Analysis"
```

## On this page

```
{:.no_toc}
- TOC
{:toc}
## Product Geolocation Analysis : Self-Managed
自社の製品が世界中のどこで使われているかを理解することは、顧客や製品のグローバル展開、関連するロケーション
インサイトをより深く理解するための重要なステップです。
現在、GitLabのお客様の大半は、「GitLabセルフマネージドインスタンスをダウンロードし、インストールし、ホストす
る」(/handbook/marketing/strategic-marketing/dot-com-vs-self-managed/#why-you-probably-want-gitlabcom)
ことを選択しており、そのために私たちは「素晴らしいセルフマネージドカスタマーエクスペリエンスを提供することに
重点を置いている」(/direction/#strategic-challenges)。
セルフマネージ製品のライフサイクルや投資すべき機能について、データに基づいた適切な判断を下すために、[セル
フマネージのお客様](/is-it-any-good/)は、[セルフマネージのインスタンスでusage pingを有効にする
[(https://docs.gitlab.com/ee/user/admin_area/settings/usage_statistics.html#instance-level-
statistics)ことで、またはカスタマーサクセスチームと値を共有することで、インスタンスレベルで毎週GitLabに
[usage ping](/handbook/customer-success/tam/usage-ping-faq/)を送信しています。
このインスタンスレベルのデータにより、GitLabは国レベルの統計や、インスタンスの採用率、バージョン採用率、イ
ンスタンスのライフサイクルの傾向を把握することができます。
**このページの目的: **。
 * Product Geolocation dashboard] の操作方法の理解を助ける
(https://app.periscopedata.com/app/gitlab/731086/Product-Geolocation-:-Self-Managed-(doesn't-
include-SaaS))
 * Product Geolocation Dashboardの作成に使用されたデータモデルの理解を深めることができます。
 * GitLabでのあなたの役割に最も適した認定を受けることで、あなたの理解度を評価してください。
   * Product Geolocationダッシュボードの使用方法については、[Dashboard user
certification](https://docs.google.com/forms/d/e/1FAIpQLScflqXOnU-
W6kz24b5qD715uw9_1s6tfKF34qf1tqvFgguIVw/viewform)をご利用ください。
   * Sisense ダッシュボードの開発について詳しく知りたい方は、【ダッシュボード開発者認定】(
https://docs.google.com/forms/d/e/1FAIpQLSeqicaMfWVUfFsex9 o6GTkWJKobYBT8qucz9YNmyDm5ZKqiA/viewform) &
受けてみてください。
?usp=sf link)となっています。
   * SQLを使ってデータを深く理解するためには、[SQL Certification] (
https://docs.google.com/forms/d/e/1FAIpQLScH9CkiACQ1worzldjUi6cUWUL03tXrLEEaZALABabZPV7GuQ/viewform) &
受講してください。
?usp=sf_link)となっています。
 * そして、全体的にみんなが貢献できるように。
### Quick Links
<div class="flex-row" markdown="0" style="height:80px">

 <a href="https://app.periscopedata.com/app/gitlab/731086/Product-Geolocation-:-Self-Managed-(does</pre>
not-include-SaaS)" class="btn btn-purple"
style="width:33%;height:100%;margin:5px;float:left;display:flex;justify-content:center;align-
items:center;">製品のジオロケーション(セルフマネジメント・ダッシュボード)</a>。
 <a href="https://www.youtube.com/watch?v=F4FwRcKb95w&feature=youtu.be" class="btn btn-purple"</pre>
style="width:33%;height:100%;margin:5px;float:left;display:flex;justify-content:center;align-
items:center;">Sisenseをはじめよう</a
</div> </div>
<br> <br>>
```

#### ### Getting Started

そのためには、お客様に理解していただく必要があります。

1.このダッシュボードを使ってどのようなKPI/PIがサポートされるのか

merged.md

2021/11/8

1.メトリクスをどのように説明するかを説明するキーワード

- 1. ダッシュボードのデータソース
- 1. さらに探求するために、Sisenseでビジュアルや分析を自分で作成することができます。Sisenseディスカバリーツールを使用して始めるのが良い方法です。Sisenseを始めてみたい方はこちらをご覧ください。

(https://about.gitlab.com/handbook/business-technology/data-team/direction/self-service/#selfservice-dashboard-developer)

1.さらに深く掘り下げるために、snowflakeでデータを探索することができます。Snowflakeでの探索の利点は、追加の情報(つまり他のデータソース)に結合できることです。Snowflakeでの探索に関する追加情報は[こちら](https://about.gitlab.com/handbook/business-technology/data-team/direction/self-service/#self-service-sql-developer)をご覧ください。

Product Geolocation Dashboardは、以下のKPIやPIの把握に役立ちます。

- [アクティブホスト](/handbook/product/performance-indicators/#active-hosts)
- [ロスト・インスタンス](/handbook/product/performance-indicators/#lost-instances)
- [有料ユーザー](/handbook/product/performance-indicators/#paid-user)
- [有料UMAU](/handbook/product/performance-indicators/#paid-umau)
- [月間ユニーク・アクティブ・ユーザー数 UMAU](/handbook/product/performance-indicators/#unique-monthly-active-users-umau)

ダッシュボードをご覧になる際には、以下の点をご理解ください。

- [アカウント](/handbook/sales/#additional-customer-definitions-for-internal-reporting)
- アカウントインスタンス 1つのアカウントにマッピング可能なレポートされたインスタンスの合計数
- [ホスト](https://docs.gitlab.com/ee/development/product\_analytics/event\_dictionary.html)
- 「インスタンス](https://docs.gitlab.com/ee/development/product analytics/event dictionary.html)
- Instance User Count インスタンス上のユーザーの合計数
- [有料ユーザー](/handbook/product/performance-indicators/#paid-user)
- [製品階層](/handbook/marketing/strategic-marketing/tiers/#overview)
- [使用方法 Ping](https://docs.gitlab.com/ee/development/product\_analytics/event\_dictionary.html)

注:製品ジオロケーション分析をサポートする追加データは、「オレンジ

](/handbook/engineering/security/data-classification-standard.html#orange)または[イエロー](/handbook/engineering/security/data-classification-standard.html#yellow)に分類されます。これには、アカウントからのORANGEの顧客メタデータ、SalesforceやZuoraからのコンタクトデータ、GitLabのNon公開財務情報などが含まれ、これらはすべて公開されるべきものではありません。このダッシュボードからデータを共有する際には、その詳細が GitLab の組織内に留まるようにし、外部に共有する際には適切な承認を得られるように注意しなければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保存しないよう、常に注意を払う必要があります。これらのデータは、[Snowflake](/handbook/business-technology/data-team/platform/#data-warehouse)と[Sisense](/handbook/business-technology/data-team/platform/#data-warehouse)とに表すべきであり、特別な承認がない限り、これらのアプリケーションを通じてのみ共有されることが理想的です。

\*\*ORANGE\*\*。

- 説明行またはレコードレベルの顧客および個人データ。
- オブジェクトです。
  - dim\_billing\_accounts`.
  - dim\_crm\_accounts`。
  - `usage\_ping\_mart`

### ソリューション・オーナーシップ

- ソース・システム・オーナー
  - バージョンアップしました。`@jeromezng`
  - セールスフォース: `@jbrennan1`.
  - Zuora: `@andrew\_murray` (英語)

- 場所(IPアドレス):`@m\_walker`さん
- ソースシステムのサブジェクト・マター・エキスパート。
  - バージョンアップしました。`@jeromezng`

- セールスフォース: `@jbrennan1`.
  - Zuora: `@andrew\_murray` (英語)
  - 場所 (IPアドレス) : `@m\_walker`さん
- データチーム・サブジェクト・マター・エキスパート: `@mpeychet ` `@m walker`

#### ### データソース

Product Geolocation Dashboardは、`usage\_ping\_mart`データモデルを使って構築されています。

- データは主に使用率のpingデータから得ており、顧客セグメントのフィールドはZuoraとSalesforceから得ています。
  - 分析は、自己管理インスタンスから暦月に受信した最後の使用pingを考慮して標準化されており、
- `is\_last\_ping\_in\_month`として利用できます。これにより、使用状況のメトリクスがインスタンス間で重複しないようになります。
- Geolocationフィールドは、インスタンスの\*ホスト\*のIPアドレスに由来するもので、必ずしもインスタンスの物理的な位置を示すものではありません。
- usage\_ping\_martを作成する際にどのようなテーブルが考慮されたかを理解するには、Usage Ping Mart Entity Relationship Diagramをご覧ください。
- ERDを理解することで、異なるデータソースがどのように結合するかを理解し、必要に応じてモデルを修正することができます。

### エンティティ・リレーションシップ・ダイアグラム

| [Usage Ping Mart](https://app.lucidchart.com/documents/view/be5f5dc8-8ad5-4586-af53-93ff5e00f720/0\_0)| usage\_ping\_id | 使用状況のpingと関連する顧客セグメンテーションメトリクスを探索するためのマート | #### Sisenseのデータディスカバリー機能

## ##### Sisenseでさらに探究する

Sisenseは私たちが視覚化のために使っているツールです。ダッシュボードデザイナーはSQLに精通している必要があります。しかし、Sisenseのデータディスカバリー機能を使えば、ドラッグ&ドロップでチャートを作成することができるので、制限を感じることはありません(SQLクエリは必要ありません)。

```
<figure class="video_container">
```

<iframe src="https://www.youtube.com/embed/h\_b9A8F7Ic8" frameborder="0" allowfullscreen="true">
</iframe>
</iframe>
</irrame>

</figure>

詳細はこちら【SisenseのData Discoveryページ】(

https://dtdocs.sisense.com/article/data-discovery) 。

## #### Sisense Snippets

スニペットは、異なるチャートで再利用可能なSQLコードの文字列を作成するために使用されます。詳しくは、[Sisense SQL Snippetsページ](https://dtdocs.sisense.com/article/snippets)をご覧ください。

##### usage\_pings\_mart スニペット

このスニペットは現在、Sisenseに[usage\_pings\_mart]という名前で存在しています(

https://app.periscopedata.com/app/gitlab/snippet/usage\_pings\_mart/553a4fc6bf004b749eb60a144d722ccc/ 編集)。

```
````sql
WITH pings AS (
```

SELECT \*

FROM analytics.usage\_ping\_mart
WHERE ping\_source = 'SelfManaged'
 AND is\_last\_ping\_in\_month =
 TRUE AND date\_id >= 20191101

```
AND [ping_product_tier=product_tier]とする。
AND [ping_country_name=Usage_Ping_Country].
```

### リファレンスSQL

### 1ヶ月間に使用量のpingを送信したアカウントの合計数

```
[use_pings_mart]

SELECT
  ping_monthです。
  COUNT(DISTINCT account_id) AS
  total_accounts FROM pings
GROUP BY 1
```

#### 国ごとの使用量pingを送信するインスタンスの月間数

```
[use_pings_mart]

SELECT
ping_monthです。
ping_country_ nameAS country_name,
COUNT(DISTINCT uuid) AS instances_reporting
FROM pings
GROUP BY 1,2
```

## データプラットフォームソリューション

## データの系統

- dbtモデルの系統図
- ◆ ジオロケーションへのIPアドレスのマッピングは、無料のgeolite2Maxmindデータベースから取得しています。
- ◆ 位置情報もMaxmindのデータベースから取得していますが、iso3の国コードフィールドを除いては、リポジトリの ZuoraCountryCSVから取得しています。

## DBTソリューション

- テーブル間の大きな結合や、IPアドレスとジオロケーションのマッピングがless-than/greater-than結合節で構成されることを 避けるために、今回のマージリクエストで最初に実装されたように、IPアドレスは他のモデルとは別にジオロケーション に段階的にマッピングされています。
- この方法では、SisenseでIPアドレスを不明瞭にすることができますが、異なるデータベーステーブル間でIPアドレスを照合する機能は維持されます。

## トラステッドデータソリューション

トラステッドデータフレームワーク

EDM エンタープライズ ディメンション モデル バリデーション

◆ エンタープライズ次元モデル検証ダッシュボード

RAWソースデータパイプラインの検証

#### バージョン

- 信頼できるデータダッシュボード
  - o tdfタグを含むすべてのバージョンデータテストを報告します。

#### ライセンス

- 信頼できるデータダッシュボード
  - o tdfタグを含むすべてのライセンスデータテストを報告します。

#### カスタマーDB

- ◆ 信頼できるデータダッシュボード
  - tdfタグを含むすべてのCustomerDBデータテストのレポート。

### 手動でのデータ検証

◆ 手動使用のPing検証ダッシュボード

layout: handbook-page-toc title:"製品利用データ"説明"製品利用データは、お客様が製品としてのGitLabをどのように、いつ、どこで利用しているかを定量的に測定したものです。"

## このページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

## 製品使用データ

**製品使用データ**は、お客様がいつ、どこで、どのようにしてGitLabを製品として使用しているかを定量的に測定したもので、GitLabチームがより良い製品を作り、お客様の導入を促進し、お客様の維持を改善するために使用します。製品利用データのページでは、GitLabチームのメンバーが製品利用データを調査し、顧客のインサイトを開発するために使用できる情報とツールを提供します。

顧客の製品使用データのビジネスユースケース。

- 顧客のGitLabインスタンスにおけるユーザー導入/オンボーディングの洞察
- ◆ 特定のDevOpsステージにおける機能の利用率の増減傾向の監視 アクティブユーザーの減少の監視
- 最新バージョンが遅れているアカウントの特定 デプロ
- ◆ イに使用されたインストールの種類の特定 GitLab
- ◆ Runnersの数の増減の可視化
- GitLabの実装アーキテクチャの可視化、インフラのスケールアップ/ダウン
- 社内のチャンピオンと利用や採用パターンについて話し合い、新規ユーザーの導入やアカウントの全体的な価値実現の ために、どこにエネルギーを投入すべきか戦略を練る

将来的には、以下のようなユースケースへの対応を予定しています。

- 顧客と最終的なチームによるステージの採用をROIと成熟度で追跡する。これは、カスタマーサクセスプランの目標に対する進捗状況の把握にもなります。
- ◆ 活用度、顧客ライフサイクルのフェーズ、エンゲージメント (=Time-to-Value) に基づいて、デジタルエンゲージメント (=アクションやコンテンツ) を促進する。

## 製品使用データのソース

## ソースキーユースケースソースからEDWへのデータフロー

#### ソースキーユースケースソースからEDWへのデータフロー

使用状況Ping	Gainsight 製品使用率、xMAU、 推定MAU	セルフマネージドインスタンスから送信されたJSONペイロード -> version.gitlab.com -> VersionPostgresDatabase <- pgp -> snowflake.raw.version_db
シートリンク	Gainsight製品使用状況	顧客ポータル→顧客Postgresデータベース <- pgp -> snowflake.raw.tap_postgres.customers_db_license_seat_links
バージ ョンチ	Noneversion -> VersionPostgresDataba	ase <- pgp -> snowflake.raw.version_db.version_checks
ェック GitLab.com	ProductAdoptionDashboardです。 Gainsight製品の使い方(近日公開	gitlab.com -> replicas/clones <- pgp -> snowflake.raw.tap_postgres

## セルフサービスの機能

このデータソリューションは、3つのSelf-Service Data機能を提供します。

- 1. **Gainsightユーザーの皆様**ヘセルフマネージド製品の使用状況データがGainsightで利用できるようになり、Gainsightユーザーは特定のワークフローの作成、トレンドの可視化、顧客の健全性スコアカードの作成、ユースケースの採用戦略の検討などが可能になりました。Using Product Usage Data in Gainsight a full guide(製品使用状況データをGainsightで使用する方法)をご覧ください。
- 2. **ダッシュボード開発者**。新しいダッシュボードを構築し、既存のダッシュボードを顧客製品採用データにリンクさせる ための完全な次元モデルコンポーネントを含む新しい**Sisense**データモデルです。
- 3. **SQL開発者**。A EnterpriseDimensionalModel」をご覧ください。R2A Objects」タブを参照してください。

データプラットフォームコンポーネント

データプラットフォーム技術の観点からも、このソリューションは有効です。

- 1. Gainsight Data Pump EDWからGainsight、GainsightからEDWへ
- 2. 製品使用状況データのためのEnterprise Dimensional Modelの拡張
- 3. Data Pipeline Healthダッシュボードのテストとデータ検証の拡張機能
- 4. ERD、dbtモデル、および関連するプラットフォームコンポーネント

クイックリンク

Gainsightの製品使用状況データ カスタマーサクセスでのGainsightの活用

WIPお客様の製品採用ダッシュボード

WIP製品使用データ - 知識評価

Sisense Discoveryを使い始める

セルフサービス・ウォークスルー・ビデオ

## データセキュリティの分類

Product Usage Dataに含まれるデータやそれを支えるデータの多くは、オレンジまたはイエローです。これには、オレンジ色のアカウントの顧客メタデータ、SalesforceやZuoraのコンタクトデータ、GitLabの非公開財務情報などが含まれますが、これらはすべて公開すべきものではありません。このダッシュボードからデータを共有する際には、その詳細が GitLab の組織内に留まるようにし、外部に共有する際には適切な承認を得られるように注意しなければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保存しないよう、常に注意を払う必要があります。このデータはSnowflakeとSisenseに残すべきであり、特に承認されない限り、これらのアプリケーションを通じてのみ共有されることが理想的です。

## ORANGE

- 説明行またはレコードレベルの顧客および個人データ。オブジェ
- クトです。
  - o dim\_crm\_account
  - dim\_billing\_account
  - o dim\_ip\_to\_geo
  - o dim\_location

#### **YELLOW**

- 説明します。GitLab 財務データ。集約や合計を含む。オブジェクトです。
- o dim subscriptions
  - o prep\_recurring\_charge

### ソリューション・オーナーシップ

- ソース システム オーナー:
  - 使用法 Ping@jfarris
  - セールスフォース@jbrennan1
- ソースシステムのサブジェクト・マタ
  - o ー・エキスパート: 使い方
  - Ping@jfarris
  - Gainsight:jbeaumont セールス フォースJBRENAN1
- データチームのサブジェクト・マター・エキスパートrparker @snalamaru

### 主な用語

- 1. お客様
- 2. 使用方法 Ping
- 3. GitLabセルフマネジメント・サブスクリプション
- 4. GitLab SaaSサブスクリプション
- 5. シートリンク
- 6. 製品カテゴリー、製品階層、配送
- 7. バージョンチェック
- 8. ビラブルメンバー。API, 定義, EDM フィールド名: billable\_user\_count
- 9. アクティブなユーザーCustomerDocs, MetricDictionary, EDM フィールド名: active\_user\_count

## キーメトリクス、KPI、PI

- 1. データマート メトリクスの定義
- 2. イベントベースのメトリクス
- 3. ユーザーベースのメトリクス
- 4. ステージおよびグループのパフォーマンス指標
- 5. ARR

## メトリックフォーマット

- 1. ユーザーベースのメトリクスアクション/イベントを実行した
  - ユーザーの数 過去28日間
  - 過去7日間
    - 例"# of users who completed a merge request in last 28 days"
- 2. イベントベースのメトリクス。実行されたアクション]の数
- 3. イベント/属性/ユーザー]のカウント数を合計します。
  - 1. 例"# of runners" または "# of auto\_devops\_enabled projects"
- 4. 指標となるメトリクス。ある属性が真か偽か
  - 1. 例"シェアードランナーが有効かどうか"
- 5. **パワーユーザーのメトリクス**。このアクション、このアクション、そしてこのアクションを行ったユーザーの数

## セルフサービス・データ・ソリューション

セルフサービス・ダッシュボード開発者

Sisenseでのチャート構築を始めるには、Sisenseの10分間のDataOnboardingVideoを見るのが良いでしょう。ダッシュボードを作成した後は、そのダッシュボードを簡単に見つけられるようにしたいものです。トピックは、ダッシュボードを1つの場所に整理し、簡単に見つけることができる素晴らしい方法です。ダッシュボードの右上にあるトピックへの追加アイコンをクリックすると、トピックを追加できます。ダッシュボードは、関連性のある複数のトピックに追加することができます。トピックには、Finance(財務)、Marketing(マーケティング)、Sales(販売)、Product(製品)、Engineering(エンジニアリング)、Growth(成長)などがあります。

セルフサービスのSOL開発者

## キーフィールドとビジネスロジック

- GitLab SaaSとGitLab Self-Managedの顧客デプロイメントから得られた製品使用状況データは、定期的にエンタープライズデータウェアハウスに投入され、GainsightとSalesforceで利用されます。
- 私たちはUsage Pingを利用して、自己管理された顧客の利用データを導き出します。自己管理された顧客の製品利用データは、 自己完結型のUsage Pingパケットに大部分が含まれています。
- ◆ SaaS顧客製品利用データは、ソースデータベースのテーブルを使用して再構築されます。
- シートリンクデータは、自営かSaaSかに関わらず、すべてのお客様のライセンス使用率データを含みます。集約
- ◆ されたメトリクスは7日間と28日間の期間で収集され、Usage Ping ペイロードのトップレベルキーである counts\_weekly と counts\_monthly のサブキーに aggregated\_metrics が追加されます。
- ◆ すべてのタイムフレームの集約されたメトリクスは、名前に aggregate\_ のプレフィックスが追加された count トップレベ
- ◆ ルキーに存在します。Gainsightの消費量の基礎となるテーブルは、Self-Managedの料金プランに関連付けられたすべての Zuoraサブスクリプションのセットに基づいて構築されています。このテーブルに含まれる事実のセットを構築するために 、Customers DBからのSeat Linkデータが優先度の高いUsage Pingメトリクスと組み合わされています。各月の最新の受信 および最新のUsage Ping(指定されたサブスクリプションIDのcreated\_at dateによる)およびSeat Link(dim\_subscription\_idによる)ペイロードが報告されます。

## エンティティ・リレーションシップ・ダイアグラム

## Diagram/EntityGrainPurposeKeywords

製品使用デー タ**ERD**  以下 のす べて カラム名、カラムのデータタイプ、カラムの制約、主キ 一、外部キー、テーブル間の関係など、すべてのテーブ

ル構造を表示します。

顧客、使用状況Ping、サブスクリプション、シートリンク、セルフマネージド、SaaS、プロダクト、デリバリー、アカウント

リファレンスSQL

スニペットライブ ラリ説明

データプラットフォーム ソリューション

Gainsightデータポンプ

データは、GitLab SaaSおよびGitLab Self-Managedの顧客デプロイメントから得られたものです。

## EDWからGainsight Data Pumpへ。

データチームはGainsightのネイティブ機能を利用してSnowflake Enterprise Data Warehouseからデータを読み込みました。データチームは、Gainsightがアクセスできる読み取り専用のマートレベルのテーブルを構築し、そこには現在利用可能なすべてのデータが含まれています。時間の経過とともにデータチームがメトリクスや顧客セグメントを追加すると、このテーブルは自動的に追加データで更新されます。この「インターフェース」をGainsight Data Pumpと呼びます。

#### GainsightからSnowflake Data Pipelineへ。

データチームは、GainsightからSnowflakeへの新しいソースデータパイプラインを開発し、新しいカスタムオブジェクトやGainsightで作成されたデータを含めることで、使用pingの一致率を高めるなどの改善を行いました。

ProductUsagedatadevelopmentementalStreams」の図は、Self-ManagedおよびSaaS Product UsageをGainsightに提供するための開発アプローチを示しています。

データの系統

dbtソリューションは、Source ModelsからGainsightまでの各テーブルを介したデータの流れを表す次元モデルを生成します。

• dbtモデルの完全なデータ系統は、DataFlowDiagramに掲載されています。

## トラステッドデータソリューション

TrustedDataFrameworkの概要を見る

詳細なテストの実施方法については、dbtガイドの例を参照してください。

製品使用状況 信頼性のあるデータ ダッシュボード

dbtテスト、ソースモデルの新しさ、レコード数、最終実行日、ゴールデンレコードの検証などを表示する詳細なダッシュボードです。 最新のEnterprise Dimensional Modelのテストと実行に関するレポートです。

(WIP) 製品使用状況 信頼性の高いデータダッシュボード

RAWソースデータパイプラインの検証

データパイプライン・ヘルスバリデーション

layout: handbook-page-toc title: "SaaS製品イベントデータ" 説明"製品利用データは、有料・無料ユーザーがGitLabを製品としてどのように、いつ、どこで利用しているかを定量的に測定するものです。"

{::options parse\_block\_html="true" /}<sub>o</sub>

このページについて

{:.no\_toc .hidden-md .hidden-lg}<sub>o</sub>

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}.

## SaaS製品のイベントデータ

**SaaS製品イベント・データ・パイプライン**は、ビジネス・ユーザー(プロダクト・アナリスト、成長担当PM、コア・プロダクトPMなど)に、当社のSaaS製品(Gitlab.com)のユーザーと名前空間の活動を記録したファクト・テーブルとマート・テーブルのセットを提供します。

このテーブルの汎用性により、ビジネス関係者は以下の例のように様々な分析を行うことができます。

## 実行可能な分析の種類

- 1. ユーザージャーニー分析。同じネームスペースで、異なる製品ステージがどれだけ頻繁に使用されているか。どのステージが組み合わせて使用されているかを確認します。
- 2. ネームスペース・ステージの採用。新しいネームスペースが使用開始後数日で「Create」や「Verify」などのステージをどの程度採用しているかを評価する。
- 3. 組織ごとのステージ。ネームスペースが最初の数日間でどのようにステージを採用し、それが有料コンバージョンや長期的なエンゲージメントとどのように相関しているかを確認します。

ERDとデータフロー、テーブルの記述

データチームは、これらのファクトテーブルを管理しており、これを利用して、:

- fct\_event\_400: イベントレベルの粒度でのファクトテーブル。つまり、このテーブルではイベントごとに1つの行が見つかります。イベントとは、例えば、issueの作成やci\_pipelineの開始などです。fct\_daily\_event\_400: 近日中に実装予定の
- テーブルで、進捗状況はこのイシューで確認できます。
- fct\_daily\_xmau\_400: 近日中に実装される予定のテーブル、進捗状況はこの課題で確認できます。

これらのモデルは、prep\_eventという名前のdbtモデルを介して作成されます。

prep\_eventモデルでは、月ごとに分割されたテーブルを作成します。Snowflakeを見ると、月ごとに1つのスキーマがあります。

つまり、1ヶ月に1つのprep\_eventテーブルがあるということです。dbtでは、スキーマの定義がdbt\_project.yml.

クエリの中には、データを月ごとのサブセットに限定するWHERE文もあります。以下は、

prep eventテーブルのERDです。

これらのテーブルで現在追跡されているイベントのリストは、ここで入手できます。いくつかの新しいイベントは非常に簡単に追加できます。

## このテーブルに新しいイベントを追加するには?

テーブル全体の複雑さはprep eventテーブルに含まれる(dbtドキュメントはこちら

このモデルは少し長く、多くの異なるテーブルからのデータを集約/統合しています。とはいえ、新しいイベントを追加するのは、割と簡単な操作です。以下にその手順を示します。

#### Gitlabデータチームのワークフロー

すべては課題から始まるので、まず最初のステップとして、Data Projectに課題とMRを開設します。MRで作成したブランチに入ったら、以下の手順で作業を進めてください。

#### データを格納するCTEの追加

テーブルに含めたいデータが準備中のテーブルに(フィルタリングの必要なく)完全に取り込まれている場合、これはかなり簡単に実行できるアクションです。ここで示したシンプルなCTEマクロで、新しいタプルを追加するだけでよいのです。

```
{{ simple_cte([
    ('prep_ci_pipeline', 'prep_ci_pipeline'),
    ('prep_action', 'prep_action'),
    ('prep_ci_build', 'prep_ci_build'),
    ('prep_deployment', 'prep_deployment'),
```

```
('prep_issue', 'prep_issue'),
  ('prep_merge_request', 'prep_merge_request'),
  ('prep_note', 'prep_note'),
  ('dim_project', 'dim_project'),
  ('dim_namespace', 'dim_namespace'),
  ('prep_user', 'prep_user')
])}}
```

取り込みたいデータが既存のテーブルのサブセットである場合は、操作が少し複雑になります。

1. まず、上記のようにsimple\_cteマクロに、キャプチャしたいイベントに対応するテーブルを追加します。
2. 次に、最初のアクションで呼び出したテーブルの行をフィルタリングするための別のCTEを以下のように作成します。

例えば、成功したCIパイプラインだけを選択したいとします。このイベントはprep\_ci\_buildで取得されています。 テーブルで、failure reasonカラムでフィルタリングする必要があります。

まず、prep\_ci\_build テーブルを simple\_cte マクロに追加します(まだ追加されていない場合)。次に、

failure reason列でフィルタリングする新しいCTEを追加する必要があります。

```
successful_ci_pipelines AS (
    SELECT * )
    FROM prep_ci_pipeline
    WHERE failure_reason IS NULL
)
```

## JSONにエントリを追加する event\_ctes

JSONはこのようになります。

```
{
    「event_name": イベントの名前。issue_creationなどのアクション名をつけることをお勧めします。
    "source_cte_name": 必要なデータを含むCTE名で、先に定義したものです。 "user_column_name": ユーザーID を含むカラム。
この情報を得るためには、dbtを参照する必要があります。
    「project_column_name": プロジェクトのIDを含むカラム、この情報を得るためにはdbtを参照する必要があります。
    "primary_key": CTEのプライマリキー
}
```

#### MRのテスト

変更をテストする最も簡単な方法は、CIパイプラインの機能を使って行うことです。最初のステップは、本番のデータベースをクローンすることです(以下の例)。これには10分から15分かかります。

そして、以下の小さなビデオに示すように、dbtジョブを実行す

る必要があります。 Dgif

dbt job specify\_1\_modelを®をクリックして実行し、以下の変数を追加する。

- ◆ キーを使用しています。DBT\_MODELS
- ◆ 値:プリップ\_イベント



ジョブが無事に実行されたら、R&Dチームのメンバーにレビューを依頼します。

MRに関する質問は、Slackの#dataチャンネルやMerge Requestを通じて、データチームのメンバーに遠慮なく質問してください。

#### データの埋め戻し

データが統合されたら、prep\_eventデータをバックフィルする必要があります。そのためには、データエンジニアに t\_prep\_dotcom\_usage\_events\_backfillというAirflow Dagを起動してもらい、過去3年分のデータのバックフィルを実行してください。

## layout: handbook-page-toc title: "SaaS Service Ping Automation" $\subset$

のページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

## サービスPingの概要

以前はUsage Pingと呼ばれていたサービスPingは、GitLabインスタンスで毎週実行されるバックグラウンドプロセスで、分析目的に有用な一連のメトリクスを収集、集約、パッケージ化する役割を担っています。メトリクスは簡単に拡張することができ、新しいメトリクスはメトリクス・ライフサイクルに従って定期的に変更されます。メトリクスの完全なセットは、Metric Dictionaryで定義されています。インスタンスからメトリクスが収集されると、それらはJSONペイロード(「ping」)にまとめられ、GitLab Versions Appに投稿され、Snowflakeなどのダウンストリームプロセスに同期されます。以下は、ServicePingペイロードの例です。サービスPingには、セルフマネージドサービスPingとSaaSサービスPingという2つの主要なバリエーションがあります。セルフマネージドサービスPingはシングルテナント型のGitLabに対して実行され、SaaSサービスPingはマルチテナント型のGitLab.comに対して実行されます。

## サービスPingの使用例

サービスPingのデータは、GitLabがどのように使われているかをプロダクト、サポート、セールスの各チームが理解するためのインサイトを提供します。例えば、このデータは以下のことに役立ちます。

- 1. サポート GitLabxMAUKPIs KPI
- 2. ステージの月間アクティブユーザー数 (SMAU) を算出することで、ステージや機能の成功を測ることができます。
- 3. どの機能が使われているのか、あるいは使われていないのかを理解し、顧客がGitLabの豊富な機能を活用できるようにガイダンスを提供する。
- **4.** 月ごと(または週ごと)にカウントを比較し、インスタンスが異なる製品機能をどのように使用しているかを大まかに把握することができます。
- 5. GitLabインストールの分類と理解に役立つその他の事実の収集

### セルフマネージドサービスPing

セルフマネージドのお客様は、サービスPingをセットアップして実行し、GitLabの独自のデプロイメント(インスタンス)の分析を行います。お客様はオプションでサービスPingを無効にすることができ、その場合はメトリクスがGitLabに送信されません。お客様は、コンソールアプリケーションを通じてサービスPingのデータにアクセスできます。

## SaaSサービスのPing

GitLab.com (またはGitLab SaaS) は、基本的にGitLabがホストするマルチテナント版のGitLabです。手動で生成されたサービスPing (Manual SaaS Service Ping) がSaaS用に実装されており、Self-Managedインスタンスで実現しているものと同等の分析カバーをSaaSでも提供しています。

しかし、現在のプロセスには、大きく分けて2つの問題があります。

## パフォーマンスの問題

- プロセスにエラーが発生しやすく、お客様のライブ活動とリソースの競合が発生
  - 59/ 243

merged.md

2021/11/8

する。プロセスに時間がかかり、個々の指標の問い合わせが定期的に失敗する。

- プロセスはオフピーク時に手動で実行されるように実装されてい
- るプロセスはエンドツーエンドで管理するためのスタッフを必要とする

#### データカバレッジの不足

また、手動のSaaSサービスPingは、インスタンスレベル(サイト全体)のデータしか生成できないため、セールスやカスタマーサクセスなど、GitLab.comのお客様の個々の採用状況を測定するためにネームスペースレベルのより詳細なデータを必要とする人たちのニーズをすべて満たすものではありません。

この2つの主要な問題を解決するために、データチームは「Automated SaaS Service Ping」を開発しています。Automated SaaS Service Pingは、ビッグデータ、自動化、スケールのために設計されたシステムであるEnterprise Data Platformでネイティブに実行されるプログラムのセットです。Automated SaaS Service Pingが完全に運用されると、Manual SaaS Service Pingは廃止されます。Automated SaaS Service Pingを構成する主なサブプログラムは2つあります。

- SaaS Instance Service Ping GitLab.comインスタンスのサービスPingを毎週自動生成。
- SaaS Namespace Service Ping すべてのGitLab.comインスタンス->名前空間のサービスPingを毎週自動生成。

## 4種類のサービス pingプロセスの実行環境 3種類の環境

サービスPingには、運用中のものと開発中のものがあり、全部で4種類あります。

	1. セルフマネージドサー ビスPing	2. マニュアルSaaSサー ビスPing	3. Automated SaaS Instance Service Ping	4. 自 動 化 さ れ た SaaS名前空間サービ スのPing
Where Run	環境1:お客様のセルフマ ネージドインスタンスの 場合	環境2:GitLab.comのイ ンフラ内	環境3:データプラッ トフォーム基盤	環境3:データプラットフォ ーム基盤
Run リジェンス	ByGitLab (Automatically)	プロダクト・インテ	エアフロー(自動)	エアフロー(自動)
		(手動)		
				頻度週刊
誌週刊誌週刊誌コード		オーナープロダクト	インテリジェンスプログ	ダクトインテ
リジェンスデ	ータチームデータチーム			
	ソースコードRuby,	SQLRuby,	SQLPython, SQL,	dbtPython, SQL, dbt
データ				InstanceInstanceNamespace
Granularity				リティ・インスタンス・イン
グラニュラ スタンス・ネ	ニームスペース			
開発			## N. ### N.	HH 2/V
ステータス			実写化実写化	開発中

## (自動) SaaSサービスのPing実施

実写化実写化

## プロセス概要

(自動化された) SaaSサービスPingは、AirflowでオーケストレーションされたPythonプログラムとdbtプロセスのコレクションで、Enterprise Data Platform内で毎週実行されるようにスケジュールされています。自動化されたSaaSサービスPingプロジェクトは、すべてのソースコードと構成ファイルを保存します。このプログラムは、2つの主要なデータソースに依存しています。RedisベースのカウンターとSQLベースのpostgresテーブルです。この2つのデータソースは、Snowflakeの自動データパイプラインとして実装されており、SaaS Service Pingの実装プロセスとは独立して実行されるようになっています。

- SaaSからのSQLベースのpostgresデータは、pgp経由で同期され、RAW.SAAS\_USAGE\_PINGスキーマで利用できるよう
- ◆ になります。 Redisデータはプログラムの実行時にアクセスされ、RAW.SAAS\_USAGE\_PINGスキーマにも保存されます

0

自動化されたSaaSサービスPingは、2つの主要なデータ処理フェーズで構成されています。

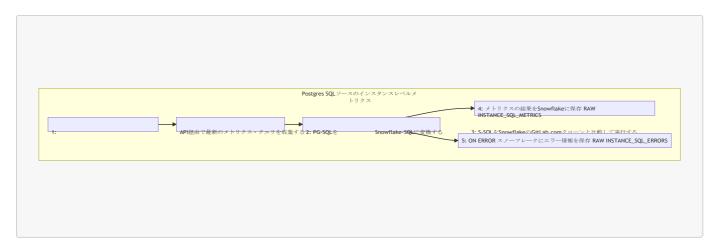
1. フェーズ1では、MetricDictionaryで定義されているメトリクスの収集とジェンター化を行います。 2. フェーズ2では、メトリクスをTrusted Data Model (FCTとDIMテーブル) 形式に変換します。

フェーズ1:測定基準の収集と生成

### SaaS Instance Service Ping

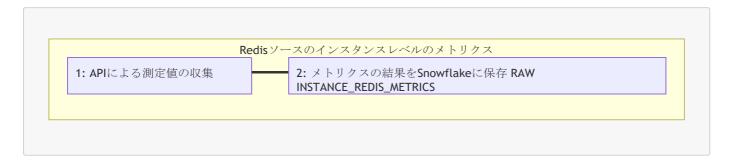
SaaS Instance Service Pingは、「プロセスの概要」で説明したとおりに実行されます。

## インスタンスのSQLベースのデータフロー



データフローの詳細については、ServicepingeのREADME.mdファイルをご覧ください。

#### Redisベースのデータフロー



## SaaS 名前空間サービス Ping

SaaS Namespace Service Ping は、SaaS Instance Service Ping よりも細かいレベルのメトリクスを生成します。このプロセスでは、GitLab.comのすべてのネームスペースのリストにアクセスし、各ネームスペースをループして最後親のネームスペースレベルの使用率メトリクスを生成します。名前空間テーブルはプログラムへの入力となり、効率化のために、従来の1対1の名前空間ループではなく、SQL SET OPERATIONとSQL GROUP BY namespace-idが使用されています。最終的なメトリクスの出力は、最終的な親のネームスペース・レベルに格納されます。名前空間サービスPingの欠点は、現在のところSQLソースのメトリクスしか利用できず、analytics unique visits.g analytics contributionのようなRedisソースのメトリクスは現在利用できないことです。



## メトリクス収集・生成プロセス 擬似コード

- 1. GitLab.comのPostgresソースデータのパイプラインが稼働しており、Snowflakeでは新鮮な最新のデータがそれぞれ、RAW.SAAS\_USAGE\_PINGスキーマとPREP.SAAS\_USAGE\_PINGスキーマです。
- 2. サービスPingのパイソンプログラムの開始
  - 1. **SQLベースの**メトリクス生成の開始
    - 1. PostgresのSQLソースのメトリクスを開始

- 1. 最新のPostgres SQL-sourced (PG-SQL) メトリクスクエリのセットをMetrics Dictionary API Query Endpointから取得します。
- 2. pythonトランスフォーマーを使用して、インスタンスレベルのPG-SQLをSnowflake SQL (S-SQL) に変換します。
- 3. Snowflakeデータウェアハウスで利用可能なSaaS GitLab.comクローンデータに対してS-SQLを実行し、 その結果をRAW.SAAS\_USAGE\_PING.INSTANCE\_SQL\_METRICSに格納します。エラーが発生した場合、データはRAW.SAAS\_USAGE\_PING.INSTANCE\_SQL\_ERRORSテーブルに格納されます。
- 2. Redisベースのメトリクスの開始
  - 1. Redis APIの呼び出し
  - 2. データはJSON形式でピックアップして保存されます。おおよそのサイズは2k行程度で、通常は1回のロードで1つのファイルになります(現時点では週1回のロードです)。データは

RAW.SAAS\_USAGE\_PING.INSTANCE\_REDIS\_METRICSに保存されます。

- 2. 名前空間レベルのメトリクス生成開始
  - 1. NamespaceQueriesJSONから最新のメトリクスクエリを取得
  - 2. Snowflakeデータウェアハウスで利用可能なSaaS GitLab.comクローンデータに対して名前空間クエリを実行し、結果をRAW.SAAS\_USAGE\_PING.GITLAB\_DOTCOM\_NAMESPACEに格納する。

フェーズ2:メトリクスの信頼性の高いデータモデルへの変換

すべてのソースメトリクスがSnowflake RAW. SAAS\_USAGE\_PINGスキーマで利用可能になると、データをTrusted Data Modelフォーマットに変換するためのdbt処理を開始します。

- **SQL**ベースのメトリクスのdbt処理
- ◆ Redisベースのメトリクスのdbt処理
- 名前空間レベルのdbt処理

## 既知の制限事項/改善事項

- ◆ 名前空間レベルのRedisソース・メトリクスはまだ利用できません。
- Snowflakeには冗長な "レガシー"サービスpingプロセスがあり、これを廃止する必要があります。

サービスPingメトリクスの種類

Service Pingでは、主に2種類のメトリクスをサポートしていま

- す。SQLメトリクス: Postgresテーブルから取得されたメ トリクス
- ◆ Redis metrics: Redisベースのカウンターから取得したメトリクス

#### SQL Metricsの実装

SQLベースのメトリクスのワークフローは、最も複雑なフローです。SQLベースのメトリクスは、実際には、インスタンスの Postgres SQLデータベースに対して実行されるSQLクエリによって作成されます。大規模なテーブルの場合、これらのクエリの実行には非常に時間がかかります。例えば、Counts.ci\_builds メトリックは、最大級の ci\_builds に対して COUNT(\*) を実行しています (10 億行以上の dbt テーブルを参照)。このモジュールの目的は、SaaSインスタンスのpostgres SQLデータベースの代わりに、Snowflakeデータベースに対して実行することです。

プロダクト・インテリジェンス・チームは、メトリクスを計算するために実行するすべてのSQLクエリを取得できるAPIエンドポイントを作成しました。以下はそのファイル例です。

APIエンドポイントに関するテクニカルドキュメントはこちらか

らどうぞ。 それでは、JSONレスポンスで受け取ったいくつかの

```
"counts "です。{
   "assignee_lists"。 "SELECT COUNT(\sqrt{-}-\sqrt{-}") FROM \sqrt{-}" WHERE \sqrt{-}-\sqrt{-}"
3",
  "boards""SELECT COUNT(~~-- ~~) FROM ~~~~"
  = 'Ci::Build'"となります。
   "ci internal pipelines"。"SELECT COUNT(ci pipelines\.id\) FROM `ci pipelines\ WHERE
(ci_pipelines\.source\ IN (1, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13) OR"
i.s.....NULL)
  "ci_external_pipelines"。"SELECT COUNT(CI_PIPELINES\.ID\) FROM 'CI_PIPELINES\ WHERE"
"ci_pipeline_config_auto_devops "です。"SELECT COUNT(ci_pipelines\.id\) FROM ci_pipelines\ WHERE
ci_pipelines\.config_source\ = 2".
   "ci_pipeline_config_repository "です。"SELECT COUNT(ci_pipelines\.id\) FROM `ci_pipelines\ WHERE `
ci_pipelines\.config_source\ = 1 "です。
  "ci_runners":"SELECT
```

\"ci\_triggers\"",

そのためには、Snowflakeのテーブル (GitLab Saasから同期されたもの) に対して実行できるようにすることが目標です。その

merged.md

2021/11/8

ためには、PostgresのSQLテーブルと同じカラム名、同じ粒度のテーブルを用意する必要があります。

現在、SnowflakeでPostgresのデータを変換している様子は以下の通りです。



重複を避け、妥協のない正確なデータを提供するためには、デデュープ・レイヤーが有効です。

次に、SQLベースのメトリクス・クエリを実行できるテーブルを特定しました。次に、これらのテーブルに対してクエリを実行するために、SQLステートメントを変換する必要があります。

このSOL文を変換するスクリプトを実行しています。

をこのSQL文に追加します。

そして、これらのクエリをすべて実行し、その結果をjsonに格納して、それをテーブルのraw.saas\_usage\_ping.instance\_sql\_metrics.このテーブルには以下の列があります。

- ◆ query\_map: すべてのクエリを格納する
- ◆ run\_results: 結果を格納する ping\_date:
- クエリが実行された日付 run id: 処理の一
- ◆ 意な識別子

エラーが発生した場合、データはRAW.SAAS\_USAGE\_PING.INSTANCE\_SQL\_ERRORSテーブルに保存されます。このテーブルは、データをjsonデータタイプで保持します。このテーブルの上には、エラーをTrusted Data Modelに取り込み、SQL処理の不具合を簡単に可視化する仕組みが構築されています。

### Redisメトリクスの実装

Redis カウンターは、GitLab で起きている任意の状況のうち、データベースに永久的な記録が残らない高頻度の発生を記録するために使われます。このような場合、バックエンドエンジニアは、例えば navigation\_sidebar\_opened のように状況を表す名前を定義し、それが発生する瞬間を(既存の実行経路に専用のコードを追加することで)任意に決定します。

プロダクト・インテリジェンス・チームは、データ・チームがいつでも好きな時にすべてのRedisメトリクス値を取得できるAPIエンドポイントを作成しました。JSON レスポンスの例は UsageDataNonSqlMetricsAPI のページにあります。なお、-3 はメトリクスが Redis でないことを意味しており、API はその値を取得しません。JSONレスポンスを受信すると、Snowflakeの RAW.SAAS\_USAGE\_PING.INSTANCE\_REDIS\_METRICSテーブルに保存されます。APIエンドポイントに関する追加の技術文書はこちらを ご覧ください。このテーブルには、以下の列があります。

- jsontext: 実行されたすべてのクエリを
- 格納するもの ping\_date: クエリが実行
- ◆ された日付 run\_id: 処理の一意な識別子

エアフローの設定

毎週土曜日に実行されるAirflow dag saas-instance-usage-pingを作成し、以下のすべての操作を実行しました。API Endpoint

◆ からのクエリの取得

- クエリを変換してSnowflakeのデデュープ層で実行できるようにする クエリの
- 実行
- ◆ 結果をSnowflakeに保存

## RAWからPRODデータベースとSisenseへ

現在は、RAWで保存されたデータの変換を限定的に行っています。将来的には、このようなデータフローになる予定です。

つまり、作成されたデータセットは、モデルprep\_usage\_data\_flattenedで現在のデータパイプラインとUNIONされます。

## layout: handbook-page-toc title: "Sales Funnel" $\subset$

のページについて

TOC {:toc}

{::options parse\_block\_html="true" /}。

## セールスファネル

セールスファネルは、ギットラボのLead to Cashプロセスの中核をなすもので、潜在顧客からの問い合わせを、獲得した機会のクローズまで追跡します。このプロセスでは、顧客はリードから、マーケティング適格リード、営業許可されたオポチュニティ、そして最終的にはクローズしたオポチュニティ、または失われたオポチュニティへと移行します。このハンドブックページでは、GitLabチームメンバーがセールスファネル全体を分析し、インサイトを生み出すための分析ツールやリソースを提供しています。

このページの目的は

- TD: Sales Funnel Targets vs. Actualダッシュボードとその派生ダッシュボードのナビゲーション方法の理解を助ける。TD:
- ◆ Sales Funnel Targets vs. Actuals Dashboardの作成に使用されたデータモデルの理解に役立ちます。
- GitLabでのあなたの役割に最も適した認定を受けることで、あなたの理解度を評価してください。
  - o ダッシュボードの使い方を詳しく知りたい方は、ダッシュボードのユーザー認定を受けてください。
  - o Sisenseダッシュボードの開発について詳しく知りたい方は、DashboardDeveloper認定を受けてください。
- ◆ そして、全体的にみんなが貢献できるように。

## Release Train Cadence:

- 1. 要件収集期限: Sales Key Meetingの1週間後。
- 2. UAT用にリリースされたのは、Sales Key Meetingの1週間前。

メンテナンスの予定です。

1.必要に応じて、2週間に1度、金曜日の午前9時から午前11時(米国東部時間)に定期メンテナンスを実施します。

クイックリンク

TD: Sales Funnel - Targets vs. Actuals ダッシュボード セールスファネルの次元モデルとダッシュボー ドのトレーニング

**TD**: Sales Funnel Management View

- Preloaded Cuts

TD: スタンダードブッキングのカット

Sisense Discoveryの使用開始について

## はじめに

そのためには、お客様に理解していただく必要があります。

- ◆ このダッシュボードでサポートされるKPI/PIは?
- ◆ メトリクスをどのように説明するかを説明するキーワー
- ▶ ドダッシュボードを支えるデータソース
- ◆ さらに探求するために、Sisenseでビジュアルや分析を自分で作成することができます。Sisenseディスカバリー・ツールを 使用するのが良い方法です。Sisenseを始めたい方はこちらをご覧ください。
- さらに深く掘り下げるために、snowflakeでデータを探索することができます。Snowflakeでの探索の利点は、追加の情報 (つまり他のデータソース) に結合できることです。Snowflakeでの探索に関する追加情報はこちらをご覧ください。

# キーターム、メトリクス、KPI/PI、そしてキーフィールドとビジネスロジック

#### ▶ 主な用語

- ◆ 注文タイプ
- **販売資格のあるソース。**機会がどのように作られたか**購入チャ**
- ◆ ネル。アカウントが購入した方法。GTM戦略
- フォーカスアカウント。アカウントベースドマーケティングチームが、セールスやフィールドマーケティングと密接に連携して、特にターゲットとするアカウントを表します。GTM戦略の価値である「アカウント・セントリック」、「アカウント・ベース・ネット・ニュー」、「アカウント・ベース・エクスパンション」のうち、「フォーカス・アカウント」に相当します。
- マーケティングチャネル (イニシャルソース
- セールスヒエラルキー・ライブ販売エリア > 販売地域 > 所在地地域 > 販売セグメント
- Sales Hierarchy Stamped:セールスエリア > セールスリージョン > ロケーションリー
- ▶ ジョン > セールスセグメント Key Metrics, KPIs, and PIs
  - Marketing Qualified Lead (MQL) の略。販売機会を得ることができる適格なリードの数。現在、MQLは一度だけですが、 将来的には、コンタクトを育成する中で、MQLは複数回認定される可能性があります。MQLとは、人口統計学的、企業 統計学的、および行動学的情報に基づいて、一定の基準に達したリードのことで、当社では累積100ポイントとしていま す。MQLスコア」は、様々な行動やプロフィールデータで構成されており、プラスまたはマイナスのポイント値で重み 付けされています。パーソンスコアが更新されるたびに、LeanDataはそのレコードをフローで処理する必要があるかど うかのチェックを実行します。
  - Sales Accepted Opportunity (SAO) の略。購入権限を持つ見込み客の担当者とのミーティングが設定されているオポチュニティ、またはGitLabがソリューションとなりうるイニシアチブ、またはユースケースや潜在的なシート数、今後60日以内の次のステップを含む適合性を持つオポチュニティです。この案件は、SDRが作成した営業担当者によって検証され
  - ◆ 、受理されています。**SAO/MQLです。**1ヶ月間にSAOに変換されたすべてのMQLのスナップショット変換率。1ヶ月間の SAO数を1ヶ月間のMQL数で割ったもの。
  - **新しいロゴ。**獲得したファーストオーダーアカウントの数。
  - Closed Won IACV (Net\_ARR):クローズドな獲得案件の予約収入。
  - クローズドウォンディール。収益に結びついた機会の数。Win Rate (勝率)。解決
  - に至った案件のうち、期間中に何件の案件を獲得したかを示す。ASP on New First
  - Order Deals:初回受注案件の平均案件サイズ。
  - IACV (Net\_ARR) Created:各日に作成されたパイプラインの値です。パイプラインは時間の経過とともに変化します。これは、作成日に基づく現在のパイプラインです。
  - SAOからClosed Won Dealまでのサイクルタイム: Closed Won案件の結論を得るまでの平均的な時間。
  - SAOからClosed Wonへのコンバージョン率。月間ベースでSAOがClosed Wonの案件に転換したコンバージョン率。
  - SAOから獲得済みクローズド案件への転換率: SAOから獲得済みクローズド案件への転換率。SAOの数: SAOがクローズド案件の結論に達するまでの時間を基準にしたSAOの数。
- ▶ キーフィールドとビジネスロジック

TD: Sales Funnel - Target vs. Actual Dashboard\*\*は、Sisenseの3つの主要なマートビュー、\*\*[mart\_crm\_opportunity]\*\*、
\*\*[mart\_crm\_person]\*\*、\*\*[mart\_sales\_funnel\_target]\*\*から生成されたレポートビューから作成されます。Sisenseのビューでは、シンプルなSQLの文字列を、必要なだけのチャートで参照することができます。これにより、エンドユーザーは新規ロゴ数や
SAO(Sales Accepted Opportunity)数などの計算の裏にあるロジックを理解することができ、また、どの結合が可能かを理解するための推測作業を省くことができます。

独自のダッシュボードを作成するには、Sisenseで利用できる以下のビューを参照します。

- [mart\_crm\_opportunity] 該当するファクトとディメンションを結合して、セールスファネルのオポチュニティセクションのビューを取得するマートテーブル。
- [mart\_crm\_person] 該当するファクトとディメンションを結合し、セールスファネルのマーケティング適格リードセクションのビューを得るためのマートテーブル。
- [mart\_sales\_funnel\_target] 該当するファクトとディメンションを結合して、セールスファネルターゲットのビューを取得するマートテーブル。
- [rpt\_crm\_opportunity\_closed\_period] [mart\_crm\_opportunity]ビューをクローズド期間にピボットし、有用な日付のアグリゲーションフィールドを生成するレポートです。
- [rpt\_crm\_opportunity\_accepted\_period] [mart\_crm\_opportunity]ビューを受理された期間にピボットし、有用な日付のアグリゲーションフィールドを生成するレポートです。
- [rpt\_crm\_opportunity\_created\_period] [mart\_crm\_opportunity]ビューを作成期間でピボットし、有用な日付のアグリゲーションフィールドを生成するレポートです。
- [rpt\_sales\_funnel\_target] [mart\_sales\_funnel\_target] ビューをターゲット期間にピボットし、有用な日付のアグリゲーションフィールドを生成するレポートです。
- [rpt\_crm\_person\_mql] MQLのカウントに使用される[mart\_crm\_person] WHERE is\_mql = 1をフィルタリングするレポート。

このビューでは、ERD(Entity Relationship Diagram)に見られるようなデータモデルを考慮しています。

- ▶ データリネージの
- ▶ クエリ例
- -- MQLsのカウント SELECT DATE\_TRUNC('month',mql\_date\_first) AS mql\_month, COUNT(mql\_date\_first\_id) AS actual\_mqls FROM "PROD". "COMMON\_MART\_MARKETING". "MART\_CRM\_PERSON" WHERE is\_mql = 1 GROUP BY 1 ORDER BY 1 DESC
- --SAOのカウント SELECT DATE\_TRUNC('month',sales\_accepted\_date) AS sao\_month, COUNT(\*) AS actual\_saos FROM "PROD". "COMMON\_MART\_SALES". "MART\_CRM\_OPPORTUNITY"WHERE is\_sao = TRUE GROUP BY 1 ORDER BY 1 DESC

```
</details> </details
<br>
<style> #headerformat {
background-color:5px; text-align: center;} #6666c4; color: black; padding:5px; text-align: center;}.
<h1 id="headerformat">セールスファネル標準レポート</h1>
<詳細
<summary markdown='span'>
 Sales Funnel Management View - Preloaded Cuts
</summary>
TD: Sales Funnel Management View - Preloaded Cuts]
(https://app.periscopedata.com/app/gitlab/828239/TD:-Sales-Funnel-Management-View---Preloaded-Cuts) ダ
ッシュボードは、セールスファネルから主要な9つのKPIを追跡します。
1.ネットARR
1.新しいロゴマーク
1.パイプラインの作成
1.セールス・アクセプタンス (SAO
1.マーケティング・クオリファイド・リード (MQL) について
1.トライアル
1.平均販売価格(ASP
1.MOLからSAOへ
1. 勝率
KPIは、四半期ごとのペースでそれぞれの目標値と比較され、異なる次元で切り分けられます。これらの切り口により、
経営陣はビジネスが異なる次元でどのように実行されているか、高レベルの概要を素早く把握することができます。これ
は、ビジネスのどの分野が軌道に乗っていて、どの分野にもっと注意を払う必要があるかという疑問に答えるものです。
現在、販売管理のカットは以下のスニペットで動いています(使用方法はスニペットの説明欄に記載されています)。
1.[main_qtd_view_sales_funnel]
(https://app.periscopedata.com/app/gitlab/snippet/main_qtd_view_sales_funnel/55d49d17d8cf4cc6bf976c6
1da2d0017/edit)
1.[ratio_qtd_view_sales_funnel]
(https://app.periscopedata.com/app/gitlab/snippet/ratio_qtd_view_sales_funnel/55d5211e130f45a29a3a15
62307f95d9/edit)
販売管理pythonモジュール】(https://gitlab.com/gitlab-)
data/periscope/-tree/periscope/master/custom modules/sales_management) のスニペットと一緒に、BI層でレポ
ートを視覚化するために使用されます。
</details> </details
<詳細
<summary markdown=span>
 Standard Bookings Cuts
</summary>
```

TD: Standard Bookings Cuts](https://app.periscopedata.com/app/gitlab/831911/TD:-Standard- Bookings-Cuts) ダッシュボードでは、予約済みネットARRの前年比(Year over Year, YY)や、四半期および年間の財務目標に対するパフォーマンスをさまざまな角度から追跡します。

各カットは、「四半期ビュー」と「年度ビュー」で構成されています。

ダッシュボードを使用するには、フィルターで会計年度を選択し、また ドリルダウン "フィルタです。後者は、「Sales Segment Drilldown」フィールドの粒度を選択します。

全体 | 大規模-パブセックの内訳 | 大規模-地域の内訳 | 中間市場の内訳 | 大規模-地域の内訳 | --1 --| --| USイースト | USイースト USウェスト USウェスト I EMEA I EMEA APAC I APAC | 大きい 大きい | 大きい PubSec PubSec PubSec │ 大規模+PubSec | ファーストオーダー | チーム | テリトリー(500未 | | ミッド-マーケッ | ミッド-マーケ | ミッド-マーケッ | ミッド-マーケット | ット 1 SMB SMB SMB I SMB WW WW WW WW

#### キーフィールドとビジネスロジック

- \* [ATR (Available To Renew)](https://about.gitlab.com/handbook/sales/sales-term-glossary/#available-to-renew-atr)
- \* ATR:特定のカテゴリーの正味ARR/小計ATR。例:セグメントがSMB、成長タイプがContractionの場合、%ATRを計算するには、`SMBのNet ARR | Contraction` / `SMBのATR`となります。
- \* TRX: チャンスの数
- \* %MIX(ARR)。特定のカテゴリーの純ARR/小計ARR。例)セグメントがSMBで、Deal Sizeが5-25kの場合の %MIX(ARR)を計算する場合。SMBのネットARR | 5-25k` / SMBのARR`。
- \* 予約数の割合特定のカテゴリーの純ARR / 特定の四半期の総ARR
- \* 総予約数に対する割合特定のカテゴリーの純ARR / 特定のカテゴリーの総ARR 小計
- \* ProServ #: proserv\_amountが0と異なるオポチュニティの数
- \* [A R (Attach Rate)](https://about.gitlab.com/handbook/sales/performance-indicators/#proserve-deal-and-dollar-attach-Rate)
- \* 売上セグメントのドリルダウンフィールドに表示される「米国東」の行は、「東」と「東」の両方で構成されています。 LATAM」地域
- \* 販売セグメントのドリルダウン = APACは、「East」、「West」、「LATAM」、「EMEA」以外のすべての地域で構成されます。
- \* セールス・セグメント・ドリルダウン=SMBは、セグメントSMBに加えて、すべての 大規模 | と「中規模 | について
- \* Channel Type = Fulfilled」は、「Channel Type Fulfilled」に「NULL/Missing Channel Type」を加えたものです。

</details> </details

```
<style> #headerformat {
background-color:6666c4; color: black; padding:5px; text-align: center;
}
</style>
<h1 id="headerformat">追加のリソース</h1>。
```

<詳細

<summary markdown='span'>
トラステッドデータソリュー
ション

</summary>

セールスファネルモデルでは、Trusted Dataのテストとその結果に`sales\_funnel`タグを使用しています。これは [Trusted Data Dashboard](https://app.periscopedata.com/app/gitlab/756199/Trusted-Data-Dashboard)で 最も簡単に見ることができます。

概要は[Trusted Data Framework](https://about.gitlab.com/handbook/business-technology/data-team/platform/#tdf)をご覧ください。

[dbtガイドの例](https://about.gitlab.com/handbook/business-technology/data-team/platform/dbt-guide/#trusted-data-framework)を参照してください。

さらなるテストの実施に関する詳細と例 </details> </details

<詳細

<summary markdown='span'>

EDM エンタープライズ ディメンション モデル バリデーション

```
</summary>
(WIP) Enterprise Dimensional Model Validation Dashboard]
(https://app.periscopedata.com/app/gitlab/760445/WIP:-Enterprise-Dimensional-Model-Validation-
Dashboard)は、最新のEnterprise Dimensional Modelのテストと実行をレポートします。
</details> </details
<詳細
<summary markdown='span'>
   RAWソースデータパイプラインの検証
</summary>
[Data Pipeline Health Validations](https://app.periscopedata.com/app/gitlab/715938/Data-Pipeline-
Health-Dashboard)
</details> </details
<詳細
<summary markdown='span'>
   データセキュリティの分類
</summary>
セールスファネル分析のデータの多くは、[オレンジ] (/handbook/engineering/security/data-classification-
standard.html#orange)または[イエロー] (/handbook/engineering/security/data-classification-
standard.html#yellow)です。これには、ORANGEのアカウントの顧客メタデータや、GitLabのNon publicの財務情報
が含まれており、これらはすべて一般に公開されるべきではありません。このダッシュボードからデータを共有する際に
は、その詳細が GitLab の組織内に留まるようにし、外部に共有する場合には適切な承認を得られるように注意しなけ
ればなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを
保存しないよう、常に注意を払う必要があります。これらのデータは、[Snowflake](/handbook/business-
\texttt{technology/data-team/platform/\#data-warehouse)} \ \\ \succeq \texttt{[Sisense](/handbook/business-technology/data-team/platform/\#data-warehouse)} \\ \succeq \texttt{[Sisense](/handbook/business-technology/data-team/platform/#data-warehouse)} \\ \succeq \texttt{[Sisense](/handbook/business-technology/data-team/platform/#data-warehouse)} \\ \succeq \texttt{[Sisense](/handbook/business-technology/data-team/platform/#data-warehouse)} \\ \succeq \texttt{[Sisense](/handbook/business-team/platform/#data-team
team/platform/periscope/)に残すべきであり、特に承認されない限り、これらのアプリケーションを通じてのみ共有さ
れることが理想的です。
**ORANGE**
- 説明行またはレコードレベルの顧客および個人データ。
- オブジェクトです。
    - dim crm person
   - dim_crm_account`.
YELLOW** (イエロー
- 説明します。GitLabの財務データで、集約や合計を含みます。
- オブジェクトです。
   - `fct crm person`
    - `fct_crm_opportunity` の略。
</details> </details
く詳細
<summary markdown='span'>
    ソリューションの所有権
</summary>
* ソース・システム・オーナー
    * セールスフォース: `@jbrennan1`.
* ソースシステムのサブジェクト・マター・エキスパート。
    * セールスフォース: `@jbrennan1`.
* データチーム・サブジェクト・マター・エキスパート:`@paul armstrong` `@jeanpeguero` `@jjstark` `@iweeks`
</details> </details
レイアウト:ハンドブック-ページ-トック
のタイトルである。 "Trusted Data Solutions Dashboard"
## On this page
{:.no_toc}
```

- TOC

```
{:toc}
Coming Soon!
layout: handbook-page-toc
title:"ウェブトラフィック分析"
## On this page
{:.no_toc}
- TOC
{:toc}
# # Google Analyticsによるウェブトラフィック解析
ウェブトラフィック分析は、訪問者やウェブサイト上での行動についての洞察を得て、より良い結果をもたらすのに役
立ちます。BigQueryとGoogle Analyticsは、ウェブサイトを改善するために必要な情報を提供し、ウェブサイトを最
高のものにしてくれます。
**以下のサブドメインが含まれます:**。
- about.gitlab.com _(マーケティング分析のためのプライマリ)_。
- docs.gitlab.com
- forum.gitlab.com
- customers.gitlab.com
- learn.gitlab.com
### Googleアナリティクスデータの解析になぜBigQueryデータが使われるのか?
- サンプリングデータがなくなる
- ディメンション数は無制限
- 異なるスコープを1つのレポートにまとめる_(例:セッションとヒット数)_。
- カスタムチャンネルグルーピングの構築とデータエラーの修正
- Google Analyticsのデータをサードパーティのデータソースと組み合わせる (例: Salesforceのトライアルキ
ャンペーンメンバーとトライアルページビュー) (例: Salesforceのトライアルキャンペーンメンバーとトライア
ルページビュー)
### Quick Links
- Googleアナリティクスの基礎を学ぶ。[アナリティクス アカデミー オンライン講座
- BigQueryデータのフォーマットとスキーマを理解する。[BigQuery Export Schema -
Columns] (https://support.google.com/analytics/answer/3437719?hl=en)を参照してくだ
- GitLab公開サイトへのGoogleアナリティクス導入情報 [Marketing Web Analytics - Inbound
Mktg](https://about.gitlab.com/handbook/marketing/growth-marketing/inbound-
marketing/analytics/#dimensions-vs-metrics)
- BigQueryデータを利用したSaaS Trial Sisenseのダッシュボード。[SaaS Trial
Dashboard] (https://app.periscopedata.com/app/gitlab/697554/SaaS-Trial-
Dashboard)
### Data Caveats and Constraints
- ユーザー数 (Nbr of users) は、BigQueryとGoogle Analytics 360 UIの間で常に約0.1~0.2%の相違がありま
す。これは、Google Analyticsではユーザー数が概算で表示され、BigQueryでは正確に表示されるためです。
- ヒットレベルのデータ (例:ページビュー)を報告する場合は、次のようにリンクする必要があります。
Google_analytics_360.session_hit`テーブルを `analytics.ga360_session_xf` に変更し、
session_date でフィルタリングできるようにしました。これは現在、ヒットレベルのデータでは利用できない
ためです。
### Data Security Classification
```

Google Analyticsでは、クライアントIDや訪問者IDなどの固有の識別子が保存されますが、これには \*\*変更\*\* の分類が必要です。

### ソリューション・オーナーシップ

- ソースシステムオーナー:`@shanerice`さん
  - マーケティング戦略とパフォーマンスのサブジェクト・マター・エキスパート: `@vjuhasz`.
  - Data Team Subject Matter Expert: `@paul\_armstrong`

## `@jeanpeguero` ### Key Terms

- \*\*セッション\*\* (別名:訪問):一定の時間内に行われるユーザーのウェブサイトとのやり取りのこと。例えば、1 つのセッションには、複数のページビュー、イベント、ソーシャルインタラクションを含めることができます。セッションは午前0時に失効します。また、GitLabのサブドメインでは、一定時間の活動がないと、\*\*60分\*\*で失効します
- \*\*ページビュー\*\*: Google Analyticsのトラッキングコードが挿入されたページをユーザーが閲覧した回数です。これはすべてのページビューを対象としています。したがって、ユーザーがページを更新したり、ページからナビゲートして戻ってきたりした場合、これらはすべて追加のページビューとしてカウントされます。
- \*\*ユニーク・ページビュー\*\*。ユニークページビュー数は、同一ユーザーが同一セッション中に複数回閲覧したページを1回だけカウントすることで算出されます。
- \*\*新規訪問者\*\*。新規訪問者/ユーザーとは、指定された期間以前にサイトを訪れたことのないユーザーのことです。
- \*\*リターニング・ビジター\*\*:リターニングビジター/ユーザーは、以前にサイトの少なくとも1つのページに少なくとも1回の訪問をしています。これは、Google Analyticsが以前の訪問を示すクッキーを検出できるかどうかによって決まります。Googleがクッキーを検出できない場合、ユーザーが個人のブラウザ設定でクッキーを無効にしていない限り、今後の記録のためにクッキーが設定されます。
- \*\*直帰率\*\*: ユーザーがどのようにサイトにアクセスしたか、またどのくらいの時間そのページに滞在したかに関わらず、1ページだけでサイトを離れた訪問数をパーセントで表したものです。
- \*\*ヒット\*\*:ページビュー、イベント、トランザクションを含む、訪問中のインタラクション。
- イベント\*\*: イベントは、トラッキングコードでは記録されないアクティビティを測定するのに役立ちます。ユーザーが動画を再生したとき、CTAをクリックしたときなどがこれに該当します。
- \*\*カスタムディメンション\*\*: ブログ記事のタグ、ブログ記事のタイプ(ブログ記事の作成者)など、追加の非標準的なデータをGoogle Analyticsに記録することができます。

## ### Key Metrics, KPIs, and PIs

- \*\*セッション\*\*: SisenseのBigQueryウェブトラフィックデータをGoogle Analytics UIで適切にクエリしてマッチさせるためには、ミッドナイトスプリットをオーバーライドするセッションIDを作成する必要があります。
- `CONCAT(visitor\_id, CAST(visit\_start\_time AS STRING)`)
- \*\*Users\*\*:ユーザーの計算は、visitor\_id `COUNT(DISTINCT visitor\_id)`に基づいて行う必要があります
- 。Google Analyticsでは、ユーザーを計算する際にすべての行を考慮します(セッションとは異なります)。
- \*\*新規ユーザー\*\*:ユーザーの計算は、訪問者のIDに基づいて行う必要があります `COUNT(DISTINCT visitor\_id)` where total\_new\_visits is not null`.
- \*\*ページビュー\*\*: この指標はヒットレベルのデータに基づいており、すべての行が `COUNT(1)` where `HIT\_TYPE = 'PAGE'` とカウントされる必要があります。
- \*\*ユニークなページビュー\*\*: この指標もヒットレベルのデータに基づいており、セッション内のユニークなページビューをカウントするためには、セッションIDを連結する必要があります(詳細は
- \_Sessions KPI\_) with `PAGE\_PATH` where `HIT\_TYPE = 'PAGE'`.全ページの結果は、Google AnalyticsのUIと比較して、約-0.2%の差がある場合があります。
- \*\*バウンスレート\*\*。近日公開予定

## #### キーフィールドとビジネスロジック

- session\_date`: 日付を含むすべてのレポートは、`session\_date`に基づいていなければなりません。 ビジットスタートタイム
- Visitor id`: セッションIDの作成とユーザー数の算出に使用されます。
- total\_visits`: セッションの数です。この値は、インタラクションイベントがあるセッションでは `1` です。 セッション内にインタラクションイベントがない場合は、`null`となります。Google Analyticsではインタラクティブなセッションのみが表示されるため、nbr of sessionsを検索する際にはこの点を考慮する必要があります。
- `hit\_type`: ヒットの種類で、`PAGE`または`EVENT`のいずれかになります。
- `page\_path`:ページのURLパス、\_e.g:/フリートライアル/\_
- ホスト名`: これは、異なるサブドメイン\_(about.gitlab.com、docs.gitlab.comなど)でフィルタリングするために使用されます。

#### スニペット

近日公開予定

#### リファレンスSQL

```
以下のSQLクエリは、Google Analytics UIの結果と一致する各主要指標の計算を定義しています(All Data
- Raw No Filtersビュー)。
2020年9月の日別セッション数、ユーザー数、新規ユーザー数**。
````sql
SELECT
 DATE_TRUNC('day',session_date)::date
AS visit_date,
 COUNT(DISTINCT IFF(total_visits=1, CONCAT(visitor_id, CAST(visit_start_time AS STRING)), NULL)AS
interactive Sessions,
 COUNT(DISTINCT visitor_id)
AS users,
 COUNT(DISTINCT IFF(total_new_visits IS NOT NULL, visitor_id,
NULL))AS new_users
from legacy.ga360_session_xf
WHERE DATE_TRUNC('day',session_date)::date >= '09/01/2020'
 AND DATE_TRUNC('day', session_date)::date <= '09/30/2020'
GROUP BY 1
ORDER BY 1
```

## 2020年2月以降の日別総ページビュー数およびトライアルページビュー数

```
このサブクエリーは、一意のセッションIDごとにsession_dateを取得し、それをリンクさせるために必要です。
ヒットレベル表へ
WITH ga_unique_sessions AS (
SELECT
DISTINCT(CONCAT(visitor id, CAST(visit start time AS STRING)))AS session id,
SESSION DATE
from legacy.ga360_session_xf
総ページビュー数とトライアルページビュー数
), ga_pageview_hits AS (
SELECT
CONCAT(visitor_id, CAST(visit_start_time AS
                                                            Δς
                                                            session id,
STRING)) page_path,
 CASE
 WHEN page path = '/free-
   trial/' THEN COUNT(1)
 else 0
 end
                                                            AS trial_pageviews,
 count(1)
                                                            AS total_pageviews
from legacy.ga360_session_hit
where hit_type = 'page'
GROUP BY 1,2
)、ファイナルAS(
SELECT
 SESSION_
                                                            DATEAS visit_date,
SUM(trial_pageviews)
                                                            AS trial_pageviews,
                                                            AS total_pageviews,
SUM(total_pageviews)
COUNT(DISTINCT ga_pageview_hits.session_id,page_path )
                                                           AS unique_pageviews
FROM ga_pageview_hits
LEFT JOIN ga_unique_sessions
ON ga_pageview_hits.session_id=ga_unique_sessions.session_id
WHERE DATE_TRUNC('day', SESSION_DATE)::date >= '2/1/2020'
GROUP BY 1
```

SELECT \*
FROM final

エンティティ・リレーションシップ・ダイアグラム

## 近日公開

データプラットフォームソリューション

データの系統

Google Analytics 360/BigQueryのデータは、ロサンゼルスのタイムゾーン(GTM-08:00)を使用している「All Data - Raw No Filters」ビューに基づいています。Google Analytics 360/BigQueryのデータは、Sisenseでは2020年2月19日から利用可能です。これらは、Google Analytics 360のUIでクエリ結果を比較する際に知っておく必要があります。

Google Analytics BigQueryの各行は、1つのセッションを表しています。Google AnalyticsのUIには多くのディメンションやメトリクスがあるため、このスキーマはデータベースレベルでは少し圧倒されるかもしれません。下の図は2つのセッションを表しており、BigQueryのデータ構造をよりよく理解するのに役立ちます。

image{: .shadow}.

BigQueryのエクスポートスキーマについては、以下のインタラクティブツリーマップを参照して

ください。 **BigQuery** 

ソース

Webトラフィック分析を行う上で重要となるBigQueryのテーブルは大きく分けて3つあります。

- 1. analytics.ga360\_session\_xf: セッション数、ユーザー数、ユニークビジター数の算出に使用されるセッションレベルのデータ。
- 2. google\_analytics\_360.session\_hit: ページビュー、イベントなどを計算するためのヒットレベルデータ。
- 3. ga360\_session\_custom\_dimension\_xf: GitLab固有のディメンションのレポート用(これは、各カスタムディメンションに定義されたスコープに依存します。

DBTソリュー ション

近日公開

トラステッドデータソリューション

近日公開

EDM エンタープライズ ディメンション モデル バリデーション

近日公開

RAWソースデータパイプラインの検証

近日公開

layout: handbook-page-toc title:"xMAU分析" 説明"xMAUは、月間アクティブ使用量(MAU)を把握するための様々なレベルを把握するための単一の用語です。"

このページについて

{:.no\_toc}

TOC {:toc}

# xMAU分析

xMAUは、月間アクティブユーザー数(MAU)を把握するための様々なレベルを表す単一の用語であり、アクション(AMAU)、グループ(GMAU)、ステージ(SMAU)、トータルCMAUを含みます。製品グループに有用な単一の指標を提供し、製品全体の重要業績評価指標にうまくマッピングするために、xMAUの各指標は上述の順序で上に向かって連鎖していきます。

xMAUメトリクスは、Usage Ping(セルフマネージドインスタンスレベルの粒度)とGitLab.com(SaaSネームスペースレベルの 粒度)から得られます。このアナリティクス・ワークフローは、様々な顧客セグメントにおける各レベルのxMAUメトリクスの 分析を可能にし、報告、推定、予測メトリクスの報告の基礎となります。

## このページの目的は

- ◆ Product Adoption Dashboardの操作方法の理解 Product Adoption Dashboard
- に使用されるデータモデルの理解
- ◆ GitLabでのあなたの役割に最も適した認定を受けて、理解度を確認してください。ダッシュボ
  - o ード・ソリューションへの理解度を測るには、ダッシュボード・ユーザー認定を受けて
  - ください。ダッシュボード開発者認定資格は、今後取得する予定です。
  - o さらに、推定xMAUを記録するための方法を理解することをお勧めします。

クイックリンク

製品採用ダッシュ ボード

SiSense xMAU入門 ダッシュボード例 推定アルゴリズムのページ

# はじめに

- 1. 重要な用語、測定基準、KPI/PIの理解
- 2. データモデルの理解

## 主な用語

- アカウント
- ホストイン
- スタンス
- Instance User Count インスタンス上のユーザーの合計数
- 有料ユーザ
- ◆ 一製品階層
- 使用量 Ping
- バージョン

キーメトリクス、KPI、PI

以下の指標の説明は、Product Team Performance Indicatorのページにあります。

- Action Monthly Active Users (AMAU)
- Stage Monthly Active Users (SMAU)
- セクションの月間アクティブユーザー数(セクショ
- ンMAU) セクションの総月間アクティブユーザー数(
- セクションCMAU) 複合月間アクティブユーザー数( CMAU

各指標には、3つの異なるバージョン(記録、推定、予測)があり、「製品チームのパ

◆ フォーマンス指標」のページで説明されています。

• the Sisense Style Guide 現在、使用量を特定した記録されたpingメトリクスは、Centralized Dashboardにチャート化されていますが、今回のEstimated valuesの最初のバージョンに取り組んでいます。

#### xMAUの算出方法は?

xMAUは主にサービスPingデータソースによって計算されます。プロジェクトの開始時、ステージとグループのプロダクトマネージャーは、xMAUチャートを作成する特定のサービスPingメトリックを選択しました。

- ◆ 使用されたGMAUメトリ
- クスのリスト 使用された SMAUメトリクスのリスト

現在、xMAUメトリクスのSSOTはこのスプレッドシートで、Sheetload経由でデータハウスにインポートされています。つまり、特定のメトリクスのGMAU、SMAU列を更新すると、その変更はハンドブックで更新されたxMAUチャートに伝搬するということになります。

このSSOTをこのスプレッドシートから辞書のYAMLファイルに移行する計画がありますが、その作業は本号で行います。

メトリクスの定義についてさらにご質問がある場合は、プロダクト・インテリジェンス・チームにお尋ねください。彼らは現在、ここで利用可能なメトリクス辞書を管理しています。また、メトリクスがデータベースの計算である場合は、メトリクスの値を生成するために実行したSQLクエリを提供することができます。

### 日付の範囲

すべてのインスタンス (自己管理型とSaas/GitLab.com) について、その月に生成された最後のpingを使用してxMAUを計算しています。インスタンスには、サービスのpingを生成する曜日がランダムに割り当てられますが、その割り当ては時間の経過とともに変化します。例えば、あるインスタンスが火曜日にPingを生成するように割り当てられた場合、そのインスタンスは常に火曜日にPingを生成します。pingが生成される曜日はインスタンスごとに異なるため、28日カウンタに記録される正確な日付範囲も異なります。

現在、「今月の最後のping」は、週の初め(月曜日)の日付で決定されています。(SQLではDATE\_TRUNC('month', ping\_created\_at\_week))を使用しています。そのため、「その月の最後のping」が技術的には測定月に起こらないこともあります。ここでは、2021年8月の例を紹介します。

Ping 作成日	Ping作成週間	<b>2021年8</b> 月のxMAUに使用され ています。
2021-08-05	2021-08-02	FALSE
2021-08-12	2021-08-09	FALSE
2021-08-19	2021-08-16	FALSE
2021-08-26	2021-08-23	FALSE
2021-09-02	2021-08-30	TRUE

## xMAU」と「有料xMAU」の違い

## 有料xMAUの定義

上記の各指標は、xMAUと有料xMAUについて計算されます。有料xMAUは現在、「28日間のローリング期間中に、Usage Pingデータを介してセルフマネージドの有料インスタンスにロールアップしたユーザー、*または*GitLab.com Postgres Database Imports を介してSaaSの有料ネームスペースにロールアップしたユーザー」を月間アクティブユーザーと定義しています。(例として「PaidStageMonthlyActiveUsers-PaidSMAU」を参照してください)。

現在、GitLab ビジネス機能は、どの名前空間やインスタンスが OSS、EDU、内部プロジェクト、または有料プランタイプを持つ が ARR を貢献しないその他のサブスクリプションに属するかを識別する標準的な方法を持っていないため、xMAU の現在の実装では、これらのサブスクリプションに関連するユーザーを「有料」として含めています。

xMAUとpaid xMAUを計算するために、バージョンアプリとGitlab.comのPostgresデータベースの2つの主なデータソースがあります。下の表は、これらの計算に使用されるデータソースをまとめたものです。

デリバリー	xMAU	有料xMAU
SaaS	バージョンア ップのお知ら せ	Gitlab.com Postgres Table**。
セルフマネジ メント	バージョンア プリ	バージョンアップしました 。

#### メモ

\*:SaaSのxMAUでは、gitlab.com用に生成されたペイロードを使用しています。これらのペイロードは、uuid = ea8bf810-1d6f-4a6a-b4fd-93e8cbd8b57f のインスタンスにリンクされているので、簡単に識別できます。

\*\*:SaaSの有料XMAUでは、Gitlab.comのpostgresテーブルを使用する必要があります。Gitlab.com インスタンスから生成された Usage Ping ペイロードは、インスタンスレベルでのハイレベルな統計情報を提供します。これは集約された数値であり、例えば 製品の階層、プランの種類、ネームスペースなど、さらに細かく分類することはできません。

- SaaSの有料xMAUを生成するためには、データウェアハウスに格納されているGitlab.comデータベースのレプリカにxMAUカウンターを複製する必要があります。
  - バッチカウンタは、シンプルなSQLで生成されたカウンタです。カウンタの生成に使用されたSQLクエリはアクセス可能で、簡単に再現できます。
  - 残念ながら、これはすべてのカウンターに対してできることではありません。Redisのカウンタは、SQLで生成された カウンタではありません。Redisカウンタは、ページビューやフロントエンドのインタラクションなど、Postgres DBにはないアクションも追跡します。
- そのため、Gitlab.com の Postgres レプリカを使って再現できるのは一部のメトリクスだけです。つまり、現時点では Monitor StageのようなSaaS PaidのxMAUメトリクスの一部を計算することができません。

\*\*\*:セルフマネージの有償 xMAU を計算するために、使用量Pingペイロードのエディション・フィールドを使用し、EEP、EES、および EEUエディションの使用量 Ping のみを選択します。エディションの値はライセンスが生成れた時点で、licenseDotデータベースのライセンステーブルのプラン列から得られ。つまり、現在はEDU/OSSのサブスクリプションを有料のxMAUの計算から除外していません。

## データソース

上記の説明に基づいて、私たちが使用している2つの主要なデータソース

- ◆ があります。サービスPingデータソース
- Gitlab.comのデータソース

## データマート

私たちは、ユーザーがさまざまな製品データソースを探索できる一連のデータマートを構築しました。

## マートサービス Ping使用量 Pingデータ

mart\_service\_ping\_product\_usage\_dataは、すべてのService Pingデータのための最も包括的なデータモデルです。このデータ・モデルは、Service Pingデータを、サブスクリプション、CRM Account...などの財務およびGTMデータ・ソースと結合します。

これにより、ここで定義されているすべての月別およびすべての時間別指標の使用データを取得することができます。

## マート推定xMAU

mart mart\_estimated\_xmauは、推定XMAU PIを簡単に生成するために構築されています。このモデルには、このページで説明されているすべての推定ロジックが含まれています。

エンドユーザーは、非常にシンプルなチャートを使って、xMAUの推定値を作成することができます。

```
SELECT
Report_month,
Product_tier,
SUM(推定月次指標値_sum)
FROM common_mart_product.mart_estimated_xmau
WHERE xmau_level = 'SMAU' AND stage_name =
'create' GROUP BY 1,2```。
```

mart\_paid\_estimated\_xmauデータマートが作成され、特別に支払われたxMAUのチャートが作成されました。

サービスPingデータパイプライン

# プロダクトマネージャーのためのSisense Snippets

プロダクトマネージャーのために作成されたスニペットをすべて網羅した別のページを作成しました。このトピックについては、このページをお読みになることをお勧めします。

また、これらのスニペットをどのように使用するかの例を満載したダッシュボード回も作成しました。

これらのスニペットを使って、特定のxMAUチャートのターゲットを更新するには?

## ダイナミック・ターゲット・バリューの導入

Usage Ping Metrics の SiSense チャートに Monthly Target Value を実装することができるようになりました。これを実現するには、パフォーマンス指標(xMAU、またはUsage Pingから得られる他のPI)のあるymlファイルを更新し、以下に説明する例に従う必要があります。

----

- name: Configure:Configure - Adoption PI - GitLab Managed Terraform Stateを持つプロジェクトの数 base\_path:"/handbook/product/performance-indicators/ops-section-performance-indicators/" definition:のGitLab Managed Terraform Stateを使用するプロジェクト数のローリングカウントです。

## 過去28日間

目標: FY22.1Q末までに6000人

ターゲット名: Projects with terraform states 月

間記録されたターゲット。

"2021-04-20":6000

"2020-11-01":3000

"2020-07-01":700

毎月の推定目標値"2021-05-

01":6000

"2020-11-01":3000

"2020-07-01":700

org:オペレーションセクション

## ターゲットを正しく埋めるには

ymlセクションに追加される日付は、各日付間隔の終了日です。つまり、例えば以下の例では

毎月の推定目標値

"2021-05-01":6000

"2020-11-01":3000

"2020-07-01":700

そして、次のような目標を掲げます。

- start\_dateから2020-07-01までは700 2020-08-01
- から2020-11-01までは3000
- ◆ 2020-12-01から2020-05-01まで、目標は6000

そのため、yml定義の中に2つの異なるセクションがあり、それを使用することができます。

## ▶ 月間記録目標

PIが推定PIの場合(xMAUのチャートなど)。ダイナミックターゲットを使ってチャートを可視化するには、このスニペットを使用す

る必要があります。このスニペットの詳細については、こちらをご覧ください。 できることの例

# ▶ 月間推定目標値

ターゲットを設定したいPIが記録された値に基づいている場合(推定値を含まない)。ダイナミックターゲットを設定したチャートを表示するには、次のようなコードが必要です。

このスニペットの詳細はこちらをご覧ください 何ができるかの例

#### スニペットを使って静的なターゲットを取得する

当社の標準化されたスニペットとビジュアライゼーションを使用したすべての埋め込みPI/xMAUチャートでは、ターゲットの更新は非常に簡単な4ステップのプロセスです。

- 1. 興味のあるPIにアクセスし、sisenseのリンク pi\_pageをクリックします。
- 2. チャート 編集」をクリックします。
- 3. 最後の値を変更する(これが目標値です。10進数でやれば一定の割合で増えます。完全な数値を入力すると、その数値が目標ラインになります。)  $\square$  ターゲット
- 4. 保存ボタンをクリックすると、 保存されます。

# エンティティ・リレーションシップ・ダイアグラム

私たちの目標のひとつは、分析に必要なすべてのデータを簡単に提供できる1モデルを作ることです。しかし、ソリューションの 開発を続けていると、このモデルでは得られない情報が出てくることがあります。ここで、エンティティ・リレーションシップ・ダイアグラムの理解が役立ちます。このモデルは、アクセスしているレイヤーを作成するためにどのテーブルが結合されているかを示しています。このモデルは、より深く掘り下げて、さらなる洞察を得ようとするときに役立ちます。

## Diagram/EntityGrainPurpose

ファクト 月間使用量Pingデータ ERDとデータフロー

host\_id, instance\_id, reporting\_month, metrics\_path さまざまなインスタンスおよびサブスクリプションの ディメンション別に、機能の使用状況を把握すること ができます。

## ▶ データの分類

xMAU分析をサポートするデータの中には、[オレンジ](/handbook/engineering/security/data-classification-standard.html#orange) または[イエロー](/handbook/engineering/security/data-classification-standard.html#yellow)に分類されるものがあります。これには、アカウントからのORANGEの顧客メタデータ、SalesforceやZuoraからのコンタクトデータ、GitLabのNon公開の財務情報などが含まれ、これらはすべて公開されるべきではありません。このダッシュボードからデータを共有する際には、その詳細が GitLabの組織内に留まるようにし、外部に共有する際には適切な承認を得られるように注意しなければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保存しないよう、常に注意を払う必要があります。これらのデータは、[Snowflake](/handbook/business-ops/data-team/platform/#data-warehouse)と

[Sisense](/handbook/business-ops/data-team/platform/periscope/)に保存され、特に承認されない限り、これらのアプリケーションを通じてのみ共有されることが理想的です。

## ORANGE

- 説明行またはレコードレベルの顧客および個人データ。オブジェク
- トです。
  - o dim\_billing\_accounts

0

dim\_crm\_accounts

## o ユーズドピングマート

## ▶ ### ソリューション・オーナーシップ

セルフサービス・ダッシュボード・ソリューション

	ダッシュボード目的
エグゼクティブ・オーバービュー - TBD	このダッシュボードは、すべてのメトリクスの現在のステータスをエグ ゼクティブ・オーバービューとして表示します。
Dev Section Analysis - TBD	このダッシュボードでは、関連するメトリクスの現状をセクション ごとに紹介しています。
DRAFT: ハンドブック更新のための一元化された ダッシュボード	
トラステッドデータソリューション	
トラステッドデータフレームワーク	
手動でのデータ検証	
• 手動での使用方法PingValidationDashboard	

{:toc}

# xMAU分析

xMAUは、月間アクティブユーザー数 (MAU) を把握するための様々なレベルを表す単一の用語であり、アクション (AMAU)、グループ (GMAU)、ステージ (SMAU)、トータルCMAUを含みます。製品グループに有用な単一の指標を提供し、製品全体の重要業績評価指標にうまくマッピングするために、xMAUの各指標は上述の順序で上に向かって連鎖していきます。

xMAUメトリクスは、Usage Ping(セルフマネージドインスタンスレベルの粒度)とGitLab.com(SaaSネームスペースレベルの粒度)から得られます。このアナリティクス・ワークフローは、様々な顧客セグメントにおける各レベルのxMAUメトリクスの分析を可能にし、報告、推定、予測メトリクスの報告の基礎となります。

このページの目的は

- Product Adoption Dashboardの操作方法の理解 Product Adoption Dashboardに
- 使用されるデータモデルの理解

Version Trusted Data Dashboard

- GitLabでのあなたの役割に最も適した認定を受けて、理解度を確認してください。ダッシュボード・
  - o ソリューションへの理解度を測るには、ダッシュボード・ユーザー認定を受けてください。ダッシ
  - ュボード開発者認定資格は、今後取得する予定です。
  - さらに、推定xMAUを記録するための方法を理解することをお勧めします。

クイックリンク

 製品採用ダッシュ
 SiSenseを始めよ
 xMAUのダッシ
 推定アルゴリズムのページ

 ボード
 う
 ュボード例

# はじめに

- 1. 重要な用語、測定基準、KPI/PIの理解
- 2. データモデルの理解

## 主な用語

- アカウント
- ホストイン
- スタンス
- Instance User Count インスタンス上のユーザーの合計数 Paid
- User
- 製品階層 使
- 用状況 Ping
- バージョン

キーメトリクス、KPI、PI

以下の指標の説明は、「Product Team Performance Indicator」のページにあります。アクション月間

- アクティブユーザー数 (AMAU
- ステージ・マンスリー・アクティブ・ユーザー (SMAU
- セクション 月間アクティブユーザー数(セクション
- MAU) セクション合計月間アクティブユーザー数(セ
- クションCMAU) 合計月間アクティブユーザー数 ( CMAU

各指標には3つの異なるバージョン(記録、推定、予測)があり、その説明は

- Product Team Performance Indicatorページ
- the Sisense Style Guide 現在、使用量を特定した記録されたpingメトリクスは、Centralized Dashboardにチャート表示されていますが、今回のEstimated valuesの最初のバージョンに取り組んでいます。

xMAU」と「有料xMAU」の違い

## 有料xMAUの定義

上記の各指標は、xMAUと有料xMAUについて計算されます。有料xMAUは現在、「28日間のローリング期間中に、Usage Pingデータを介してセルフマネージドの有料インスタンスにロールアップしたユーザー、*または*GitLab.com Postgres Database Importsを介してSaaSの有料ネームスペースにロールアップしたユーザー」を月間アクティブユーザーと定義しています。(例として「PaidStageMonthlyActiveUsers-PaidSMAU」を参照してください)。

現在、GitLab ビジネス機能は、どの名前空間やインスタンスが OSS、EDU、内部プロジェクト、または有料プランタイプを 持つが ARR を貢献しないその他のサブスクリプションに属するかを識別する標準的な方法を持っていないため、xMAU の 現在の実装では、これらのサブスクリプションに関連するユーザーを「有料」として含めています。

xMAUおよびPaid xMAUの算出について

xMAUとpaid xMAUを計算するための主なデータソースは、Version AppとGitLab.comのPostgresデータベースの2つです。下の表は、これらの計算にどちらのデータソースが使われているかをまとめたものです。

デリバリー	xMAU	有料xMAU
SaaS	バージョンア ップのお知ら せ	GitLab.com Postgres Table**。
セルフマネジメ ント	バージョンア プリ	バージョンアップしました。

メモ

\*:SaaSのxMAUでは、gitlab.com用に生成されたペイロードを使用しています。これらのペイロードは、uuid = ea8bf810-1d6f-4a6a-b4fd-93e8cbd8b57f のインスタンスにリンクされているので、簡単に識別できます。

\*\*:SaaSの有料XMAUでは、GitLab.comのpostgresテーブルを使用する必要があります。GitLab.com インスタンスから生成された Usage Ping ペイロードは、インスタンスレベルでのハイレベルな統計情報を提供します。これは集約された数値であり、製品の階層やプランの種類、ネームスペースなど、さらに細かく分類することはできません。

• SaaSの有料xMAUを生成するためには、データウェアハウスに格納されているGitLab.comデータベースのレプリカに xMAUカウンターを複製する必要があります。

- バッチカウンタは、シンプルなSQLで生成されたカウンタです。カウンタの生成に使用されたSQLクエリはアクセス可能で、簡単に再現できます。
- 残念ながら、これはすべてのカウンターに対してできることではありません。Redisのカウンタは、SQLで生成されたカウンタではありません。Redisカウンタは、ページビューやフロントエンドのインタラクションなど、Postgres DBにはないアクションも追跡します。
- そのため、GitLab.comのPostgres Replicaを使って再現できるのは一部のメトリクスのみです。つまり、今のところ、 Monitor StageのようなSaaS Paid xMAUメトリクスの一部を計算することはできません。

\*\*\*:セルフマネージの有償 xMAU を計算するために、使用量Pingペイロードのエディション・フィールドを使用し、EEP、EES、およびEEUエディションの使用量 Ping のみを選択します。エディションの値はライセンスが生成され時点で、licenseDotデータベースのライセンステーブルのプラン列から得られ

データモデルとスニペット

2種類のデータソリューションを作成しました。

- xmau\_202011は、PMが自分のXMAU指標を素早くグラフ化し、グループや個人のダッシュボードに表示したり、ダッシュボードに埋め込んだりすることができるスニペットです。詳細はこちら
- mart\_monthly\_product\_usageテーブルです。PMがデータをさらに調査したい場合は、トップユーザーが誰であるかを理解したり、他の次元(案件の規模、会社の規模、製品の階層...)で結果を分解したりすることができます。

Product Adoption Dashboardは、1つのメインスニペットを使って作成されます。

## xmau 202011

注: このリンクは、このスニペットを使ったSisenseビューにリダイレクトされます。RUN SQLボタンを押すと、データセットを見ることができます。

以下の説明には、いくつかの例があるとより理解しやすいでしょう。

- 開発部門のCMAU推定値
- リリース管理グループの有料GMAU推定値 プランステージの
- SMAU推定値(エディション別

このスニペットは、主に推定および記録されたXMAUに推定値を与えるために使用され

ます。使用するには4つのパラメータを宣言する必要があります。

- xmau\_type: すべてのXMAUデータを表示するか、有料データのみを表示するかを決定します。許容され
- る値は「All」または「paid」です。 xmau\_level: これらの値の中から1つを選択します。cmau'、'umau'
   'smau' 'gmau'
- フィルター:これはxmau\_levelに選ばれた値に依存します。
  - CMAU: Allを宣言するか、セクション名を選択することができます(キャメルケース、可能な値は'dev'、 'enablement'、'ops'、'secure\_protect')。
  - 。 SMAU:  $セクション名(キャメルケース、可能な値は'dev'、'enablement'、'ops'、'secure_protect'です)、またはステージ名(キャメルケース)のいずれかを選択します。ステージ名の一覧を以下の表に示します。$
  - GMAU: ステージネーム (キャメルケース、ステージネームのリストはこちら) またはグループネーム (camel\_case) のどちらかを選択できます。グループ名のリストはこちらの表にあります
- target: チャートにターゲットラインを設定することができます。空欄のままだと何も表示されません。ターゲットとして設定したい値を入力することで静的なターゲットラインを作成することもできますし、毎月の成長目標(例えば10%)を10進数の値(この例では0.1)で入力することで動的なターゲットラインを作成することもできます。

このスニペットは、限られた寸法のコンパクトなテーブルを返します。

- created\_month: レポート作成月
- product\_tier: こちらの定義を参照してください。これは、インスタンスにリンクされたサブスクリプションの product\_tierではなく、インスタンスが属するproduct\_tierを表します。('All'、'target'も潜在的な値です)

- の配信を行っています。SaaSまたはセルフマネージド(「All」、「target」も潜在的な値です。
- の内訳です。配信='SaaS'にはSaaS、配信='セルフマネージド'には録音されたセルフマネージドとに分かれます。 推定自己管理。これにより、当社の設計基準に従ってXMAUチャートを作成することができます(例はこちら)。

なお、「All」や「target」も潜在的な値です。

• 版である。CE、EE、SaaS (「All」、「target」も可能性のあ

る値です) そして、1つのメジャー。

• mau value: 特定のディメンションのアクティブ・ユーザーの数

推定アルゴリズムの詳細については、現在の方法論と私たちのビジョンについての詳細をこちらでご覧いただけます。

#### 月間製品使用量の推移

このモデルでは、使用状況のpingデータと、ライセンス、salesforce、zuoraのデータを結合します。つまり、販売/財務と製品のデータを同時に扱うことができるのです。例えば、次のような質問に答えられるようになります。

- XMAU計画に最も貢献している大学上位10校(課題を作成したユーザー数
- 特定のステージの月別および製品層別の採用率は、SMAUを業界
- 別に分割して検証する

dbtモデルは、各カラムの定義がしっかりと書かれています。また、この新しいモデルを使って答えられる質問を示すダッシュボードも作成しました。

課題を作成したユーザー数の月別推移をプロットする基本的なSQLクエリ。

```
SELECT
report_month,
main_edition,
delivery,
SUM(月次指標の値)
FROM legacy.mart_monthly_product_usage
WHERE metrics_path = 'usage_activity_by_stage_monthly.plan.issues' GROUP
BY 1,2,3
```

CIパイプライン機能を使って、EDU/OSSサブスクリプションのトップ10を抽出する基本的なクエリです。

```
SELECT
 host_name,
 license_id,
 ping_id。
 サブスクリプション名 slugify,
 ping_product_tier,
 main_edition,
 monthly_metric_value
FROM legacy.mart_monthly_product_usage
WHERE metrics_path = 'usage_activity_by_stage_monthly.verify.ci_pipelines'
 -- スニペット[last_month]は、最後の暦月のみから値を取得するために使用されます AND
 reporting_month = [last_month]
 AND DELIVATION = 'Self-Managed'
 is program subscription は、EDU/OSS プログラムの一部であるすべてのサブスクリプションに対して TRUE
に設定されるブーリアンフラグです。
 AND is_program_subscription
ORDER BY monthly_metric_value
DESC LIMIT 50
```

有料のセルフマネージドインスタンスでデプロイメントを作成したユーザ数をSalesforceの業種別に分割して月ごとにプロットする基本的なクエリです。

2021/11/8 merged.md

```
SELECT
  reporting_month,
 IFF(ultimate_parent_industry IS NOT NULL, ultimate_parent_industry, 'Unknown') AS
ultimate_parent_industry,
  SUM(月次指標の値)
FROM legacy.mart monthly product usage
WHERE metrics_path = 'usage_activity_by_stage_monthly.release.deployments' AND
  is_paid_subscription
GROUP BY 1,2
```

# エンティティ・リレーションシップ・ダイアグラム

私たちの目標のひとつは、分析に必要なすべてのデータを簡単に提供できる1モデルを作ることです。しかし、ソリューション の開発を続けていると、このモデルでは得られない情報が出てくることがあります。ここで、エンティティ・リレーションシップ ・ダイアグラムの理解が役立ちます。このモデルは、アクセスしているレイヤーを作成するためにどのテーブルが結合されて いるかを示しています。このモデルは、より深く掘り下げて、さらなる洞察を得ようとするときに役立ちます。

#### Diagram/EntityGrainPurpose さまざまなインスタンスおよびサブスクリプシ ファクト 月間使用量Pingデータ host\_id, instance\_id, ERDとデータフロー ョンのディメンション別に、機能の使用状況を reporting\_month, metrics\_path 把握することができます。

▶ データの分類

xMAU分析をサポートするデータの中には、[オレンジ](/handbook/engineering/security/data-classificationstandard.html#orange)または[イエロー](/handbook/engineering/security/data-classification-standard.html#yellow)に分類され るものがあります。これには、アカウントからのORANGEの顧客メタデータ、SalesforceやZuoraからのコンタクトデータ、 GitLabのNon公開の財務情報などが含まれ、これらはすべて公開されるべきではありません。このダッシュボードからデータを 共有する際には、その詳細が GitLab の組織内に留まるようにし、外部に共有する際には適切な承認を得られるように注意しな ければなりません。また、行やレコードレベルの顧客メタデータを扱う際には、個人のデバイスやラップトップにデータを保 存しないよう、常に注意を払う必要があります。これらのデータは、[Snowflake](/handbook/business-technology/datateam/platform/#data-warehouse)と[Sisense](/handbook/business-technology/data-team/platform/periscope/)に残すべき であり、特別な承認がない限り、これらのアプリケーションを通してのみ共有されるのが理想的です。

## **ORANGE**

- 説明行またはレコードレベルの顧客および個人データ。オブジェク
- トです。
  - dim\_billing\_accounts
  - o dim\_crm\_accounts
  - o usage\_ping\_mart

## ▶ ### ソリューション・オーナーシップ

セルフサービス・ダッシュボード・ソリューション

	ダッシュボード目的
エグゼクティブ・オーバービュー - TBD	このダッシュボードは、すべてのメトリクスの現在のステータスをエグ ゼクティブ・オーバービューとして表示します。
Dev Section Analysis - TBD	このダッシュボードでは、関連するメトリクスの現状をセクション ごとに紹介しています。
DRAFT: ハンドブック更新のための一元化された ダッシュボード	このダッシュボードのチャートは、ハンドブックのエンベッドに使用されます。

これらのスニペットを使って、特定のxMAUチャートのターゲットを更新するには?

当社の標準化されたスニペットとビジュアライゼーションを使用したすべての埋め込みPI/xMAUチャートでは、ターゲットの更新は非常に簡単な4ステップのプロセスです。

- 1. 興味のあるPIにアクセスし、sisenseのリンク pi\_pageをクリックします。
- 2. チャート 編集」をクリックします。
- 3. 最後の値を変更する(これが目標値です。10進数でやれば一定の割合で増えます。完全な数値を入力すると、その数値が目標ラインになります。) シターゲット
- 4. 保存ボタンをクリックすると、 保存されます。

トラステッドデータソリューション

トラステッドデータフレームワーク

手動でのデータ検証

- ◆ 手動使用Ping検証ダッシュボードバージョン
- 信頼できるデータ ダッシュボード

layout: handbook-page-toc title:"予測されるXMAUのアルゴリズム"

予測アルゴリズムの説明

パフォーマンス指標のページによると、現在xMAUには3つのバージョンがあります。予測されるxMAUは、現在の成長率で3年後に使用量がどのようになるかをプロダクト・リーダーシップに示すためのものです。

最初に提案された解決策は、次の式を適用して、M月の予測xMAUを計算することです。

予測xMAU(M月) = 予測xMAU(先月) + (現在のMoM成長量 x month\_difference(先月とM月の間))

MoM growth Amountは、前四半期の平均月間絶対成長量として算出されます。つまり、仮定すると、過去4

ヶ月間のEstimated Plan SMAUは以下のような数字になります。

月	エスティメイト
	SMAU
M-3	94
M-2	95
M-1	97.5
М	100

平均MoM成長量(M - (M-3) / 3)が2で、現在の値が100です。12ヶ月後のSMAU予測値を算出したいと思います。

予測プランSMAU(先月+12 $_{7}$ 月)=推定プランSMAU(先月)+(2 $_{12}$ )=100+24 = 124

上記の計算式により、Plan SMAUの12ヶ月後の予測値は124となります。

SaaS、Self-Managed CE、Self-Managed EEの成長率の調整

SaaS、CE、EEの間の成長傾向は大きく異なる可能性があります。これを考慮して、アルゴリズムでは、あるグループのCE、EE、SaaSの利用量を考慮し、それぞれの配信タイプの平均成長率に基づいて成長率を調整しています。

そこで、「Predicted SMAU for Plan in 12 Months」を計算してみましょう。以下のような前提条件があります。

配信と版	SMAU	前月比成長率
CE	50	2%
EE	20	4%
SaaS	30	3%
合計	100	

そのため、CE、EE、SaaSについては、上記で説明した同じ計算式を適用して、それぞれ異なるPredicted SMAUを算出します。

```
予測SMAU、CE = 50 x (1+ (0.02 x 12)) = 50 x 1.24 = 62
予測SMAU、EE = 20 x (1+ (0.04 x 12)) = 20 x 1.48 = 30
予測SMAU、CE = 30 x (1+ (0.03 x 12)) = 30 x 1.63 = 41
予測SMAU = 予測SMAU,CE + 予測SMAU,EE + 予測SMAU,SaaS = 62 + 30 + 41 = 133
```

そのため、12ヶ月後のSMAU予測値は133となります。この

ロジックを使ったWIPダッシュボードは、こちらでご覧い

ただけます。

次のステップ

私たちは、この問題に取り組むためのさまざまな方法を検討しています。私たちは、さまざまな選択肢をまとめたダッシュボードを作成しました。

追加のオプションを検討。

- 一定のMoM成長率を使用:つまり、毎月同じ成長率を適用します。つまり、Predicted(Month +1) =予測値(月)×MoM成長率
- Prophetという外部の予測用Pythonライブラリを使って、前月のデータに基づいて成長を予測する。

layout: handbook-page-toc title:"予測されるXMAUのアルゴリズム"

xMAU バリエーション - Trusted Data Framework

xMAUのKPIは、ビジネスの健全性とその成長を評価する非常に重要な指標です。KPIは、月に1度開催されるプロダクト・キー・レビューで精査されます。また、ハンドブックの製品パフォーマンス指標のページにも掲載されています。

これらの理由から、データチームは一連のツール(ダッシュボード、アラート、カスタムテスト)を作成し、KPIの直接の所有者がKPIの変動を容易に監視し、極端な傾向や疑わしい傾向がある場合には警告を受けることができるようにしました。

SiSenseのアラート

SiSenseにはアラート機能があり、クエリの結果が定義された値/閾値よりも高い、同じ、または低い場合に、ユーザーに通知を送ることができます。

KPIオーナーがメトリクスの変動を把握するのに役立つxMAUメトリクス用のSQLアラートセットを作成しました。

- [td\_xmau] Est.CMAU変動監視。このアラートは、CMAU KPIの月次変動が±10%を超えた場合に発生します。
- [td\_xmau] Est.UMAU変動監視。このアラートは、UMAU KPIの月次変動が±10%を超えた場合に発生します。

• [td\_xmau] Est.SMAU Variations Monitoringです。このアラートは、いずれかのStageでMoMの絶対値の変化が10% を超えた場合にトリガされます。SMAU値が1000未満のStageは除外されます(これらのStageはSMAUメトリクスが新たに計測されており、当初は極端な変動を受けると想定しています)。

• [td\_xmau] Est.GMAU変動の監視。このアラートは、グループのいずれかがMoMの絶対的変化が10%を超えた場合に トリガされます。SMAU値が1000未満のグループは除外されます(これらのステージでは、SMAUメトリクスが新た に計測されており、当初は極端な変動を受けるものと想定しています)。

これらのアラートは現在、毎週月曜日と毎月1日に送信されています。

SiSenseのダッシュボード

データチームが作成したTDxMAUダッシュボードは、迅速な理解に役立つ一連のビジュアルをまとめたものです。 これらのKPIの健康状態を

# \$515**=**N5**=** ページが見つかりませ お探しのページは存在しないか、移動さ れています。

## 次のステップ

現在、アラートは1人のチームメンバー(Mathieu Peychet)にのみ送信されています。アラートを受信するユーザーのリストを作成する必要があります(私たちの提案は、Product Intelligence PMとProduct Analytics Teamに送信することです)。

また、このアラートシステムをより良く、よりスマートにするために、他のソリューションも試しています。

- 現在、推定xMAUの計算をdbtに移行しています。
- **GreatExpectations**は、カスタムデータテストを作成するために設計されたオープンソースのソリューションで、試してみる予定です。

layout: handbook-page-toc title:"データ開発" 説明"このページでは、データ開発のライフサイクルを定義しています"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .hidden-md .hidden-lg}。

## GitLabでのデータ開発

GitLabでは、インサイトとビジネス上の意思決定を促進するデータソリューションを構築するために、3つの異なる、しかし相互に関連するアプローチを展開しています。これらのアプローチは互いに補完し合い、ビジネス、解決すべき問題、問われている質問に応じて、必要なスピード、品質、信頼性のレベルで結果を提供することに重点を置いています。これらのアプローチは相互に補完し合い、進化していくものであり、必要に応じて初期の段階で開発したものを後の段階で活用することができます。初期の段階で開発されたデータソリューションは、ビジネス上の必要性が十分にあれば、後の段階で改善・強化することができます。すべての分析は、確立されたデータ分析プロセスに従って行われます。

この3つのアプローチ「Ad-Hoc」「Business Insights」「Trusted Data」を紹介します。

	アドホック	ビジネスインサイト	信頼できるデー タ
方向性/緊急性 使うとき		ルーチン/オペレーション	ミッションクリ ティカル
	分析	分析	分析
<b>SiSense</b> によるビジュアライゼーション	任意	必須	必須
<b>玉ジ</b> タープライズディメンションを利用して構築	オプション	任意	必須
データカタログ	への登録 任意	必須	必須
<b>嗜賴性</b> あ高いデータの開発を ション	フォローする オプ	任意	必須
<b>管頼性の高い</b> データ を実施 オプション	を使用してテスト	任意	必須
ソース クによる監査可能性(	システムへのリン オプション	任意	必須
信頼できるデータ」	ブランド化N/A	N/A	必須
Business losights' <sub>夕開発</sub>	ブランド <b>N/A</b>	必須	N/A

アドホックは、分析作業の典型的な最初のステップであり、1回限りまたは限定的に使用するレポートやダッシュボードを提供することになります。アドホック開発は、既存のデータソリューションでは質問に答えられない場合に行われます。その場限りの分析のために開発されたコードは、長期的なソリューションに活用するために書かれたものではなく、結果を迅速に提供するためのものです。その場限りの分析を行うために、Analystは通常、Enterprise Data Warehouseに対してSQLクエリを書いて実行したり、SisenseやPythonなどのツールを使って分析するためのデータを抽出したり、あるいは既存のダッシュボードを活用したりします。場合によっては、テキストファイルやスプレッドシート、その他のデータソースから新しいデータを入手する必要があります。

ほとんどの場合、アドホックレポートは当面のビジネスニーズを解決し、それ以上の対応は必要ありません。しかし、アドホック分析の結果、追加の分析を必要とする結果が得られることもあります。また、アドホック分析の結果が、より信頼性の高いソリューションに発展させるのに十分なほど重要な場合もあります。このような場合には、ビジネスインサイトソリューションまたはトラステッドデータソリューションの作成が決定されます。

## ビジネスインサイト データ開発

ビジネスインサイトは、安定した信頼性の高いレポートが必要とされるものの、構造化されたエンタープライズ次元モデルが まだ利用できないソリューションの大半を占めています。ビジネスインサイトソリューションは、それぞれのメトリクスのSSOT としての役割を果たし、レポート全体の中で重要な役割を果たしています。

ビジネスインサイトソリューションは、データテスト、コードレビュー、データカタログへの登録などの品質検証が含まれる点で、アドホックレポートとは異なります。ビジネスインサイトソリューションは、EDMの一部を利用することができますが、完全にEDMに基づいているわけではありません。しかし、Trusted Data Solutionと比較すると、Business Insightsソリューションは、完全なテストカバレッジとEDMカバレッジがありません。

信頼されるデータ開発

Trusted Dataは、企業が利用可能な最も完全で、信頼性が高く、正確な分析を提供します。組織が成熟し、アナリティクスの価値が高まるにつれ、Trusted Dataは進化し、開発の厳しさも進化していきますが、基本的なステップは一貫しており、要件収集、設計、反復的なワイヤーフレーム、テスト、運用監視などが含まれます。

☞データチーム開発\_プロセス 信頼

されるデータソリューション基準

すべてのTrusted Dataソリューションは、以下の基準を満たす必要があります。

- 1. ビジネス上の問題が定義され、明確な収益への影響が確立されている。
- 2. 開発を管理するためにデータプロジェクトエピックが作成される
- 3. 要件と成功基準がEpicに取り込まれ、追跡される
- 4. v1.0およびv1.1のスコープが定義されており、リリースサイクルが前もって設定されている(例:毎週、隔週、毎月)。
- 5. ダッシュボードのワイヤーフレームをLucidまたはSisenseで作成し、ユーザーと共有し、"最終ドラフト"に向けて繰り返し作業を行います。
- 6. 次元データモデルは、エンタープライズ次元モデルバスマトリックスに設計・統合されています。
- 7. 信頼できるデータテストの作成と展開
- 8. ソリューションは、ソースシステムへのデータ検証を含むユーザーアクセプタンステストの段階に入ります。
- 9. ソリューションがデータカタログに登録されていること
- 10. 必要なトレーニングやユーザーイネーブルメントを含む、ソリューションの展開

layout: handbook-page-toc title: "Data Team Data Management Page" 説明"データ管理ページは、エンタープライズデータプラットフォームと関連する活動の管理、セキュリティ、および統治に関する内容をカバーしています。"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .hidden-md .hidden-lg}。

データマネジメントページへようこそ

このページでは、Enterprise Data Platformおよびそれに関連する活動の管理、セキュリティ、およびガバナンスに関するプラクティスとポリシーについて説明します。Enterprise Data Platformの技術的なコンポーネントは、GitLab Tech Stackに記載されています。

データセキュリティの実践

エンタープライズデータプラットフォームは、多くのシステムから収集したデータを取り込み、処理し、保存します。すべてのデータが同じ重要性を持っているわけではありません。私たちはCriticalSystemTierフレームワークとDataClassificationStandardを使用して、どのデータが最も重要で、どのようにセキュリティを確保するのが最適かを判断しています。

Sisense

Sisenseでロールベースのデータアクセススキームを導入してい

ます。ユーザーアクセスはOktaで管理 ロールとスペースでデータアクセスを管理

各ユーザーにはSisense Roleが割り当てられ、ダッシュボードやレポートへのデータアクセスが可能になります。

• SisenseスキーマはSnowflakeデータアクセススキーマと相互に作用し、ユーザーがどちらのシステムからもデータへの「バックドア」を持つことがないようにします。

2021/11/8

•

٠

٠

追加のコントロールは以下の通りです。

- 未使用のダッシュボードをアーカイブ化
- システムへのアクセスはAPIキーで管理

## スノーフレーク

役割ベースのデータアクセススキームをSnowflakeに導入します。

- ユーザーアクセスはOktaで、アクセスリクエストはGitLabで管理されています。
- 各ユーザーには、さらに1つ(職能に応じた役割)[/handbook/business-technology/data-team/platform/#snowflakepermissions-paradigm]が割り当てられ、この設定はPermifrostで管理されます
- SnowflakeスキームはSisense Data Accessスキームと相互に作用し、ユーザーがどちらのシステムからもデータへの「バックドア」を持つことがないようにします。

追加のコントロールは以下の通りです。

- データの分類基準に基づき、データはデータベースとスキーマで管理される すべてのクエリ/
- ユーザー/プロセスには、事前に定義されたウェアハウス、またはコンピュートリソースが割り
- 当てられる(パスワードはローテーションされる) [/handbook/business-technology/datateam/platform/#passwords].

# 一般的なデータ・セキュリティ・コントロール

- データコントロールを定義する目的では、Enterprise Data PlatformはTier1システムです。
- 重要: お客様のプライベートREDデータは、Enterprise Data Platformに永久保存することはできません。

コントロール	RED	ORANGE	YELLOW
一般的なデータ管理			
データ登録リスト	必須	必須	
休息時の暗号化	必須	必須	必須
トランジットでの暗号化	必須	必須	必須
プライバシーの見直し	必須	推奨	必要なし
データ保持の手順	必須	推奨	必要なし
データインフラコントロール			
マルチファクタ認証	必須	必須	必須
役割ベースのアクセス	必須	必須	必須
アクセスロギング	必須	必須	推奨
データウェアハウスコントロール			
四半期ごとのSnowflakeユーザー監査	必須	必須	必須
四半期ごとのSiSenseユーザー監査	必須	必須	必須
四半期ごとの変更管理レビュー	必須	推奨	必要なし
四半期ごとのREDデータスキャナー	必須	N/A	N/A
エンドポイントデバイス			
マルウェア対策	必須	必須	必須
フルディスク・エンクリプション	必須	必須	
四半期ごとのデータパージ	必須	必須	

データインフラ:データウェアハウスの一部としてデータにアクセスまたは処理し、エンドユーザーがデータを利用できるようにするためのあらゆるシステムが含まれます。

- データウェアハウスの制御。エンタープライズ・データ・ウェアハウスは、ティア1システムです。
- エンドポイント・デバイス。データウェアハウスにアクセスするすべてのエンドポイントは、Tier 1に分類されます。

四半期ごとのデータヘルス・セキュリティ監査

四半期ごとの監査では、適切な担当者が正しいデータアクセス設定を行い、データパイプラインが正しく実行されているかなど、システムのセキュリティを検証します。

このプロセスは、「四半期データの健全性とセキュリティ問題のテンプレートによっ

てサポートされています。ここでは、活動のチェックリストのサンプルを紹介します

ロスノーフレーク

- Snowflakeから離脱した従業員の活動を停止させる
- オフボーディングされたGitLabチームメンバーのすべてのSnowflakeアカウントは、オフボーディングされた日から無効化されるべきです。このアクティビティは、オフボーディングされたGitLabチームメンバーのアクティブなアカウントがあるかどうかをチェックします。その後、アクティブなアカウントはすべて非アクティブになります。
- 監査を実施した時点から過去60日以内にログインしていないアカウントをSnowflakeから無効にします。
- 。 60日以上ログインしていない指定ユーザーのSnowflakeアカウントは無効化されます。非活性化後、ユーザーには その旨が通知されます。GitLabチームメンバーが再びアクセス権の付与を希望する場合は、通常のARを作成する 必要があります。マネージャーの承認後、アカウントが有効になります。
- すべてのユーザーアカウントに多要素認証が必要であることを確認する。

• Sisense

- Sisenseからオフボーディングされた社員を無効にする。
- オフボーディングされたGitLabチームメンバーのすべてのSisenseアカウントは、オフボーディングされた日から無効化されるべきです。このアクティビティは、オフボーディングされたGitLabチームメンバーのアクティブなアカウントがあるかどうかをチェックします。その後、アクティブなアカウントはすべて非アクティブになります。
- 監査を実施した時点から過去60日以内にログインしていないアカウントをSisenseから無効にします。
- 。 **60**日以上ログインしていない**Sisense**アカウントは無効化されます。無効化された後は、ユーザーに通知されます。 **GitLab**チームメンバーが再びアクセス権の付与を希望する場合は、通常の**AR**を作成する必要があります。マネージャーが承認すると、アカウントが有効になります。
- すべてのユーザーアカウントに多要素認証が必要であることを確認する。
- □信頼できるデータ
  - Golden Record TDテストをすべて確認し、合格していることを確
  - 認する。データサイレンを確認し、REDデータの存在を確認する。
  - TD: Sales Funnel ダッシュボードへのビジネスロジック変更のレポートを四半期ごとに作成する。新しいディメンション、新しいファクト、新しいマートの追加、結合の変更、新しい計算フィールドの追加などのビジネスロジックです。

layout: handbook-page-toc title: "Data Quality" 説明 "GitLab Data Quality Programは、GitLabの生産性や効率性に影響を与えるデータ品質の問題を特定、監視、修正することを目的としています。"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .hidden-md .hidden-lg}。

データクオリティープログラム

GitLab Data Quality Programは、データの継続的な改善を通じて、GitLabの生産性と効率性を向上させることを目的としています。

このプログラムは、データの問題を特定し、監視し、修正するための信頼性の高いプロセスを構築することで、この目標を達成します。データクオリティーの範囲はGitLabが管理するすべてのデータあり企業内データの利用可能性にのみ制限さます

EDWは、大規模かつ複数のソースシステムにまたがる大量のデータをスキャンし、データ品質の問題を検出する機能を備えているため、データウェアハウス。

データ品質に関する問題をお持ちですか? Data Quality Issue Workflowをご覧ください。{: .alert .alert-success}。

## データ品質問題の種類

従来のデータ品質プログラムでは、品質問題を、完全性、正確性、一貫性、妥当性、一意性、完全性など、いくつかのタイプに分類しています。このようなニュアンスは、データ品質の専門家ではない人と接するときに混乱を招くことがあります。この問題を単純化するために、GitLab データクオリティプログラムでは、以下のデータクオリティ問題の種類を認識しています。

- 不正確なデータInaccurate Dataとは、表すべき実世界の値と一致しないデータのことです。不正確なデータの例としては、3桁の米国郵便番号があります。
- 欠損データ欠損データとは、本来存在するはずのNULLまたは空のフィールドやレコードのことです。欠損データの例としては、住所録の郵便番号がNULLであることが挙げられます。
- 重複したデータ重複データとは、同じデータが繰り返されるべきでないときに繰り返されるデータのことです。重複 データは、データの報告方法に基づいて自然に発生することがあるため、特定するのが難しい場合があります。重複 データの例としては、CUSTOMERマスターテーブルにある2つの(ほぼ)同一の顧客レコードが、両方とも1人の「現 実世界」の顧客にリンクされている場合が挙げられます。

## データ品質システムの構成要素

データ品質システムは、時系列で問題点を把握するスコアカードと、データの既知の問題点を特定する検出ルールで構成されています。

- データ品質スコアカード データ品質スコアカードは、データカスタマおよびデータクリエータが使用するダッシュボードです。ダッシュボードには、サブジェクトエリアの個々の検出ルールのステータスによって測定された、サブジェクトエリアの全体的な品質が表示されます。特定の目的のために、特定の独立したデータ品質スコアカードを作成することができ、今後も作成される予定です。例えば、当社は「データ品質スコアカード-製品使用データ」を積極的に開発しており、Zuora課金システムの品質を測定するために、別の「データ品質スコアカード-Zuora」を開発する予定です。
- データ品質検出ルール データ品質検出ルールは、フィールドや行のデータの品質を事前に定義された条件と比較してチェックするSQLベースのテストです。検出ルールを実行するには、データがすでにEnterprise Data Warehouseに存在している必要があります。検出ルールは列挙され、1つのSQLステートメントには1つのテストしか表現されません。検知ルールの例は以下の通りです。
  - 検知ルール1: 不正確なデータ Account Location レコードのStateフィール
  - ド 検知ルール2:重複したデータ アカウントマスターレコードのアカウント名
  - 検出ルール3:データの欠落-新しいUsage Pingの送信にはライセンスキーが必要です。

### 運用プロセス

毎週、検出ルール「バッチ」が実行され、出力が永続的なテーブルに保存される。永続的テーブルには、実行日、検出ルール 識別子、トランザクションIDが含まれ、ソースシステムへのリンクが可能になる。この永続的テーブルは、スコアカードを作 成する際のベースとなります。

### データ品質の問題を解決する

リメディエーションとは、品質問題を解決、修正、または排除するプロセスのことです。修正は「データ作成者」、つまりソースデータの作成に責任を持つ人やチームが担当します。品質問題を特定すること、または特定を支援することは、「データカスタマー」の責任である。

製品データ品質スコアカード

目的 - Product Data Quality Scorecard は、製品使用データに関するデータ品質の問題を定量化します。スコアカード

・ダッシュボードには、以下の情報を表示するビジュアライゼーションがあります。

• 製品のデータ品質検出ルールのそれぞれの合否の割合。合格したレコードの割合は、条件またはデータ品質検出ルールを満たしたレコードの総数の割合を取ることで計算されます。計算に使用される式は以下の通りです。
\*((passed\_record\_count/processed\_record\_count)100)

同様に、「失敗したレコードの割合」は、条件やデータ品質検出ルールを満たすことができなかったレコードの総数に対する割合を取ることで算出されます。\*((failed\_record\_count/processed\_record\_count)100)

検知ルールの合格と不合格は、合格したレコードの割合としきい値制限の割合を比較して決定されます。現時点では、しきい値は50に設定されています。正確なしきい値はDRIが決定する必要があります。

通過したレコードの割合がしきい値を超えた場合、検出ルールのステータスは緑になります。例えば、通過したレコードの割合が72%(50%以上)の場合、72%のレコードがデータ検出ルール/条件を満たしていることを意味します。

通過したレコードの割合がしきい値以下の場合、検出ルールのステータスは「赤」になります。例えば、合格したレコードの割合が40%(50%未満)の場合、60%のレコードがデータ検出ルール/条件を満たしていないことを意味します。これらのレコードには注意が必要で、データはソースチームによって修正される必要があります。

• トレンド分析チャートは、1週間の間の各データ品質検出ルールの合否率の変化を示します。

トレンド分析チャートのデータポイントをダブルクリックすると、**TD: Product Data Quality Detailed Dashboard V1.0**が表示され、各検出ルールで影響を受けた**ARR**の合計値とともに、データ行と集計値の詳細情報が表示されます。

• 各日の集計結果には、Daata品質検出ルールのそれぞれについて処理された行の総数と、ルール/条件を満たす (パスする ) 行の数、ルール/条件を満たさない (フェイルする) 行の数が、ルール実行日によって追跡される各日について表示されます。

ダッシュボード「TD: Product Data Quality Scorecard - Overview V1.0」と「TD: Product Data Quality Detailed Dashboard V1.0」のデータは、毎日更新されます。

クイックリンク

**TD: Product Data Quality Scorecard** 

- Overview V1.0

**TD: Product Data Quality Detailed** 

**Dashboard V1.0** 

現在、製品の使用状況データについて確認されているデータ品質検出ルールは以下の通りです。

検出ルールID	ルール	説明DRI
1	ホスト名のインスタンスタイプがない	
2	サブスクリプションIDがないライセンス	
3	ライセンスが見つからないサブスクリプション	
4	ライセンスの開始日が将来になるセルフマネージドプランのサブスクリプション	
5	ライセンス開始日がライセンス終了日よりも長いセルフマネージドプランのサブスクリプション	
6	過去に購読終了日を迎えた期限切れのライセンスID	
7	名前空間IDがないSaaSサブスクリプション	

追加リソース

ガイド&ブック

- データの前に立つ
- 非侵襲的データガバナンス
- データ・ライフサイクル・マネジメント (DLM) または同等の堅牢なアプローチ。

#### SaaSツール

FivetranもStitchも、マネージドサービスであるため、独自のデータ品質チェックを行っています。これらのデータの抽出時 に問題が発生した場合は、ベンダーのサポートチームにご相談ください。

#### カスタム

サードパーティのツールからの抽出については、独自のAPI接続やAirflowで管理しています。これらのすべてが、データ品質チェックの ためのメソッドを備えているわけではありません。

#### **BambooHR**

BambooHRのデータはカスタムコードで抽出しています。データの品質チェックには、API からの 200 レスポンスの確認や、JSON データに最低数のレコードが存在することが含まれます。

#### Postgres パイプライン

私たちの Postgres\_Pipeline (gitlab.com、customers.gitlab.com、license.gitlab.com、version.gitlab.com からのデータを扱う) は、送信元と送信先のデータベースの行数が一致しているかどうかをチェックします。

トランスフォーメーション・データ・クオリティ

私たちは、倉庫内のすべての変換にdbtを使用しています。私たちは、すべての新しいdbtモデルのテストを要求し、必要に応じて定期的にテストを更新します。これらのテストは、抽出テストやチェックと同様に、上述のデータ品質哲学に沿って書かれるべきです。

## データ パイプライン ヘルス ダッシュボード

## 課題を見る

第1回目のイテレーションでは、以下を中心に追加しました。

- キーテーブルの日々のレコードの挿入と更新の速度をテストする SQL ステートメント (行数テス
- ト) キーテーブルのキーフィールドの集計値をテストする **SQL** ステートメント (列値テスト
- キーテーブルのキーレコードの存在をテストするSQL文(ゴールデンレコードテスト
- これらの結果をシンプルに可視化するワイヤーフレームのダッシュボード

layout: handbook-page-toc title:"Data Team Services" description:"Data Team Services"

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

スペシャルイベントのオンコール

Special Event On-Callは、データチームが提供するサービスで、エグゼクティブ・プランニング・セッションなどの短期間の特別なイベントに必要な迅速な対応能力を提供します。オンコールサービスは、サポートに多大な労力と調整が必要なため、ほとんど利用されていません。データチームは、既存の仕事やスタッフの空き状況により、オンコールのリクエストに応えられない場合があります。

どのような仕組みになっているのでしょうか?

オンコールのリクエストが確認されると、データチームはイベント専用のSlackチャンネルを設定します。イベント参加者は、このチャンネルを使ってリクエストを送信し、レスポンスを受け取ります。解決までの時間は、リクエストの複雑さによって異なります。データチームは新しいリクエストに対して、(a)リクエスト完了までの推定時間、(b)トラッキング用の課題へのリンクを返信します。データチームがリクエストの解決に長時間を要する場合は、リクエスト者に通知します。

イベント終了後、データチームは次のことを行います。

- 作成された課題のリストとそのステータスを含むイベントサマリーレポートの提供
- → フォローアップが必要な課題の特定
- Slackチャンネルの閉鎖

オンコールサービスの依頼

- 1. イベントの10営業日前までに、DataTeamプロジェクトで課題を作成します。
- 2. 以下の内容を「説明」に追加します。
  - 1. イベントの開始日と終了日
  - 2. 毎日の開始時刻と終了時刻
  - 3. 取り組むべき質問の種類と主題カテゴリ(財務、販売、成長など)について
  - 4. 想定される質問の総数
  - 5. Slackチャンネルを使用しているリクエスターの推定数
- 3. 詳細が確定したら、データチームがスタッフのモデルを作成します。
- 4. データチームは、オリジナルの問題でお客様に回答します。
  - 1. コミュニケーションに使用するSlackチャンネル
  - 2. スタッフの配置計画(名前とカバー時間を含む)(個人にDMを送らず、すべてのコミュニケーションを確立されたSlackチャンネル内で行ってください。

社内データチームのプロセス

- 1. データチームのマネージャーは、オンコールのリクエストを確認し、既存の仕事の都合やスタッフの空き状況を考慮して、サポートできるかどうかを判断します。
- 2. リクエストに対応できない場合は、データチームのマネージャーがリクエスト者に対応し、代替のサポートプランを 共同で検討します。
- 3. リクエストがサポートされる場合は、データチームのマネージャーが行います。
  - 1. 名前とカバー期間を含むスタッフモデルの作成
  - 2. 影響を受けるデータチームのメンバーに対して、人員配置計画とカバー率の要件を調整する
  - 3. 影響を受けるプロジェクトの関係者に、それぞれのプロジェクトへの潜在的な影響を通知する
  - 4. 新しいSlackチャンネルを「data-oncall--」という形式で作成します。
  - 5. オリジナルのIssueをSlackチャンネル名とスタッフのモデルで更新する
  - 6. スタッフ配置計画とローテーションスケジュール(名前を含む)をslackチャンネルにピン留めする
  - 7. 原稿依頼者への機種・プランの確認

オンコール・イベントサポートの流れ

データチームのオンコールサポートプロセスでは、コミュニケーションを合理化し、責任を明確にするためにティアリング方式を採用して います。

- ティア1メンバーは最初の対応者であり、イベントSlackチャンネルの監視、リクエストの確認、リクエストへの対応、課題の作成、可能であれば課題の解決を行います。
- Tier2メンバーは、Tier1メンバーからの要請に応じて、Tier1メンバーが対応できないリクエストをサポートします

ティア1とティア2のメンバーは、オンコールイベント中も定期的に連絡を取り合う必要があります。休息が必要なTier 1チームメンバーは、Tier 1の責任を引き継ぐ資格のあるチームメンバーを見つける責任があります。

## ティア1の責任

- 1. 事前に設定した取材期間中に、イベントSlackチャンネルをモニターします。
  - 1. 新規のリクエストに対しては、投稿後30分以内に「Request Received, We Are Working On It」のような返答をする。

- 2. 最初のリクエストから1時間以内に、より正式な回答を調整して伝え、以下を含む。
- 3. リクエスト者に代わって作成される新しい課題へのリンク(まだ作成されていない場合
- 4. 問題解決に要する時間の目安 (例:「1時間以内に回答します」など
- **5.** 問題が特に複雑であったり、解決に数時間を要する可能性がある場合には、明確なメッセージを表示します(例: 「このリクエストはすぐに対応することが難しく、解決には時間がかかります」)。
- 2. 受信したリクエストは可能な限りすべて解決してください。それができない場合は、Tier2のメンバーに接続して助けを求めてください。
- 3. 定期的に、Issueコメントを使ってリクエスト者に最新情報を提供し、特に緊急性が高い場合は、イベントのSlackチャンネルに最新情報を提供します。

#### Tier2の責任

- 1. ティア1メンバーのサポートができること。
- 2. 定期的にイベントのslackチャンネルをチェックし、全体の状況を把握してください。
- 3. 長時間の休息が必要な場合は、Tier1メンバーにアラートを出す。

layout: handbook-page-toc title:"データチームの方向性"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

このページには、将来の見通しに関する内容が含まれており、現在の状態または予定されている機能セットや能力を正確に反映していない可能性があります。

## 四半期目標

## FY22-Q1

FY22-Q1は、データチームにとって非常に重要な四半期となりました。これは、フュージョンチームによる新しい組織モデル、Trusted DataソリューションとSales Funnel Dashboardの最初のメジャーリリース、製品使用状況データのGainsightへの統合など、FY21の初期に開始された多くのイニシアチブをチームが完了したためです。当チームは、当四半期にすべての主要なコミットメントを成功裏に達成し、第2四半期に向けて成功と規模拡大のための体制を整えることができました。

- FY22-Q1Retrospective
- FY22-Q1ReportCard
- FY22-Q1 Handbook Jamboree

		ThemeObjectiveNotes and Key Results
データプロ グラムの成 熟	目的1: FY22フュージョンチームの定義、Project Compassのサポート、アンケートによる社内関係者の満足度の報告、優先順位と成果物の整合性の確認。	Fusion Teams、XMAU Handbook Page、Gainsight SM Usage Data、Marketing Data Martを展開しました。
レベル2-コ ーポレート ・インテリ	目的2:GTMパフォーマンス管理の効率化	Shipped TD: Sales Funnel dashboard, Shipped Analytics Hub MVC
レベル <b>2</b> - プラットフ ォームの成 熟度	目的 <b>3</b> :エンタープライズデータプラットフォームの信頼性と 安全性の向上	四半期監査のv1を作成、Data Pump v1を出荷

### FY22-Q2

# FY22-Q2DataTeamHandbookUpdates

- \* FY22-Q2 Retrospective Coming Soon
- FY22-Q2 Report Card Coming Soon

FY22-Q2 Handbook Jamboree Coming Soon

#### テーマ目的

Level 2-Corporate IntelligenceObjective 1: Accelerate GTM Teams and Project Compass

**Level 2-Corporate** 

IntelligenceObjective 2: Enable GitLab to be Public

#### Company

Level 2-Platform MaturityObjective 3: Improve Product Decision Quality and Accuracy with Trusted Data(信頼できるデータによる製品決定の質と精度の向上

FY22-Q3

テーマ目的

Level 2-Corporate IntelligenceObjective 1: Accelerate GTM Teams and Project Compass Level 2-

Corporate

IntelligenceObjective 2: Public Company Readiness

Level 2-Product & Customer IntelligenceObjective 3: Accelerate R&D Teams with Trusted Data

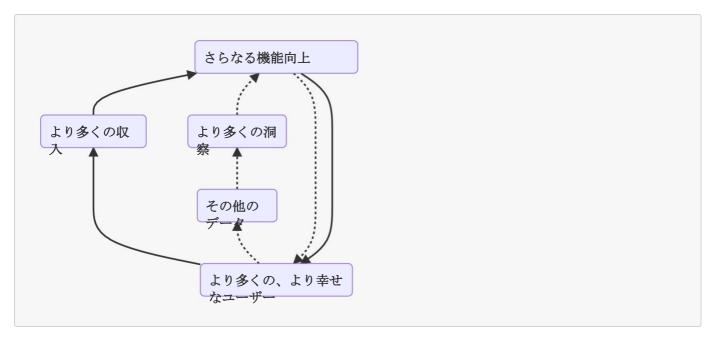
## 戦略

ミッションを達成し、責任を果たし、GitLabが成功した上場企業になるための重要なステップとして、私たちはEDP(Enterprise Data Platform)を構築しています。EDPは、単一の統一されたデータおよび分析スタックであり、セルフサーブデータやデータクオリティーなどの幅広いデータプログラムを備えています。EDPは、GitLabのKPI、部門横断的なレポートや分析、さらにはチームメンバー全員が信頼できるデータを使ってより良い意思決定を行えるようにします。EDPは、データパブリッシングやプロダクトなどの機能により、GitLabの分析能力をさらに加速させます。これらの機能は、ビジネスシステムやGitLab製品に統合され、お客様が使用するために強化され、集約されたデータを提供します。この加速は、GitLabのオープンコアやデベロップメントスペンドのフライホイールのように、「データフライホイール」の開発によって行われます。

データフライホイール

カスタマー&プロダクトインテリジェンス フライホイール

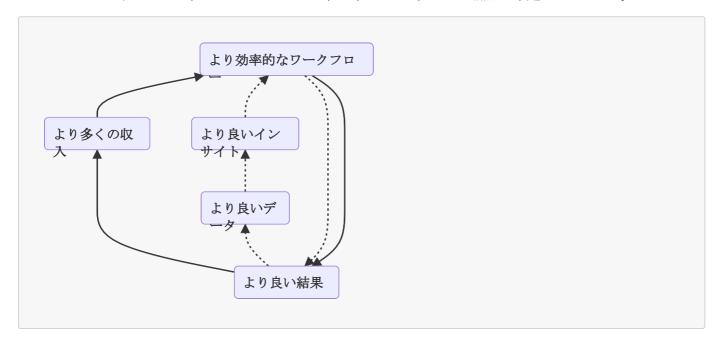
カスタマー&プロダクトインテリジェンスフライホイールは、カスタマーエクスペリエンスの向上に重点を置き、ユーザーと 製品のインタラクション、カスタマーユースケース、製品開発、製品採用、およびカスタマージャーニーのほとんどの側面に 関わるデータと分析を網羅しています。



コーポレート・インテリジェンス

コーポレート・インテリジェンス・フライホイールは、(社内の)ビジネス効率の向上に焦点を当てており、これはビジネスワークフローの計測、監視、および改善によって達成されます。コーポレート・インテリジェンス・チームの一般的なアウトプットは以下の通りです。

パフォーマンス・ダッシュボード、バランス・スコアカード、KPI、MBOなど、データを活用した関連フレームワーク。



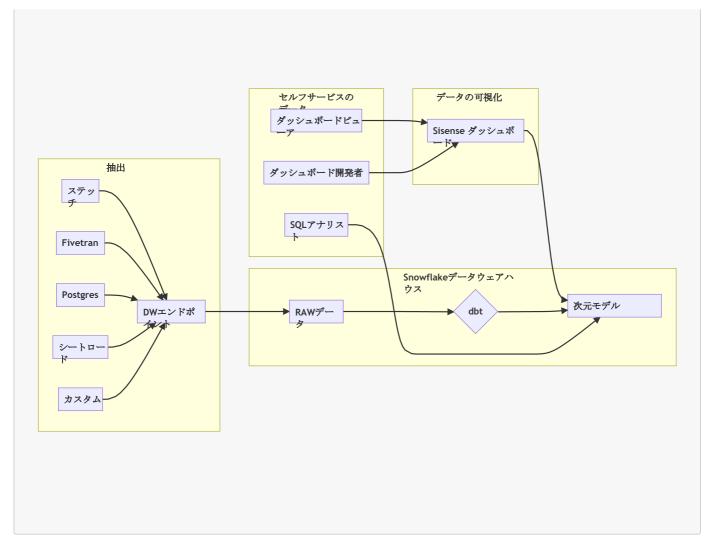
### 短期的な方向性

私たちの短期的な目標は、GitLabのKPIをサポートし、GitLabの最も重要な2つのビジネスプロセス(リードからキャッシュまでのサイクルと製品リリースから採用までのサイクル)をカバーするSelf-Serve Dataを提供することで、データケイパビリティモデルのレベル2に到達し、GitLabを公開企業にすることです。

優先順位の高い順に、私たちが積極的に構築しているEDPの要素は以下の通りです。

- 1. SisenseのダッシュボードやSnowflakeでのSQL分析など、レポートや分析のための唯一の情報源となるエンタープライズ次元モデルです。
- 2. 信頼できるデータフレームワークにより、すべてのデータソリューションが徹底的にテストされ、自動化された継続 的な検証が行われます。
- 3. KPIを含むパフォーマンス分析をサポートするためのデータ・ビジュアライゼーション。
- **4.** セルフサービスデータプログラムは、すべてのGitLabチームが信頼できるデータに確実にアクセスできるようにする ためのプログラムで、データチームの関与は必要ありません。セルフサービスデータプログラムには、次の**3**つのオ プションがあります。
  - 1. ダッシュボード・ビューアー GitLabチームのメンバーがSisenseポータルにログインし、事前に作成されたダッシュボードにアクセスできます。
    - 1. ダッシュボード開発者 Sisenseの認定を受けたGitLabチームメンバーがSisenseポータルにログインし、エンタープライズディメンショナルモデルで利用可能な信頼できるデータをソースとした独自のダッシュボードを構築します。
    - 2. SQL Developer SQL Analysisの資格を持つGitLabチームメンバーが、Snowflakeでホストされている Enterprise Dimensional Modelから信頼できるデータを照会するための独自のSQLを作成します。

merged.md	2021/11/8



### 長期的な方向性

当社の長期的な方向性は、マスターデータ管理、データレイク、高度な分析など、成熟したエンタープライズデータプラットフォームに見られる機能をEDPに追加することで、「年単位」で測定されます。また、レベル2に到達した後は、次のようになります。

- オープンソースのデータプロジェクトに貢献する方法を探したい
- データのパイプラインや処理コンポーネントとしてMeltanoを使用したい。
- EDPの一部をGitLab.comと統合し、GitLabの顧客に詳細な分析機能を提供したい OktaのBusiness@Workのような
- DevOps業界のベンチマーク・レポートを提供したい
- レベル3に到達するために必要な要素を確保するために、データ技術スタック全体を再検討する。

## 成功の測定

短期的な方向性に対する進捗状況を以下の方法で測定します。

- 1. データチームのKPI
- 2. データバリューピラミッドに沿った当社の結果のビジネスインパクト
- 3. 当社が提供するデータ機能は、データ機能モデルに対応しています。
- 4. The Data Team Quarterly Report Card

長期的な進捗を測る基準はまだ定めていません。データチームのKPI

- 1. 過去最高のセルフサービス・データ顧客数を実現
- 2. 月間 アクティブなセルフサービス・ダッシュボード開発者の数
- 3. 月間アクティブなセルフサービス型SQL開発者の数
- 4. ダッシュボードのトラフィックに占めるユーザー生成コンテンツの割合(月間

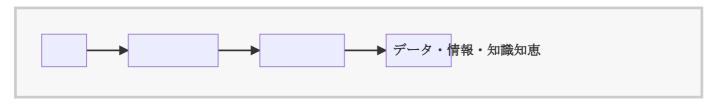
# 完全なエンタープライズデータプラットフォーム

次の表は、大企業が必要とする広範なデータと分析を解決できる、成熟したエンタープライズデータプラットフォームの機能を表しています。記載されているすべての機能が、GitLabの短期的なニーズや既知の長期的なニーズを満たすために必要なわけではありません。特定の機能を実装するかどうかは、明確なビジネスニーズによって決定され、最終的な結果は参照例とは大きく異なる可能性があります。

データアーキテクチャ	データセキュリティ	データ品質
記述式	診断	アドバンスド・アナリティク ス
報告	ダッシュボード	セルフサービス
運用データストア	データウェアハウス	データレイク
データモデル規格	エンタープライズ次元モデル	データマート
リファレンスデータ管理	データエンリッチメント	マスターデータマネジメント
データパイプライン	データトランスフォーメーション	リアルタイムデータ
データのエクスポート	データ出版	データ製品
データタクソノミー	データカタログ	データポータル

## データバリューピラミッド

私たちは、すべてのGitLabチームがデータバリューピラミッドを上(下の図では左から右)に移動し、基本的なメトリクスやカウントを知恵に変えて、お客様のためにより良い製品を作り、ビジネスをより効率的に運営し、ビジネスモデルに新しい機能を追加するのを支援したいと考えています。データバリューピラミッドでは、現在、主に「データと情報」の段階で仕事をしています。



## データケイパビリティモデル

データケイパビリティモデルは、ギットラボの戦略をサポートするための目標状態の要件を特定するために使用されます。

GitLabが上場企業になるためには、Lead-to-CashとPublic-Facingの指標がケイパビリティモデルのレベル2に達する必要があります。

		レベル特性ベネフィット
(5)Prescriptive	製品に組み込まれたリアルタイムの複合的な分析、行動や認識の形成 、データ分析は戦略的な差別化要因である。	新しいデータ製品、意思決定 のROI向上
(4)Predictive	データサイエンス」 起こりそうなことを洞察し、広めていく とエフォートレスアナリティクスの生産、エンタープライズデー タの品質とガバナンス	信頼できる顧客生涯価値。 拡大・解約予測、製品組み込み 型アナリティクス
(3)Strategic	広範囲かつ容易にドリルダウン可能な分析、ドリルダウン可能 なクロスファンクショナル・スコアカード、ダッシュボード、エ ンタープライズ・データ・ウェアハウス	カスタマー <b>360</b> &ヘルススコア 予測可能で信頼できるデータレポーティング、堅牢なセルフサービス、データ・アット・スケール
(2)Advanced: 参考文献	運用可能な自動化されたレポートとダッシュボード、自動化され たテストによる信頼性と検証されたデータ、手動および自動化さ れた統合の混合、一部のデータサイロを含むコア統合データ	信頼できるデータ、セルフサービスのデータ、重要なパフォーマンス指標、拡張の ための安定したプラットフォーム

レベル	/特性べ	ネフ	1	ッ	ኑ

(1)リアクテ	静的なリストやレポート、履歴やラグ(過去30日、90日、365日)	ヒストリカル・タブラル・レポート、データ
ィブ	に強い関心がある、予測不可能な速度、最小限のクロスファンクショナル分析、データサイロ	・ビジュアライゼーション
(0)なし	レポート作成に一貫性がない、結果が広く信頼されていない、 安定した分析インフラがない	

layout: handbook-page-toc title:"データチャンピオンプログラム" 説明"Discover how GitLab uses a Data Champion Program in concert with the Data Team to promote data literacy and acumen"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

このページには、将来の見通しに関する内容が含まれており、現在の状態または予定されている機能セットや能力を正確に反映していない可能性があります。

## 概要

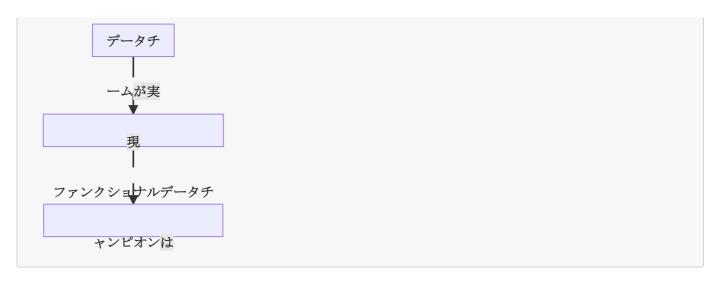
データチャンピオン・プログラムは、既存のチームメンバーで組織全体のデータ能力を開発するための最新のソリューションです。すべてのデータニーズを中央のデータチームに依存するのではなく、データチャンピオンプログラムは、中央のデータチームの専門家が作成した専門知識、コンテンツ、プログラムを活用して、機能別チーム内でデータを開発するインサイドアウト、トレーナー育成のアプローチを採用しています。

データ・チャンピオン・プログラムの目標は以下の通りです。

- データを会社の重要な資産として活用し、全員がデータを使って最善の意思決定ができるようにする。
- セルフサービス・ダッシュボードの開発やSQL分析など、機能チームがデータを自給自足できるようにするためのパスの作成
- データ駆動型の技術やアプローチの導入を促進するデータの質とデータの説
- 明責任を向上させる
- 活発なDataコミュニティの発展に貢献する

## アプローチ

データチャンピオン・アプローチでは、各機能別チームの中から、(a)そのチームのデータDRIとして、(b)データチームの安定したカウンターパートとしての役割を果たす個人を特定します。データチームのリーダーとデータチャンピオンは共同でデータサクセスプランを策定し、ビジョン、責任、目標を明確にし、機能チーム独自のデータニーズ、適性、既存のデータ知識を考慮します。この計画では、GitLabデータ資産(レポート、ダッシュボード、データモデル)、セルフサービス機能、データロードマップをカバーするトレーニングスケジュールを確立します。データチャンピオンは、チーム内でのデータ能力の向上を支援するために、地域のデータコンピテンシーセンターを通じてデータプラクティスを促進する責任があります。



ファンクショナルチームメンバー

# データチャンピオンの責任

- 1. データチームとの連携によるデータサクセスプランの策定
- 2. 機能別チームのデータDRIとして、データ品質問題を含むすべてのデータ要求を把握する。
- 3. データイシューとエピックの優先順位付けと積み重ねを行い、データイシューボードを最新の状態に保つ。
- 4. イシューボードを通じてデータチームに優先事項を伝える
- 5. プロダクトマネージャーのためのデータ」など、各チームの対象者に合わせたデータハンドブックページの作成
- 6. データに関するニュースやプラクティスを定期的に伝え、カスケードしていく
- 7. ソースシステムのデータ品質向上に貢献
- 8. ソースシステムの変更を監視し、必要な更新をEnterprise Data Warehouseに組み込む手助けをする

### 要求事項

- 週に1~2時間程度のプログラムへの参加(取り組みの進展に応じて時間を調整
- ) 個人およびチームのデータコンピテンシーの向上への取り組み 中級SQLの経
- 脳
- データビジュアライゼーションの経験
- データモデリングの経験; Zuoraの例、Salesforceの例

# プログラムターゲット

	エリア現在の状態	ターゲット
の状態データ 機能チームの 優先順位付け	Data TeamData Champion, supp	ported with Value Calculator
ファンク		
ショナル		
チームの	データチーム	データチャンピオンとデータチームの維持
データイ	) - 9 ) - 4	/ 一クリヤンにオンと/ 一クリームの無行
シューボ		_
ード		
機能チーム	Does Not Exist(	存在しない)データチャンピオンは、所属する機能別チームのデータDRIと
のデータDRI	して活動します。	行在しなり、ケーケンドンとなるでは、別層りの機能加入。これのケーケーにと
wan -		
<u>機能チームハ</u>		
ンドブックの		
データ	ページが存在し	ないデータチャンピオンとデータチームの維持

2021/11/8 merged.md

	エリア現在の状態	ターゲット
の状態通信する 機能チームへ のデータニュ ース	定期的なコミュニケーション 席し、データを渡します。	データチャンピオンは、データチームのミーティングに定期的に はありません 機能チームに関連するニュース
Sisense トレーニ ング	オンラインチュートリアルとデー タチームによる定期的なライブ トレーニング	データチャンピオンが機能チームのためにカスタマイズし、必要なトレーニングを提供します。
データの収集と処理	データチームは、すべての Handbookドキュメントを管理し 、dbtドキュメントの更新をサ	データチャンピオンは、ソースシステムの品質やモデルの変更に関する最新情報を入手し、必要な更新情報をデータチームに伝達して実装する(ハンドブックの維持を含む
データのプラ イバシーと保 護	ポートします。 プライバシー審査を担当するデ ータ、法務、セキュリティ、コ ンプライアンスの各チーム	データチャンピオンは、データ、法務、セキュリティ、コンプライアンス の各チームが、プライバシー評価を必要とする可能性のあるデータ要素を 理解し、適切な文書を維持できるように支援します。

layout: handbook-page-toc title:"データプログラムレベル2リファレンスソ

リューション"このページについて

{:.no\_toc}

TOC {:toc}

目的

このページには、将来の見通しに関する内容が含まれており、現在の状態または予定されている機能セットや能力を正確に反 映していない可能性があります。

公開企業は、主要な財務、顧客、成長の指標を確実かつ予測可能な形で共有するとともに、ビジネスパフォーマンスを継続的に改 善するために、リードタイムからキャッシュタイム、製品のアイデアから採用までのプロセスを分析する必要があります。これらの 活動は、データ機能モデルのレベル2で定義された機能によって支えられています。このページでは、現実的な例を示し、今後の開 発の参考とするために、「製品ジオロケーション分析」のレベル2データソリューションを紹介します。

ソリューション概要 - 製品ジオロケーション分析

自社の製品が世界中のどこで使われているかを理解することは、顧客や製品のグローバル展開をより深く理解するための 重要なステップであり、ロケーションを意識したインサイトを構築することにつながります。

このデータソリューションは、3つのSelf-Service Data機能を提供します。

- 1. ダッシュボードビューアー 国、地域、月、年ごとのGitLabデプロイメントを視覚化する新しいSisenseダッシュボード(セ ルフサービスダッシュボード)です。
- 2. ダッシュボード開発者 新しいダッシュボードを構築し、既存のダッシュボードをジオロケーションデータにリンクさ せるための完全な次元モデルコンポーネントを含む新しいSisenseデータモデルです。
- 3. SQL Developer 新しいEnterprise Dimensional Modelの対象領域

データプラットフォームの観点から、このソリューションは実現します

- 1. GeoLocation分析のためのEnterprise Dimensional Modelの拡張機能です。
- 2. 新しいdim\_countryテーブル
- 3. Data Pipeline Healthダッシュボードのテストおよびデータ検証拡張機能
- 4. ERD、dbtモデル、および関連するプラットフォームコンポーネント

最後に、これは過去1年間に完了した以下のようないくつかのアドホックな問題に対する長期的な自動化ソリューションで

す。顧客とユーザーのデータを国別に抽出

• 各国ユーザーデータ

• SheetloadでGeoLite2フリーデータベースをアップロード

知識評価と証明書

Self-Service Data 認定プログラムは、Learning and Development 認定プログラムに基づいています。Self-Service Dataプログラムでは、対象となるDashboard DeveloperまたはSQL Developer Knowledge Assessmentを完了した場合に、個別の証明書が発行されます。ナレッジ・アセスメントへのリンクは、以下の各セクションにあります。

データの分類

#### **ORANGE**

• IPアドレス

ソリューション・オーナーシップ

- ソースシステムのオーナー@rparker2
- ソースシステムのサブジェクト・マター・エキスパート。@jeromezng
- データチームのサブジェクト・マター・エキスパート。@rparker2

#### 主な用語

- 地域 データはGitLab Sales Territoriesを使って可視化されます。AMER, APAC, EMEA
- Country データはISOCountryによって可視化されます。
- IP-ジオロケーションマッピング IPアドレスを地図上の地理的位置にマッピングすること。

キーメトリクス、KPI、PI

- 国または地域別、月別、年別の使用Ping数
- 国または地域別の GitLab.com ページビュー数 (月別、年別) KPIs 定義な
- [
- PI 定義されていない

セルフサービス・データ・ソリューション

セルフサービス・ダッシュボード・ビューアー

ダッシュボード目的

世界の製品の 成長国、地域、時間ごとのGitLabの採用状況を可視化します。

Geolocation DataのData Health Dashboardこのソリューションをサポートする
ために使用されるGeolocationデータの データヘルス。

セルフサービス・ダッシュボード開発者
データスペース

Global

以下に示す「製品ジオロケーション分析」モデルと1-1の関係を持つデータモデルを含みます。

セルフサービス・ダッシュボード開発者認定証

証明書を取得するには、Self-Service Dashboard Developer Knowledge Assessmentで100%の評価を得て、新しいダッシュボードのスクリーンショットをアップロードする必要があります。Knowledge Assessmentを完了すると、回答がメールで送信され、このメールが証明書となります。

セルフサービスのSQL分析

- データはusage-pingとsnowplowから取得しています。
- Usage-pingは、ホストのユーザーがどこから来るかではなく、インスタンスがどこにHOSTされているかについての情
- 報を含みます。Snowplow は、ユーザーがどこから来ているかについての最良の情報を含んでいます。ユーザーの中にはVPNや同様のプロキシ・ソリューションを使ってGitLabを導入している人もいるので、これらのソースから得られる IPとジオロケーションのマッピングは正確ではないかもしれません。しかし、Snowplowは私たちが持っている最高のデータソースであり、私たちの最高の真実の情報源でもあります。

エンティティ・リレーションシップ・ダイアグラム

it. Cualu	Diagram/Ent	目的	キーワード
ityGrain 製品	Activity By	分析に使用できるディメンションとファクト	dim_date, dim_country,
<b>分析</b> ロケーション	目	国別、地域別、時間別の <b>GitLab</b> 利用状況。	fct_country_activity_by_day
	DIM_DATED	すべての日付の中心となるディメンションです。	
AY		<b>CountryCentral</b> すべての国と地域のためのディメン	ションです。
	dim_countryIS ISC	O-3166とGitLab Sales Territoriesをソースとしています。	

リファレンスSQL

SisenseやdbtのすべてのプロダクションSQLは、読みやすさと保守性のためにSQLスタイルガイドを遵守しなければなりません。

NORAMの日別国別ページビュー

### 2020年の国別トップ100ネームスペース

#### セルフサービス型SQL開発者認定証

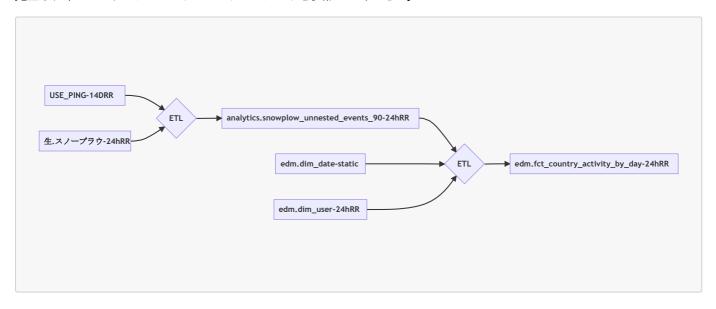
証明書を取得するには、Self-ServiceSQLDeveloperKnowledgeAssessmentで100%の評価を得て、新しいSQL文のスクリーンショットをアップロードする必要があります。知識評価が完了すると、回答がメールで送信され、このメールが証明書となります。

データプラットフォームソリューション

全体的なソリューションは、当社のEnterprise Dimensional Modelガイドラインに準拠しています。

データの系統

完全なリネージグラフについてはdbtのドキュメントを参照してください。



DBTソリューション

dbtソリューションは、RAWのソースデータから次元モデルを生成します。

トラステッドデータソリューション

エンタープライズ次元モデルの検証

## バリデーション 予想される結果

3		別)	۰
2	APACからのトラフィックの割合がAMERを超えないこと。		
1	マップされた国の総数が300を超えないこと。		

ソースデータパイプラインの検証

## バリデーション 予想される結果

1	直近14日間で、新しいUsage_pingデータがアップロードされました。
2	<b>Usage_ping</b> データで表されるアカウント数の合計は、 <i>予想される結果</i> よりも多い。
3	〜 除雪機のデータで表されるアカウント数の合計は、 <i>予想される結果を</i> 上回ります。

layout: handbook-page-toc title: "Self-Service Data"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}.

# 概要

このページには、将来の見通しに関する内容が含まれており、現在の状態または予定されている機能セットや能力を正確に反映していない可能性があります。

データ民主化は、データチームにとって共通の目標ですが管理すべきデータの多様性、量、速度、正確を考えると、達成するのは難しいかもしれません最終的には、効果的なデータ民主化ソリューションはすべて、データ顧客に焦点を当て、データを見つけやすく、理解しやすく、実用的なものにするソリューションを提供する必要があります。

- 簡単に見つけられること。データ お客様が自分に関連するデータを簡単に見つけられること
- が*理解しやすいこと。*データは、ビジネスに適した用語や概念を用いて、一貫した方法でラベル付けされ、整理 されている必要があります。
- アクション可能であること。意思決定や行動を支援し、結果を出すためのインサイトを提供する、または支援するデータであること。

GitLabセルフサービス・データプログラム

GitLab Self-Service Dataプログラムは、信頼できるデータを使ってGitLabがより速く行動できるようにするとともに、データチームが拡張できるように、4つの機能を提供します。

- データカタログ データカタログは、データカスタマーが対象分野に基づいてデータ定義、ダッシュボード、モデルを見つけるのに役立ちます。
- ダッシュボード・ビューアー すべてのGitLabチームメンバーは、Oktaからアクセスできる常時接続のポータルを通じて、事前に構築されたダッシュボードに*アクセスできます。*
- ダッシュボード開発者 GitLabチームメンバーで、独自のデータ可視化チャートやダッシュボードを作りたい人向け
- SQL開発者 GitLabチームメンバーで、SQLに精通しており、SQLベースの分析を行いたい人向け

## セルフサービス機能の概要

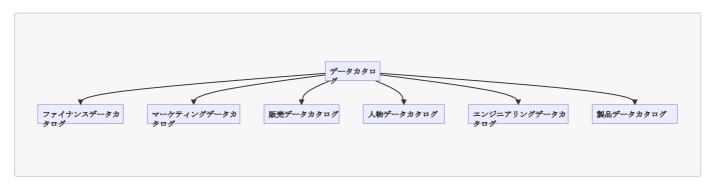
	ダッシュボードビ ューア	ダッシュボード開発者	SQLデベロッパ ー
アクセスデータカタログ	はい。	はい。	はい。
既成の <b>Sisense</b> ダッシュボードにアクセ ス	はい。	はい。	はい。
新しいSisenseダッシュボードの構築	いいえ	はい。	いいえ
SQLを書いてデータを調べる	いいえ	いいえ	はい。
認証が必要	いいえ	はい。	はい。

### セルフサービス・データ

## データカタログ

データカタログはハンドブックに掲載されており、機能ごとに整理されたデータダッシュボード、モデル、定義のインデックスが含まれています。ハンドブックの一部として、データカタログは、共通の基準と単一のソース・オブ・トゥルース(真実の)データへの調整を促進します。データチームのハンドブックセクションの一部として、データカタログは関連するコンテンツで定期的に更新されます。

#### データカタログ構造



データカタログハンドブック ページ内容



セルフサービス・ダッシュボード・ビューアー

セルフサービス・ダッシュボード・ビューアーは、GitLabチームのメンバーがSisenseで公開されている事前に作成されたダッシュボードにアクセスできるようにします。前提条件として

None

### **Access Requests:**

• □必須ではありません - すべてのGitLabチームメンバーは、オンボーディング時にアクセス権が与えられます。

セルフサービス・ダッシュボード開発者

セルフサービス ダッシュボード デベロッパーは、認定GitLabチームメンバーに、Enterprise Dimensional Modelを利用して明確に定義され、検証されたKimball標準のデータモデルに基づいて、Sisense Data Discoveryを使用して独自のダッシュボードを構築する機能を提供します。

#### 前提条件

- □データモデルの読解力
- ◆ □ レポートやチャートのデータビジュアライゼーションを開発した経験がある
- □方 レビュー Sisense トレーニング
- ロハンドブックに掲載されている「Self-Service-ready Data Subject Areas」の中から、「Self-Service Dashboard Developer Knowledge Assessment」に合格する。

アクセスリクエスト。

- □ 新規アクセス・リクエストを開き、シングル・パーソン・アクセス・リクエストを使用する
- □完成した知識評価へのリンクを含める Snowflakeへのア
- □ クセスを要求する

セルフサービスのSQL開発者

Self-Service SQL Developerは、認定GitLabチームメンバーにSnowflakeデータウェアハウスへのSQLプロンプトを提供します。現在、すべてのデータが標準的な形式でモデル化されているわけではないため、調査を行ったり、「RAW」データやモデル化されていないデータを読んだり、経験豊富なチームから学んだりするために、追加の時間が必要になる場合があります。データチームがEnterprise Dimensional Modelを構築し、SQLの例やデータモデルを含むレベル2ソリューションを追加で提供することで、SQL開発者の生産性が向上し、SQL分析がより簡単かつ迅速に行えるようになります。

## 前提条件

- □中級SQL
- □データディクショナリをナビゲートする能
- □力 データモデルを読んで**SQL**を書く能力
- ◆ □ Data Team SQL Style Guideの遵守
- ◆ □ハンドブックに掲載されている「Self-Service-ready Data Subject Areas」の中から、「Self-Service SQL Developer Knowledge Assessment」に合格する。

アクセスリクエスト。

- ■新規アクセス・リクエストを開き、シングル・パーソン・アクセス・リクエストを使用する
- □完成した知識評価へのリンクを含める Snowflakeへのア
- □クセスを要求する

layout: handbook-page-toc title:"データチームドキュメンテーションガイド" 説明"このドキュメンテーション方法は、現在データチームでは使用されていませんが、異なるタイプのドキュメンテーションをどのように考えるかの良いガイドとなっています。"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

このドキュメント作成方法は、現在データチームでは使用していませんが、さまざまなタイプのドキュメントをどのように考えればよいのか、よい指針となります。

ペルソナ

データチームのハンドブックには、たくさんの情報が詰まっています。すべての情報が、すべての人のために書かれているわけではありません。社内でも出身地が違えば、求めている情報も違います。

読者の皆様に参考にしていただくために、大多数のユーザータイプを想定した**4**つのペルソナを作成しました。ペルソナは以下のとおりです。

ビルダー祭

{:#builder}には

ビルダーの主な属性

- データのクエリと変換の両方で強力なSQL能力を持っている優れたデータ運用
- 能力を持っている
- データへの直接アクセス (Snowflakeの
- Rawデータ) 強いドメイン知識
- ドキュメンテーションや分析の貢献者/協力者 典型的な役割

データエンジニア コアデータアナリ

・スト

ユーザー風を

{:#user}になります。

ユーザーの主な属性

データを照会するための低~中レベルのSQL能力

- \* 基本的なデータの流暢性
- \* データへの直接アクセス (Snowflakeの
- \* Rawデータ) 強いドメイン知識

典型的な役割です。

ディストリビューション・ア

- \* ナリスト/エンジニア・プロダ
- クトマネージャー

• エンジニアリング・マネージャー

## コンシューマー

{:#consumer}です。

### 消費者の主な属性

- 限られたSQL能力限られた
- データ流暢性 強いドメイン
- ◆ 知識
- データへのアクセス制限 (Sisenseの

### み) 代表的な役割

- ピープルマネージャー ビ
- ジネスパートナー
- Non-data / technical individual contributor

### チャンピオン[!!!]

{:#champion}になります。

チャンピオンの主な属性

- 限られたSQL能力と限られた時間
- ドメイン知識は豊富だが、分析的洞察は直属の部下に頼るデータへのア
- クセスが限られている

### 典型的な役割です。

- 取締役/VP/CXO インベ
- ・スター

## みんなの

{:#everyone}になります。

これは、すべてのペルソナに関連するものがある場合のための、包括的なグループ分けです。

ドキュメントの種類

TheDocumentationSystem」では、ここでは説明しきれないほど詳しく説明されていますが、一般的な考え方として、ドキュメントには4つのタイプがあるとされています。これらはそれぞれ異なっており、異なる目的を持って書かれています。それは以下の通りです。

チュートリアル

### {:#tutorials}

のチュート

リアルです

ね。\*

は学習志向

初めての方でも安心して始められる

ように、レッスンでは

例え:小さな子供に料理を教えること

ハウツーガイド**≥**を

{:#howto}を参照してください。

ハウツーガイドです。

目標を持っている

•

•

•

- 特定の問題を解決する方法は、一連の
- ステップであることを示しています。

類語:料理本のレシピ

説明はこちら♀。

### {:#explanation}

の説明です。

- ・ は、理解重視の説明です。
- 背景や文脈を提供する Analogy: 料理

の社会史に関する記事 Reference 🗐

{:#reference}になります。

リファレンスガイドです。

- は、情報を重視し、機械
- ◆ が正確で完全であること
- を説明します。

類語:参考となる百科事典の記事

## 使用方法

ペルソナタイプとドキュメントタイプの両方のヘッダーの横に絵文字があるのがわかります。この絵文字は、ハンドブックのデータ チームの部分で使用しており、読者が各セクションの対象者や目的をすぐに理解できるようにしています。例えば、次のような ものです。

- **大**図 データチームの誰かのためのハウツーガイドであることを示します。例としては、**Snowflake**でのプロビ ジョニングの方法などです。
- 口歌 これがユーザーペルソナの参考資料であることを示します。例としては、Sisenseで作業するためのヒントやコツのセクションです。
- \$ ♀ 全てのペルソナについての説明であることを示します。例としては、データチームの憲章が挙げられます。

絵文字は、そのセクションに興味を持つすべての人を網羅しているわけではないことにご注意ください。私たちの目的は、読者が何を読もうとしているのかを説明するためのクイックガイドを提供することです。

データチームハンドブックへの寄稿

データチームのすべてのチームメンバーは、ハンドブックに貢献することが奨励されています。@data-teamのみんなは、データチームのハンドブックページのコードオーナーですが、義務があります。

- 1. すべてのハンドブックの変更はピアレビューされ、ピアが承認を与える前にマージされることはありません。
- 2. ハンドブック変更のためのすべてのMRには適切な名前がついていますが、これはHandbook Changelogを通じて変更を追跡するために重要です)。)データチームのハンドブックの変更であることを示してください。
- 3. MRが直接影響を与える可能性のあるチームメンバーにタグを付け、チームメンバーがレビューするために24~48時間の時間を確保する。ハンドブックへのわずかな変更であっても、少なくとも1人のチームメンバーに確認してもらうのが良い方法です。
- 4. データチームのデモミーティングでは、あなたのハンドブックへの貢献について自由に議論・デモすることができます。

layout: handbook-page-toc title: 'GitLab Experimentation Best Practices' description:"実験によって、私たちは学び、正しい体験を顧客に提供し、顧客とGitLabにとってより良い価値を生み出すことができます。"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .hidden-md .hidden-lg}.

なぜエクスペリメント・ベスト・プラクティスなのか?

実験によって、私たちは学び、お客様に正しい体験を提供し、お客様とGitLabにとってより良い価値を生み出すことができます。実験は簡単なように見えますが、組織がベストプラクティスに従わなければ、不正確な結論を出すリスクがあまりにも高くなります。

例えば、実験を行う際によくあるミスステップは、不正確な結論や判断につながる可能性がありますので、以下にご紹介します

- 1. ピーキング(記事はこちら)。テスト期間を前もってロックしておかないと、A/Bテストが終わる前に結果を確認して 行動を起こすことで副作用が発生するピーキング問題が発生してしまいます。A/Bテストの中間結果を意思決定の準備 をした状態で見る回数が多いほど、基準がないのに統計的に有意な差を示してしまう確率が高くなります。
  - 1. 2つのピーキングケースはp値を2倍にします。
  - 2. 5回の覗き見でp値が3.2倍になる。
- 2. シンプソンのパラドックス: テスト飛行の途中でテストグループの配分を不均衡に変更した場合に発生します。テストグループの配分比率を変更すると、テストグループ内の潜在的なセグメントの比率が変わり、結果に誤差が生じます。より正式には、シンプソンのパラドックスとは、ある母集団における2つの変数の間の関連性が、その母集団を小集団に分けたときに現れたり、消えたり、逆転したりする統計現象のことです。

実験方法から得られる価値を最大化し、GitLabでの決定の不正確さを減らすために、GitLabでのすべての実験においてベストプラクティスに従うことを推奨します。このドキュメントでは、Gitlabで採用すべきベストプラクティスの概要を説明します。

## 実験の計画

仮説を立てる。アップサイドの定義

### 実験を始める前に

- 改善しようとしているビジネスチャンスの定義
- どのようにしてビジネスチャンスを拡大しようとしているのか、なぜそのようなアプローチを取りたいのかを説明してください。

裏付けとなるデータや分析があれば、それを文書化する。

コントロールのUIとフローを文書化し、ワイヤーフレーム/コンセプトをテスト する

目標。目標がなければテスト測定は意味がない

- 実験で動かしたい主要な指標を定義します。これは、ロールアウトのシナリオを定義するために使用する指標です。
  - 理想的には、1つの指標だけであることが望ましい(最大2つの指標)。
  - 理想的には、ビジネス上のKPI (例: ARRの増加、有料ユーザーのサインアップなど)
  - 。 に結びつけるべきです。 サンプルサイズを計算するために、予想される変化の大きさ と方向を定義します。
  - ・ テストのワンテール/ツーテールと期間を決定 → サンプルサイズに基づき、KPIが変化すると予想される方向性に基づく。
  - 例
    - テストコンセプト新規ユーザー登録時のステップ数の削減
    - 主な指標はまた、有料ユーザー数などの収益につながるKPIに結びつけることが理想的です。
- 副次的なメトリクスを定義する。一次指標がなぜ、どのように影響を受けたのかを検証し、詳細に理解するための 指標があります。
  - 二次指標は、一次指標の動きが本物かどうかをクロスバリデーションするのに役立ちます。
  - 5%や10%という高いアルファ値を使用しているため、一次指標への影響が偶然のものである可能性があり、二次 指標は一次指標の変化が本物かどうかを検証するのに役立ちます。テストのケイデンスが改善されれば、この点 はより重要になるでしょう。

#### ■ 例

- テストコンセプト。新しいユーザー登録フローに、より目立つトライアル登録ボタンを導入する。主
- な目標D90コンバージョン率の向上。新規ユーザー登録の際にトライアルに申し込む人が増えることで、製品の価値をより深く理解してもらい、D90コンバージョン率を5%向上させる\*\*ノイズの可能性
- :主要な指標が5%向上しても、トライアルに申し込む新規ユーザーの数が増えなかった場合、**D90**コンバージョンの向上はノイズである可能性があり、さらなる検証が必要です。
- \*\*ガードレールの指標を定義する。これは、短期的な利益のために長期的にはビジネスを害することがないようにするための ものです。

\*\*

- これらは、実験がもたらす長期的な影響を追跡する指標に基づいています。例
- 新規ユーザーの登録フローでトライアルを増やす実験をする場合、D90コンバージョンやトライアルから有料へのコンバージョンをガードレール指標として設定するとよいでしょう。
  - もうひとつの例は、あるソーシャルメディア企業で、「通知無効率」をガードレール指標として定義しています。プッシュ通知の増加によりセッション/DAUが1%上昇するごとに、「通知無効率」の増加はX%以内に収まるべきだというガードレールを設定しています。
  - 別の例。あるSaaS企業では、ガードレールの指標として12ヶ月間の解約率を設定し、12ヶ月後に分析を再検討することがあります。

プレアナリシス:前もって計算しておく

- データを用いて問題、機会、アプローチを検証する(例:記述的分析
  - 例特定のメールタイプのメールへッダーコピー変更をテストしても、このメールのボリュームが少なければ、有意な価値はないかもしれません。有意性を達成することはできないでしょう。このような場合には、事前・事後テストを行うか、よりトラフィックの多いメールタイプを選んでテストを行うことを検討してください。
- 片側/両側、α(5%)、検出力(80%)を用いて、必要なサンプルサイズを算出する。実験期間(どのくらいの期間 実験を行うか)を前もって決めておく。
  - o ロールアウトの決定に高いアルファ値を使用する必要があるかどうかを事前に判断します。
- 実験を行うための費用対効果の分析を行う。組織内で実験のコストが高く、技術リソースが限られている場合、正しい アイデアを優先させるために、それぞれのテストから得られる潜在的なアップサイドを理解できなければなりません。
  - 会社のKPIの1つとして、この変更による年間のアップサイドの可能性を計算する。それが不可能な場合は、 KPIへのアップサイドをログ値で定義し(0.01%, 0.01%, 0.1%, 1%, 10%, 100%インパクトのように)、異なるアイデアの間で大まかな比較ができるようにします。
  - 段階的な機能変更と段階的な変更の推進

実験のセットアッププラン。テストのセットアップ方法

- 実験から得たいビジネス上の疑問に答えるために必要なバリエーションの数を決定する (例: A/B/nまたはMVT (多変量解析)のセットアップ
- どのような仮説の質問に答えられるか、答えられないかを含む分析計画を文書化します。
- 必要な機器の定義(既に利用可能なデータに加えて、追加で必要なトラッキングデータ)。
  - 開発者が新しいトラッキングについて議論し、労力やパフォーマンスへの影響を考慮して実装できるように 、トラッキングのニーズを「重要」または「あったらいいな」と指定する。
- セグメンテーションや除外基準を含む、テスト対象となる母集団の定義
- いつ実験を行うかを実験カレンダーに記載し、競合をチェック クリティカルパスで実験を行う場合や、大規
- 模なユーザーに影響を与える場合には、立ち上げ時の重みや立ち上げプランを定義する
  - 重要な分野でリスクの高いテストを行う際には、1%から5%/10%、そして50%へと段階を踏んでいきます。
    - シンプソンのパラドックスに注意して、重量配分が比例して比較可能なテスト期間のみを分析してください。
  - ランピングの成功基準の指定
    - 例プライマリー指標に大きな影響はなく、5%程度の変化でも検出できる(=感度)。

リスクとクリティカル・パスに基づくテストグループのサンプル配分に関する一般的なガイダンス

失敗のリスク、機能や影響を受ける人々の重要性を考慮した実験打ち上げの重み 付けマトリックス クリティカル・パス・ペー ジまたはクリティカル・プ ロダクト・フィーチャー

LMH

故障時の潜在的なコードリスク

L50

10%

5%

М

10%

10%

5%

Н

10%

5%

1%

# 実験検証。

Experiment Validation プレローンチ。すべてが整っているか確認

- 開発チームとデータエンジニアチームが機能を構築した後、開発環境または本番前の環境で、レポートデータと**UI**が意図したとおりに動作するかどうかを検証します。
  - 推奨: 開発、QA、アナリストのうちいずれか2名にトラッキングとレポーティングを検証して
  - もらう 推奨: PM、QA、アナリストのうちいずれか2名にUI機能を検証してもらう
- コントロールとテストエクスペリエンスのスクリーンショットを、今後の参考のために実験ドキュメントに貼り付けます。

Experiment Validation ポストローンチ。期待通りの結果が得られているかを確認する

- 本番開始から1~2日後、実験の報告データが有効かどうかを確認する母
  - 集団の割り当てに偏りがないかどうかを確認する
- 1%の重さでスタートした場合、2~3日後に最初の読み込みを行い、事前に取り決めた計画通りに5%や10%の重さに設定します。
  - 1%の実験の目的は、あくまでも「モノが壊れない」ことであり、結果を読み取ることではありません。
- 5%や10%で立ち上げた場合は、当初の計画通り50%にウェイトを変更する
  - 平日と週末では顧客層が異なる可能性があるため、週単位の季節性がある場合は1週間単位、月単位の季節性がある場合は1ヶ月単位で分析してください。
  - 1%のテストは、実験的な経験によって生み出された影響を読み取ることを目的としていないため、上記は 1%には適用されません。

# 実験分析とコミュニケーション。

- 課題や実験データに偏りがないことを確認する
  - 可能であれば、これらの検証を自動化する 割り当ての母集団バイアス 統計的有意性を
- 用いたテスト結果の分析
  - セカンダリーメトリクスを活用して、テストが成功した理由や方法を検証し、より理解を深めます。
  - 展開シナリオを理解し、人口コホートによるパフォーマンスの有意差を見つけるために、主要なディメンション 別に結果を分析する。
    - 多重比較を行うと、統計的ノイズが増えるので注意が必要です。そのため、ディメンションレベルの分割では 分析前に仮説を立てるか、修正P値を使用してください。
- このテストと同時に実施された、事前に特定された相互作用実験との相互作用結果を分析します。
- 統計的に信頼できる結果を伝える。

- 有意な影響があった場合、信頼区間(@80%)をつけて結果を伝える。
- 有意な影響がない場合は、一般の人々がどの程度の観察された変化率で、統計的に有意な結果と結論付けられるかを理解できるように、感度を持って結果を伝える。
- 結果をハンドブックに記載し、テストの詳細な説明、ランプの概要、主要な結果、メトリクスのスナップショット、インサイト、次のステップを記載します。
- 結果をパートナー(技術部門およびビジネス部門)および分析仲間に提示し、追加の洞察を集め、学習内容を他の人に 伝えます。フォローアップの分析と洞察を文書化する。
- PMのために最終的な結果に基づいて、実験がロールアウトされたり、コードから破棄されたりすることを確認します。

## 実験から得られた知見を制度化する。

- ビジネスパートナーと一緒に、詳細な結果をより広い組織のグループで発表し、他のチームがどのように学習から利益 を得ることができるか、価値を最大化するために他のチームがどのような行動を取ることができるか、協力することが できるかについての会話を促進します。
- 失敗した同じアイデアを繰り返さないために、すべての結果を文書化し、機能や製品タグで検索できるようにする。これ により、複数の実験からメタアナリシスを行い、より幅広い洞察を得ることができます。
  - 例)類似企業において、過去の20数例のデータを活用して、通知タイプの増加と通知無効化の関係を把握した
  - 例比較対象となる企業では、過去の実績(~50件)のメタデータを用いて、異なるメールタイプのボリュームを増 やすことによるセッションの上昇を把握しました。
  - 例:同業他社では、過去8回の分析結果を活用して、クロスマーチャンダイジングの追加スポットの増分価値を把握しました。
  - o Growthチームはすでに、実験の結果を一元化されたハンドブックのページに記録しています。このやり方を Gitlab 全体で採用し、一元化したページを作るべきです。

## 実験のガバナンス。

- 自動通報システム。
  - 自動化された警告システムを有効にして、実験による主要メトリクスへの著しいマイナスの影響を監視します。これにより、結果を覗き見する必要がなくなります。
    - P <= 0.01またはP <= 0.5を2日間連続で使用しています。
  - 実験の結果、一貫して有意なポジティブな結果がx日間継続した場合、自動アラートを有効にする(ランダムウィンを減らすために3日間とした)。
- 長期にわたる実験を定期的に見直し、製品のパフォーマンスに影響が出ないようにする。
- セットアップ、トラッキング/データの失敗、誤った実装、コンフリクト/実験の相互作用などによる実験の失敗 を記録する。これにより、実験プラットフォームの健全性を監視することができます。
- 実験カレンダーを作成することで、コンフリクト管理を可能にし、実験の量と速度を把握することができます。

## オポチュニティ・アクション・アイテムを検討する。

注) 短期とは、数時間/数日分の努力を意味します。長期的には、数週間、数ヶ月分の努力が必要です。

現在の実験プラットフォームの機能に関する初期のフィードバックと観察に基づいて、これらは私たちが検討することを推奨 するアクションアイテムのリストです。ギャップがまだ存在し、その優先順位を確認するために、影響を受ける関係者と協議 する必要があります。

# サンプルサイズの要求を検出力計算で標準化する。(優先度H)

- 短期的にはアルファ(5%または10%)、片側/両側、検出力80%の標準的なサンプルサイズ計算機を定義する。
- 結果のダッシュボードに有意性と信頼度を確実に反映させる(優先度H
  - 短期的には**Z**検定の公式や**Python**の関数を使って、**Sisense**実験分析フレームワークに直接、有意差計算と信頼区間を 組み込む。
- P値のしきい値を調整して、テスト結果を次元別に分割する機能。(優先度M)
  - 長期的には実験のメトリック/ファネル/ディメンションを自由に評価し、ディメンションやフィルターによる分割分析を可能にするための代替ツールやデータトラッキングを検討する。
- サンプルサイズの問題を解決する(優先度H
  - 。 短期的には実験のサンプルサイズが小さくても対応できるように、実験のサンプルサイズ分析を片側または両側 で定義する。

0

短期的にはセカンダリー指標を活用して、プライマリー指標での方向性の読み取りに自信を持つ。

- •
- 0
- •
- 0
- \_
- 0
- - 0
    - 0

- 短期的には十分なサンプルを収集した後に実験を中断し、後にガードレール・メトリックへの影響を測定することで、ガードレール・メトリックを長期的な影響測定に使用する。
- 長期的。ベイズ実験解析フレームワークを活用し、サンプルサイズが小さい問題に対するデータの不確実性を 伝える。
- テスト結果のバイアス検証の自動化。(優先度M)
  - 。 短期的には実験結果に重みや母集団の偏りがある場合に、自動的にフラグを立てる方法の検討 長期的結果を
  - 分析する際に、自動化されたバイアス警告を導入する。
- 勝った実験と負けた実験の自動アラート。(優先度L)
  - 短期的には負けている実験や常に勝っている実験について、P値や結果が有意に残っている日数に基づいて自動 化されたアラートを作成する方法を検討する。
- 実験内および実験間の無作為化に偏りがないか分析する。(優先度M)
  - 短期的には既存のシステムの限界に基づいて、実験内および実験間の無作為化バイアスを検証する必要があるかど うかを検討し、分析する。
  - 長期的なもの実験内および実験間のバイアスを継続的に監視する自動化システムを構築する。
- 実験カレンダー。(プライオリティーM)
  - 短期的には現在の実験カレンダーを理解し、実際の実験割り当てデータと比較することで、テストカレンダーの 精度を検証する自動化された方法を模索する。
  - 長期的には組織全体で実行されているすべての実験と、実験の期間を記録する単一のシステムがある。実際のテスト割り当てデータに対する自動検証を可能にする。

## 分析用に作られたツール。

- 検出力と信頼度の変化を利用したサンプルサイズ計算機
- WIPSisense実験ダッシュボードへの有意差検定の組み込み(概念実証

## リファレンス

- Growthの実験ページ。Growthはすでに実験結果を一元化したページに記録しています。
- 成長実験ダッシュボード (Daveによる) は、発売後に実験データの精度を迅速に検証するのに役立ちます。7ステッ
- プのA/Bテストプロセス 計画、設計、QAに時間を費やす必要がある
- 低トラフィックサイトのA/Bテスト
- テスト計画の重要性

# 用語の説明。\_

- デフォルトの体験。実験に割り当てられていないときに参加者が受ける経験。これはほとんどの場合、対照体験と同じである。
- 実験カレンダー。過去の実験、現在の実験、今後の実験を、実験対象者、ウェイト、実験期間とともに表示するシステム。これにより、コンフリクトの計画を立てたり、相互に作用する実験間の相互作用を分析したりすることができます。
- 感度。与えられた母集団とその指標に対して、有意な閾値を検出できる最小の変化は何か。

layout: handbook-page-toc title: "Data Team Calendar - Meetings" description: "GitLab Data Team Calendar"

## このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc.toc-list-icons.hidden-md.hidden-lg}。

{::options parse\_block\_html="true" /}。

# ミーティング

データチームのGoogleカレンダーは、チームミーティングのためのSSOTです。また、データスペースの関連イベントも含まれています。このカレンダーには誰でもイベントを追加することができます。このカレンダーに掲載されているイベントの多くは、毎月のキーレビューを含め、出席を必要としない、参考のためのイベントです。データチーム全体のイベントを作成する場合は、Googleカレンダーで勤務時間を確認し、事前に勤務時間外のミーティングについて話し合っておくとよいでしょう。就業時間外のミーティングが必要な場合は、誰がミーティングをするかを交代することを検討してください。

ノーミーティング・フライデー

金曜日に絶対にミーティングをしないとは約束できませんが、チームは「No Meetings Fridays」を採用しています。

メーカーのスケジュールの偏り

データチームのメンバーのほとんどは、マネージャーのスケジュールよりもメーカーのスケジュールに強く共感しており、私たちはこれを支持しています。

- 個々のコントリビューターが会議に参加する時間を20%以下に制限すること
- ミーティングを一日や一週間に分散させるのではなく、連続した単位にまとめ、定期的に必要なミーティ
- ングのほとんどを火曜と木曜に行う。

ミーティングプロトコル

チームはGitLabミーティングのプラクティスとスタンダードに従っています。

- ミーティングには必ずアジェンダがあります。
- 会議のテーマは優先順位の高い順に並んでいます。
- アジェンダには、トピックや質問に沿ってメモが追加されます。
- すべてのアジェンダは、データチームミーティングの共有ドライ ブに保存されます。

定期的なミーティング

デイリースタンドアップ

データチームのメンバーは、日々のスタンドアップにギークボットを使用しています。これらのスタンドアップノートは、#data-dailyに投稿されます。Geekbotが「今日は何に取り組もうと思っていますか?と聞かれたら、チームメイトが積極的にブロックを解除できるように、具体的な内容を答えてみましょう。Salesforceの作業」ではなく、「sfdc\_opportunity\_xfモデルにオポチュニティのオーナーを追加する」と考えてみてください。Geekbotがメッセージを送ってきたら、すぐに返事をしなければならないというプレッシャーはありません。あなたがその日に取り組んでいることがチームに伝わるような応答をGeekbotにしてください。そうすれば、優先順位が変わったり、必要な追加情報があったりした場合に、チームがあなたの理解を助けることができます。

火曜日 企画・結果発表会

毎週火曜日の14:00-17:00 (UTC) には、ビジネスに特化したフュージョンチームとプラットフォームエンハンスメントチームの成果を出すことに焦点を当てています。データチームのリーダーは、ほとんどのセッションに定期的に出席し、すべてのミーティングの議題を確認することが求められます。私たちは2週間のマイルストーン間隔で仕事をしているため、この火曜日のミーティングの焦点は週ごとに変わります。

時間(CET/UTC/ET/PT) ミーティング

1600/1400/1000/700 GTM Fusionチームミーティン

グ1700/1500/1100/800 ProductFusionチームミーティ

ング

1800/1600/1200/900 データ・エンジニアリング・チーム・ミーティング

DEMOミーティング

毎週木曜日は、データスタックに焦点を当て、開発中のソリューションを紹介しています。

• デモミーティング。ソフトウェアや関連コンテンツのデモを行うことは、ソリューションへの関心を高め、知識を深める

のに最適な方法です。デモには、SQLの例、ハンドブックのページ、dbtフロー、ERD、ダッシュボード、その他の関連コンテンツが含まれます。画面共有やウォークスルーも推奨します。必要に応じて、ビジネス関係者もデモに招待します。

時間 (CET/UTC/ET/PT)	ミーティング
1700/1500/1100/800	データチームのデ

モ

#### 月例会議

時間(CET/UTC/ET/PT)	デイ	ミーティング
Varies	最初のフルウィ ーク	オールハンズ&レトロスペク ティブ
Varies	先週の金曜日	Handbook Jamboree
Varies	Varies	月刊ピザパーティー

layout: handbook-page-toc title:"Merge Request Roles and Responsibilities" description:"GitLab データチーム MR の責務"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg},

{::options parse\_block\_html="true" /}。

マージリクエストの役割と責任

レビュアー

GitLab データチームのすべてのメンバーは、同僚やコミュニティの貢献者のマージリクエストに対してコードレビューを行う ことができますし、そうすることが推奨されています。マージリクエストをレビューしたい場合は、誰かがあなたに割り当てて くれるまで待つこともできますが、公開されているマージリクエストの一覧を見て、フィードバックや質問を残すことも大歓迎 です。

なお、すべてのチームメンバーはすべてのマージ要求を確認できますが、マージ要求*を受け入れる*ことができるのはメンテナに限られています。

### レビュアーの責任は

- 技術的な実装を確認するために
- コードがビジネス目標を達成するために
- 作成したデータモデルのデータ品質をチェックする

## ために Codeowner

コードオーナーシップとは、プロジェクトメンバーとプロジェクト内の特定のフォルダやファイルを結びつけるGitLabの機能です。これは、「このコードについて誰に聞けばいいのか」「このコードの変更を誰がレビューすべきか」という質問に答えるためのものです。

コードオーナーになることは、プロジェクトメンテナーになるための過程の一部です。もしあなたがファイル(例えば新しい **dbt**モデルセット)の唯一の作成者であれば、あなたはそのファイルの事実上のコードオーナーとなります。オーナーシップの範囲を広げたい場合は、以下の手順に従ってください。

- 1. CODEOWNERSファイルに、変更したいオーナーシップを記載したMRを作成します。
- 2. 参加したいエリアをすでにカバーしている他のコードオーナーと協力して、担当したいコードを変更する少なくとも5つのMRでペアを組む。
  - MRは、コードベースの小さな変更だけでなく、アーキテクチャ的なものや、完全に機能する追加機能を生み出す ものも反映させる必要があります。この基準に達するには、5つ以上のMRが必要になる場合があります。
  - あなたは、コードオーナーに割り当てられたすべてのMRの一次レビュー者と
    - 145/ 243

なります あなたは、マージする権限があるかのようにMRをレビューします

○ MRの準備が整い、あなたがそれをマージしたと感じたら、それをコードオーナーに割り当て、あなたがMRをマージしただろうとコメントします。

o コードオーナーは必要に応じてフォローアップレビューを行い、MRをマージするか、あなたのためにメンテナー に割り当てます。

- 作成したオリジナルのMRでドキュメントを作成する
- 3.5MRの閾値に達した時点で、コードオーナーは申請者のマネージャーと協力して
- 4. 拒否された場合、MRを終了し、3ヶ月以上の期間、上司と協力して再応募する
- **5.** 承認された場合、そのMRをメンテナに割り当ててマージします。マージされると、コードオーナーの承認リストに自動的に 追加されます。

#### メンテナ

データチームのプロジェクトにおけるメンテナーは、どの職種とも同義ではありません。ここでは、データチームはメンテナの責任についてエンジニアリング部門が示した前例を参考にしています。データチームのすべてのプロジェクトには少なくとも1人のメンテナがいますが、ほとんどのプロジェクトには複数のメンテナがいて、アナリティクスのようにdbtとオーケストレーションに別々のメンテナがいるプロジェクトもあります。

メンテナの責任は次のことを保証することです。

- データチームのプロセスに従うこ
- と MRでの最終確認

データチームのメンテナーになるには

メンテナになるためのガイドラインはありますが、具体的なルールはありません。メンテナは、GitLab Data プロジェクトのコードベースについて高度な知識を持っている必要があります。あるプロジェクトのメンテナになる前には、そのコードベースをよく理解し、一つ以上のドメインの専門知識を身につけ、そして私たちのコーディング標準を深く理解しなければなりません。これらの条件が両方とも満たされたとき、あなたはメンテナになる準備ができたと言えるでしょう。

- 1. あなたがレビューしたMRは、追加で必要な大幅な変更をすることなく、一貫してメンテナのレビューを通過しています。
- 2. 作成したMRが、大幅な変更を要求されることなく、レビュアーやメンテナのレビューを一貫して通過する。

それらの主観的な要件が満たされていれば、自分をメンテナとして追加するプロセスとなります。

- 1. 関連するプロジェクトに "Add as project maintainer"というタイトルで課題を作成します。
- 2. 以下のように、課題にドキュメントを追加します。
  - なぜあなたがメンテナーを引き受けることができるのかを説明してください。
  - メンテナリングの範囲を説明する(プロジェクト全体、dbt、pythonなど
  - あなたが作成し、レビューした最近のMRで、あなたの準備ができていると思われるもの。
- 3. 課題が作成されたら、正式なレビューのためにペアになりたいメンテナをタグ付けしてください。
- 4. メンテナは、少なくとも10個のマージ・リクエストであなたとペアを組ませます。
  - MRは、アーキテクチャの変更や、多くのファイルやディレクトリに変更を加えた完全に機能する機能リリースなど、多様な範囲を示している必要があります。
  - あなたは、これらの10個のMRの主なレビュアーとなります。
  - 融合する力を持っているかのように、MRの見直しを行う。
  - MRの準備が整い、あなたならそれをマージしただろうと思ったら、それをメンテナに割り当て、あなたならMR をマージしただろうとコメントします。
  - メンテナは必要に応じてフォローアップレビューを行い、あなたのためにMRをマ
  - ージします あなたが作成した課題のMRをドキュメント化します
- 5. 10MRの閾値に達した時点で、メンターは申請者のマネージャーと協力して決定します
- 6. 拒否された場合、問題を解決し、3ヶ月以上の期間、上司と協力してから再度申請してください。
- 7. 承認されたら、MRを作成してチームページのエントリーにメンテナスを追加します。
- 8. MRを上司に割り当て、関連するプロジェクト(インフラ、アナリティクスなど)やエリア(dbt、エアフローなど)の既存のメンテナに言及します。
- 9. 関連するグループ (例: dbt) の既存のメンテナが重大な異議を唱えず、少なくとも半数のメンテナがレビュアーの準備ができていることに同意すれば、私たちは自分たち自身で新しいメンテナを獲得したことになります。
- 10. プロジェクトのオーナーが、プロジェクトでのあなたの特権を増やします

マージリクエストのワークフロー

データチームは、データアナリティクスチームとデータエンジニアリングチームの間にある1つのチームとして運営されています。そのため、各MRには最低でも3人の関係者がいることが予想されます。以下のシナリオをご覧ください。

- 1. DA著者、DAレビュー、DEメンテナのマージ
- 2. DEの著者、DAのレビュー、メンテナのマージ
- 3. DEの著者, DEのレビュー, メンテナのマージ

データエンジニアリングチームはデータプラットフォームに責任を持っているので、すべてのMR依頼にはデータエンジニアが含まれていなければなりません。

layout: handbook-page-toc title:"新しいデータソース" 説明"新しいデータソースを 追加する方法"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}o

このページでは、データウェアハウスに新しいデータソースを追加するプロセスについて詳細に説明します。データチームは、データウェアハウスに新しいデータソースを追加するあらゆるイニシアチブを採用しています。私たちは、データが価値につながることを想定していますが、データウェアハウスに新しいデータを追加するには、コストがかからないわけではありません。

開発(データチーム、サポートに関わる他のチーム、要求者であるお客様のリソースの割り当て)と、データパイプラインの 稼働維持(ストレージ、コンピュータリソース、インシデント管理、監視など)には、時間とお金がかかります。

新しいデータソースの追加をリクエストする前に、以下の点を考慮してください。

- 妥当なビジネスケースがあるのか?
  - ビジネスケースとしては、規制上の要件を満たすことが必要な場合もあります。
  - ビジネスケースは、投資コストよりも潜在的な価値がはっきりしているため、単純明快な場合もあります。ほとん
  - どの場合、価値が不明瞭であったり、コストが不明瞭であったりするため、定量化するのは困難です。そのような場合は、データチームと率直に話し合ってみてください。私たちには経験がありますし、正当化の手助けもできます。また、これは必ずしも科学的な計算である必要はありません。
- データウェアハウスにデータが入った時点では、作業は完了していません。データは生データ層に「ただ」入っているだけで、Sisenseのデフォルトではアクセスできません。データはdbtを介してEDM(Enterprise Dimensional Model) にロードする必要があります。フォローアップが必要で、このページで説明されているプロセスの上に行われます。
  - 当社はデータプラットフォームへの貢献を歓迎しているため、下流のモデリングはビジネスチームが担当することも可能です。dbt、SQL、データモデリングに関する幅広い知識が必要です。
  - 下流のモデリングはデータチームが担当することもあります。計画を立て、会社の優先順位に沿って優先順位を設定する必要があります。これは、データウェアハウスに新しいデータソースを追加する範囲ではないため、後から手配する必要がある。
- データウェアハウスに新しいデータソースを追加することは、1回きりの作業ではありません。データがデータウェアハウスに抽出されると同時に、定期的に(週に1回、1日に1回、1日に複数回など)データが更新されます。これは、実装後に何かが起こったり、うまくいかないことがあることを意味します。このプロセスを必要に応じてサポートするために、ソース側(ビジネスと技術)からDRIが必要になります。
- データはEDMの外でも使うことができます。例えば、ファンクションアナリストが生データ層で使用することができます。 生データにアクセスするには、ARを上げる必要があります。
- データがEDMに入ってしまうと、Sisenseではダッシュボードを作成して作業を行うことになります。そのためには、いくつかの技術的な知識が必要です。
- データプラットフォームは、ソースシステムにアクセスする必要があります。当たり前のことのように聞こえますが、これは必ずしも簡単なことではありません。

新しいデータソース

新しいデータソースを追加する場合のワークフローは次のようになります (注:これは私たちがすべての作業に使用している 既存のワークフローであり、ここでは新しいデータソースを追加するという文脈で明示しています)。

新しいデータソースを追加するためのプロセス。

新しいデータソースのステージ (ラベル) サマリー

ワークフロー	すべての初期情報は要求者によって提供され、データチームによって評価されます。
::1 - トリア	7 CONDINERROS ANTECES O CIENCIANO 7 77 MICS O CITIMICAVS 76
ージ	人似不用或不上上不少儿。
ワークフロー	今後の開発のためのソリューションデザインと完全な作業内訳が作成されます。
::2 - バリデ	四半期に一度、データチームのOKRが定義されます。作業量や優先順位に応じて、次の四半期の
ーション	OKRに新しいデータソースを組み込むことができる。
ワークフロー	新しいデータソースの追加は、次の四半期のOKRのリストに入っています。すべての作業内訳活
::3 - スケジ	動が添付されたエピックが作成されます。
ューリング	開発が進んでいる
ワークフロー	paper and the second se
::4 - スケジ	開発は検討中
ュール	
ワークフロー	開発が阻害される
::5 - 開発	
ワークフロー	
::6 - レビュ	
_	

ワークフロー

::X - ブロッ

キング

ワークフロー::1 - トリアージ

新しいデータソースのリクエストには、必ず問題が発生します。これは、すでに抽出されているが、拡張が必要なデータソースや、まったく新しいデータソースにも当てはまります。新規データソースのテンプレートをご利用ください。すべての初期情報はリクエスト者によって提供され、データチームによって評価されます。データがデータウェアハウスにまだ存在していない場合は、クロスチェックが行われます。

ワークフロー::2 - バリデーション

新しいデータソースに関するすべての詳細は、完全な作業内訳を作成することを目的として、表示されます。デ

#### ータの範囲

要件に基づいて、抽出する必要のあるデータポイント (どのテーブルかなど) が決定される。総データ範囲が与えられなければならない。

#### 抽出液

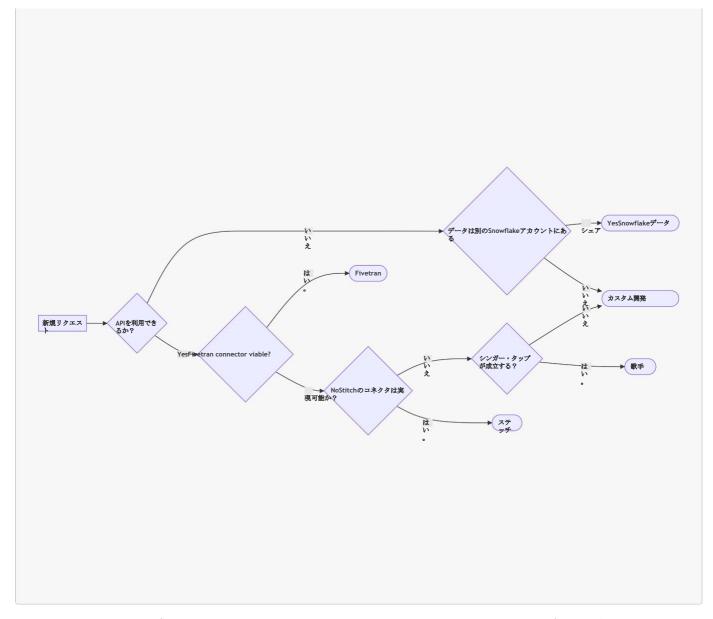
データチームには、ソースシステムからデータを抽出するためのさまざまな機器が用意されています(ランダムに発

- 注されます)。カスタムメイド (Python経由
- フィベラ
- ン・メルターノ
- スノーフレークデー
- タの共有 ステッチ

どの楽器を使うかの判断は、常に以下のような組み合わせで行われます。

- 1. 実装のための努力。
- 2. メンテナンスに力を入れています。
- 3. 伸び縮みする能力

どの機器を使用するかは、データチームが決定します。適切な抽出ソリューションを決定するために、次のような決定図が使用されます。



カスタム開発とは、GitLabデータチームが設計・開発したソリューションのことです。その例として、現在のPGPやZuora Rev Proの抽出が挙げられます。

## アクセスリクエスト

データチームにソースシステムへのアクセス権を提供しておくと便利ですが、今すぐアクセスリクエストを出す必要はありません。

### 作業内訳

検証の最終目標は、ソリューションデザインと完全な作業内訳、DRIとサイズの見積もりを添付することである。実行しなければならないすべての作業が記述されています。私たちは、すべての作業が課題に変換され、T字型のサイズがM/5/8以上にならないようなレベルを目指しています。

### MNPIデータ

データ・ソースにMNPIデータがあり、このデータがデータ・ウェアハウスに向けて抽出されようとしているかどうかを判断する必要があります。データ・ソースにMNPIが含まれており、このデータが抽出される場合は、課題のラベルを新しいデータ・ソースMNPIに変更します。

## ワークフロー::3 - スケジューリング

ビジネスケース、実装にかかる労力、データチームの作業量、データチームの優先順位に基づいて、実装のスケジュールが決定されます。スケジューリングについては、GitLab Data Team Planning Drumbeatに従います。これは、四半期ごとに

データチームは、新しいデータソースのリクエストがいつピックアップされるかを決定します。新しいデータソースのリクエストが 残っている場合

スケジューリングは、データチームが注目していないということではありません。それは、まだスケジュールが組まれていないということなのです。

- 1. 次の四半期のOKRの定義は行われませんでした。四半期に一度、データチームは次の四半期のOKRを設定します。
- 2. ビジネスケース、導入の手間、データチームの作業量、データチームの優先順位などの理由から、次の四半期のOKRには適合しませんでした。

ワークフロー::4 - スケジュール

新しいデータソースを実装する要求が、次の四半期のOKRの範囲内にある場合、その要求はスケジュールされます。課題が添付されたエピックが作成されます。作業内訳に基づいて、すべての課題が作成され、対応する課題のウェイト、ワークフロー「4 - スケジュール」のラベル、適切なマイルストーン(わかっている場合)が割り当てられ、エピックに添付されます。

ワークフロー::5 - 開発

作業の実行が始まると、その課題は開発中であり、通常の開発ライフサイクルに従います。

開発期間中、データエンジニアはデータへのアクセスについて、すべての関係者と調整を行います。データへのアクセスは、データやユースケースに応じて、生のスキーマに提供することができます。

ワークフロー::6 - レビュー

作業の実行が終了すると、課題はレビューに入り、通常の開発ライフサイクルに沿ったものになります。

ワークフロー::X - ブロッキング

外部からの介入を必要とするために実行が継続できなかった場合、問題はブロック化されます。明確な問題提起がなされ、適切な人材が最短でアサインされなければならない。

layout: handbook-page-toc title:"Data Team - Planning Drumbeat" 説明文"GitLab Data Team OKR and Milestone planning process"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /},

データチーム・プランニング・ドラムビート

データチームのPlanning Drumbeatは、四半期ごとに行われる事前設定のシーケンスです。プランニングドラムビートは、GitLabの会計年度と四半期の構造に従っています。データチームのプランニングは以下のようになっています。

プランニング・ドラ ムビートによるオリジ

ナルへのリンク

マイルストーンの命名規則:マイルストーンFYxx-Qxx-MSxx。例: FY22-Q02-MS01 LT=リーダーシップチーム (@rparker2, @iweeks, @dvanrooijen2, @mlaanen)。

DT = データチーム (プロダクトフュージョンチーム、GTMフュージョンチーム、データ

エンジニアリングチーム)全体的な野心は、なぜ、いつ、何をするのかを構造化する

ことです。

• OKRは「なぜ」につながる

- ワークブレイクダウンで何をするか
- マイルストーン計画でいつ何をするか

このようなトップダウンのアプローチは、私たちの仕事のやり方にも通じるものがあります。

### OKRプランニング

データチームのOKRは、ビジネス・テクノロジーのOKR、財務部門のOKR、CEOのOKRと連携し、サポートする部門のOKRと一致することを目指している。エンタープライズ・データ・ウェアハウスの開発と運用に必要な技術とデータ・インフラストラクチャーの作業の性質上、これは必ずしもそうではありません。

FQの終わりまでに、データチームは自分たちのKRを成功させるため、また他のチームがKRの成功を測定するのを支援するために必要なすべての行動の概要を示す。これはイシューを使って行います。イシューでは、関連するアクションのそれぞれについて、全員が予想されるすべてのステップを説明します。これは、チームが予測している障害や懸念を提起する絶好の機会です。これらは将来の参考のために記録しておきます。

これらのOKRは、セントラルデータチームが四半期に行う作業の60%を占めています。残りの時間は、緊急に発生した問題と、アドホック/探索的な分析に分けられます。

OKRアライメントの実例は以下の通りです。

- 1. FY21-Q2CEO Objective 2、データチームのObjective 2と連携し、成長戦略をサポートするための自動化されたデータ パイプラインを提供する。
- 2. FY21-Q2CEOの目標3、KR4、KR5は、データチームの目標1に合わせて、財務ARRデータマートを提供する。

OKRレビュー

OKRレビューでは、データリーダーシップチームが、現在のFQ KRの現状をレビューします。OKRの結果は評価され、必要に応じて次のFQ OKRに伝達されます。

作業内訳

作業内訳の成果は、次の四半期に実施する作業、成果物、責任の詳細な説明であり、以下のように決定されます。

- 1. 今後のOKRの定義
- 2. OKRレビュー
- 3. 新しい/将来を見据えたインサイト
- 4. チームの可用性
- 5. チームメンバーの抱負

仕事の内訳はチームワークであり、全員が貢献することが奨励されています。T

シャツのサイジングの考え方

当社では、新しい課題や大作、長期的な構想を実現するために必要な作業量を迅速に把握するため、T-シャツサイジングアプローチを採用しています。T-shirtサイジングアプローチは、詳細な作業計画を作成するための作業分解をサポートするように設計されていますが、全体のスコープを管理するために十分なレベルの詳細を提供することも目的としています。

| サイズ | 専任者の時間 | 重さ(課題点) | 例 | | :--:| :--:| :--: | XS | < 半日 | 1 | 既存のハンドブックページの更新#データの調査・対応。新しい信頼できるデータテスト。S | 1日 | 2-3 | 新しいハンドブックページ; 典型的なトリアージ問題; 既存モデルの上に新しいダッシュボードを追加。M | 1週間 | 5-8 | 新しいモデルを必要とする新しいダッシュボード。StitchやFivetranを使った新しいデータソース。L | 2~3週間 | 13 | 新しいファクトテーブルの実装とテスト。完全なXMAUソリューション。XL | 1~2ヶ月

|新しいシステムへの新しいデータポンプ。複雑なソースAPIを持つ新しいデータソース。|XXL | 2-4ヶ月 | 52+ |新しいデータソースを使った新しい次元のモデル対象領域。

マイルストーン計画

データチームは、マイルストーンと呼ばれる2週間の間隔で作業を行います。1四半期に1つのマイルストーンは、1四半期が13週間であるため、3週間かかります。3週間かかるマイルストーンとは、大型連休を挟む場合や、チームの大半が休暇やContributeに参加している場合などです。マイルストーンは、水曜日に始まり、火曜日に終わります。これにより、金曜日に土壇場でマージすることを抑制し、チームはマイルストーンの先頭でマイルストーン計画会議を行うことができます。

マイルストーン計画では、以下の点を考慮する必要が

• あります:休暇のスケジュール

- 会議のスケジュール チー
- ◆ ムメンバーの空き状況
- チームメンバーの仕事の好み(スペシャリティは好みとは異なるマイルスト

ーン計画のタイムラインは以下の通りです。

- ミーティングの準備 責任者Milestone Planner 未解決の問題を調
  - 査し、具体化する。
  - チームロードマップとの整合性に基づき、マイルストーンに課題を割り当てます。
  - 注:この段階では、必要な場合を除き、課題が個人に割り当てられることはありません。

	日現在のマイルストーン	次のマイルストーン
0 - 第1水曜	マイルストーンスタート	-
目	ロールマイルストーン	
	中間点	
<b>7</b> - 第 <b>1</b> 火曜日	マイルストーンからずれ てしまう恐れのある問題 は、アサイニーが提起し なければならない	-
	MRの最終提出日について	
10 - 第2 金曜日	MRは、ドキュメントとテストを含めて、すぐにマージできるようにする必要があります。  金曜日(「家族と友人の日」の場合は木曜日)には、 MRの合流を行わない。 マイルストーンの最終日	マイルストーンはほぼ最終段階 マイルストーンプランナーは、次のマイルストーンに向けて課題の 優先度とチームのキャパシティを確認します。
13 - 第2	準備ができたMRはマージ	_
- 月曜日	されることができる 会議	
	の日	マイルストーン計画
<b>14</b> - 第 <b>2</b> 火曜日	未完成の課題は、マイルストーンから削除するか、次のマイルストーンに回す必要があります。	Sync-meetingでは、現在のマイルストーンのレトロパースペクティブを実行し、作成されたマイルストーン計画に従って次のマイルストーンの調整/開始を行います。すべての未完成の課題は、マイルストーンから削除するか、次のマイルストーンにロールオーバーする必要があります。

プランニング・ドラムビートの目的

このプロセスの目的は、ベロシティの理解を深めることで、作業の計画と見積もりの能力を向上させることにあります。マイルストーン計画では、課題点を設定しますが、これはマイルストーン間で平均値を共有するため、一貫性を測るのに適しています。そして、完成した仕事にコミットできると考えられる時期に基づいて、マイルストーンに引き込みます。そして、優先順位に応じて課題の優先順位を決定します。

一度に担当する課題は2つ以下になるでしょう。このアプローチ

には、以下のような多くの利点があります。

- 1. 最優先のプロジェクトを確実に完了させることができる
  - 2. リーダーシップが阻害されている問題を特定することができます。

3. データ機能の外部から優先順位を設定されている専門アナリストを含むデータチームの作業にリーダーシップの視点 を提供します。

- 4. チームメンバーの安定したスループットを促進する
- 5. ステークホルダーに対して、自分の要求がどこに優先しているのかを明確にする。
- 6. チーム全員に、自分の任務がどこにあるのか、優先順位を明確にする。
- 7. 次のマイルストーンを計画する際に、課題がすでにランク付けされているため、プレッシ

ャーを軽減することができます。計画の鼓動と実行は、チームの努力の賜物です。

layout: handbook-page-toc title:"Data Triage Guide" 説明"GitLab データ・トリアージ・ガイド"

このページについて

{:.no toc.hidden-md.hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg},

{::options parse\_block\_html="true" /}。

データトリアージ

Triager

データチームでは、各国のタイムゾーンを考慮して、以下のようなトリアージスケジュールを実施しています。

UTCデー	データアナリスト	データエンジニア
日曜日	@ken_aguilar	-
月曜日	@chrissharp	VEDPRAKASH2021
火曜日	@chrissharp	@paul_armstrong
水曜日	@michellecooper	rbacovic
木曜日	@ken_aguilar	ラドゥーラ
金曜日	@jeanpeguero	@jjstark

休暇中や優先度の高いプロジェクトに従事しているチームメンバーは、カバー範囲を探し、カバー範囲を引き継ぐチームに連絡する責任があります。これはデータチームのGoogleカレンダーで更新する必要があります。毎週同じ人に*月曜日の仕事を*押し付けないように、データエンジニアは時々交代でトリアージの日を設け、その場限りのコラボレーションを行います。

チームに専任のトリアジャーがいることで、傍観者の影響にも対応できます。スケジュールは日々のオーナーシップ情報を明確に共有していますが、オンコールのポジションではありません。オーナーシップを明確にすることで、チームの他の全員が1日の大半を深い仕事に費やす余地を作ることができます。トライアングルは、このような追加の時間的要求があってもうまく達成できるような仕事の種類に合わせて一日の計画を立てるようにします。

Data Triagerは、Dataチームへのリクエストや問題に最初に対応する人です。

- Data Analyst triagerは、GitLabチームメンバーからのリクエストに応えることを主な仕事としています。#data-triage に投稿されるissueを介して、または#dataのslackで直接リクエストに応えます。
- データエンジニアのトリアージ担当者は、主に当社のデータプラットフォームに関する問題を解決する責任があります。 これらの問題は、#analytics-pipelines スラックチャンネルか、DE-TriageErrors ボードに投稿されます。これらのエラーから 作成された課題は、DETriageErrorsの課題テンプレートを使用する必要があります。
  - 。 割り当てられたトリアージの日には、データエンジニアは主にアクティブな問題や、DE-TriageErrorsボードに 表示された問題に集中しなければなりません。
- データアナリストとデータエンジニアのトライアングルは、それぞれプラットフォームの問題とGitLabチームメンバーのリクエストを担当します。

• トリアージの責任についての詳細は、データエンジニアのトリアージトレーニングセッションのビデオをご覧ください。

他の GitLab チームメンバーからデータチームのプロジェクトに入ってきた問題の多くは、理解、見積もり、優先順位付けのために追加の情報やコンテキストが必要です。このような質問をして、問題を早期に解決することがトリアジャーの優先事項です。

### 注: Data Analyst triagerの場合

- は、すべての質問に対する答えを知っているとは限らない。
- 対象分野の専門家 (SME) や機能的なDRIである他のチームメンバーにcc (言及) してフィードバックを求めるべきで
- ある。 躓いた問題を文書化し、チームメンバー全員に知識を広めるのを助けるべきである。

データトリアージの方法

DataTeamプロジェクトで課題を作成します。タスクと任務は課題テンプレートに記載され

ています。よくある質問と一般的な問題をお読みください。

トリアージ・ボット

トリアージの一部は、Analyticsプロジェクトに設置されたGitLab Triage Botによって支援されます。このボットは1時間ごとに実行され、ポリシーファイルで定義されたルールに基づいてアクションを起こします。GitLabTriageのREADMEには、ルールのフォーマットや定義に関するすべてのドキュメントが含まれています。

triage botのポリシーファイルの変更は、CIジョブ "dry-run:triage "を実行して、ログ出力を確認することで、MRでテストする必要があります。このCIジョブはドライランです。つまり、プロジェクト内で実際には何もアクションを起こさず、ポリシーが実際に実行された場合に何が起こるかを出力します。

エンドオブデイ・ラップアップ

日々のトリアージをより良く、より効率的に行うために、その日のうちに作業をまとめます。以下の情報は、データアナリストとデータエンジニアが毎日提供しています。

- タイムトラッキング。その日のトリアージにかかった時間と、どのような活動が行われたかを記録すること。一般的な 考え方としては、どこに作業負荷がかかっているかを理解し、それらの活動を最適化することです。分析目的に役立つ ように作業内容を記載してください。例えば
  - データインジェストの問題解決に2時間、データ
  - 変換の問題解決にソース×1時間
  - 壊れたSisenseダッシュボードの解決に1時間
  - トリアージと(新しい)問題の再ルーティングに1時間
  - Slackの質問への回答に0.5時間
- グラウンドホッグの問題定期的に発生する問題は、私たちを悩ませ、無駄な時間を費やしています。何度も出てくる 問題をリストアップしてみてください。これは、現在の風景の中で対処すべきスポットを見つけるためです。

トリアージラウンドアップは、各マイルストーンの終わりに、データリーダーシップチームによって行われ、マイルストーンのトリアージの取り組みを集約します。情報提供の目的は、それを有用なものとし、Triageを向上させることであることを念頭に置いてください。

GitLab.comのDB構造の変更

定期的に変更される最も重要なデータソースのひとつが、GitLab.comのデータベースです。日々の業務を中断させないためには、データベースへの変更を追跡・確認する必要があります。GitLab.comデータベースへの変更はすべて、db/structure.sqlファイルに行われます。データチームは、db/structure.sqlに変更が加えられた場合、MRにラベルを適用することで、Danger Botを介して通知を受けます。

Data Warehouse::Impact CheckというラベルがDanger Botによってデータチームの行動を促すために追加されます。

トリアージでは、TriagerはData Warehouse::Impact CheckというラベルでMRをチェックします。

以下のアクションは、データチーム・トライジャーによって実

行されます。すべてのMRが判定される

。 **SQL**ファイルの変更によって操作が中断されていない場合は、ラベルはData Warehouse::Not Impacted に変更されます。

- SQLファイルの変更により操作が中断された場合。ラベルは
  - Data Warehouse::Impactedに変更されます。
  - 新しい課題がGitLab Data Teamプロジェクトに開設され、正しいDRIに割り当てられ、元のMRにリンクされます。
  - 影響は課題の中で決定されます。
  - 任意のMRは、ロードの問題、下流のdbt処理、Sisenseの使用を克服するために作成されます。
  - GitLab.comのMRのMergeによると、Mergeは計画的に行われます。
  - すべての関係者に通知されます。

判定マトリクス。\*\*

## ChangeCall to Actionが必要\*。 X 新しいテーブルの作成 $\overline{\checkmark}$ 表削除 $\overline{\checkmark}$ テーブル名の変更 フィールド追加 X $\overline{\mathbf{V}}$ 削除されたフィールド フィールド名の変更 $\overline{\checkmark}$ フィールドのデータタイプの変更 ? 制約条件の変更

\*デフォルトでは、すべてのテーブルやカラムをロードしていません。そのため、新しいテーブルやカラムが追加されても、特定のビジネス上の要求があった場合にのみ、これらのテーブルを読み込むことになります。現在の構造に変更を加えることで、操作が中断される可能性があるかどうかを判断する必要があります。

判断材料は豊富ではありません。すべてのMRを注意深くチェックする必要があります。

GitLab Postgres データベースにアクセスできない

gitlabにクローンされたPostgresデータベースにアクセスできない場合、airflowのタスクログに以下のエラーが表示されます。

 $sqlalchemy.exc. Operational Error: (psycopg 2. Operational Error) \ FATAL: \ The \ database \ system \ is \ starting \ up \ n$ 

b'FATAL: The database system is starting up\'D

#### 以下の手順で行います。

- 1. DE Triageテンプレートを使用して課題を開きます。
- **2.** gitlab\_com\_data\_reconciliation\_extract\_load、gitlab\_com\_db\_extract、gitlab\_com\_db\_incremental\_backfill 、gitlab\_com\_scd\_db\_sync という名前の gitlab.com の DAG をすべて一時停止します。その理由は、警告を抑えて不要なリソースを使わないためです。
- 3. アラートチャンネルを見て、「GitLab Job has failed」を検索してアラートを見つけます。サンプルのアラートは以下のような内容です。

Firing 1 - GitLab Job has failed

GitLabジョブの "clone "リソース "zlonk.datalytics.dailyx "が失敗しました。

:chart:プロメテウスのグラフを見る:label:ラベ

ルです。アラートネーム: JobFailed

Alert\_type: Symptom (症状)

環境: GPRD 環境: GPRD

Fqdn: blackbox-01-inf-gprd.c.gitlab-production.internal

Job: clone モニタ:default プロバイダ:gcp 地域:us-east

リソース: zlonk.datalytics.dailyx

Severity: s3 Shard: default Stage: main Tier: db タイプ:

zlonk.postgres 少ない

表示

4. 問題を調査するために @sre-oncall の slack ハンドルに連絡し、incidentdeclare を使用してインシデントリクエストを上げてください。これにより、SREオンコールチームが対応するための本番のインシデント問題が作成されます。また、インシデントをより広く認知してもらうために、@gitlab-data/engineersにも連絡してください。

- 5. インフラの課題と、提起されたトリアージの課題をリンクします。
- 6. 問題が解決するか、@sre-oncallの担当者またはDBREチームの誰かから確認されたら、接続の安定性を検証するために、失敗したタスクの1つを単独でクリアして再実行してみてください。
- 7. DAGの場合、gitlab\_com\_scd\_db\_sync、gitlab\_com\_data\_reconciliation\_extract\_load、 gitlab\_com\_db\_incremental\_backfillは、24時間以内に1回しか実行されないので、失敗したタスクをクリアして、実行するようにします。
- 8. DBTがソースの更新に間に合わない日に実行された場合は、ソースの更新が遅れていることをトリアージテンプレートを使って#dataチャンネルに通知します。

Zuora Stitch Integration シングルまたはセットのテーブルレベルリセット

テーブルを完全に埋め戻すために、StitchでZuoraデータパイプラインのテーブルをリセットすることは、どのような場合にも起こり得ます(例:ソースに新しいカラムが追加された、テクニカルエラーが発生したなど)。現在、ZuoraStitch統合はテーブルレベルのリセットを提供してないため、統合内のすべてのテーブルのリセットを実行する必要がますこの余分なコストとリスクが発生ます

このため、以下の手順でテーブルレベルのリセットに成功しています。この例では、Zuora購読テーブルを使用していますが、Stitch Zuoraデータパイプラインの他のテーブルにも適用できます。

ステップ1:-既存のテーブルの名前を変更し、バックアップを識別するために日付の接尾辞を付けます(推奨フォーマットはYYYYMMDD)。

alter table "raw". "zuora\_stitch". "subscription "の名前を "raw". "zuora\_stitch". "subscription\_20210903 "に変更します。

Step 2:-定期的な統合を一時停止する。

☑通常の統合を一時停止する

ステップ3: - Stitchで新しい統合Zuora-Subscriptionを作成します。

この設定では、抽出頻度を30分に、抽出日を2012年1月1日に設定し、すべてのデータが抽出されるようにしました。

複製するのはサブスクリプションテーブルの

み ステップ4:-新しく作成した統合を実行しま

す。

新しく作成した統合機能を手動で実行し、完了するのを待ってください。完了すると、ホームページに正常に表示されます。完了したら、次のステップに進む間、データがずれないようにするため、新しく作成した統合タスクを一時停止します。

### Step 5:-記録を確認する

新規に作成したテーブル "RAW". "ZUORASUBSCRIPTION". "SUBSCRIPTION"で、統合UIでstitchに読み込まれたと表示される行数と テーブルに読み込まれた行数が同じかどうかをクロスチェックします。

Step 6:-メインスキーマにテーブルを作成します。

新たに読み込んだデータをZUORA\_STITCHスキーマに移動します。なぜなら、新しい統合ではテーブルが画像の上に記載されているZUORASUBSCRIPTIONです。

create table "raw". "zuora\_stitch". "subscription" clone "raw". "zuorasubscription". "subscription";

注:テーブルのポストクローンに主キーが存在するかどうかをチェックし、存在しない場合はリンクに主キーをチェックし、それらのカラムに制約を追加します。

ステップ7:-レコード数のチェックを行い、新しいテーブルのレコード数が少なくなっていないかを確認します。

select count(\*) from "RAW". "ZUORA\_STITCH". "SUBSCRIPTION\_20210903" where deleted =
'FALSE'; select count(\*) from "RAW". "ZUORA\_STITCH". "SUBSCRIPTION";

ステップ8:-新しいスキーマをドロップする

drop schema "raw". "zuorasubscription" cascade ;

ステップ9:一時的なZuora-Subscription統合を削除し、通常の統合を有効にする ステ

ップ10:-通常の統合を実行し、検証する

これは、テーブルに以前発生したエラーがなくなり、データがテーブルに入力されていることを確認するためです。2つの異なる抽出器による重複したIDをチェックし、データが正しくテーブルに入力されていることを確認します。

select id, count(\*) from "RAW". "ZUORA\_STITCH". "SUBSCRIPTION" group by id カウント(\*)が1の場合

メモ 詳細はMRを参照してください。

## トリアージFAO

Data Triageは24時間365日のサポート、もしくは24時間サポートしなければならないシフトですか? 通常の勤務時間内に、トリアージテンプレートに記載されているトリアージ当日のタスクリストを実行する必要があります。

何か問題が見つかった場合、直接本番で修正するのか、それともインシデントの一部として捉え、定義された時間内に解決するのか。

トリアージの日には、その場にいたデータチームのメンバーが、すべての失敗、質問、エラーを探します。

- Slackチャンネル; #data-prom-alerts #analytics-pipelines and #data 新たに追
- 加された課題
- Sisense の TDF ダッシュボード

これには、最後の人がサインオフしてからのすべての失敗が含まれ、それ以降、その人がサインオフするまでのすべての失敗 について問題が発生します。データのパイプラインが破損し、データの読み込みや更新に遅延が予想される場合。その場合

関係するチームに[トリアージテンプレート]を使って通知しなければなりません(https://gitlab.com/gitlab-data/analytics/-/issues/new?

異なる種類の問題については、ETAがありますか?

パイプラインが壊れている場合は、修正する必要があります。現在、私たちはデータ資産のSLOの定義に取り組んでいます。データ抽出パイプラインについては、こちらに包括的な概要が掲載されています。

トリアージの日に通常の勤務時間、つまり米国時間の午前11時まで働いたとします。通常の勤務時間後にパイプラインが 切断され、データの入手が遅れた場合はどうなりますか?

そうですね、私たちの存在のメリットは、時間のオーバー幅が広いことです。トリアージに参加している人が米国のタイムラインに先行していれば、問題をタイムリーに解決できるというメリットがあります。デメリットは、米国のタイムラインをその日のうちにフルカバーできないことです。これは将来に向けての注意点です。

一般的な問題のトリアージ

このセクションでは、一般的な問題とその解決策について説明します。

エアフロータスクの失敗!?

DAG gitlab com db extract

タスク gitlab-com-dbt-incremental-source-freshness

背景を説明します。このエキスは、GitLab.com環境のコピー(レプリケーション)データベースに依存しています。これがレプリケーションの遅延の原因となっている可能性が高いです。

セットアップの詳細はこちら

考えられる手順、解決策、対処法-レプリケーション・ラグの確認

- 必要に応じてDAGを一時停止する
- データギャップの確認
- 埋め戻し作業
- DAGのリスケジュール

注: GitLab.comのデータソースは非常に重要なデータソースであり、一般的に使用されています。ビジネス関係者に適宜お知らせください。

便利な正規表現

これらの用語が存在しない行にマッチ

^(?!.\*(<一番目に探す用語>)/.\*\*

例: Airflowのログをきれいにするため。

^(?!.\*(テストでの失敗|データベースのエラー)).\*\$

layout: handbook-page-toc title:"Data Team Learning and Resources" 説明"GitLab データチームライブラリ"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

## Powered by GitLab Team Members

- How Data Teams Do More With Less By Adopting Software Engineering Best Practices Thomas's talk at 2018
   DataEngConf in NYC
- ナッシュビルで開催された「2019 Music City Tech Conference」でテイラーが行った同様の講演のスライド

- LocallyOptimisticAMA:August2019withTaylorMurphy 才一
- プンソースアナリティクスの力を示す4つの例 GitLabClを
- 使った初めてのdbtプロジェクトのデプロイメント
- クラウドネイティブな世界でのデータオプス
- データ分析の3つのレベル-データ組織の成熟度を評価するためのフレームワーク GitLabを
- 使ったDataOpsの導入方法
- GitLab Dataチームの運営で得た教訓
- Views on Vue Podcast with Jacob Schatz and Taylor Murphy
- HowtodoDataOpswithGitLab-CustomerCall (GitLab Internal)
- GitLabforML-CustomerCall (GitLab Internal)

## おすすめの読み物、聴き物、見物

- AIのニーズの階層 データ
- ・メタ・メトリクス
- エンジニアはETLを書くべきではない
- スタートアップ創業者のためのアナリティクスガイド
- Functional Data Engineering バッチデータ処理のための最新パラダイム
- KeepitSQLStupid、DataEngConfSF'18でFishtownAnalyticsのConnorMcArthurがAIのためのdbt DevOpsを
- 説明した講演です。
- データサイエンティストはDevOpsから何を学ぶか
- 一人のアナリストが語る「GoodtoGreat」への道
- データの価値とは?第1回、第2回、第3回
- データプラクティスの構築
- スタートアップのデータチームにデータエンジニアは必要か?
- データサイエンスは今までとは違う(注: GitLabにまだデータサイエンティストがいないの
- はこのためです) WhyDon'NeedDataScientists
- リソースWrittenbydbtCommunityMembers あな
- たの会社はデータドリブンには向いていません
- か?あなたにとって「セルフサービス」アナリティクスとは
- 組織内のデータサイエンスチームを統合するためのモデル
- 成熟したアナリティクス・ワークフローの構築(注:「アナリティクスはソフトウェア工学のサブフィールド」という 前提を説明しています。)
- YouTubeのDataOpsプレイリスト

## データニュースレター

- サルゴス&エ
- シックス カロ ージカ
- The Carpentries
- DataEng Weekly
- DataElixir
- Data is Plural
- データサイエンス・ラウンドアップ・
- ニュースレター データサイエンス・
- ・ ウィークリー
- LantrnsAnalytics (プロダクトアナリティ
- クス) ミュージック&テック
- Normcore Tech NumLock News

**OneShotLearning** 

**SFData** 

# データブログ

Airbnb AskGoodQuestions BufferBlog

- •
- •
- •
- •
- •
- •
- •
- •
- ٠
- ٠
- •

- カロージカ
- FishtownAnalyticsBlog
- GoDataDriven
- MBAMondays
- ModeAnalyticsBlog
- Multithreaded
- Sisense Data Blog
- Silota
- WesMcKinneyBlog デー
- タ運用
- レティナAIブログ
- StitchFixAlgorithmsBlog
- Five Thirty Eight
- Data.gov

# データビジュアライゼーションリソース

- データを使ったストーリ
- ーテリング

## **DataRevelations**

- データ・ビジュアライゼーシ
- ョン・カタログ EagerEyes
- FiveThirtyEight'sDataLab
- FlowingData
- データからViz
- Gravy Anecdote
- ~ JunkCharts
- MakeaPowerfulPoint
- Makeover Monday
- Perceptual Edge
- PolicyViz
- TheFunctionalArt
- The Pudding
- Visualising Data
- VizWiz
- WTF ビジュアライゼーション

## データ スラック コミュニティ

- データビズ協会
- データサイエンス・コミュ
- 🚦 ニティ dbt
- 。大いなる期待 Locally
- Optimistic Measure
- ・メルタノ
- 。 オープンデータコミュニティ
- Pachyderm
- プリフェクトパ
- \_ イカロリナス
- R for Data Analysis
- SoftwareEngineeringDaily
- **TheDataSchool**

# テクニカル・ラーニング・リソース

- クリス・アルボン
- Mode SQL チュートリアル

- dbt Tutorial
- Technically.devPostonSQL
- ElementsofDataScience
- MachineLearningResources (GitLab Internal)
- Codecademy
- DataQuest
- KhanAcademy
- HackerRank(Exercises)
- Udacity
- StanfordUniversityMini-Courses
- The Data School by Chartio
- W3Schools

layout: handbook-page-toc title:"Data Team Organization" description:"GitLab データチームの組織"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

## データチームの編成

データチームは、卓越した技術と主題を持つセンターで構成されており、フュージョンチームと呼ばれるデータソリューションを提供するビジネスに特化したチームとビジネスに関与するチームによって補完されています。データプラットフォーム&エンジニアリングチームは、データソース、パイプライン、分析ツールに加えて、高性能なコンピュートとストレージのレイヤーを提供することで、データフュージョンチームを補完しています。そして、データチームを補完し、高度な分析機能を提供するのが、データサイエンスチームです。チームの構成や、GitLabの他のチームとの連携については、How Data Works at GitLabをご覧ください。

データ&アナリティクス・デモ

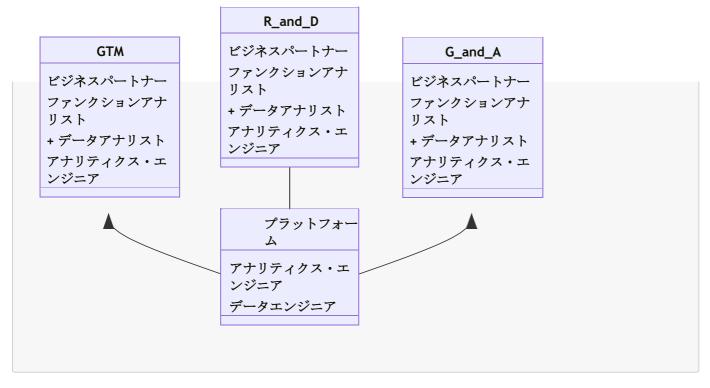
データ&アナリティクス・デモは、データプログラムに関わるすべての人が、進歩や革新を共有し、協力し合い、ただ楽しむための素晴らしい方法です。データ&アナリティクス・デモ毎週木曜日され、録画GitLabUnfilderedDataTeamのプレイリストに投稿れ。

データプログラム採用動画

この度、GitLabDataRecruitingのビデオを作成しました。ぜひご覧ください。

データフュージョンチームの組織

データフュージョンチームは、ビジネスとデータチームの間のチームメンバーで構成されています。現在のデータフュージョンチームについては、フロントページでご紹介しています。



フュージョン・チーム・アサイメント

		GTMR&DG&A
iweeks Lead/DRI	@mlaanen Lead/DRI	@iweeks DRI
	@ken_aguilar@snalamaru	@pempey $\mathcal{Y} - \mathcal{F}$
	Jeanpeguero@chrissharp	@jjstark データプラットフォームチーム 安定したカウンターパート
@michellecooper	@jeanpeguero GTM Fusion Team 安定したカウンターパート	
@paul_armstrong データプラットフォ ームチーム 安定したカウンターパート	@rbacovic データプラットフ: カウンターパート	オームチーム 安定した

データプラットフォームチーム 安定したカウンターパート

GitLab Stable Counterpartの原則に従い、すべてのFusion TeamにはData Platform Team Stable Counterpartが割り当てられています。データプラットフォームチームのステイブルカウンターパートは、データプラットフォームチームとフュージョンチームの間で、自分の時間、仕事、優先順位を分けています(一般的に平均50%ずつ)。安定したカウンターパートは、フュージョンチームの方向性や優先順位を把握しており、必要に応じてデータプラットフォームチームと議論します。例えば、安定したカウンターパートが割り当てられた時間内に処理できる以上の需要がある場合や、アーキテクチャの方向性を変更する必要がある場合などです。安定したカウンターパートは、これを認識し、フラグを立て、該当するステークホルダー(一般的には、データプラットフォームチームとFusionチームのリード/DRI)に対処します。

フュージョン・チーム・オペレーション

データフュージョンチームを成功させるためには、以下の要素が不可欠です。

- **1. DataforGTMSeries**」や「**DataforR&DSeries**」を通じて、ビジネスパートナーやデータチャンピオンとの定期的かつ透明性の高いエンゲージメントを行います。
- 2. データチームによる作業計画の立案 プランニング・ドラムビート
- 3. 定期的なCSAT調査を実施し、継続的な改善を目指してデータフュージョンチームにフィードバックを提供する

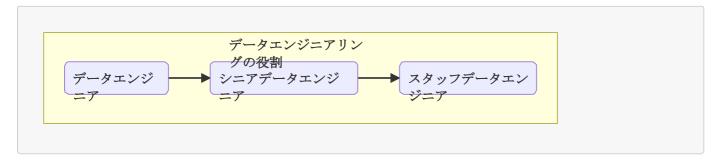
ステークホルダーの皆様には、私たちのイシューボードに沿って作業範囲を理解していただくことをお勧めします。

#### 1. **GTM**

- 2. 研究開発
- 3. G&A: 人、エンジニアリング・アナリティクスは近日公開
- 4. プラットフォーム

データの役割とキャリア開発

データエンジニアリング



## データ分析

```
データアナリスト

の役割

ジュニア・データ・

インターン

フィックアナ
リスト

フィックアナリスト

フィックアナリスト
フィックアナリスト
フィックアナリスト
フィックアナリスト
```

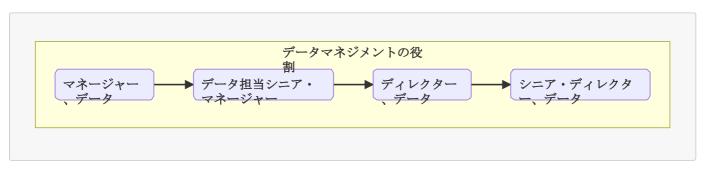
### データサイエンス

```
データサイエンス
の役割
データサイエ
ンティスト
シニアデータサイエ
ンティスト
```

## アナリティクス・エンジニアリング



## データ管理



layout: handbook-page-toc title:"Data Analytics Handbook" description:"GitLabデータ分析チームハンドブック"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# GitLabでのデータ分析

データ・アナリティクス・チームのミッションは、信頼性と拡張性のあるデータ・ソリューションによって、ビジネス上の意 思決定と戦略のインパクトを最大化することです。

私たちは、GitLabバリューとデータチーム・プリンシプルによって、すべてのGitLabチームがデータ・バリュー・ピラミッドの上に到達するのを支援します。

データ分析の責任

Data Team's Responsibilitiesのうち、Data Analytics Teamが直接担当するのは以下の通りです。

- 会社の主要業績評価指標の定義、データベース、およびデータの可視化の管理と運営データの可視化、データモデリング、デー
- タ品質、およびデータ統合を含むカスタマイズ可能なデータサービスの提供データ&アナリティクスシステムに関連する会社のガ
- バナンス、リスク、およびコンプライアンスプログラムのサポート

さらに、データ・アナリティクス・チームは以下の責務をサポートします。

- データリーダーシップチームと
  - 。 GitLabのデータ資産の価値を最大化するためのデータ戦略の策定と公開 データの成果物、進行中
  - のイニシアチブ、今後の計画に関する定期的なアップデートの配信
- データエンジニアリングチームと
  - 。 すべてのGitLabチームのレポーティング、分析、ディメンショナルモデリング、データ開発をサポートするための、会社の中心となるエンタープライズデータウェアハウスの構築と維持
  - マスターデータ、リファレンスデータ、データ品質、データカタログ、データパブリッシングなどのデータマネジメント機能の開発
  - o 誰もがデータと分析を活用できるよう、セルフサービスのデータ機能を提供する
  - GitLabデータシステムのデータ品質に関するプラクティスとプログラムの定義と推進の支援

データ・アナリティクス・チームの最初のハンドブック

GitLabでは、Handbook Firstを掲げ、これを推進しています。

- 1. ビジネスパートナー(すべてのGitLabチームメンバー)がハンドブックを最新かつ正確な情報で更新するのを支援する。
- 2. 内容を見直し、データの機密性やセキュリティを考慮した上で、ハンドブックの情報がデータを最もよく表していること を確認する。
- 3. データチームのハンドブックに、データアナ
  - リティクスに関するエキサイティング な内容が追加されました。
    - **今後のコンテンツの例**データビジュアライゼーションのヒントとコツ」「データストーリーテリング
  - データチームのプロセスやコードベースの変更を反映させる(より良いプロセスやコードに興奮するのは当然のことです。

アナリティクスとは何か、そしてなぜそれを行うのか?

アナリティクスとは、データを情報、知識、そして知恵に変えていくことです。

- 1. 効率化とコスト削減
  - 1. Human Capital Location Factor, Offer Accept Rate, Cost Per Hire
  - 2. ファイナンス ARR、新規パイプラインの作成、CC障害率
  - 3. セールスチャネル-勝率、サービス添付率。
  - 4. マーケティング 機会の創出、製品のダウンロード
- 2. ビジネスモデルの改善
  - 1. プロダクト・インテリジェンス サインアップの増加、新機能の提供、トップ機能
  - 2. 顧客分析 企業統計、採用傾向、利用とサブスクリプションの比較
  - 3. お客様の購買動向 リニューアル、アップグレード、ダウングレード
- 3. カスタマー・エクスペリエンスの向上
  - 1. データ製品 業界ベンチマーク、データAPI、アルゴリズム
  - 2. データに基づいたプロセスの改善
  - 3. インプロダクトインサイト

## 短期的な方向性

21年度から下期にかけての短期的な目標は、GitLabをデータ・キャパシティ・モデルの「レベル (1) リアクティブ」から「レベル (2) アドバンスド」に引き上げることです。

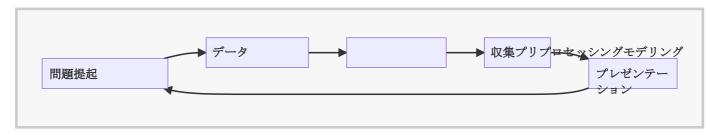
これを実現するために、データアナリティクスチームは

- データフュージョンチームの中で、同じビジネス目標に向かって、同じペースでデ
- ータ分析プロセスを強化することができる。
- GitLabのデータ・ビジュアライゼーションとデータ・ストーリーテリングのスキルを高める

データ分析プロセス

データ分析」「高度なデータ分析」「データサイエンス」は、すべて「データ分析プロセス」から始ま

ります。データ分析プロセスは5つのステップで構成されています。



データ分析では、コンテキストがすべてです。コンテキストは、データアナリストが問題やデータをどのように捉え、どのような方法でデータからの洞察を得るかを導きます。

### 問題提起

問題記述」のステップでは、ビジネスパートナーと一緒にビジネス上の問題を定義することが重要です。データアナリストは、ビジネスパートナーがデータの洞察を求めているビジネス上の問題を明確に定義することで、その問題が利用可能なデータで技術的に解決可能であることを保証し、その問題に対して最大限のビジネス上の洞察を提供するための成功基準を設定することができます。

#### データ収集

データ収集のステップでは、データアナリストは、分析の影響や一般化を制限する可能性のある「データの禁忌」を調べて理解し、「データの偏り」や「データの懸念」を評価することが重要です。各ビジネスシステムでどのようにデータが収集されているかを理解することで、正確なデータを適切な方法で報告することができます。このステップでは、信頼性の高い変換を行うために、データタイプが正確であることを確認することが重要です。

したがって、このステップはソースシステムの所有者とのフィードバックループの一部となります。データチームのメンバーが、あるデータフィールドが欠けていたり、正確に把握されていなかったりすることを発見した場合、GitLabの誰もがソースシステムのオーナーに連絡を取り、データフィールドの更新や追加を依頼する責任があります。

#### 前処理

前処理ステップでは、データアナリストがデータを準備し、クレンジングし、データの品質を検証して、スケーラブルなレポート作成に最適なデータ構造を設計します。このステップでは、モデルをレビューし、ビジネスに適したモデルであることを確認することが重要です。

### モデリング

モデリングのステップでは、GitLabチームメンバー全員のために、Enterprise Dimensional Model形式の新しいデータモデルを作成することを目指します。まず、新しいデータ構造が正確なビジネスプロセスを反映していることを確認するために、エンティティリレーションシップダイアグラム(ERD)から始めます。すべてのデータモデルは、データエンジニアリングチームによってレビューされます。

GitLabがデータ能力モデルのレベル(4)予測に達した場合、このステップでは探索的データ分析とデータフィーチャリングが行われます。

### プレゼンテーション

プレゼンテーションのステップでは、データアナリストは、強力なデータ・ビジュアライゼーションによるデータ・ストーリーテリングのスキルを発揮します。データアナリストは、利害関係者と効果的にコミュニケーションをとるために、実用的なビジネスインサイトをパッケージ化して表示することが重要であると認識しています。実際、ビジネスパートナー(GitLabチームメンバー)にデータを信頼してもらうためには、洞察の背景にある詳細な理由を示すドリルダウン機能を備えたビジネスインテリジェンス(BI) ダッシュボードなどのデータソリューションに、強固なデータ品質チェックを設けることと同じくらい重要です。

データ分析の素晴らしいところは、それぞれの洞察が、解決すべき新たなビジネス上の疑問の波を伝播させ、それによってデータ 分析プロセスが再び循環することです。

layout: handbook-page-toc title:"データサイエンスハンドブック" 説明"GitLab データサイエンスチームハンドブック"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# ギットラボのデータサイエンス

データサイエンスチームのミッションは、モデルベースのインサイトにより、ビジネス、顧客、製品をよりよく理解することです。このチームは、機械学習モデルを構築し、データを分析して、その結果を信頼できるスケーラブルな製品に変換することで、このミッションを達成します。

## プロジェクト

### 買い方の傾向

- ◆ 開始日:2021-06-30 プロジェ
- クト
- データサイエンスの仕事の内訳 ス
- ラックチャンネル

#### プロジェクトバックログ

- ・ 購入(アップセル) 予測ARR
- ゴールデンジャーニー グロースチーム Hila
- Qu PtB.first order グロースチーム Hila Qu

• PtB.churn - カスタマー・サクセス・チーム - David Sakamoto

2021/11/8 merged.md

- ユーザーセグメンテーション/ペルソナ活動、ユースケース 製品チーム -
- Hila Qu コミュニティセンチメント分析/twitter、facebook 製品チーム
- GitLab MLOps 製品開発 製品チーム Taylor
- Feature \$ARR uplift prediction Product Team Anoop Dawar

データサイエンスロードマップ - エグゼクティブサマリー

			タイムライン実現	記可能な価値 ホーリス ボーリス ボーリス イン れる 成果
<b>FY-22 Q3 /</b> イテレーテ ィブ	拡大の傾向( PtE	営業チームがARRを増加させる機会を獲得できるようにする。	セールスチーム( Jake Bielecki	TBD
<b>FY-22 Q4 /</b> イテレーテ <del>イブ</del>	<del>契約傾向(PtC</del>	<del>- 営業チームがARRの減少を防ぐことができ</del> る。	セールスチーム( Jake Bielecki	TBD
FY-23 Q1 / Q2 イテレーテ ィブ	ユーザーセグメンテ ーション(US	お客様	のプロファイ グを安定する ーム (Hila Qu	TBD

- データサイエンスとは一体何なのか!?こちらのビデオやスライドで概要をご紹介しま
- す。データサイエンティストの責任については、こちらをご覧ください。
- さらに、データサイエンスチームは以下のような責務をサポートします。
  - o データリーダーシップチームと
    - 分析能力を拡大するためのデータサイエンスのロードマップの策定と公開 データの成果物、進行
    - 中のイニシアチブ、今後の計画についての定期的なアップデートの配信
  - データエンジニアリングチームと
    - GitLabデータシステムのデータ品質に関するプラクティスとプログラムの定義と推進の支援

# プロジェクトの構成

データサイエンティストは、CRISP-DM (Cross-IndustryStandardProcessforData Mining) に従い、6つのフェーズで構成される プロジェクトを遂行します。プロジェクト構造の詳細な内訳とコード例はこちらをご覧ください。

#### ビジネス理解

最初のステップは、プロジェクトのエンドユーザーのビジネスニーズと成功基準を定義することです。このフェーズには 、要件の収集、ステークホルダーへのインタビュー、ビジョンの定義、製品のユーザーストーリー、モデルのユースケー スなどが含まれる。このフェーズは、GitLabのOKRや会社の価値観であるイテレーションと密接に関連しています。

#### データの理解

データの理解は、ビジネスの理解と密接に関連しています。この段階では、ビジネスがすでに導入しているデータソースと、そう でないデータソースを定義する必要があります。この段階では、データサイエンティストは、データエンジニアやデータアナ リストと密接に連携し、不一致やリスクを定義します。外部データを導入する必要がある場合は、このフェーズでリクエス トを出す必要があります。

## データの準備

ビジネスニーズとデータソースを定義した後は、いよいよデータに飛び込んでいきます。この段階では、「2回目のデータ品 質チェック」と「探索的データ分析」を行う必要があります。この段階では、データサイエンティストは、次の段階であるモデ

リングに役立つデータ、その種類、相関関係、分布についての理解を深めます。

モデリング

モデリングとは、機械学習モデルを構築することであり、モデルのトレーニング、テスト、バリデーションから構成されます。 ビジネス理解の段階では、最初のアルゴリズムを選択し、前処理を定義する必要があります。このフェーズでは、すべて のアルゴリズムがすべてのデータタイプで動作するわけではないため、データ準備の成果が必要となります。

評価

このフェーズでは、モデルのパフォーマンスを測定し(ビジネス・アンダースタンディングで必要な場合)、その結果をステークホルダーやエンドユーザーに提示してフィードバックを得る。このフェーズは、GitLabのバリューの1つである「反復」と強く結びついている。フィードバックの後、データサイエンティストは、ビジネスゴールが達成される限り、上記のフェーズを反復する。

展開• 生産性向上

実行可能な予測モデルが作成されたら、次はそのモデルをTBDデータサイエンスのプロダクションパイプラインに展開します。このプロセスでは、新しいデータが利用可能になると自動的にモデルが更新され(「再学習」)、定期的にすべての対象レコードの予測が生成されます。

# 成功基準

GitLabのデータサイエンスプロジェクトは、すべて成功基準で始まり、成

功基準で終わります。成功基準は以下のように測定されます。

value\_provided = future\_state - current\_state

私たちは、データサイエンスプロジェクトから得られるビジネス上の成果として、直接的な影響と間接的な影響の2種類を認識しています。

ダイレクトインパクト

直接的な影響を与えるプロジェクトでは、プロジェクトの成果を直接金銭的価値に置き換えることができます。これは通常、プロセスにかかる時間を短縮したり、企業の収益を増加させたりすることを目的とした最適化プロジェクトによって達成されます。以下は、ダイレクトインパクトの成功基準の例です。

節約できた時間

"私たちのプロジェクトは、プロセスXに費やす時間を10,000時間、500kドルの価値で削減することに貢献しました。"

稼いだお金

"私たちのプロジェクトは、新しい有名な顧客をもたらし、年間収益に50百万ドルを追加しました"

間接的な影響

間接的な影響を与えるプロジェクトは、直接金銭的価値には結びつきませんが、他の意思決定者に金銭的価値をもたらすプラットフォームを提供します。これらのプロジェクトは通常、重要な洞察が拡張可能な製品としてステークホルダーに提供され、意思決定者の能力向上に利用されるものである。以下は、間接的な影響を与える場合の成功基準の例です。

視認性の向上

"私たちのプロジェクトにより、ステークホルダーの意思決定能力が高まり、販売戦略が10%改善されました。"

アンロックインサイト

「私たちのプロジェクトによって、製品チームは主要な顧客が誰であるか、そして彼らがどのような行動をとるかを最終的に理解することができました。この知識によって、ARR100万ドル相当の新規顧客を獲得することができました。"

プロジェクトの洞察/成果の記述

グーグルの公式「XYZ」に慣れること

☑ を実行することで、[Y]によって測定されたように、[Z]を実行します。"

X=何を達成したかを表す動詞に置き換える必要があります。

Y = は、達成度の尺度に置き換えてください。例えば、収益の増加やコストの削減があった場合、金額や割合を示すことができます。

**Z** = は、説明したい結果をどのように達成したかを表す動詞に置き換えてください。例えば、"詳細な分析を行うことで"**ABC**を達成した。

フルセンテンスの例は以下の通りです。

評細な分析を行うことで、10%の削減を実現します。

上記の式は、プロジェクトの成果をより効率的かつ簡潔に表現するための提案であることにご留意ください。

# Gitlabのデータサイエンスツ

ール

- 設定済みのJuypterLabイメージ。データサイエンスチームは、一般的なpythonモジュール(pandas、numpyなど)、ネイティブのSnowflake接続、gitサポートなどがあらかじめ設定されたJupyterLabを使用しています。共通のフレームワークで作業することで、モデルの作成やインサイトの導出を迅速に行うことができます。このセットアップは誰でも自由に使用することができます。詳細はJupyterGuideをご覧ください。
- gitlabds パイソンツール。一般的なデータ前処理(ダミーコーディング、外れ値検出、変数削減など)やモデリング作業 (モデルの性能評価など)を自動化するための関数です。pip install gitlabdsで直接インストールするか、上の JupyterLabイメージの一部として使用します.
- モデリングテンプレート (Coming Soon!)

# データサイエンスと機械学習に関する有用なリソース

- ジェイク・バンダープラス著「*Pythonデータサイエンスハンドブック」。*Pythonを学び、データサイエンスに足を踏み入れようとしている初心者の方に最適です。
- Sebastian Raschka & Vahid Mirjaliliによる *Python Machine Learning。* pythonの基礎レベルを前提とした、より高度なトピックを扱っています。
- The Elements of Stastical Learning, Data Mining, Inference, and Prediction」(Trevor Hastie、Robert Tibshirani、
  Jerome Friedman著)。一般的に使用されている予測技術の背景にある統計や論理を深く掘り下げて解説しています。ただし、統計学や数学の知識が必要な場合もあります。

# データサイエンスの共通用語

以下は、チームで使用される一般的なデータサイエンス用語です。

データサイエンス (DS) - コンピュータサイエンス、統計的手法、分野の専門知識を用いて、データから洞察を引き出す学際的な分野

機械学習(ML) - データのパターンを決定するために、明示的にプログラムされていないアルゴリズムの使用と開発

アルゴリズム - 特定の問題を解決するために使用される、コンピュータに実装可能な命令のシーケンス

特徴 - 国、年齢など、分析に使用できるデータセット内の単一列。変数や属性と呼ばれることもある

Feature Engineering - データを選択し、組み合わせ、特徴に変換するプロセス Weight - 特徴の強

さを決定するために割り当てられる数値

モデル - 入力データに基づいて決定された重みを持つ、適用されたアル

ゴリズム インピュテーション - 欠損したデータを置換した値で置き換え

るプロセストレーニング - 入力データに基づいてモデルに重みを割り

当てること

テスト - 学習したモデルの予測値と実際の値を比較する 分類 - 各データポイ

ントのカテゴリーを予測するプロセス

回帰 - 各データポイントについて連続的な数値変数を予測するプロセス クラスタリング -

データセット内のグループ化を定義するプロセス

# PythonとSQLの比較

SQL (Structured Query Language) に慣れている方は、Python (ライブラリpandasを使用) で分析を行う際に句を使用したいと思うかもしれません。以下に、一般的なSQLコマンドに相当するPythonのコマンドを示します。

なお、以下の方法は、期待通りの結果を得るための唯一の方法ではなく、同じ結果を得るために**Python**には多くの方法があることをご理解ください。

## **SELECT \* FROM**

#データセットを'data'として読み込んだと仮定すると、Jupyterのセルdataにデータフレームの名前を入れるだけです。

## SELECT col1, col2 FROM

```
columns = ['col1',

'col2'] data[columns].

上記は次のように1つのコードにまとめることができます: data[['col1','col2']].
```

## **COUNT**

#データフレームの長さを定義する(行数に相当) len(data)

上記と同様ですが、さらに列数を表示します data.size

### LIMIT

```
#最初の5行をリストア
ップ data.head()
```

#上位10行を表示したい場合は、head()の中に数字を入れることで実現できます function data.head(10)

### WHERE

```
#column1が1以上の値を持つすべての行をリストアップ condition
= data['column1'] > 1
データ[条件]
```

Pythonは複合条件を受け入れることに注意してください。PythonのANDに相当するものは&、ORに相当するものは | です。

## GROUP BY (集計機能付き

```
#列ごとの行数をリストアップ 1列目
data.groupby(['col1']).size()

#列1ごとの行の平均値を列挙する
data.groupby(['col1']).mean()

#列1ごとの行の最小値を列挙する
data.groupby(['col1']).min()

#列1ごとの行の最大値を列挙する
data.groupby(['col1']).max()
```

layout: handbook-page-toc title:"Data Engineering Handbook" description:"GitLab データエンジニアリングチームハンドブック"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# ギットラボのデータエンジニアリング

データエンジニアリングチームのミッションは、安全で信頼できるデータプラットフォームを構築し、誰もがアナリストになれるようにして、*唯一の*制限がデータやアナリスト自身になるようにすることです。私たちは、GitLabの価値観とデータチームの原則によってこれを実現します。

データエンジニアリングの責任

Data Team's Responsibilitiesのうち、データエンジニアリングチームは以下を直接担当しています。

- すべてのGitLabチームのレポーティング、分析、ディメンショナルモデリング、データ開発をサポートするための、会社の中心となるエンタープライズデータウェアハウスの構築と維持
- 新しいデータソースを統合し、対象分野、活動、プロセスの分析を可能にする。
- Single Source of Truthを実現するためのEnterprise Dimensional Modelの構築と維持
- マスターデータ、リファレンスデータ、データ品質、データカタログ、データパブリッシングなどのデータマネジメント機能の開発
- 誰もがデータと分析を活用できるよう、セルフサービスのデータ機能を提供する
- GitLabデータシステムのためのデータクオリティープラクティスとプログラムの定義と推進の支援

• データモデリング、データ品質、データ統合など、カスタマイズ可能なデータサービスの提供

データ・エンジニアリング・チームの最初のハンドブック

GitLabではHandbook Firstを掲げています。コードベースやプロセスに変更を加える場合は、作業を開始する前に、あるいはコードを変更するためのMRと同時に、ハンドブックのMRを作成しなければなりません。

例えば、以下のようなものがあります。データソースを追加する場合、まず、ハンドブックの抽出とロード、システム図のセクションを更新するために、マージリクエストを作成します。このMRは、変更が作成されると同時に更新され、変更と一緒にマージされます。

# データエンジニアリングチームの役割

チームとして効率的に活動するために、私たちはデータエンジニアリングチームのメンバーに以下の責任を指定し、割り当てます。これらの割り当ては永続的なものではありませんが、これらの責任についてチーム内に直接責任を負う個人が誕生します。割り当てられたエンジニアは、これらの責任に対するメンテナンスと反復的な改善の両方に責任を負います。

プロダクション・オペレーション

本番運用は、データエンジニアリングチームにとって常に最優先事項であり、本番環境とそれがサポートするSLOをサポートまたは影響するすべての技術的運用を含みます。これには、モニタリング、テスト、デプロイメント、コードレビュー、および一般的なDevOpsが含まれますが、これらに限定されるものではありません。

プロジェクトマネジメントとプランニング

プロジェクト管理・計画を担当するエンジニアは、ワークフロープロセス、マイルストーン計画、および課題のトリアージを監督、指導、改善します。作業が効率的に計画・完了され、データチームと **GitLab** の両方の目標を適切にサポートし実現することを保証します。

ユーザーエクスペリエンス

このエンジニアは、当社のデータプラットフォームのユーザビリティを担当します。データエンジニアリングのミッションにもあるように、私たちはすべての人のためのプラットフォームを構築しています。割り当てられたエンジニアは、最も使いにくい経験を念頭に置いて、データプラットフォームに関する解決策を積極的に見つけ出します。

セキュリティ、コンプライアンス、プライバシー

この役割は、当社のプラットフォームの安全性、プライバシー、およびコンプライアンスを維持し、監査に関する質問の窓口となるだけでなく、当社のデータプラットフォームのセキュリティとプライバシーに影響を与える変更を検討する主要なエンジニアとしての役割も果たします。また、当社のプラットフォームのセキュリティ機能を繰り返し改善していく必要があります。

データウェアハウスアーキテクチャ

イタレーションを安価に抑え、ウェアハウスの使い勝手を向上させるために、ディメンションウェアハウスの導入を決定しました。このエンジニアは、当社のウェアハウス・アーキテクチャの完全性を守り、アーキテクチャの改善が必要な箇所のビジョンを導く役割を担っています。これらの変更や保護はすべて、社内外のお客様を第一に考えたものでなければなりません。

layout: handbook-page-toc title:"Data Team Platform" 説明"GitLab データチーム・プラットフォーム"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse block html="true" /}

## クイックリンク

- データインフラデ
- ータパイプライン
- データCIの仕事
- dbt Guide
- エンタープライズデータウェ
- アハウス Jupyterガイド
- パーミフロスト
- パイソンガイド
- Sisense (Periscope)
- Sisense Style Guide
- Snowplow
- SQLスタイルガイド

データスタック

## エンタープライズデータプラットフォーム

アナリティクス機能の運用・管理にはGitLabを使用しています。すべてはイシューから始まります。変更はマージリクエストによって行われ、パイプライン、抽出、ロード、変換、分析の一部などが変更されます。

ステージツール

エクストラクトスティッチ、フィベトラン、カスタム

#### LoadingStitch, Fivetran, and Custom

	OrchestrationAirflow
データ	ウェアハウススノーフレ
ークトランスフォーメ	ーションdbtとPythonスクリ
プトデータビジュアラ	イゼーション Sisense For

## **Cloud Data Teams Advanced**

# Analyticsjupyter

抽出と読み込み

現在、データソースのほとんどにStitchとFivetranを使用しています。これらは既製のELTツールであり、データソースから Snowflakeデータウェアハウスへのデータの移動を構築、維持、または指揮する必要がありません。Stitch/Fivetranのすべて のデータソースのフルリフレッシュは、セキュリティ認証を変更するのと同時に行います(約90日ごと)。フル・リフレッシュを実行する前に、すべてのテーブルを削除します。

StitchとFivetranは、データパイプラインの開始を自分たちで処理します。つまり、AirflowはStitchやFivetranのスケジュールのオーケストレーションには関与していません。

ソースオーナーシップについては、TechStackApplicationsシート(社内限定)をご覧ください。

データソース

次の表は、データウェアハウスにロードしているすべてのRAWデータソースのインデックスです。開発のバックログと優先順位は、NewDataSource/PipelineProjectManagementシートで管理しており、最新のステータスと進捗管理のためにGitLabissuesへのリンクを設けています。新しいデータソースのハンドブックのページでは、データチームが新しいデータソースのリクエストをどのように処理するかを説明しています。

#### キー

• パイプライン。データを複製するための技術です。

• RF(Replication Frequency): 新しいデータや更新されたデータをどのくらいの頻度で読み込むか。

- Raw Schema。データが保存されているRAWデータベースのスキーマ。
- PREP Schema。ソースモデルが実体化しているPREPデータベースのスキーマ。オーディエ
- ンス。データの主要ユーザー。
- SLO: Service Level Objectiveの略。私たちのSLOは、リアルタイムからデータビジュアライゼーションで表示される分析結果までの時間を
- ・ツール
  - xは未定義または実行されていないことを示す

データソース	パイプライン	生スキーマ	プリップスキーマ	········· ,	MNPI SLO	
アダプティブ	メルタノ	tap_adaptive		ceYes		Finan
エアフロー	ステッチ	エアフロー・ス テッチ	エアフロー	データチーム	24h / 24h	No
BambooHR	エアフロー	バンプーアール	センシティブ	People	12h / 24h	No
クリアビット	х	х	х	х/		хNо
カスタマードッ ト	pgp	tap_postgres	お客様	製品	24h / x	No
ERD	スノーフレ	デマンドベース	デマンドベース	Marketing	24h /	No
<del>ブマンドベー</del>	タスク				x	
ジッター	х	ジッター	x	х	х	いい え
GitLab. です。	compgp		tap_postgresgitla	製品 b_dotcom	6h / xNo	
			エン	<b>ノジニアリング</b>		
GitLab Ops DB	pgp	6h /	tap_postgresgit	Lab_	opsEng	ineering
			INO		х	
GitLab Profiler DB					ххххх /	хNо
<b>Google</b> アナリティク	ファイブトラン goo	gle analytics 360 fivetran	google analytics 360	Marketing	6h /	No
ス 360				<u>-</u>	32h	
Google Cloud ビリング		xgcp_billinggcp_		billingEnginee	ringx /	хNо
グラファイトAI		エンジニアリング_エクストラ	クト	xエンジニアリ	ング	24h /
	No				48h	
グリーンハウス	<b>z</b>	シートロード	greenhousegreenh	nousePeople	24h / 48h	No
ハンドブッ クの <b>YAML</b>						ブック G ログデーク
ファイル						
ハンドブッ ク <b>MR</b> デ						
ータ						

エアフロー 24h エアフロ		gitlab_data_yamlgitlab_data_	yamlMultiple 	ドブッ <b>お</b> り ト/ 24h	クハ ツク <b>24</b> <b>No</b>
	エアフロー	ハンドブックハンドブックマルチプル		ノー1m	1
ライセン ドット El	ス オートマチック プロセス RD 品	ライセンス_0	dblicense_db 製	24 h / 48 h	No
マルケ	・ト ファイブトラン	marketo	oxMarketing	24h / 24h	No

データソース		PipelineRaw	SchemaPrep	SchemaAudienceRF /	です。 SLO	MNPI
Netsuite ファ	・イブトラン		netsuite_	fivetrannetsuiteF 6h / 以下はその例	です。	はい。
	PMGxpmgpm	ıqxx / xNo			24h	
	rwoxpmgpm	ISAX / XIVO			7	
PTOby	RootsSnowpipept	ogitlab_pto		生 人	, /x	daysN
クオルトリ	ックス	エアフロー		QUALITRICS		
	QUALITRICSMAF	RKETING	12H /	NO	401-	
					48h	
SaaSの利用	エアフロー		Abil II		1	weekN
	saas_usage_pin	gsaas_usage_ping	製品Ping		/ x	
セールスファ	ャース stitchsfdcSales	スティッチ 6h / 1枚		salesfor	ce_	はい。
	stitchsfucsales				24h	<del>- はい。</del>
シートロート	<b>SheetLoad</b>	sheetloadsheetload	24h / dMultiple		48h	
スノープラウ	ウスノープラウ 15m		スノープラ	・ウスノープラウ	製品	<del>* ۱۱۱ / Ž</del>
					/ 24h	
サノスレーク	スノーフ タスク	Ĭ	prometheus , xNo	prometheus	24 h エン:	<b>ジニアリン</b>
	ナートラチ		,		24 h	
バージョン DB	オートマチ ック プロセス 品		バージョン_	_dbversion_db 製	/ 48 h	No
Xactly	メルタノ	tap_		xactlyN/A		
				•		はい。
					N/ A	いいえ
Zendesk	スティッチ		zendesk_			
	stitchzendeskSu	pport	6h /		24h	いいえ
ズーム	メルタノ	tap_		zoomN/AF		はい。
					/ N/ A	
ズーラ	スティッチ		zuora_			
	· · ·				stit chz uor	
					uor aFi nan ce	
					ce 6h /	
					, 24h	48
<b>Zuora</b> の収益	エアフロー	24h /	zuora_reven	uezuora_ revenueFin	ance	ア ナ



新しいデータフィールドの追加

BambooHRの抽出物に新しいフィールドを追加する。

- AnalyticsプロジェクトでBambooHRテンプレートを使用して新しい課題を作成します。
- Data Team ManagerおよびCompensation and Benefits Managerの承認を得る。
- 承認されたら、Compensation and Benefits Managerと、抽出物を検証するData Engineerに割り当てます。

データチームによるデータソースへのアクセス

新しいデータソースをデータウェアハウスに統合するために、データチームの特定のメンバーは、UIとAPIの両方で、データソースへの管理者レベルのアクセスが必要になります。適切な分析を構築するために必要なすべてのデータを引き出すためには、APIを通じてこの管理者レベルのアクセスが必要であり、準備された分析の結果をUIと比較するためには、UIを通じてのアクセスが必要となります。

機密性の高いデータソースは、必要なレポートを作成するためのアクセス権を持つデータエンジニアとデータアナリストを 1名以上に限定することができます。場合によっては、2人のデータエンジニアに限定することもあります。自動抽出プロセ スのために、追加のアカウントを要求することもあります。

機密データは以下のセキュリティパラダイムによってロックダウンされます。Sisenseはデフォルトではいかなるデータにもアクセスできないため、機密データにアクセスすることはありません。Sisenseのアクセスは常に明示的に許可されます。

#### **DataSiren**

データチームがデータウェアハウスのどこに機密データがあるかを完全に把握し、Sisenseが機密データにアクセスできないようにするために、dbtと社内で開発されたツールのライブラリdatasirenを使ってデータウェアハウスの定期的なスキャンを行っています。このスキャンは現在、週1回のペースで実行されています。細かい結果は、SnowflakeのPREP.DATASIRENスキーマに格納され、機密性の理由からPeriscopeでは利用できません。高レベルの結果はPeriscopeで公開されており、ここにあるシンプルなダッシュボードもその一つです。

Qualtrics Mailing List Data Pump / Qualtrics SheetLoad

Qualtricsのメーリングリストデータポンププロセスは、コード上ではQualtrics SheetLoadとも呼ばれ、データウェアハウスからQualtricsにメールをアップロードする際に、チームメンバーのマシンにダウンロードする必要がありません。このプロセスの名前がSheetLoadと同じなのは、Google Sheetsを使ってqualtrics\_mailing\_listで始まる名前のファイルを探すためです。見つかったファイルのうち、1列目にid列があるものについては、そのファイルをSnowflakeにアップロードします。そしてできあがったテーブルをGitLabのユーザーテーブルと結合し、メールアドレスを取得します。その結果を新しいメーリングリストとしてQualtricsにアップロードします。

プロセス中、Google シートはプロセスの状態を反映して更新されます。最初の列の名前は、プロセスの開始時に処理中に設定され、その後、メーリングリストと連絡先がQualtricsにアップロードされたときに処理中に設定されます。カラム名を変更することで、依頼者にプロセスの状態を知らせ、デバッグを支援し、各スプレッドシートでメーリングリストが一度しか作成されないようにします。

エンドユーザーエクスペリエンスについては、UX Qualtricsの

## ページに記載されています。デバッグ

スプレッドシートにエラーが発生し、リクエストファイルに明らかな問題がない場合は、通常、スプレッドシートの再処理を試みることが最初の行動となります。過去には、新しい GitLab プラン名が gitlab\_api\_formatted\_contacts dbt モデルに追加されたときや、Airflow タスクがファイルの処理中にハングアップしたときに、再処理が必要になりました。この処理は、スプレッドシートの所有者が、この処理によって作成された部分的なメーリングリストを使用していないことを確認し、スプレッドシートに追加の変更を加えないようにするために、協調して、または要求に応じてのみ実行する必要があります。

Qualtrics Mailing Listリクエストファイルを再処理するには、以下の手順で行います。1.AirflowでQualtrics Sheetload DAGを無効にします。1.エラーとなったスプレッドシートから作成されたQualtricsのメーリングリストをすべて削除します。Qualtrics - API のユーザー認証情報を使用してQualtricsにログインし、メーリングリストを削除することができるはずです。メーリングリストの名前は、qualtrics\_mailing\_listの後のスプレッドシートファイルの名前に対応します。この名前は、スプレッドシートファイルのタブの名前と同じである必要があります。1.Airflowで Qualtrics Sheetload DAG を再度有効にし、Airflowのタスクログを注意深く監視しながら実行します。

スノープラウ・インフラ

セットアップの詳細については、「スノープラウ・インフラストラクチャー」のページを参照してください。

データソースの概要

- カスタマーサクセスダッシュボード
- Netsuite
  - Netsuiteとキャンペーンデー
- タのバージョン(ping)
  - なお、2019年10月までは、データチームはバージョンデータソース全体を「ping」と呼んでいました。しかし、Usage pingはバージョンデータソースの1つのサブセットに過ぎないため、現在では、version.gitlab.comデータソースを参照するために「version」または「version app」を使用し、バージョンデータソースの特定の使用データ機能を参照するために「Usagedata」または「Usagepings」または「Pings」を使用しています。データの抽出において、Servicepingデータの取り込みに関しては、Service pingページまたはService pingのReadme.mdページに詳細が記載されています。
- Salesforce
- Zendesk

## オーケストレーション

オーケストレーションにはAirflow on Kubernetesを使用しています。私たちの具体的なセットアップ/実装方法はこちらをご覧ください。また、データ基盤のページにも詳細が記載されています。

## データウェアハウス

現在、データウェアハウスとしてSnowflakeを使用しています。エンタープライズデータウェアハウス(EDW)は、GitLabの企業データ、パフォーマンス分析、主要業績評価指標などの企業全体のデータのための唯一の情報源です。EDWは、レポーティング、ダッシュボード、分析のための共通のプラットフォームとフレームワークをすべてのチームに提供することで、GitLabのデータドリブンな取り組みをサポートします。ポイントツーポイントのアプリケーション統合を除いて、現在および将来のデータプロジェクトはすべてEDWから推進されます。EDWは、GitLabの様々なソースシステムからのデータを受信することで、データ品質のベストプラクティス、対策、改善策を伝え、推進し、可能な限り最高のデータを使ってすべての意思決定がなされるようにします。

Snowflakeサポートポータルへのアクセス

Snowflakeサポートポータルにアクセスするには、以下の手順を実行してください。

- qitlabのメールアドレスをコミュニティポータルに登録
- この登録により、gitlab のメールに Welcome to the Snowflake Community という件名の歓迎メールが送信されます。このメールでは、登録を完了するように求められ、その一環としてコミュニティポータルのパスワードを設定するように求められます。
- 終わったら、snowflakeのコミュニティアカウントに再度ログインし、ホームページで「ケースの送信」をクリックします。初めて、アクセス権を持たないユーザーがsnowflakeでケースを送信することができます。アクセスのためのフォームを入力するように求められます。
- フォームでは、すでにスノーフレークのお客様のアクセスを選択します。次のページでは、Account Name(アカウント名)、Cloud Name(クラウド名)、Region Name(地域名)の情報が求められます。以下は、snowflakeコンソールからこの情報を引き出す1つの方法です。
  - アカウント名 CURRENT\_ACCOUNT()を選択します。
  - 地域名 CURRENT\_REGION()を選択します。
  - クラウド名 リージョン名の値に基づいて、クラウド名を特定できます。
- 申し込みが完了すると、「Request received] Case# instant」という件名の確認メールが届きます。万が一、メールが届かない場合は、フォームを再送信してください。
- 24時間以内に、件名が「Case# -Self Register Enable Case access」のリクエスト確認メールが届くことを掲示する。

倉庫へのアクセス

データウェアハウスにアクセスするために

• アクセス要求プロジェクトで課題を作成し、必要なアクセスレベルを文書化します。共

• 有アカウントをリクエストしてはいけません。各アカウントはユーザーに関連付けられている必要があります。

- 私たちは、このブログ記事で説明されているユーザーの許可に関するパラダイムにゆるやかに従っています。
- 既存のアカウントのミラーリングを依頼する際には、制限されたSAFEデータへのアクセスはプロビジョニング/ミラーリングされないことに注意してください(現在はrestricted\_safeロールを通じて提供されています)。

## SAFEデータへのアクセス

すべてのSAFEデータは、Snowflakeの別々のデータベーススキーマ内のテーブルに保存されます。1つのテーブルにアクセスすると、すべてのSAFEテーブルにアクセスできます。SAFEデータへのアクセスには

- 1. あなたの直属の上司の承認
- 2. 部門のVP (またはそれに相当する) の承認を得ていること。
- 3. GitLab Dashboard Indexに定義されているSAFE Dashboard Space Ownerの承認。

#### SAFEデータにアクセスするには

- 1. アクセスリクエストを作成し、お客様のニーズと意図をお伝えください。
- 2. 直属の上司、部門担当副社長(またはそれに相当する人物)、およびGitLabダッシュボードのインデックスヘッダーに定義されているSAFEスペースオーナーに承認を求める。SisenseでSAFEデータへのアクセスが60日以内に承認されていれば、承認は必要ない。このステップをスキップして、特定のARにリンクする。
- 3. リクエストが承認されると、Snowflakeのプロビジョナーにタグを付け、プロビジョナーがリクエストを処理します。
- 4. 処理が完了すると、SnowflakeでSAFEデータ (スキーマ) にアクセスできるようになります。

## スノーフレーク パーミッション パラダイム

Snowflakeのパーミッション管理にはPermifrostを使用しています。私たちのSnowflakeインスタンスの設定ファイルは、この roles.ymlファイルに格納されています。また、Permifrostに関するハンドブックのページもあります。

私たちは、このような一般的な戦略に基づいて役割管理を行っています。

- すべてのユーザーは、関連するユーザーロールを持っています。
- ◆ 機能的な役割は、共通の権限セットを表すために存在します(analyst\_finance、data\_manager、product\_manager
- ) データの論理的なグループには、それぞれのオブジェクトの役割があります。
- オブジェクトの役割は主に機能的な役割に割り当てられる
- より高い権限を持つロール (accountadmin、securityadmin、useradmin、sysadmin) がユーザーに直接割り当てら
- れるサービスアカウントには同じ名前のロールが割り当てられる
- 追加の役割は、用途やニーズに応じて、サービスアカウントの役割またはサービスアカウント自体に割り当てることができます。
- テーブル&ビューの粒度で個別の権限を付与可能
- Warehouseの使用は、必要に応じてどのような役割にも付与することができますが、機能的な役割に付与

## することをお勧めします。 ユーザーの役割

すべてのユーザーは、ユーザー名と一致する独自のユーザーロールを持ちます。オブジェクトレベルのパーミッション(データベース、スキーマ。

テーブル)は、Snowflakeではロールにのみ付与することができます。ロールはユーザーまたは他のロールに付与することができます。私たちは、ユーザーがデータベースと対話するために1つのロールを使用する必要があるように、すべての権限がユーザーロールを介して流れるように努めています。例外として、accountadmin、securityadmin、useradmin、sysadmin などの特権ロールがあります。これらのロールは、より高いアクセス権を与えるため、使用する際には意図的に選択する必要があります。

#### 機能的役割

機能的な役割とは、一般的にジョブファミリーにマッピングされる権限と役割付与のグループを表します。主な例外は、アナリストロールです。アナリストロールには、組織のさまざまな領域に対応するいくつかのバリエーションがあります。 analyst\_core、analyst\_finance、analyst\_peopleなどがあります。アナリストは関連するロールに割り当てられ、必要なスキーマへのアクセスが明示的に許可されます。

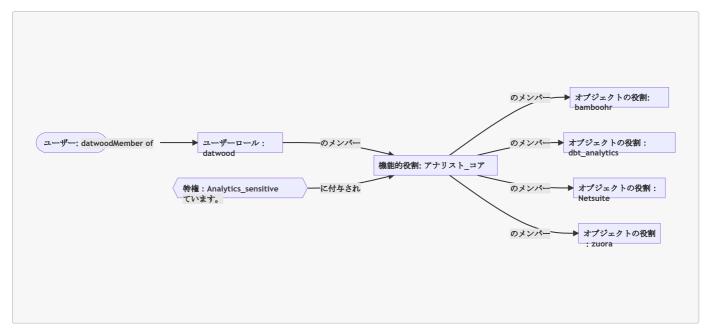
機能的な役割はいつでも作成することができます。これは、非常に似通ったジョブファミリーと権限を持つ複数の人がいる 場合に最も意味のあるものです。

オブジェクトの役割

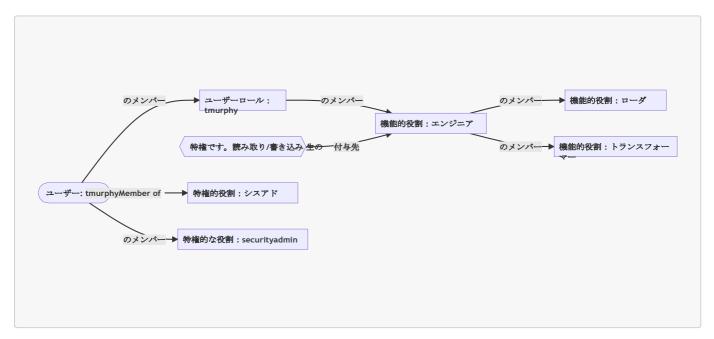
オブジェクトロールは、一連のデータへのアクセスを管理するためのものです。一般的には、特定のソースのすべてのデータを表します。zuoraオブジェクトロールはその一例です。このロールは、Stitchから送られてくる生のZuoraデータと、prep.zuoraスキーマのソースモデルへのアクセスを許可します。ユーザーがZuoraデータへのアクセスを必要とする場合、そのユーザーのユーザー・ロールにzuoraロールを付与するのが最も簡単な解決策です。何らかの理由でオブジェクト・ロールへのアクセスが意味をなさない場合は、テーブルの粒度で個別の権限を付与することができます。

例

## これは、Data Analyst, Coreの役割階層の例です。



これは、データエンジニアとアカウントアドミニストレーターの役割階層の例です。



これは、セキュリティオペレーションエンジニアの役割階層の例です。



\*\*Snowflakeのロールを管理する\*\* {: .panel-heading}。

ここでは、新しいユーザーとユーザーロールをプロビジョニングするための適切な手順を説明します。

- プロビジョニングラベルが適用されたオリジナルのリクエストをリンクするGitLab Data Teamプロジェクトの課題を確認します。
- Snowflakeにログインし、securityadminロールに切り替えます。
  - o すべてのロールがsecurityadminの所有下にあること
- user\_provision.sqlスクリプトをコピーし、初期ブロックのemail、firstname、lastnameの値を置き換えます。 パ
- スワードが必要な場合は、PasswordGeneratorを使って作成します。
  - ユーザー名とパスワードの認証情報をワンタイムシークレットまたはSlack経由でユーザーに送信
- Snowflake roles.yml permifrost設定ファイル内のドキュメント(このファイルは毎日12:00a.m. (UTC) に自動的に読
  - み込まれます) 作成したユーザーとユーザーロールを追加する
  - 新規ユーザーへのユーザーロール
  - の割り当てユーザーへの追加ロールの割り当て
- ユーザーがOktaでアプリケーションを割り当てられていることを確認する
- ◆ ユーザーがokta-snowflake-usersのGoogleグループに割り当てられているこ

とを確認する 既存のユーザーをデプロビジョニングするための適切な手順を説明

します。

- Snowflakeのデプロビジョンは、オフボーディングの問題またはアクセスリクエストの問題を介して行う必要があります。
- GitLab Data Team プロジェクトで、オリジナルのソースリクエストにリンクしているイシューがあることを確認します。 デプロビジョニングラベルが適用されています。
- Snowflakeにログインし、securityadminロールに切り替えます。
  - すべてのロールはsecurityadminの所有下にある必要があります。
- user\_deprovision.sqlスクリプトをコピーして、USER\_NAMEを置き換えます。ユーザーを削除せずにsnowflakeに残し、 disabled = TRUEを設定する理由は、ユーザーがいつアクセスできなくなったかを記録するためです。
- okta-snowflake-usersのGoogleグループからユーザーを削除します。
- Snowflake roles.yml permifrost設定ファイルのユーザーレコードを削除する(このファイルは毎日12:00a.m. UTCに 自動的に読み込まれます

詳しくは、このペアリングセッションの録画をご覧ください(GitLab Unfilteredとして視聴する必要があります)。

コンピューティングリソース

Snowflakeのコンピュートリソースは「ウェアハウス」と呼ばれています。クレジット消費量をよりよく追跡・監視するために、倉庫にアクセスする人に応じていくつかの倉庫を作成しました。倉庫の名前には、そのサイズが付加されています(analyst\_xsはextra small)。

(分) レかります )

ウェアハウス・パーポゼッション・マックスクエリ

		<b>(ガ) となりまり。</b> )
admin	パーミッション・ボットやその他の管理作業のためのものです。	10
airflow_testing_l	局所的な気流のテスト用	30
これらは	、データアナリストがデータベースへの問い合わせやモデリングを行う際に	 :使用します。
アナリスト _*さん	データ	30
	これらは、データエンジニアやマネージャーが、データを照会する際に使	用します。
エンシー/_*さん	データベースやモデリングデータ	30
fivetran_warehouse	これはFivetranが使用するための専用のものです。	30
gitlab_postgres	これは、 <b>GitLab</b> 内部の <b>Postgres</b> データベースから抽出するジョブのための です。	もの 10
グラファナ	これは、Grafanaが使用するための専用	60
ローディング	これは、当社のExtract and Loadigal ヹ゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚	60
menge request *	これらけ、マージリクェスト内のdhtジョブのGitlah Clにスコープされて)	ハキ 60

ウェアハウス・パーポゼッション・マックスクエリ

		(分) となります。)
<b>報</b> 告	これはBIツールでのクエリ用です。なお、Sisenseは4分を強制的に のタイムアウトが発生します。	30
ステッチ	これはスティッチが使用するための専用のものです	30
ターゲット_スノーフ レーク	これは、MeltanoチームがSnowflakeローダーをテストするためのものです。	5
transforming_*。	これらはプロダクションDBTジョブ用です。	60

クエリの時間が制限されている場合は、より大きな倉庫の使用を検討してください。

データストレージ

私たちは、raw、prep、prodという3つの主要なデータベースを使用しています。rawデータベースは、Snowflakeに最初にデータが読み込まれる場所で、その他のデータベースは、分析の準備が整った(または整いつつある)データ用です。

スノーフレークデータベースというものがあり、そこには**GitLab**インスタンス全体の情報が含まれています。これには、すべてのテーブル、ビュー、クエリ、ユーザーなどが含まれます。

covid19データベースがあり、これはSnowflake Data Exchangeで管理されている共有データベースです。

testing\_dbデータベースがあり、Permifrostのテストに使用されています。

roles.yml Permifrostファイルで定義されていないすべてのデータベースは、週単位で削除されます。

Sisense	rawNo.の
DatabaseVi	ewable
	PREPNO
	prodYes

生

このデータベースは、Sisenseではクエリできません。このデータに対するdbtモデルは存在しません。

- ◆ 生データには機密データが含まれている可能性があるため、権限を慎重に管理す
- る必要があるデータはソースに応じて異なるスキーマに保存される
- スキーマやテーブルごとにユーザーのアクセスを制御できる

プリップ

このデータベースは、Sisenseでは照会できません。

- ソースモデルは、データソースに対応する論理スキーマで構築されます(例:sfdc、zuora)
- PREPARATION これは、dbtモデルが構築されるデフォルトのスキーマです。
- SENSITIVE

## Prod

このデータベースと、その中のすべてのスキーマとテーブルは、Sisenseが問い合わせ可能です。

public と boneyard を除いて、すべてのスキーマは dbt で制御されています。詳細はdbtガイドを参照してください。分

析プロジェクトのフォルダ構造

以下の表は、アナリティクス・プロジェクトのmodels/ディレクトリ内のフォルダに格納されたモデルが、データウェアハウスでどのように実体化されるかを示すマッピングです。

これは、dbt\_project.ymlという設定ファイルに記載されています。

スノーフレークのフォ ルダー dbt/models/	Sis	sense ${\mathcal O}$	db. schem aDetailsQ ueryable
d.common	common/pro	ファクトのトップレベルフォルダと の寸法を表示します。ここには モデルを入れないでください。	はい。
	common/prep/prep.prep.	prepモデルを使用して、 作成する。 No facts/d	を lims
	common/sensitive/prep.	sensitiveFacts/dimsには、 機密デー	
		せん。	ま 
レート/prod.curateYes			コモン/キュ
<b>F/</b>	ディムテーブル。はい。	common/prod/prod.commonProd	uctionのファク
	common_mapping/prod.common_mapping	マッピング、ブリッジ、または を含み まず。	はい。
common_mapping/prep/	prod.common_mapping	の <b>準備</b> テーブル。 マッピング、ブリッジ、ルック アップテーブル	はい。
みます。	marts/prod. marts	マートレベルの	 データを含 
νρά ) .	prep/prep.準備	生産のための一般的な準備モデル	いいえ。
	legacy/prod.legacy	非次元的に構築されたモデルを 含みます。	はい。
	source/prep.source	ソースモデルを含んでいます。 スキーマはデータソースに <b>基</b> づく	いいえ
ワークスペース/	prod.workspace_workspace	SQLの対象とならないワーク スペースモデルや dbt規格。	はい。
共通/制限付き	prod.restricted_domain_common	制限されたファクトとディメンションのトップレベルフォルダ。 レギュラーの同等品 共通のスキーマではなく 、限定されたデータを対 象としています。	はい。
common_mapping/resticted	prod.restricted_domain_common_mapping	制限付きのマッピング、ブリッジ、 またはルックアップテーブルを含み ます。 レギュラーのエクイップメント 共通のマッピングスキーマを使用し ていますが、制限されたデータを 対象としています。	はい。
マート/制限付き	prod.restricted_domain commonmarts	はい。	
レガシー/リストリクテッド	prod.restricted_domain_legacy_43	無次元的に構築された制限付きモ デルを含む。通常のレガシーに相当 する スキーマではなく、制限	はい。

スキーマではなく、制限

**dbt**で自動的に更新することなく、ユーザーのためにデータを保存する必要があるデータウェアハウスのユースケースでは、 STATICデータベースを使用します。これにより、アナリストや他のユーザーが独自のデータ・リソース(テーブル、ビュー、その 他)を作成することもできます。

一時的なテーブル)を使用しています。)静的データベースの中には、機密データ用のスキーマがあります。静的データベースのユースケースで、機密データの使用または保存が必要な場合は、データエンジニアのために課題を作成してください。

タイムゾーン

ウェアハウス内のすべてのタイムスタンプデータは、UTCで保存する必要があります。SnowflakeセッションのデフォルトのタイムゾーンはPTですが、これをオーバーライドしてUTCがデフォルトとなっています。これは、current\_timestamp()がクエリされたときに、結果がUTCで返されることを意味します。

Stitchはタイムスタンプ明示的にUTCに変換します。Fivetranも同様に行っています(サポートチャットで確認済み)。

このルールの唯一の例外は、ファクトテーブルの date\_id の作成に太平洋時間を使用することです。これは常に get\_date\_pt\_id dbt マクロで作成し、\_pt\_id サフィックスを付ける必要があります。

スナップショット

{:#snapshots-definition}になります。

データチームハンドブックの複数の箇所でスナップショットという言葉を使用していますが、この言葉は文脈によっては混乱を招く恐れがあります。辞書で定義されているスナップショットとは、「ある時点での記憶場所やデータファイルの内容を記録したもの」のことです。私たちは、この言葉を使うときは常にこの定義を使うように努めています。

#### dbt

最も一般的な使い方は、dbtsnapshotsに関するものです。dbt snapshotsが実行されると、ユーザーが指定したクエリに基づいてデータの状態を取得し、データの状態の全履歴を含むテーブルを更新します。このスナップショットには、特定のスナップショットが有効な期間を示す valid\_to と valid\_from フィールドがあります。より詳しい技術情報については、dbtガイドの「dbt snapshots」セクションをご覧ください。

dbt snapshotsで生成・管理されるテーブルは、生のヒストリカル・スナップショット・テーブルです。これらの生のヒストリカル・スナップショットの上にダウンストリーム・モデルを構築し、さらにクエリを行うことになります。snapshotsフォルダには、dbtのモデルが格納されています。これは、24時間の間に複数のスナップショットが取得されている場合に便利です。また、生の履歴テーブルから最新のスナップショットを返すモデルも作成します。

## その他の用途

グリーンハウスのデータはスナップショットのようなものだと考えられています。グリーンハウスから毎日データベースダンプが送られてきて、それをSnowflakeに読み込みます。もしこれらのテーブルのdbtスナップショットを取り始めたら、グリーンハウスデータのヒストリカル・スナップショットを作成することになります。

一部の yamlファイルや BambooHR で行っている抽出は、スナップショットと考えることもできます。この抽出では、ファイルやテーブルの全体を取得し、それをタイムスタンプ付きの行としてウェアハウスに保存します。つまり、これらのファイルやテーブルの過去のスナップショットがあるということですが、これらはdbtと同じ種類のスナップショットではありません。同じvalid\_toとvalid\_fromの動作を得るためには、追加の変換を行わなければなりません。

## 言語

- スナップショット 特定の時点でのデータの状態
- Take a snapshot 現在のデータの状態を取得して保存するジョブを実行します。dbtのコンテキストで使用できます。yamlの抽出ジョブを参照することはお勧めしません これらは "run the extract"となります。
- ヒストリカル・スナップショット 特定のソース・テーブルの複数の時点でのデータを含むテーブル。dbtが生成した スナップショットテーブルを参照するために最もよく使用されます。また、yaml抽出テーブルの参照にも使用できます。
- 最新のスナップショット 保存されているデータの最新の状態です。dbtのスナップショットでは、valid\_toがnullのレコードが該当します。BambooHR と yaml の抽出では、抽出ジョブが最後に実行されたときの状態です。Greenhouse rawの場合、これはウェアハウスにあるデータを表しています。Greenhouse データのスナップショットを開始する場合、発言者は raw テーブルを意味しているのか、Historical snapshots テーブルの最新レコードを意味しているのかを明確にする必要があります。

また、堅牢性を高めるために、倉庫のデータをGCS(Google Cloud Storage)にバックアップしています。ジョブの実行には dbtのrun-operation機能を使用しています。現在、すべてのスナップショットを毎日バックアップし、 $1_{7}$ 月間保存しています。 Calogica社のブログで紹介されている基本的な手順を実行しました。

## アドミン

Snowflakeを稼働させるために、管理作業を行います。保存場所の作成

データウェアハウスにデータをロードするためには、通常、データはストレージのバケットから読み出されます。バケットから読み込むためには、そのバケットをSnowflakeの許可リストの一部として追加し、ステージを作成する必要があります。

まず、現在のストレージロケーションをすべて選択します。property\_valueでproperty=STORAGE\_ALLOWED\_LOCATIONSとなっている値をコピーします。

```
desc integration gcs_integration;
```

<<<c\_paste\_here\_>>>> + 新しいバケットの位置の値を上にして、以下のクエリに値を貼り付けます。値は, で区切る必要があります。

```
統合を変更する gcs_integration
SET STORAGE_ALLOWED_LOCATIONS = ('<<<_paste_here_>>>>')
```

そして、新たなステージでは、収納場所を追加して作成することができます。

```
ステージ "raw "を作成します。"PTO".pto_load
STORAGE_INTEGRATION = GCS_INTEGRATION URL = 'bucket location';
```

## トランスフォーメーション

当社では、すべての変換にdbtを使用しています。このツールを使う理由と方法については、dbtガイドをご覧ください。

トラステッドデータフレームワーク

{:#tdf}になります。

データ顧客は、重要な意思決定を行うために、信頼できるデータをデータチームが提供することを期待しています。そしてデータチームは、提供するデータの品質に自信を持つ必要があります。しかし、これは解決が難しい問題です。エンタープライズデータプラットフォームは複雑で、複数の段階でデータの処理や変換が行われ、数十人から数百人の開発者やエンドユーザーが24時間体制でデータの変更や問い合わせを行っています。Trusted Data Framework(TDF)は、技術チーム*やビジネスチーム*がアクセス可能な、データ処理段階におけるデータのテストとモニタリングのための標準フレームワークを定義することで、これらの品質と信頼性のニーズをサポートします。TDFは、既存のデータ処理技術から独立したスタンドアローンのモジュールとして実装されており、独立したデータ監視ソリューションのニーズを満たしています。

- アナリストやエンジニアだけでなく、誰もが信頼できるデータに貢献できるようにな
- るデータ処理の上から下まで、すべての段階でデータを検証できるようになるソー
- スシステムのデータパイプラインからデータを検証できるようになる
- 次元モデルへのデータ変換を検証する 重要な企業デー
- タを検証する
- 中央のデータ処理技術から独立して展開可能

主な用語

• アサーションまたはテストケース - 個々のテストであり、実行可能なテストの最小単位である。TDFでは、テストケースはSQLステートメントとして、またはSQLコンパイルツールdbt内のYAMLコンフィギュレーションを介して表現されます。

- データスキーマ データ対象領域を構成するテーブル、カラム、ビュー、その他の構造要素で、SQLデータ定義言語(DDL)を用いて作成される。
- ゴールデンデータ ゴールデンデータとは、ビジネスにとって重要な1つのフィールドまたは複数のフィールド
- のグループから得られるデータの定数です。モニタリング テストケースの結果を追跡することで、データが使用可能な状態であることを確認します。

信頼できるデータコンポーネント

TDFの主な要素は以下の通りです。

- 1. 新しいデータソリューションから問題解決に至るまで、日々のデータ開発の中で品質を当たり前のように組み込む 好循環テストサイクル。
- 2. テストケースをSQLやYAMLで表現し、誰でも開発できるようにしました。
- 3. Trusted Data Schemaは、ビジネスプロセスやデータプラットフォームのパフォーマンスに関する知恵を開発するために、テスト結果をモニタリングやアラート、長期的な分析のために保存します。
- **4. Schema-to-Golden Record Coverage**: スキーマから重要な "**Golden** "データまで、データウェアハウス領域を幅広くカバーします。
- 5. Trusted Data Dashboard」は、ビジネスに役立つダッシュボードで、全体のテストカバレッジ、成功、失敗を可視化します。
- 6. テスト実行とは、テストケースを実行することです。
- 7. ソースシステムとSnowflakeのVirtuous Test Cycle間の行数を再計算するRow Countテ

#### スト

TDFは、ビジネスユーザーを信頼できるデータを確立するための*最も重要な参加者*として受け入れ、シンプルでアクセス可能なテストモデルを使用しています。テストエージェントとしてSQLとYAMLを使用することで、幅広いグループの人々がテストケースを提供することができます。テストフォーマットはシンプルなPASS/FAILの結果と4つのテストケースタイプで構成されています。TDFが価値を示すことで、採用は急速に拡大します。

- データ・カスタマーとビジネス・ユーザーがテスト・フレームワークを学び、自らテストを作成
- する チームは、テストを最後の手段としてではなく、常に組み込むべき貴重な活動として受け入れる
- データチームは、プロダクションダウン・レトロスペクティブの一環として新しいテストを追加し、大きな問題になる前に問題をより迅速に特定できるようになりました。
- チームは、継続的に新しいテストを開発し、テストカバレッジを拡大するための運用リズムを開発します。

時間が経てば、何百ものテストケースを開発して毎日のように実行し、データの品質を継続的に検証することも珍しくありません。

SQLとYAMLで表現されるテストケース

**SQL**はデータベースの世界共通言語であり、データを扱うほとんどの人が、ある程度の**SQL**の能力を持っています。しかし、すべての人が**SQL**に精通しているわけではなく、それによって貢献できる人が制限されてしまうことは避けたいと考えています。私たちはdbtを使ってTDFをサポートしており、**SQL**と**YAML**を使ってテストを定義することができます。

信頼性の高いデータスキーマ

すべてのテストは**dbt**で実行されるので、テスト結果の保存は簡単です。すべてのテストの実行結果をデータウェアハウスに保存します。テスト結果を保存することで、以下のような様々な価値ある機能が可能になります。

- データの可視化とパターン分析 テスト結果(日付ごとの総テスト数、分野ごとの PASS/FAIL 率など) データの対象や
- スキーマに対するテストカバレッジの測定(分野ごとのテスト数)。
- 時間経過によるシステム品質の向上(PASS率の向上)の測定 テスト結果に基づくアラ
- ートシステムの開発

これらのテスト結果は解析され、Sisenseでの問い合わせに利用できます。

すべてのテスト結果を保存するスキーマは

workspace\_data」です。注:このスキーマには、ビュ

ーしか含まれていません。

スキーマからゴールデンレコードカバレッジ

データウェアハウスの環境は急速に変化する可能性がありますが、TDFはデータウェアハウス内で最も変化する可能性の高い領域をテストでカバーすることで、予測可能性、安定性、品質をサポートします。

- 1. スキーマの整合性を検証するスキーマテスト
- 2. 列のデータ値が事前に定義されたしきい値やリテラルに一致するかどうかを判断する列値テスト
- 3. Rowcount は、事前に定義された期間におけるテーブルの行数が、事前に定義されたしきい値またはリテラルに一致するかどうかを判定するテストです。
- 4. Golden Dataテストでは、あらかじめ定義された価値の高いデータがテーブル

に存在するかどうかを判断します。 これらのテストの実装の詳細は、dbtガイドに記

載されています。

信頼できるデータダッシュボード

SisenseのTrusted Data Dashboardはこちらをご覧くださ

い。 Test Run

まだまだ続きま す。

行数テスト

行数テストでは、ソース データベースとターゲット データベース間の行数を調整するために、以下のデータを抽出します。 ソースのDBテーブルをSnowflakeテーブルにロードし、Snowflakeから類似の統計情報を抽出してソースとターゲットの比較 を行います。ソースとターゲットを完全に一致させるのは難しいです。

- タイミングの違いがあります。
- ◆ データウェアハウスでは履歴が残ることが
- あります。削除はソースデータベース上で 行われます。

シナリオによっては、行数を最高(テーブル)レベルではなく、より低い粒度のレベルでチェックすることをお勧めします。これは、論理的な分布を持つ1つまたは複数のフィールドである可能性がありますが、依然として集約されたレベルであると言えます。例えば、テーブルの挿入日や更新日などが挙げられます。

ソースの行数とターゲット(**Snowflake**データウェアハウス)の行数に基づいて、すべての行がデータウェアハウスにロードされているかどうかを判断するために、リコンシリエーションを行うことができます。

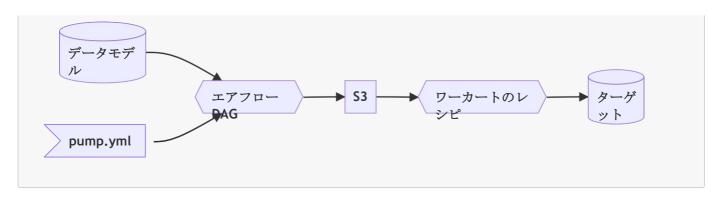
## 行数テスト PGP

このフレームワークは、テストを実行するためのあらゆる種類のクエリの実行を扱うように設計されています。現在のアーキテクチャでは、すべてのクエリが1つのKubernetesポッドを作成するため、1つのクエリにグループ化することで、Kubernetesポッドの作成数を減らすことができます。postgres DBとsnowflakeの間のレコード数とデータの実際のテストでは、ボリュームの少ないソーステーブルをグループ化し、ボリュームの大きいソーステーブルは個別のタスクとして実行するアプローチをとった。

すべてのタイプの調整を行うことを想定した新しい yaml ファイルが作成されます(そのため、既存の yaml 抽出マニフェストには組み込まれていません)。マニフェストファイルは、ボリュームの少ないテーブルのグループをまとめ、ボリュームの大きいテーブルを個々のタスクとして結合します。Postgres と snowflake からの行数テストの比較は、"PROD". "WORKSPACE DATA". "PGP SNOWFLAKE COUNTS" という名前の snowflake テーブルに保存されます。

データポンプ

merged.md	2021/11/8



誰でも簡単にSnowflakeからGitLabテックスタックの他のアプリケーションにデータを送れるようにするために、私たちは Enterprise Applications Integration Engineeringチームと提携して、Data Pumpと呼ぶデータ統合フレームワークを作りました

これらはすべて、ポンプを実行するData Pump AirflowDAGでオーケストレーションされ、毎日05:00 UTCに1回実行されるように 設定されています。

データポンプの追加

ステップ1dbtを使用して、/marts/pumps (モデルにREDまたはORANGEデータが含まれる場合は/marts/pumps\_sensitive) にデータモデルを作成しdbtモデル変更テンプレートを使用してMRを作成しこれがマージされ、SnowflakeでPROD.PUMPSまたはPROD.PUMPS\_SENSITIVEに表示されると、ステップ2と3の準備が整います。

Step 2: 「Pump Changes」MRテンプレートを使って、pumps.ymlに以下の属性を持つモデルを追加します。

- model dbtおよびsnowflakeにおけるモデルの名前です。
- timestamp\_column データのバッチ処理に使用される列の名前(列がなく、テーブルが小さい場合はnull)。
- sensitive このモデルが機密データを含み、pumps\_sensitiveディレクトリとスキーマオーナーにある場合は
- true、あなた(またはビジネスDRI)のgitlabハンドルです。

ステップ3: インテグレーションチームがデータをターゲットアプリケーションにマッピングして統合できるように、「New Data Pump」課題テンプレートを使用して、インテグレーションプロジェクトに統合課題を作成します。

データスピゴット

データスピゴットは、外部システムに制御された方法でSnowflakeのデータへのアクセスを与えるための概念/方法論です。外部システムがSnowflakeにアクセスできるようにするには、以下のようなコントロールが必要です。

- サービス専用のアカウントです。
- ◆ 必要最小限のデータのみを表示する専用のビュー (複数可)。個人を特定できる情報 (PII) は公開されません。
- 指定されたテーブル/ビューのみにアクセス可能な専用ロール(または同等のも
- の)。コストを制限・監視するための専用のXS倉庫。

新しいData Spigotを設定する手順は以下の通りです。

- 1. 上記のように、実施されているコントロールに従うこと。
- 2. 下の表に新しいData Spigotsを追加します。

接続システム データの範囲 データベースビュー

Grafana 除雪機のロード時間 prod.legacy.snowplow\_page\_views\_all\_grafana\_spigot

ビジュアライゼーション

データの可視化とビジネス・インテリジェンス・ツールとしてSisenseを使用しています。アクセスをご希望の方は、アクセスリクエストを送信してください。

データチームのためのメタアナリシス

- ◆ Sisenseの使い方!✓
- Sisense アカウント最適化 4
- Sisense アカウントメンテナンス
- 🖩 dbtイベントロギング
- スノーフレークスペンド\*の

セキュリティ

パスワード

GitLabのパスワードポリシーに基づき、パスワードのみで認証されるサービスアカウントを90日ごとにローテーションしています。パスワードが変更されたシステムと更新された場所の記録は、このGoogleシートに保存されています。

また、Snowflakeユーザーのパスワードは、毎年第3月の第1日曜日(1月、4月、7月、10月)に「SnowflakeパスワードリセットDAG」でローテーションしています。

ソフトウェアのユーザープロビジョニング

データチームは、データチームが管理するツール内のユーザーをプロビジョニングする責任があります。これには、Fivetran、Stitch、Snowflakeなどのツールが含まれます。

Snowflakeについては、このページの「Snowflake Permissions Paradigm」のセクションで説明されている強固なプ

ロセスがあります。他のツールについては、UIを介してユーザーを追加し、適切なGoogleグループが存在する場合はそれに参加します。

layout: handbook-page-toc title:"Data Team Cl Jobs" description:"GitLabデータチームのClジョブズ"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg},

{::options parse\_block\_html="true" /}。

このページでは、DataTestsとAnalyticsの両プロジェクトのMerge Requestsでデータチームが使用するCIジョブについて説明します。

パイプラインが故障した場合の対応

- 週末が経過した場合は、以前に実行されたCLONEステップを再実行し、毎週土曜日に古いパイプラインデータベースが SnowFlakeから削除されます。
- masterブランチをマージします。dbtがパッケージをどのように処理するかに起因して、パッケージの不具合によりパイプラインが失敗する可能性がありますので、常に最新のブランチで処理する必要があります。
- モデル選択の構文を確認します。一般的には、変更するモデルのファイル名をそのまま使うのが最も簡単です。それでも
- 不明な場合は、@gitlab-data/data-engineers(緊急にサポートが必要な場合はエンジニア)にタグを付けて、問題解決のサポートを依頼してください。

パイプライン変数がジョブに渡されない

- 現在、GitLab CI パイプラインに問題があります。変数を渡す必要があるパイプラインでは、一度しか変数を渡すことができません。
- つまり、dbtモデルのパイプラインを使って2回目のプロセスを実行するテストをしたい場合は、新しいパイプラインを作成する必要があります。
- 最も簡単な方法は、以下のような青い run pipeline ボタンをクリックすることです。 🗾 run\_pipeline\_button.png

# アナリティクス パイプライン

ステージ

CIの仕事はステージごとに分類されています。

\*スノーフレーク

これらのジョブは .gitlab-ci.yml で定義され

ています。 clone\_prep\_specific\_schema

準備中のデータベースで利用可能な任意のスキーマのクローンが必要な場合、これを実行します。クローンを作成するスキーマを SCHEMA\_NAME変数を使用します。すでにクローンが存在する場合は、何もしません。

クローン\_開発者固有のスキーマ

prodデータベースで利用可能な任意のスキーマのクローンが必要な場合に実行します。クローンを作成するスキーマを SCHEMA\_NAME変数を使用します。クローンがすでに存在する場合は、何もし

ません。

MRの起動時に自動的に実行され、あらゆるdbtジョブを実行できるようになります。クローンが存在するかどうかを確認するだけなので、このジョブの後続の実行は高速です。これはprodとprepのデータベースの空のクローンです。

クローン・プロード・リアル

prodデータベースとprepデータベースの実際のクローンを作成する必要がある場合に実行します。両方のデータ

ベースの完全なクローンを作成します。 clone\_raw\_full

エクストラクト、フレッシュネス、スナップショットの各ジョブを実行する必要がある場合に実行します。このジョブの後続の実行 は、以下の場合にのみ検証されるため、高速になります。 クローンが存在しています。

クローン\_ロースト\_postgres\_pipeline

postgresパイプラインやマニフェストファイルの変更をテストするために、生のtap\_postgresスキーマのクローンが必要な場合に実行してください。生のクローンが既に存在する場合、これは何もしません。

clone\_raw\_sheetload

sheetloadへの変更や追加をテストするために、生のsheetloadスキーマのクローンが必要な場合に実行します。生のクローンがすでに存在する場合は、何もしません。

clone\_raw\_specific\_schema

変更や追加をテストするために、他の raw スキーマのクローンが必要な場合に実行します。クローンを作成する raw スキーマを SCHEMA\_NAME 変数で指定します。raw クローンがすでに存在する場合は、何もしません。

フォースクローン・ボース

raw、prod、prepを強制的にリフレッシュしたい場合に実行します。これは、rawの完全なクローンを作成しますが、prefとprodの浅いクローンを作成します。

血エキス

これらのジョブはextract-ci.ymlで定義されています。

ボーンヤード・シートロード

新しいBoneyard sheetloadのロードをテストしたい場合に実行します。これには本物のprodクローンとprepクローンが利用でき

ることが必要です。

2021/11/8

新しいシートロードをテストしたい場合に実行します。このジョブは、RAWのクローンに対して実行されます。clone\_rawジョブが実行されている必要があります。

#### pgp\_test

gitlab-ops DB を除いて、postgres pipeline のマニフェストを追加または更新する場合に実行します。RAWのクローンと MANIFEST\_NAME変数、SCDテーブルの場合はTASK\_INSTANCE変数が必要です。

MANIFEST\_NAME は、マニフェストのファイル名とは異なります。postfix \_db\_manifest.yaml を除いたファイル名です

## 。 gitlab\_ops\_pgp\_test

gitlab-opsのマニフェストファイルを更新する場合に実行します。RAWのクローンが必要です。変数は一切必要ありません。
SCDテーブルの場合は、TASK\_INSTANCE変数を指定することもできます。これは、データベースに接続するためにCloudSQL Proxyが動作している必要があるため、pgp\_testジョブとは別です。

#### (©) dbt Run

これらのジョブはsnowflake-dbt-ci.ymlで定義されています。

DBTモデル変更のMRの一部として、変更が本番環境で何も壊さないことをテストするためのパイプライン・ジョブをトリガーする必要があります。これらのジョブをトリガーするには、このMRの下部にある "Pipelines "タブに行き、適切なステージ (dbt runまたはdbt misc) をクリックします。

これらのジョブは ci ターゲットにスコープされています。このターゲットは、snowplowおよびversionデータセットの

データのサブセットを選択します。dbtのすべての実行ジョブについて、ジョブの成果物が得られます。これには、コンパイルされたコードと実行結果が含まれます。

これらのジョブは、プライマリRAWデータベースに対して実行されます。

ほとんどのdbt runジョブは、テストを必要とするdbtモデルを指定する変数でパラメータ化することができます。このイントロでは、変数の設定方法の例をご覧いただけます。

変数 DBT\_MODELS は、モデル選択の構文に関するdbtドキュメントの任意の例の代用となります。

data-tests プロジェクトのテストに対する変更をテストする場合は、手動ジョブにブランチ名とともに DATA\_TEST\_BRANCH を 渡すことができます。これにより、packages.yml の data-tests パッケージのブランチが更新されます。これはdbt testを実行しているすべてのジョブで有効です。

また、モデル選択の最後に--fail-fastを追加することで、最初の失敗で素早くdbtの呼び出しを終了させることができます。詳しくはdbtのドキュメントをご覧ください。

#### 十♪♪♪♪モデル

どのモデルを実行するかは、変数 DBT\_MODELS で指定します。

## **+** るスペック\_Lモデル

変数 DBT\_MODELS で L ウェアハウスを使って実行するモデルを指定

## +₱specify\_xl\_model

XL ウェアハウスを使用して実行するモデルを変数 DBT\_MODELS で指定する。

\_\_\_\_\_\_

変数 DBT\_MODELS でどのモデルを除外するかを指定します。

\_\_\_\_\_\_

変数DBT\_MODELSでL倉庫を使って除外するモデルを指定する

## -- କଳକ

XL ウェアハウスで除外するモデルを変数 DBT\_MODELS で指定します。

## **+**₩ Specify\_raw\_model

**RAW**データベースのクローンに対して**dbt**モデルを指定します。このジョブは、RAWのクローンに対して実行されます。clone\_raw が必要です。

ジョブが実行されたことを確認します。この機能は以下のような場合に有効です。

- あなたは新しいシートロードファイルをアップロードしています。これを使って、シートを追加するのと同じMRで、 sheetload dbtモデルをテストすることができます。
- 新しいgitlab.comやその他のpgpテーブルを追加しています。これを使って、テーブルを追加するのと同じMRでdbtモデルをテストすることができます。
- dbtのスナップショットを追加して、そのスナップショットの上に構築されたモデルをテストしたい。

## 

このジョブでは、特定のシードファイルをテストします。

変数 DBT\_MODELS で種ファイルを指定します。

#### 

実行するスナップショットを変数 DBT\_MODELS で指定します。このジョブは、RAW のクローンに対して実行されます。clone\_raw が必要です。

ジョブを実行したことになります。

# + 🗗 🍩 🖏 specify\_l\_snapshot

変数 DBT\_MODELS で実行するスナップショットを指定します。このジョブは、大規模なSnowFlake倉庫を使用して、RAWのクローンに対して実行します。clone\_rawジョブが実行されている必要があります。

dbt Misc

これらのジョブはsnowflake-dbt-ci.ymlで定義されています。

すべてのテスト

すべてのテストを実行する

• 注:テストが含まれているので、dbt\_runステージジョブを実行している場合は、このジョブを実行する必要はありません。

日データテスト

データテストのみの実行

## ♀鮮度

ソースフレッシュネステストを実行します。このジョブは、RAWのクローンに対して実行されます。clone\_rawジョブが実行されている必要があります。

## Qペリスコープ クエリ

このジョブは自動的に実行され、.sqlファイルが変更されたときにのみ表示されます。最も単純な形式では、このジョブは、現在変更されているモデルのいずれかがPeriscopeで照会されているかどうかを確認します。照会されている場合、ジョブは関連するダッシュボードを確認するように通知して失敗します。照会されていない場合は、ジョブは成功します。

## 現在、この仕事の注意点は

- どのダッシュボードをチェックすべきかは教えてくれません。
- 文字列補間構文(例:retention\_[some\_variable])でクエリされたテーブルは検証できません。 テーブルがdbtでエ
  - 210/ 243

イリアスされている場合は検証できません。

説明

ここでは、periscopeのクエリの仕組みを説明します。

git clone -b periscope/master -single-branch https://gitlab.com/gitlab-data/periscope.git --depth 1

これは、periscopeプロジェクトのクローンです。

git diff origin/\$CI\_MERGE\_REQUEST\_TARGET\_BRANCH\_NAME...HEAD --name-only | grep -iEo "(.\*)\.sql"| sed -E 's/\.sql//' | awk -F '/' '{print tolower(\$NF)}' | sort | uniq > diff.txt

これは、masterブランチ(つまりtargetブランチ)から現在のコミット(HEAD)までに変更されたファイルの一覧を取得します。次に、sql ファイルだけを見つけて (grep) 、.sql を空の文字列に置き換えます (sed)。awkを使用して、ファイルの各行の最後のカラムの小文字を表示します(\$NF(フィールド数)で表されます)。フィールドの区切りにはスラッシュ(/)を使用します。出力はdirectory/directory/filenameであり、ほとんどのdbtモデルはファイル名にちなんだテーブルに書き込むと仮定しているので、これは期待通りに動作します。その後、結果をソートし、ユニークなセットを取得して、diff.txtというファイルに書き込みます。

periscope\_check.py

periscope リポジトリ全体を再帰的に検索し、現在問い合わせ可能な3つのスキーマのうち、from|join 文にマッチする文字列を探します。一致したファイルのクリーニングを行い、参照されているすべてのファイルに対応するテーブル名の辞書を作成します。その後、diff.txtを読み込んでルックアップを行い、comparison.txtに書き込み、モデル名に基づいてマッチングを行います。

if (( \$(cat comparison.txt | wc -1 | tr -d ' ') > 0 )); then echo "Check these!" && cat comparison.txt && exit 1; else echo "All good" && exit 0; fi;

ワードカウント (wc) を使用して、比較ファイルに何行あるかを確認します。**0**行以上ある場合は、その行を表示し、失敗して終了します。行がない場合は成功して終了します。

## safe\_model\_script

すべてのSAFEデータが適切なスキーマに格納れことを保証するために、MNPIデータを持つソースモデルの下流にあるすべての モデルPRODで例外タグを持つか、制限付きスキーマになっていなければならない。このCIジョブは、この状態への準拠をチェックします。MRがこのジョブに失敗した場合は、監査を受けてMNPIデータが変更されていないことを確認し、適切な例外タグを 追加するか、モデルを適切な制限付きスキーマに移行する必要があります。

スキーマ・テスト

スキーマテストのスナ

ップショットのみを実

行する

スナップショットを実行します。このジョブは、RAWのクローンに対して実行します。clone\_raw\_fullジョブが実

行されている必要があります。 specify\_tests

変数 DBT\_MODELS で指定されたモデルテストを実行します。

急パイソン

これらのジョブは .gitlab-ci.yml で定義されています。

.pyファイルが変更されたときにのみ表示されるジョブがいくつかあります。すべてのジョブは、.pyファイルが存在する新しいコミットごとに自動的に実行されます。それ以外のジョブは実行できません。他のジョブは

パーミフロストの取扱説明書

Permifrostのドライランを行うマニュアル作業。

## YAML\_VALIDATION

permissions/snowflake/roles.ymlに変更があったときにトリガーされます。YAML が正しくフォーマットされているかどうかを検証します。

## スノーフレークストップ

これらのジョブは .gitlab-ci.yml で定義され

ています。 clone\_stop

MRがマージまたはクローズされたときに自動的に実行されます。手動では実行しないでください。

# データテストパイプライン

以下はすべて、レポで提供されている変更点を使って、Prod DBに対して実行します。以下の実行にはクローニングは必要ありません。

## all tests prod

analytics & data tests repoにあるすべてのテストを実行します。

## ☐ data\_tests\_prod

アナリティクス&データテストのレポにあるすべてのデータテストを

実行します。

analytics & data tests repoにあるすべてのスキーマテストを実行します。

テストプログラムの指定

変数 DBT\_MODELS で指定されたモデルテストを実行します。

layout: handbook-page-toc title:"dbt Guide" 説明"data build tool (dbt) Guide" このページでは

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

## クイックリンク

PrimaryProject{:.btn .btn-purple-inv} dbtdocs{:.btn .btn-purple-inv}.

何を、なぜ

dbtは、DataBuildToolの略で、データウェアハウス内のデータ変換を管理するためのオープンソースプロジェクトです。dbtは、データウェアハウス内のデータ変換を管理するためのオープンソースプロジェクトです。データがウェアハウスにロードされると、チームは分析に必要なすべてのデータ変換を管理することができます。また、dbtにはテストとドキュメントが組み込まれているので、生成・分析するテーブルに高い信頼性を持たせることができます。

以下のリンクは、dbtとは何かについての優れた概要を示しています。

- とは何でしょか? これは、ツールを理解するための技術的ではない概要です。
- dbtとはより技術的な内容で、ドキュメントからの抜粋です。

しかし、なぜ私たちはdbtを使うのでしょうか?理由はいくつかあります。

第一に、オープンソースのツールであり、活発なコミュニティが存在することです。オープンソースのツールを選択することで、より大きなデータコミュニティとのコラボレーションが可能になり、独自のソリューションを採用した場合よりも早く問題を解決することができます。

**2**つ目は、バージョン管理を念頭に置いて作られていること。**GitLab**の場合、会社の構築や運営に製品を使用しているので、これは必須です。

第三に、アナリストの言語である**SQL**を話すことができます。**SQL**は多くの人の仕事に欠かせないものになってきているので、貢献できる人の数も増えています。

最後に、テストとドキュメントを最初から統合することで、チームがより速く動けるようになります

。dbtの基本についての詳細は、データアナリストのオンボーディング問題テンプレートをご覧くだ

さい。

時には、データ変換のためにdbtのパッケージを利用することがあります。パッケージ管理はdbtに組み込まれています。利用可能なパッケージの全リストはdbtHubサイトにあります。

## dbtの実行

dbtに興味をお持ちの方は、dbtのドキュメントJaffle Shopという架空の企業のデータを扱うためのセットアップに関する素晴らしいチュートリアルがあります。

dbtを使用してデータチームプロジェクトに貢献したい場合は、Snowflakeインスタンスへのアクセスを取得する必要があります。

構成

- Snowflakeインスタンスへのアクセスを確保します。
- ◆ Makeがインストールされていることを確認する (新しいMacやXCodeではインスト
- ールされているはず) ホームディレクトリに.dbtというフォルダを作成する
- ~/.dbt/フォルダの中には、以下のサンプルプロファイルのようなprofiles.ymlfileがあるはずです。
- 可能な限り小さい倉庫を環境変数として保存する必要があります。当社の**dbt**ジョブでは、倉庫を識別するための変数名としてSNOWFLAKE\_TRANSFORM\_WAREHOUSEを使用しています。この環境変数は、.bashrcまたは.zshrcファイルで以下のように設定できます。
  - o export SNOWFLAKE\_TRANSFORM\_WAREHOUSE="ANALYST\_XS"
  - ◇ より多くの計算が必要な場合は、**dbt**コマンドに--vars '{warehouse\_name: analyst\_x1}'を追加することで、変数を上書きすることができます。
- アナリティクス・プロジェ
- クトのクローンを作成しま

す。

o Dockerがインストールされていることを確認する

なお、これらのステップの多くは、新しいアナリストが実行することを推奨するオンボーディング・スクリプトで行われます。

Venvワークフロー

#### {:#Venv-workflow}

Macをお使いの方にお勧めのワークフローです。 dbtの使い

方

- DBT\_PROFILE\_PATH環境変数が設定されていることを確認します。これは、onboarding\_script.zshを使用している場合は設定されているはずですが(最新で定期的に更新されているので、これを使用することをお勧めします)、そうでない場合は、あなたの
  - .bashrcまたは.zshrcにexport DBT\_PROFILE\_PATH="/<your\_root\_dir/.dbt/"を追加するか、ローカルのターミナルセッションで同じコマンドを実行してください。
- SSH の設定が GitLab の指示通りになっていることを確認します。鍵は ~/.ssh/ にあり、鍵はパスワードなしで生成されているはずです。

○ また、メインプロジェクトのdbt depsを実行するには、このプロジェクトへのアクセスが必要になります。

- 注意:デフォルトのブラウザがクロームに設定されていることを確認してください。内蔵の**SSO**ログインはクロームでのみ動作します。
- 注意:初めてdbtを実行する前に、make-prepare-dbtを実行してください。これにより、venvがインスト
- ールされていることが確認できます。dbtコンテナを起動し、その中でシェルからコマンドを実行するには、make run-dbtを使用します。

• これにより、ローカルのprofiles.ymlやrepoファイルを含む、dbtの実行に必要なすべてのものが自動的にインポートされます。

• 現在のブランチのドキュメントを見るには、make run-dbt-docs を実行してから、ウェブブラウザーで localhost:8081 にアクセスしてください。このためには、profiles.yml で docs プロファイルが設定されている必要があることに注意してください。

Dockerのワークフロー

{:#docker-workflow}になります。

venvのワークフローは前提条件が少なく、かなり高速であるため、主にLinuxを使用しているユーザーには以下のようなワークフローを推奨します。

dbtとその依存関係をローカルで処理する際の複雑さを排除するために、メインの分析プロジェクトでは、Dockerコンテナ内でのdbtの使用をサポートしています。このコンテナはdata-imageプロジェクトから構築します。Makefileの中には、この作業を容易にするためのコマンドが用意されています。様々なmakeコマンドとその役割について質問がある場合は、make helpを使ってコマンドのリストとそれぞれの役割を確認してください。

初回実行時(およびコンテナが更新されたとき)には、以下のコマンドを実行してください。

- make update-containers
- 2. メイク・クリーンアップ

これらのコマンドは、最新バージョンのコンテナを確実に取得し、ローカルのDockerをクリーンアップします。 の環境を整え ます。

### dbtの使用

- DBT\_PROFILE\_PATH環境変数が設定されていることを確認します。これは、onboarding\_script.zsh(定期的に更新される最新のものを使用することをお勧めします)またはonboarding\_script.shを使用している場合は設定されているはずですが、そうでない場合は、.bashrcまたは.zshrcにexport DBT\_PROFILE\_PATH="/<your\_root\_dir/.dbt/"を追加するか、ローカルのターミナルセッションで同じコマンドを実行することで設定できます。
- SSH の設定が GitLab の指示通りになっていることを確認します。鍵は ~/.ssh/ にあり、鍵はパスワードなしで生成されているはずです。

また、メダンプロジェクトのdbt depsを実行するには、このプロジェクトへのアクセスが必要になります。 db◆コンテナを起動して、その中でシェルからコマンドを実行するには、make dbt-imageを使用します。

- これにより、ローカルのprofiles.ymlやrepoファイルなど、dbtの実行に必要なものがすべて自動的にインポートさ
   れます。 見つからない変数 (GIT\_BRANCH、KUBECONFIG、GOOGLE\_APPLICATION\_CREDENTIALSなど) に関する警告が表示されることがあります。 Airflowでの開発でない限り、これは問題なく、予想されることです。
- 現在のブランチのドキュメントを見るには、make dbt-docs を実行してから、ウェブブラウザーで localhost:8081 にア クセスしてください。このためには、profiles.yml で docs プロファイルが設定されている必要があることに注意してください。
- dbtコンテナに入ったら、通常通りにdbtコマンドを実行します。
- レポ内の任意のファイルに加えられた変更は、コンテナ内で自動的に更新されます。エディタでファイルを変更しても、コンテナを再起動する必要はありません。

コマンドライン・チートシ

これは、プライマリコマンドリファレンスの簡易版で、dbt固有

のものです。

- dbt clean /dbt\_modules (depsの実行時に生成される) と/targetフォルダ (モデルの実行時に生成される) を削除します。
- dbt run regular run
- モデル選択の構文(出典)。モデルを指定することで、関連性があると思われるモデルのみを実行/テストすることができ、時間を大幅に節約することができます。しかし、上流の重要な依存関係を指定し忘れるリスクがあるので、構文をし

## っかり理解しておくと良いでしょう。

- o dbt run --models modelname modelnameのみを実行します。
- o dbt run --models +modelname modelnameとすべての親を実行します。

- o dbt run --models modelname+-modelnameとすべてのchildrenを実行します。
- odbt run --models +modelname+ モデル名と、その親と子のすべてを実行します。
- odbt run --models @modelname modelname、すべての親、すべての子、およびすべての子のすべての親を実行します。
- dbt run --exclude modelname modelname以外のすべてのモデルを実行
- します。 なお、これらはすべてフォルダ選択の構文でも動作します。
  - dbt run --models folder フォルダ内のすべてのモデルを実行します。
  - dbt run --models folder.subfolder サブフォルダ内のすべてのモデルを実行します。
  - dbt run --models +folder.subfolder サブフォルダ内のすべてのモデルと、すべての親モデルを実行します。
- dbt run --full-refresh インクリメンタルモデルを更新します。
- dbt test カスタムデータテストおよびスキーマテストを実行します。
- dbt seed data-pathsディレクトリで指定されたcsvファイルをデータウェアハウスにロードします。このガイドのseedのセクションも参照してください
- dbt compile すべてのモデルをコンパイルします。このコマンドは定期的に実行する必要はありません。モデルを実行すると、dbtはモデルをコンパイルします。

オンボーディング・スクリプトを実行した場合のみ動作します。

- dbt\_run\_changed 変更されたモデルのみを実行するためにあなたのコンピュータに追加された関数です(これは dockerコンテナ内からアクセスできます)。
- cycle\_logs dbtのログを消去するためにあなたのコンピュータに追加した機能です(dockerコンテナ内からはアクセスできません)。
- make dbt-docs dbtのドキュメントをWebブラウザで表示するためのローカルコンテナを起動するコマンドです。 ローカルホスト:8081

dbtプロジェクトへの貢献のための構成

dbtへの貢献に興味をお持ちの方のために、お勧めのローカル環境の設定方法をご紹介します。

- GitHubのUIを使ってdbtプロジェクトを個人の名前空間にフォーク
- する プロジェクトをローカルにクローンする
- 以下のコマンドでdbtの仮想環境(venv)を作成します。

```
cd~。
mkdir .venv # これはあなたのルート"~"ディレクトリにあるはずです python -m venv .venv/dbt
ソース ~/.venv/dbt/bin/activate
pip install dbt
```

- 仮想環境の起動を容易にするために、.bashrcまたは.zshrcにdbt!="source ~/.venv/dbt/bin/activate "というエイリアスを追加することを検討してください。
- 同じターミナルウィンドウでdbtプロジェクトに移動します。コマンドプロンプトの最初に(dbt)と表示されているはずです。
- pip install -r editable\_requirements.txtを実行します。これにより、venvでローカルにdbtを実行したときに、あなたのマシン上のコードを使用していることが保証されます。
- どのdbtを実行してvenvを指しているか確認する
- ローカルでコードを開発し、通常通りに変更をコミットし、GitHubの自分の名前空間にプッシュする

PRのためにコードを提出する準備ができたら、相手のCLAに署名していることを確認しましょう。

スタイル&ユースガイド

### モデル構造

よりKimballスタイルのウェアハウスへの移行に伴い、ウェアハウスやプロジェクト構造でのモデルの整理方法を改善しています。以下のセクションはすべて、dbtのデフォルトであるmodelsディレクトリの下のトップレベルのディレクトリになります。この構造は、Fishtown Analyticsのプロジェクト構造を参考にしています。

\*\*Legacy Structure\*\* {: .panel-heading} {: .panel-heading},

Kimballの次元モデリングに注目する前は、CorrとStagnittoの「Agile Data Warehouse Design」で紹介されているBEAM\*アプローチのモデリングからインスピレーションを得ていました。既存のモデルの多くは、今でもそのパターンを踏襲しています。このセクションの情報は、ハンドブックの過去のイテレーションからのものです。

- 最終的な)\_xf dbtモデルのゴールは、BEAM\*テーブルでなければなりません。つまり、ビジネスイベントの分析とモデルの構造に沿って、ビジネスを測定するための、誰が、何を、どこで、いつ、何人で、なぜ、どのように、という質問の組み合わせに答えるものです。
- ベースモデル-ソーステーブルを参照する唯一の**dbt**モデル。ベースモデルは最小限の変換ロジックを持つ(通常、データインテグリティに問題のある行や、分析対象外のフラグが立っている行をフィルタリングしたり、分析しやすいようにカラム名を変更する程度)。
- エンドユーザーモデル 解析に使用されるdbtモデル。モデルの最終バージョンは、BEAM\*テーブルになることを目的としている場合、\_xfという接尾辞が付けられます。このモデルは、ビジネスイベント分析とモデルの構造に従い、ビジネスを測定するために、誰が、何を、どこで、いつ、いくつ、なぜ、どのように、という質問の組み合わせに答えなければなりません。エンドユーザーモデルは、レガシースキーマにあります。

\*\*FY21-Q4モデルの移行\*\* {: .panel-heading}。

FY21-Q4では、アナリティクス・データベースを置き換えるために、prodおよびprepデータベースを導入しました。この2つの新しいデータベースは、アナリティクス・データベースを完全に置き換えるものです。

ローカル開発も、カスタムスキーマからカスタムデータベースに変更。ソース

すべての生データは、SnowflakeのRAWデータベースに残っています。これらのRAWテーブルは、ソーステーブルまたはRAWテーブル と呼ばれています。これらは通常、元のデータソースを示すスキーマに保存されます。

ソースはdbtのsources.ymlファイルで定義します。

- dbtソースでデータベースを参照するために変数を使用し、Snowflakeクローンでの変更をテストする場合、プログラムで参照を設定できるようにしています。
- 通常の慣習に合わない名前や意味が不明瞭な名前のソーステーブルを扱う場合、元の名前が混乱していたり混乱している場合は、識別子を使ってソーステーブルの名前を上書きします。(識別子の使用に関するドキュメント)

```
# 良い
テープ
ル
- 名前: bizible_attribution_touchpoint
識別子: bizible2 bizible_attribution_touchpoint c
# 悪い
テープ
```

- name: bizible2 bizible\_attribution\_touchpoint c

ソースモデル

すべての生データの上に、非常に薄いソースレイヤーを強制的に配置しています。このディレクトリには、ソース固有の変換の大半が保存されます。これらは、生データから直接取り出し、ファクトやディメンションを作成するために必要な準備作業を行う「ベース」モデルであり、以下のこと*だけを*行う必要があります。

- フィールド名をユーザーフレンドリーな名
- 前に変更 カラムを適切なタイプにキャスト
- 当面の間、100%使えることが保証されている最小限の変換です。例えば、データが乱れていることが知られているフィールドからSalesforce IDを解析することがその例です。
- 論理的に命名されたスキーマへの配置

基礎となる生データが完璧にキャストされ、名前が付けられている場合でも、フォーマットを強制するソースモデルが存在すべきです。これは、エンドユーザーの利便性を高めるためで、見る場所が1つで済みますし、この完璧なデータが機密性の高いものである場合には、パーミッションがより明確になります。

ソースモデルでは以下のようなことをしてはいけません。

- データの削除 他のテー
- ブルへの結合
- カラムの意味を根本的に変える変換

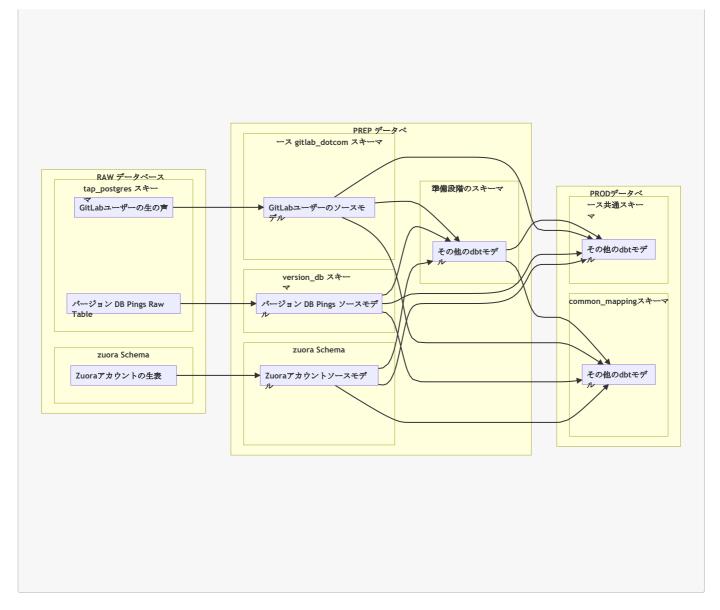
大多数のユーザーにとって、ソースモデルは「生」のデータであると考えられます。覚えておくべき重要なポイント

- これらのモデルは、データソースの種類に基づいて、準備データベースの論理的に名付けられたスキーマに書き込まれます。例えば、以下のようになります。
  - raw.zuora に保存されている Zuora データは、prep.zuora にあるソースモデルを持ちます。
  - raw.tap\_postgres.gitlab\_db\_\*に格納されたテーブルを持つGitLab.comのデータは、ソースモデルが prep.gitlab dotcom
  - o テーブルが raw.tap\_postgres.customers\_db\_\* に格納されている Customers.gitlab.com データは、ソースモデルが prep.customers\_db に格納されています。
- これらのモデルは、ソースごとに整理されるべきです。これは通常、生のデータベースのスキーマに
- 対応します。
- ソースモデルのみ、ソース/ローテーブルから選択する必要がある
- ソースモデルは、生のデータベースから直接選択すべきではありません。代わりに、以下のようにjinjaでソースを参照する必要があります。

FROM {{ source('bamboohr', 'job\_info') }}。

- ソースモデルは、/models/sources/<data\_source/>ディレクトリに配置してください。
- ソースモデルは、必要なすべてのデータタイプのキャストを行い、キャストする際には:: 構文を使用します(より少ない文字数で同じことができ、よりすっきりとした印象を与えます)。
  - の 理想的には、ソースモデルはすべてのカラムをキャストする必要があります。暗黙的ではなく明示的な方が良い。仮定のテスト
- ソースモデルでは、すべてのフィールドの命名を行い、フィールド名を標準的なフィールド命名規則に適合させる必要
- があります。 予約語を使用しているソースフィールドは、ソースモデルで名前を変更する必要があります。
- 特に大きなデータのソースモデルは、常に論理フィールド(通常は関連するタイムスタンプ)のORDER BYステートメントで終わるべきです。これは基本的にウェアハウスのクラスタキーを定義し、Snowflakeのマイクロパーティショニングを活用するのに役立ちます。

ソースモデルが生のテーブルとどのように関連しているか、また、ソースモデルが下流のすべてのモデリングのためのクリーンなレイヤーとしてどのように機能するかを視覚的に示すために、次の図を参照してください。



### 機密データ

場合によっては、生の値を公開してはいけない列があります。これには、お客様のEメールやお名前などが含まれます。しかし、このようなデータを必要とする正当な理由があり、以下のような方法でデータを保護しながらも、正当に知る必要のある人には利用できるようにしています。

与えられたモデルでは、ソースのフォーマットは上記の通りです。ソースモデルにはカラムのハッシュ化はありません。これは、セキュリティやアクセスの観点から、生データと同じように扱う必要があります。

機密性の高いカラムは、schema.ymlファイルでメタキーを使用し、sensitiveをtrueに設定することで文書化されます。例は以下の通りです。

- 名前: sfdc\_contact\_source
の説明を参照してください。SFDC Contactsカラ
ムのソースモデル。
- name: contact\_id
tests:
- not\_null
- ユニーク
- name: contact\_email
meta:
センシティブ: true
- name: contact\_name
meta:
センシティブ: true

そして、ソースモデルから2つの独立したモデル、つまり感度の高いモデルと低いモデルが作成されます。

non-sensitiveモデルでは、hash\_sensitive\_columnsというdbtマクロを使用します。これはソーステーブルを受け取り、metaフィールドにsensitive: trueを持つすべてのカラムをハッシュ化します。すべてのカラムが同じようにハッシュ化されるため、特定のジョインキーは指定されません。必要であれば、マクロの外で別のカラムをジョインキーとしてこのモデルに追加することができます。sfdc\_contactモデルは、この良い例です。2つのカラムがハッシュ化されていますが、contact\_idという追加の主キーが指定されています。

センシティブモデルでは、**dbt** マクロ nohash\_sensitive\_columns を使用してジョインキーを作成します。このマクロは、ソーステーブルとジョインキーとしてのカラムを受け取り、ジョインキーとしてハッシュ化されたカラムと、ハッシュ化されていない残りのカラムを返します。sfdc\_contact\_piiモデルは、このマクロが使われている良い例です。

すべてのハッシュにはソルトも含まれています。これらは環境変数で指定します。データの種類によって、さまざまなソルトがあります。これらは get\_saltマクロで定義され、ローカル開発用の dbt コンテナを使用する際にも設定されます。

一般的に、チーム メンバーは、Snowflake UI 内のクエリ文字列で使用されるソルトを見ることを許可されるべきではありません。テーブルモデルでは、この目標はSnowflakeの組み込みENCRYPT関数を使用することで満たされます。ビューに実体化されたモデルでは、ENCRYPT 関数は機能しないようです。代わりに、セキュアなビューを使用した回避策が利用されます。セキュアビューは、DDL の表示を所有者のみに制限し、ハッシュの可視性を制限します。セキュアなビューを作成するには、モデル構成で secure true に設定します。説明したようにハッシュ機能を利用したビューであっても、セキュアなビューとして構成されていないものは、おそらく問い合わせができません。

ステージング

Kimballモデリングを導入する前は、ほとんどのモデルがStagingカテゴリーに分類されていました。

このディレクトリには、ビジネスに特化した変換の大半が保存されます。この層のモデリングは、ソースモデルの作成よりもかなり複雑で、ビジネスの分析ニーズに合わせて高度に調整されたモデルとなります。これには

- 無関係なレコードのフィルタリング
- アナリティクスに必要なカラムの選択
- 抽象的なビジネスコンセプトを表現するためのカラム名の変更 他
- のテーブルへの結合
- ビジネスロジックの実行
- Kimballの方法論に従ったfct\_\*およびdim\_\*テーブルのモデリング

### **Dimensional Modeling**

### {:次元}

ファクトテーブルとディメンショナルテーブルは、ソースデータから抽象化され、ビジネスに関連するエンティティと プロセスを表します。なぜディメンショナル・モデリングを使用するのかという情報は、ハンドブックのEDWの部分に記載 されています。

モデリング開発プロセス

- 1. ディメンション・バス・マトリックスをEDMの青写真として拡張する。
- 2. 適切なLucidChart ERDにテーブルを追加します。
- 3. ディメンション・テーブルを展開します。各ディメンションには、未知のものを表す-1キー・バリューの共通レコード・エントリも含まれています。
- 4. ファクト・テーブルを作成します。正しいディメンジョン・キーをファクトに入力し、不明なキーには-1のキー値を使用します。

### 命名基準

セルフサービス型のデータ環境を構築する際には、命名規則をはじめとする意図的な取り組みが重要です。目標は、初心者、中級者、上級者がデータウェアハウスを簡単に操作できるようにすることです。以下のベストプラクティスに従うことで、これを可能にします。

1. FACT TABLES: fct\_<verb> 事実は、出来事や現実世界で起きているプロセスを表しています。ファクトは、アクションや「動詞」を表しているため、しばしば識別することができます。(例:セッション、トランザクション)

2. DIMENSION TABLES: dim\_<noun>=ディメンション・テーブル。ディメンションは、ファクト・レコードに記述的なコンテキストを提供します。ディメンションは、人、場所、物などの「名詞」であるため、しばしば識別することができます(例:顧客、従業員)。ディメンション属性は「形容詞」として機能します。(例:顧客タイプ、従業員部門)

- 3. 例えば、dim\_customersではなくdim\_customerのように、単数形のネーミングを使用する必要があります。
- **4.** 同じようなデータをグループ化するには、テーブル名やカラム名にプレフィックスを使用します。例えば、 dim\_geo\_location、dim\_geo\_region、dim\_geo\_sub\_regionのように、アルファベット順にソートしてもデータは論理的にグループ化されます。
- 5. 主キーおよび外部キーの命名には、ディメンションテーブル名を使用します。これにより、追加の属性を取り込むために どのテーブルを結合する必要があるかが、ユーザーに明確になります。たとえば、dim\_crm\_accountのプライマリ・キーはdim\_crm\_account\_idです。このフィールドがfct\_subscriptionにある場合、dim\_crm\_account\_idという名前になり、追加のアカウント詳細を得るためにdim\_crm\_accountに結合する必要があることが明確になります。

#### モデリング要件

- 一般的には、まず共通スキーマにディメンジョンを作成し、次にファクト・テーブルを構築するコードがこれらの共通ディメンジョンを参照して(LEFT JOIN)、共通のキーが利用できるようにします。論理的には、ファクト・テーブルにあるすべてのキーがディメンジョンで利用可能であると100%確信できる状況もあります。これは、ディメンジョンのソースと適用されるロジックに依存します。
- ファクトとディメンションの両方を、TDF (Trusted Data Framework) テストを用いてテストする必要があります。
  - ∘ ディメンションについては、-1(未知)のdimension\_idの存在や、総行数をテストすることができます。
  - ファクトについては、不適切な粒度の結合によってレコード数/行数が拡大していないかどうかをテストし、関連するファクト・レコードのディメンション・テーブルからディメンション属性を引き出して、ゴールデン・データ・レコードの期待値と比較するゴールデン・データ・テストを追加することができます。
- fct およびdim モデルは、クエリのパフォーマンスを向上させるためにテーブルとして
- 実体化する必要があります。
  - ハッシュ化されたサロゲートキーです。
  - o 自然なユニークキーです。このキーの値は、複数のカラムから構成され、一意性を生み出すことができます。
- モデルはテストされ、モデルと同じディレクトリにあるschema.ymlファイルで文書化されま
- す。すべてのファクトとディメンションには、以下の監査カラムがあります。
  - revision\_number モデルの論理的な変更を表す、手動で増やされた番号 created\_by GitLab のユー
  - 。 ザー ID。
  - ∘ updated\_by これはGitLabのユーザーIDです。
  - o model\_created\_at timestamp これは、モデルが作成されたときの静的な値です
  - o model updated at timestamp これは、モデルが誰かによって更新された最後の時間で
  - 。 す dbt\_created\_at timestamp これは、テーブルが作成されたときにdbtによって入力されます
  - **dbt\_updated\_at** タイムスタンプ データが最後に読み込まれた日付です。ほとんどのモデルでは**dbt\_created\_at** と同じになりますが、インクリメンタルモデルでは例外です。
- ディメンションテーブルをサポートするPrep(prep\_)とMapping/look-up(map\_)テーブルはcommon\_mappingに作成する必要があります。
  スキーマです。
- 追加のBridge(bdg\_)テーブルは、共通のスキーマに存在する必要があります。これらのテーブルは、2つのテーブル間の 多対多の関係を解決するための中間テーブルとして機能します。

### ERDの要件

- Lucidchartで生成されます。
- すべての関連モデルのdbtドキュメントにiframeとして埋め込ま
- れる ERDから与えられたモデルのdbtドキュメントへのクロスリ
- ンク 適切な関係の接続
- 主キーと外部キーの一覧
- テーブルの性質を示すもので、変更される可能性が低い他のカラムを少なくとも3~5個 ワーキ
- ングSQLの参照例

次元モデルは、使いやすく、ユーザーのために設計されていることを意味します。次元モデルは、ソース・モデルとは逆に非 正規化されている可能性が高く、読みやすく、解釈しやすいだけでなく、結合数を減らすことで効率的なクエリを可能にしま す。

• CTEが2つのモデルに重複している場合、別のfct\_/dim\_テーブルになることがよくあります。

マルツ

Martモデルは、ビジネスエンティティとプロセスを記述します。マーケティング、財務、製品などのビジネスユニットごとにまとめられていることが多い。

モデルがこのディレクトリにあると、データがきれいにモデリングされていて、クエリの準備ができていることをビジネス関係 者に伝えることができます。

ファクトテーブルおよびディメンションテーブルの命名規則に従い、すべてのマートは mart\_ という接頭辞

で始まります。マートは、他のマートの上に構築してはいけません。マートは、FCTおよびDIMテーブルの上

に構築する必要があります。

ワークスペース

dbtプロジェクトでは、チーム独自のコードを書くためのスペースを提供しています。これは、すべてのコーディングガイドやスタイルガイドを遵守する必要のない、チーム独自のコードです。これは、チームがより強固である必要のないソリューションを使って、より早く反復できるようにするためのものです。

プロジェクト内には/models/workspaces/フォルダがあり、チームはそこにworkspace\_<team>というスタイルのフォルダを作成してコードを保存できます。このコードはデータチームによるスタイルのレビューを受けません。マージされる前に懸念されるのは、コードが動作するかどうかと、本番稼動に影響を与えるようなひどい非効率性がコードにないかどうかだけです。

新しいスペースを追加するには

- 分析プロジェクトで課題を作成し、新しいマージリクエストを開く
- /models/workspaces/に新しいフォルダを作成する 例:

workspace\_security

• dbt project.ymlファイルに、新しいワークスペースのエントリを追加します。書き込むべきスキーマを含めます。

- フォルダに.sqlファイルを追加します。
- CODEOWNERSファイルに任意のエントリを追加
- dbtWorkspaceChanges MRテンプレートを使用し、その指示に従ってMRを提出し、レビューと最終的なマージを行います。

新たに追加されたコードがデータウェアハウスに表示されるまでには、最大で24時間かかります。

データチームは、本番のdbt実行を大幅に遅らせるようなコードを拒否する権利を持っています。もしこのようなことがあれば、ワークスペースのためだけに別のdbt実行ジョブを構築することを検討します。

一般

- モデル名はできるだけわかりやすく、可能な限り完全な単語を使用してください(例: accountsの代わりに accts.
- 新しいデータモデルのドキュメント化とテストは、データモデル作成のプロセスの一部です。新しいdbtモデルは、テス

merged.md

2021/11/8

トとドキュメントなしでは完成しません。

• 分析タイプ、データソース(複数の場合はアルファ順)、モノ、集計の命名規則に従う。

```
-- Good
retention_sfdc_zuora_customer_count.sql
悪いリテンション
.sql
```

- すべての {{ ref('....') }} ステートメントは、ファイルの先頭にある CTE の中に置く必要があります。 (これらはimport文と考えてください。)
  - これは、{{ ref('...') }} を持つすべての CTE が SELECT \* のみであるべきだということではありません。モデルにとって意味のあることであれば、ref を持つ CTE で追加の操作を行っても構いません。
- 複雑なSQLを別のモデルに分けたい場合は、DRYで理解しやすくするために絶対にそうすべきです。 materalized='ephemeral'という設定は、基本的にモデルをCTEのように扱う1つのオプションです。

### モデル構成

モデルの構成定義を提供する方法は複数あります。 **dbtdocsforconfiguringmodels**には、モデルの設定方法についての簡潔な説明があります。

モデルを構成するためのガイドラインです。

- デフォルトのマテリアライゼーションはビュー
- デフォルトのスキーマはprep.preparationです。
- ◆ モデルを無効にするには、dbt\_project.ymlの+enabled: false宣言で行う必要があります。 設定は最小の場所に適
- 用する必要があります。
  - ∘ ディレクトリ内の50%以下のモデルが同じ設定を必要とする場合、個々のモデルを設定します。
  - ディレクトリ内のモデルの**50**%以上が同じ設定を必要とする場合、以下のようなデフォルトの設定を強く検討してください。

dbt\_project.ymlを使用していますが、この設定がディレクトリ内の新しいモデルのデフォルトとして適切かどうかを考えてください。

による。

通常の使い方では、dbtは{{ ref('...') }}構文の使い方に基づいて、すべてのモデルを実行する適切な順番を知っています。しかし、モデルをいつ実行すべきかをdbtが知らない場合があります。具体的な例としては、schema\_union\_allやschema\_union\_limitマクロを使用した場合です。この場合、コンパイル時に明示的な参照が行われていないため、dbtはモデルが先に実行されると考えてしまいます。この問題に対処するために、ファイルの中で設定の後に、どのモデルに依存しているかを示すコメントを追加することができます。

```
````sql
{{config({
     "materialized": "view"
   })
}}

--に依存します。{{ ref('s snowplow_sessions') }}。

{{ schema_union_all('snowplow_', 'snowplow_sessions') }}。

````
```

dbtはrefを見て、指定されたモデルの後にこのモデルを構築します。デ

ータベースとスキーマの名前の生成

**dbt**では、カスタムのデータベース名やスキーマ名を生成することができます。これは、私たちのプロジェクトでは、以下の場所を 制御するために広く使用されています。

モデルは実体化されており、生産や開発のユースケースに応じて変化します。

データベース

デフォルトの動作は、**dbt**のドキュメントの「データベースの使用の項に記載されています。というマクロがあります。generate database nameは、書き込み先のスキーマを決定します。

このマクロの動作を独自の generate\_database\_name定義でオーバーライドします。このマクロは profiles.yml で提供された 設定 (ターゲット名とスキーマ) と model config で提供されたスキーマ設定を取り込んで、最終的なスキーマが何であるべ きかを決定します。

スキーマ

デフォルトの動作は、dbtドキュメントの「Usingcustomschemas」セクションに記載されています。というマクロがあります。generate\_schema\_nameは、書き込み先のスキーマを決定します。

このマクロの動作を独自の generate\_schema\_name定義でオーバーライドします。このマクロは profiles.yml で提供された設定 (ターゲット名とスキーマ) と model config で提供されたスキーマ設定を受け取り、最終的なスキーマが何であるべきかを決定します。

#### 開発行動

FY21-Q4では、スキーマではなく開発用のデータベースを持つように変更しました。dbtのユーザーは、TMURPHY\_PRODやTMURPHY\_PREPのように、モデルが書き込まれるスクラッチデータベースを定義する必要があります。

このスイッチは、profies.ymlファイルで定義されたターゲット名によって制御されます。ローカル開発では、決してprodやciをターゲットにしています。

マクロ

#### 命名規則

• ファイル名がマクロ名と一致していること

### 構造

- マクロの説明は、macros.ymlまたはmacros.mdのいずれかのファイルに記述してください。
- dbt-utils

私たちの**dbt**プロジェクトでは、**dbt-utils**パッケージを使用しています。これにより、一般的に便利なマクロがいくつか追加されています。重要なものを挙げてみましょう。

- group\_by このマクロは、フィールド1...Nに対してGroup Byステートメントを構築します。
- star このマクロは、except 引数にリストされている列を除いて、テーブルからすべての列を取り出します。
- surrogate\_key このマクロは、フィールド名のリストを受け取り、値のハッシュを返して一意のキーを生成します。

### 種子

{:#seeds}になります。

シードは、csvファイルからデータウェアハウス(dbtドキュメント)にデータをロードする方法です。これらのcsvファイルは dbtのリポジトリに置かれているため、バージョン管理されており、コードレビューが可能です。この方法は、頻繁に変更されない静的なデータを読み込むのに適しています。dbt seedコマンドで使用するcsvファイルは、長さが最大1k行で数キロバイト以下のものが良いでしょう。

### コラムの整理

ベースモデルを作成する際には、列に何らかの論理的な順序付けをする必要があります。私たちは、以下の4つの基本的なグルー

- プ分けを推奨します。プライマリデータ
- 外部キー

- 論理データ このグループは必要に応じてさらに細分化できる
- メタデータ

プライマリデータは、テーブルを説明する重要な情報です。主キーは、名前などの他の関連する固有の属性とともに、このグループに含まれるべきです。

外部キーとは、別のテーブルを指すすべての列のことです。

論理データは、参照しているオブジェクトを説明する追加のデータディメンションです。Salesforceのオポチュニティの場合は、オポチュニティの所有者や契約額などです。意味のあるものであれば、さらに論理的なグループ化が推奨されます。たとえば、契約額のすべてのバリエーションをグループ化することは意味があります。

任意のグループ内では、カラムはエイリアス名に基づいてアルファベット順に表示されます。

グループ化の推奨事項の例外は、定義されたマニフェストファイルによって抽出を制御する場合です。この完璧な例が gitlab.com のマニフェストで、アプリケーションデータベースからどのカラムを抽出するかを定義しています。これらのテーブル のベースモデルは、マニフェストと同じ順序で並べることができます。ファイル間の差分を比較して正確性を確保するのが簡単だからです。

• グループ内では、エイリアスのアルファベット順に並びます。

```
--良い
SELECT
 ID名
                     AS account id,
                     AS account_name,
 外部キー
 オーナーID
                    AS owner id,
                     AS parent_account_id,
 PID
                     AS zuora_id,
 ジッド
 ロジカルインフォ
 opportunity owner c AS opportunity owner,
 account_owner cAS opportunity_owner_manager,
 owner_team_o cAS opportunity_owner_team,
 メタデータ
 isdeleted ラスタ
                     AS is_deleted,
  クティヴィティデー
                     AS last_activity_date
  卜
FROMテーブル
```

• グループを介さず、エイリアスのアルファベット順に並べる

```
Less Good

SELECT

idAS account_id,
nameASアカウント名。
isdeletedAS is_deleted, lastactivitydateAS
last_activity_date, opportunity_owner c AS
opportunity_owner, account_owner cAS
opportunity_owner_manager, owner_team_o cAS
opportunity_owner_team, owneridAS owner_id,
pidAS parent_account_id,
zidAS zuora_id
FROMデープル
```

オリジナルの名前でアルファベット順に並んでいます。

タグ

dbtのタグは、プロジェクトの様々な部分にラベルを付ける方法です。これらのタグは、実行するモデル、スナップショット、シードのセットを選択する際に利用できます。

タグはYAMLファイルやモデルのコンフィグ設定で追加することができます。dbt\_project.ymlファイルには、タグの使用例がいくつかあります。Trusted Data Frameworkのタグを追加する具体的な例を以下に示します。

分析やデータテストのプロジェクトでは、すべてのタグについてSingle Source of Truth(真実の情報源)を徹底しています。どのタグが使用されているかを文書化するために、ValidTagsCSVを使用しています。マージリクエストでは、すべてのdbt CIジョブの中に検証ステップがあり、プロジェクトで使用されているすべてのタグに対してこのCSVをチェックし、不一致があった場合はジョブを失敗させます。将来的には、このCSVファイルにタグに関するより多くのメタデータを含めることを目指しています。

どのレベルでも適用されるタグは、どのテストにも適用されないことに注意してください。テスト用のタグは、schema.yml ファイル内でテストごとに明示的に適用する必要があります。

トラステッドデータフレームワーク

Trusted Data Frameworkの背景にある哲学については、PlatformページのTrusted Data Frameworkセクションを参照してください。

スキーマからゴールデンデータカバレッジ

TDF (Trusted Data Framework) テストを5つのカテゴリーに分けて実施しています。

- 1. スキーマの整合性を検証するスキーマテスト
- 2. 列のデータ値が事前に定義されたしきい値やリテラルに一致するかどうかを判断する列値テスト
- 3. Rowcount は、事前に定義された期間におけるテーブルの行数が、事前に定義されたしきい値またはリテラルに一致するかどうかを判定するテストです。
- 4. あらかじめ定義された価値の高いデータがテーブルに存在するかどうかを判断するGolden Dataテスト
- 5. Custom SQLは、上記のカテゴリーに適合しない有効なSQLをテストします。

テストは主に2つの場所に保存されます。メインプロジェクト内のYAMLファイルか、DataTestsプロジェクト内に保存されます。

スキーマとカラム値のテストは、通常、メインプロジェクトにあります。これらは、モデルと同じディレクトリにあるschema.yml とsources.ymlファイルに格納されます。

Rowcount、Golden Data、その他のカスタム SQL テストは、常に DataTests プロジェクトにあります。これは、GitLab 内部 でのみ使用されるプライベートプロジェクトです。

タギング

テストのタグ付けは、新しいテストを追加する際の重要なステップです。dbtタグでテストをラベル付けすることは、信頼できるデータダッシュボードを構築する際に、テストを解析して識別する方法です。テストをタグ付けする方法は、その種類によって2通りあります。

一つ目は、YAML 定義にタグを追加することです。これは、ソーステストの場合は YAML 定義の最上位レベルで、モデルテストの場合はカラムレベルで行います。

```
## Source Labeling in
sources.yml version: 2
のソースになります。
 - 名前: ズーラ
   のタグが表示されます。["tdf", "zuora"] です。
## Model Labeing in
schema.yml version: 2
モデルを使用しています。
 - 名前: zuora_accounting_period_source
   の説明を参照してください。Zuora会計期間のソースレイヤーで、列のクリーニングや名前の変
   更を行います。
    - name: accounting period id
      tags:["tdf", "zuora"] テ
      スト。
        - not_null
        - ユニーク
```

これらの例では、下の階層でネストされたすべてのテストにタグが適用されます。タグを

追加するもうひとつの方法は、テストファイルの先頭にある config 宣言を使用することです

0

```
{{ config({
    "tags "を使用します。["tdf", "zuora"] です。
    })
}}
WITH test AS (...)
```

### 一般的なガイダンス

- すべてのモデルはschema.ymlファイルでテストされなければなりません。
- ◆ 最低限、ユニーク・フィールド、not nullable フィールド、外部キー制約をテストする必要がある(該当する場合
- ) dbt テストの出力を MR に貼り付ける必要がある
- テストが失敗した場合は、レビューを依頼する前に修正するか説明する必

要があります。 スキーマテスト

スキーマテストは、既知のテーブル、カラム、およびその他のスキーマ構造の存在を検証するように設計されています。スキーマテストは、計画的または偶発的なスキーマの変更を特定するのに役立ちます。

すべてのスキーマテストは、PASS または FAIL ステータスになります。

スキーマテストの例

目的:このテストでは、重要なテーブルがZuora Data Pipelineに存在するかどうかを検証します。

スキーマテストをdbtのマクロとして実装しました。つまり、ユーザーは**SQL**を書く代わりに、マクロを呼び出すだけでテストを追加することができます。これは raw\_table\_existence マクロによって制御されます。

```
-- ファイル: https://gitlab.com/gitlab-data/analytics/-/blob/master/transform/snowflake-dbt/tests/sources/zuora/existence/zuora_raw_source_table_existence.sql
```

```
{{ config({
    "tags "を使用します。["tdf", "zuora"] です。
    })
}}

{{ raw_table_existence(
    'zuora_stitch',
    ['account', 'subscription', 'rateplancharge'] です。
)}}
```

### カラムの値のテスト

列値テストは、列のデータ値があらかじめ定義されたしきい値内にあるかどうか、または既知のリテラルと一致するかどうかを判定します。列値テストは、幅広い用途に使用できるため、TDFテストの中で最も一般的なタイプです。列値テストは、以下のような場面で役立ちます。

- 変更管理:リリース前とリリース後のテスト
- 重要な過去のデータの合計/集計が、過去に報告された結果と一致していること 既知
- の「承認済み」データが常に存在すること

dbtには、カラムがNULLでないこと、ユニークな値を持つこと、特定の値のみを含むこと、カラムのすべての値が他のモデルで表現されていること(参照整合性)を保証するテストがあります。

また、dbt-utilsパッケージを使用して、さらに多くのテスト機能を追加

しています。すべての列値テストは、PASSまたはFAILのステータスにな

ります。

コラムバリューテスト例1

目的: このテストは、Zuora のアカウント ID フィールドを検証します。このフィールドは常に32文字で、数字と小文字しか持ちません。

**dbt**を使用しているので、すべてのソーステーブルと下流のモデル化されたデータのほとんどについてドキュメントがあります。yamlドキュメントファイルを使って、個々のカラムにテストを追加することができます。

### 行数テスト

行数テストは、列値テストの特殊なタイプであり、その重要性と実用性のために分類されています。行数テストは、一定期間 におけるテーブルの行数が、あらかじめ定義された期待値を満たしているかどうかを判定します。正当な理由でデータ量が急 激に変化する場合は、行数テストを適切に更新する必要があります。

## 行数テスト例1

目的:このテストでは、2019年に作成された18,849件のZuora購読レコードが常にあることを検証します。

このテストはdbtのマクロとして実装されています。つまり、SQLを書く代わりに、ユーザーはマクロを呼び出すだけでテスト

merged.md

2021/11/8

を追加することができます。これは、source\_rowcount マクロによって制御されます。

```
-- https://gitlab.com/gitlab-data/data-
test/-blob/main/tests/sources/zuora/rowcount/zuora_subscription_source_rowcount_2019.sql
{{ config({
    "tags "を使用します。["tdf", "zuora"] です。
    })
}}

{{ source_rowcount(
    'zuora',
    'subscription',
    18489,
    "autorenew = 'TRUE' and createddate > '2019-01-01' and createddate < '2020-01-01'"
)}}
```

#### Rowcountテスト例2

目的:私たちは急成長しているビジネスを行っており、前日からロードされた新規購読が常に50件以上、最大で200件必要です。これは、model\_new\_reocrds\_per\_dayマクロで制御されます。

```
-- https://gitlab.com/gitlab-data/data-
tests/-blob/main/tests/sources/zuora/rowcount/zuora_subscription_source_model_new_records_per_day.s
ql
{{ config({
    "tags":["tdf", "zuora"],
    "severity":"warn",
    })
}}
{{ model_new_rows_per_day(
    'zuora_subscription_source',
    'created_date',
    50,
    200,
    "date_trunc('day',created_date) >= '2020-01-01'"
)}}
```

### ゴールデンデータテスト

トップ顧客の記録、2019年の全世界のユーザー数、100万契約を突破したときのKPIの結果など、データウェアハウスに常に存在し、決して変更してはならないほど重要なデータもあります。このようなケースの中には、新しいデータベース機能を開発することで解決できるものもありますが、これは複雑で、既存のデータ処理のワークフローとは必ずしも一致しない可能性があります。さらに、データ変換に誤ってバグが追加されたり、重要な本番テーブルに対して誤って悪いUPDATEを実行してしまうこともあります。ゴールデンデータテストは、列値テストの特殊なタイプで、既知のデータリテラルの存在を検証し、これらの問題が発生したときにキャッチするのに役立ちます。

Golden DataテストはCSVを使って実装されています。プライバシーを守り、ユーザーが必要なデータをエンコードできるようにするため、Golden DataのCSVをDataTestsプロジェクトの/dataフォルダに保存しています。これらのファイルはdbtパッケージとして本番環境にインポートされ、prep.tdfスキーマでアップロードされます。

ユーザーは、Golden Data Macrosを使用して比較を実行するテストを作成できます。

ゴールデンデータのテスト例

目的:ACMEは当社の最も重要な顧客です。このテストでは、ACMEが常にDIM\_ACCOUNTテーブルに入っていることを検証します。これは、model\_golden\_data\_comparisonマクロによって制御されます。

```
-- dim_account_golden_data
account_name, account_status, account_currency, is_deleted, crm_id
ACME, Active, USD, FALSE, 0016100001BrzkTQZY
```

```
{{ config({
    "tags":["tdf", "dim_account"]
    })
}}
{{ model_golden_data_comparison('dim_account') }}。
```

同様に、この同じデータがRAWデータベースのソーステーブルに存在する場合は、以下のようなフォーマットになります。これは、source\_golden\_data\_comparisonマクロで制御されます。

```
-- sfdc_account_raw_golden_data
name, status, currency, deleted,
id
ACME, アクティブ, USD, FALSE, 0016100001BrzkTQZY
```

```
{{ config({
    "tags":["tdf", "sfdc"] です。
    })
}}

{{ source_golden_records_comparison('sfdc','account') }}。
```

### カスタムSQL

上記のカテゴリーに当てはまらないテストを考えているかもしれません。また、任意の**SQL**をテストとして書くこともできます。留意すべき点は、行が返されなければテストは*合格ということ*です。もし、クエリから何らかの行が返された場合、テストは失敗します。

dbtのドキュメントに掲載されている例です。

```
{{ config({
    "tags":["tdf", "fct_payments"]
    })
}}

-- 返金にはマイナスの金額があるため、合計金額は常に >= 0 である必要があります。
そのため、テストを失敗させるために、これが真実ではないレコードを返します。
    order_id,
    sum(amount) AS
total_amount FROM {{
ref('fct_payments' )}}。 GROUP
BY 1
HAVING total_amount < 0
```

任意の有効なSQLを記述し、任意のdbtモデルやソーステーブルを参照することができます

。マージ要求ワークフロー

テストを追加・更新する際には、いくつかのシナリオが考えられます。

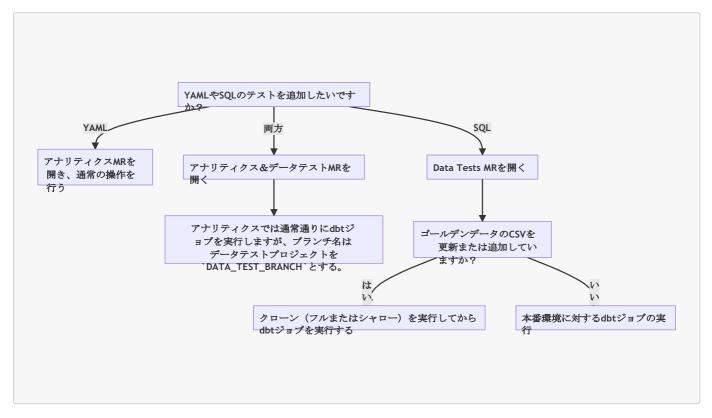
最初のシナリオは、メインプロジェクト内のYAMLファイルのテストを修正または追加することです。これは標準的なMRのワークフローに沿ったもので、何も変わりません。通常通りCIジョブを実行します。

2つ目のシナリオは、data-testsプロジェクト内のテストまたはゴールデンデータのレコードを、analyticsプロジェクト内のMRを 介して更新または追加されているテーブルに対して追加することです。これは最も一般的なシナリオです。この場合、data-testsプロジェクトのMRでパイプラインを実行する必要はありません。分析用MRの通常のdbtパイプラインは実行でき、唯一の変更点は、data-testsプロジェクトのブランチ名を環境変数DATA\_TEST\_BRANCHを介してジョブに渡す必要があることです。

3つ目のシナリオは、data-testsプロジェクトにテストが追加されているが、ゴールデンデータのCSVファイルが更新または追加されておらず、analyticsプロジェクトに対応するMRが存在しない場合です。このシナリオでは、本番データに対してテストを実行するCIジョブがいくつか表示されます。これは、マイナーな変更(構文やタグなど)が機能することを確認するのに便利です。

4つ目のシナリオは、ゴールデンデータのCSVファイルが追加・更新されているのに、対応するアナリティクスMRがない場合です。 この場合、ゴールデンデータファイルがテーブルとしてデータベースに挿入されるため、本番環境とのテストは行いたくあり ません。このシナリオでは、追加のCIジョブが表示されます。倉庫のクローンを作成するもの、Analyticsプロジェクトに格納 されているdbtモデルをクローンに対して実行するもの、クローンに対してテストを実行するものがあります。

このフローチャートは、何をすべきかの大まかな目安となります。より詳細な手順については、プロジェクト内の関連するMR テンプレートの指示に従ってください。



data-testsにマージ要求が、analyticsにもマージある場合、analyticsMRはdata-testsMRのdependencyとして設定べきですつまり、data-testsのMRがマージされる前に、analyticsのMRがマージされなければなりません。

信頼できるデータダッシュボード

Trusted Data Dashboardは、データウェアハウスの健全性を迅速に評価するために使用されます。最も重要なデータは、テストのPASSまたはFAILステータスを明確に色分けした内容で、ビジネスに適したシンプルな方法で表示されます。



### スナップショット

### {:#snapshots}になります。

スナップショットは、ソーステーブルのポイントインタイムコピーを取る方法です。 **dbt**には、スナップショットがどのように機能するかについての優れたドキュメントがあります。スナップショットについての説明や定義に注意してください。

dbt snapshotによるスナップショットテーブルの作成

スナップショットの定義は、dbtプロジェクトのsnapshotsフォルダに保存されます。発見しやすいように、異なるスナップショットをデータソースごとに整理しました。

以下は、スナップショットの実装例です。

```
{% snapshot sfdc_opportunity_snapshots %}。

{{
        config(
            unique_key='id',
            strategy='timestamp',
            updated_at='systemmodstamp',
        )
    }}

SELECT *
```

```
FROM {{ source('salesforce', 'opportunity') }}。
{% endsnapshot %}。
```

### 注目すべきキーアイテム

- データベースとスキーマはdbt\_project.ymlで設定します。データベースは環境変数で、スキーマはスナップショットに設定されています。
- *常に*ソーステーブルから選択する。重複排除が必要な場合でも、下流の**dbt**モデルからの選択は失敗しやすいため、 ソーステーブルから選択する必要がある
- スナップショットはRAWで保存されているため、スキーマやテーブルへのアクセスをロールに明示的に付与
- する必要がある ファイル名は{ソース名} {ソース テーブル名} snapshotsという命名規則に従うこと 重複排
- 除作業以外のスナップショットの変換は一切行わないことデータは常に下流側でクリーンにする
- 信頼できるupdated\_atフィールドがない場合を除いて、(チェックよりも)戦略としてtimestampを使用することをお勧めします。戦略についてのドキュメントはこちらをご覧ください。

#### テストのスナップショット

スナップショットのテストは、CIジョブのspecify\_snapshotを使ってマージリクエストで行うことができます。エンジニアは Airflow を使用してローカルでテストする必要があります。これは、適切な環境変数が git ブランチに基づいて処理されるためです。テストは master ブランチで行ってはいけません。SNOWFLAKE\_SNAPSHOT\_DATABASE環境変数を設定してローカルでテストすることは推奨されません。これは本番データを上書きしてしまうので、RAWに設定してはいけません。

### スナップショットとGDPR

GDPRを理由に、データチームがSnowflakeデータウェアハウスから個人データを削除するリクエストを受けることがあります。このような削除に対応するために、dbtマクロを使用します。マクロは、削除が必要な該当データをすべてスキャンします。これはスナップショットテーブルにも適用されます。

### フレーバーは2種類。

- 1. GDPRの削除要求は、すべてのGitLabソースに適用されるため、データウェアハウスのすべてのテーブルをチェックして更新する必要があります。マクロです。
- 2. GDPRの削除要求は、GitLab.comに関連するソースのみに適用されるため、データウェアハウス内のGitlab.comに関連するテーブルのみをチェックして更新する必要があります。マクロ

2つ目の例では、新しいスナップショットモデルを作成するときや、既存のスナップショットモデルの名前を変更するときに、dbtマクロが対象となるモデルをカバーしているかどうかを確認します。Gitlab.comに関連するGDPR削除要求があった場合に、マクロのフィルタリングが該当するスナップショットテーブルに適用されるかどうかを確認します。

スナップショット・テーブルをprodデータベースで利用可能にする

スナップショットが撮影されると、それはデータソースとなり、またデータソースとして扱われるべきです。

私たちは現在、スナップショットに基づいてモデルを生成する従来の方法を採用しています。つまり、ソースモデルはありません。 スナップショットのベースモデルは、dbtプロジェクトの/models/snapshotsフォルダにあります。注意すべき重要な項目

- スナップショットのベースモデルを書く前に、sources.ymlファイルに追加することを忘れないでください(このファイルはソートしておいてください)。
- データウェアハウスのテーブル名は、当社のデータウェアハウス設計ガイドラインと一致している必要があります。理想的には{source\_name}\_{source\_table\_name}\_snapshotsを命名規則としたいと考えています。しかし、dbtはプロジェクト内での重複したファイル名を認めていません。スナップショットとスナップショットのベースモデルが同じ名前になることを避けるために、このパターンを採用しています。
  - ベースモデルファイルの名前は、ソースのスナップショットテーブルの名前に\_baseというサフィックスを付けた ものになります。例: gitlab\_dotcom\_members\_snapshotsというスナップショットファイルがあり、このスナップショットのベースモデルをgitlab\_dotcom\_members\_snapshots\_baseとします。
  - dbtconfigのalias引数を使用して、\_baseサフィックスを削除してテーブル名を変更し、テーブル名をきれいに保ちます

0

• スナップショットされたソーステーブルに基づいて構築されたベースモデルが存在する場合は、既に記述されているクエリを再利用し、以下の修正を行ってください。

- 重複排除プロセスを削除すると、必要ありません。
- odbt\_scd\_idをスナップショットのベースモデルにプライマリキーとして追加し、名前をもっと明確なものに変更してください(スナップショットのメタフィールドに関するドキュメントはこちら)。
- dbt\_valid\_fromとdbt\_valid\_toのカラムをクエリに追加します。
- スナップショットのベースモデルとソースモデルの良い例です。 スナッ

プショットの上にモデルを構築する

場合によっては、時間的制約のある変更記録ではなく、1日あたりの記録が必要になることもあります。 dbt\_valid\_from と dbt\_valid\_to です。この場合、日付スパイニングと呼ばれる技術を使って、毎日のスナップショットを持ってデルを作成することができます。

日付スパイニングでは、スナップショット・モデルが dbt\_valid\_from と dbt\_valid\_to に基づいて日付テーブルに結合されます。ここでは、zuora\_subscription\_snapshotsのソース・モデルがdbt\_valid\_fromおよびdbt\_valid\_toに基づいてdim\_datesに結合されています。この結合により、指定された日 (snapshot\_date と呼ばれる) にアクティブだったサブスクリプションのバージョンを持つ、1 日あたりのサブスクリプションごとの 1 つのレコードが得られます。

日々の記録を生成するもう一つの方法は、dbtユーティリティー関数date\_spineを使用することです。私たちは現在、日付の詳細なソースモデルを生成するためにこの関数を使用しています。

また、便利なマクロcreate\_snapshot\_baseがあり、date\_spineを利用して、任意のスナップショットテーブルから日次レコードを持つモデルを生成します。例えば、sfdc\_opportunity\_snapshots\_baseモデルの実装を見てみましょう。

スナップショット上のインクリメンタルモデル

日ごとにレコードを生成するために日付スパイニングを使用している場合は、モデルをインクリメンタルとして実体化することを検討してください。この方法では、snapshot\_dateの条件に基づいて新しいレコードのみが追加されます。 mart\_arr\_snapshotsモデルの実装例を見てみましょう。

テストのダウンストリームへの影響

dbtモデルに加えられた変更によるダウンストリームへの影響を手動でテストする方法として、sisense\_check.pyがtransform dbtプロジェクトに含まれています。このスクリプトはPeriscopeとSisenseのセーフリポジトリをチェックし、指定された dbtモデルを参照しているビュー、スニペット、ダッシュボード、チャートを探します。このスクリプトの出力を使って、個々のビュー、スニペット、チャートを手動で評価し、カラムレベルの影響を調べることができます。

チェックするモデルのリストを提供するには、dbtlistコマンドを関連する条件で使用して、必要なモデルのセットを出力し、 to check.txtというファイルにリストをエクスポートします。

dbt list -m sfdc\_account\_source+ netsuite\_transaction\_lines\_source+ --exclude date\_details dim\_date
> to\_check.txt

このスクリプトは、PeriscopeとSisense SafeのリポジトリがAnalyticsのリポジトリと同じ親ディレクトリにチェックアウトされていることを前提としています。リポジトリがチェックアウトされていない場合、スクリプトは結果を返しません。また、リポジトリが最新でない場合、結果が不完全になる可能性があります。スクリプトは入れ子になった参照をチェックします。テーブルがスニペットで参照され、それがチャートで参照されると、チェックの出力に入れ子になった参照が追加されます。ただし、このスクリプトは、次のスキーマの直接のデータベース参照のみをチェックします。

- レガシー ボンヤ
- ードコモン\*リ
- ストリクト\_セー
- フ\* ワークスペ
- ース\*

スクリプトを実行すると、チェック用に提供されたモデルを参照しているすべてのビュー、スニペット、ダッシュボード、チャートの名前を含む ${f JSON}$ 、 ${f sisense\_elements.json}$ ファイルが出力されます。

```
"mart arr":{
      "periscope-snippets":[
          "base_pricing_customer_overview.sql",
          "mrarr_base.sql",
          "nf ptb account features.sql"
      ],
       "periscope-dashboards":{
          "Investor Relations":[
              「ライセンスユーザー -- Saas/Self-
             managed」、「トータルARPUの内訳」。
             "ARRバケットによる顧客(親顧客)カウント"
          1
      },
       "sisense-safe-dashboards":{
          "TD: ARR per Licensed User (ARPU)" を参照してください。[
              "Customer ARR by Product, Sales Segment, Account Owner
             Team", "Total ARPU Breakdown",
              "ARPU。ARPU: 顧客セグメント別自己管理」、
              「ARPU: 顧客セグメント別その他」。顧客セグ
              メント別のその他
              "ARPU:合計
          ],
      }
```

出力に含まれるビューやスニペットは角括弧[]で囲まれ、チェックしたモデルと同じレベルで表示されます。

```
"[nf_ptb_account_features]" を
参照してください。{ "sisense-
safe-dashboards":{
    "NF - WIP Sandbox "です。[
        "ARR per Periods - Extra data "です。
        "SMB - 期間別ARR - エクストラデータ",
        "当期間 - PTB For Prediction"
    ]
}
```

layout: handbook-page-toc title: "Enterprise Dimensional Model"  $\subset$ 

のページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

背景

EDM(Enterprise Dimensional Model)は、GitLabの一元化されたデータモデルで、レポーティングやアナリティクスにおいて最高レベルの精度と品質を実現し、サポートするために設計されています。このデータモデルは、バスマトリックスやエンティティリレーションシップダイアグラムなど、キンボール技法に基づいています。次元モデリングは、データ開発アプローチの第3段階(要件定義、UIワイヤーフレーム作成に続く)であり、この全体的なアプローチにより、高品質なデータソリューションを繰り返し生み出すことができます。EDMはSnowflake EnterpriseDataWarehouseに格納されており、dbtを使用して生成されています。

2021年4月現在、EDMはGo-To-Marketのファネル分析を解決し、製品の使用状況分析を解決するために積極的に拡張されています。EDMを利用したSiSenseのダッシュボードの例は以下の通りです。

- TD:セールスファネル
- TD: Customer Segmentation
- TD: Drillable Net Retention
- TD: Pricing Dashoard

### 一次元モデリングの成果物

- EnterpriseBusMatrixは、すべてのFactとDimensionのテーブルを使いやすいテーブルに統合したもので、Kimballのバスマトリックスを模しています。
- エンタープライズ・エンティティ・リレーションシップ・ダイアグラムは、ファクトテーブルとディメンションテー
- ブルのエンティティレベルの統一ビューを示します。 **DimensionalModellingDevelopmentProcess**(モデリング開発プロセス)は、命名規則を含む当社のモデリング標準をカバーしています。

エンタープライズ次元モデルの「BIG PICTURE」図

• Lucidchart社のER図テンプレートを使用して、EnterpriseEntityRelationshipDiagramのソースを構築します。

### ディメンショナル・モデリングとは?

次元モデリングは、RalphKimballによって開発されたBusiness Dimensional Lifecycle方法論の一部であり、データウェアハウスの設計に使用する一連のメソッド、テクニック、コンセプトを含んでいます。

高性能なアクセスを可能にする、標準的で直感的なフレームワークでデータを表示することを目指す論理的設計手法

Dimensional Modelingは、ビジネスプロセス指向で、4つのステップで構築することができます。

- 1. ビジネスプロセスの選択例:月次収益の把握
- 2. 穀物を宣言する 例:お客様ごとに
- 3. ディメンションの確認
- 4. 事実を確認する

ファクトテーブルとディメンションテーブル

ディメンションモデリングでは、常にファクト(測定値)とディメンション(コンテキスト)の概念を使用します。ファクトは一般的に(必ずしもそうではないが)集約可能な数値であり、ディメンションはファクトを定義する階層と記述子のグループである。

merged.md

最も単純なバージョンでは、ファクトテーブルは中心的なテーブルであり、外部キーで次元テーブルにリンクされ、スタースキーマを形成している。スタースキーマは、次元テーブルがさらに次元テーブルにリンクしているもので、スノーフレークスキーマと呼ばれ、マルチファクトテーブルスキーマは銀河と呼ばれます。

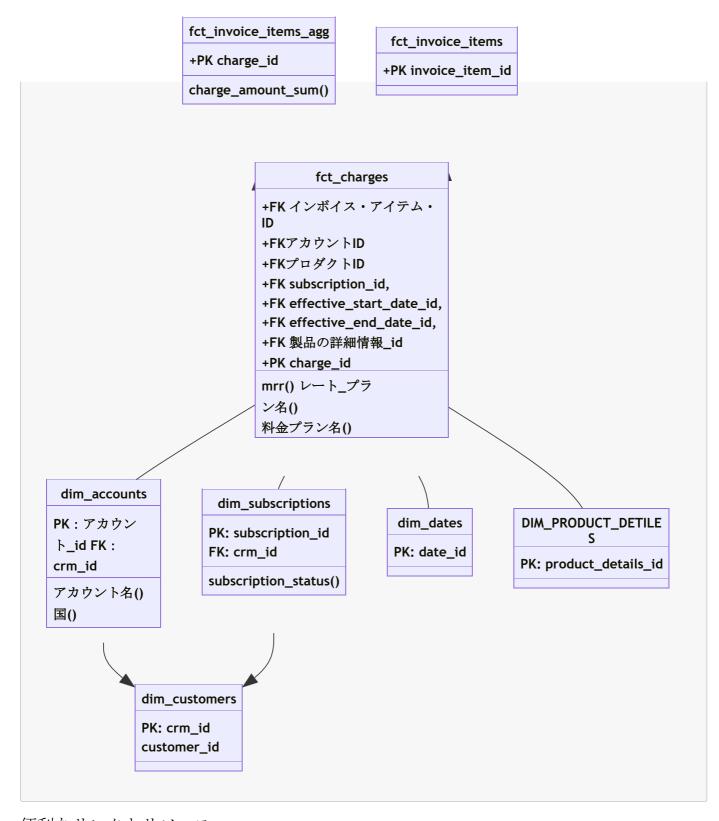
2021/11/8

なぜディメンションモデリングを使う価値があるのか

- Dimensional Modelingにはいくつかの種類がありますが、全体的なデザインは業界標準であり、何十年にもわたって使用されてきました。
- FACTとDIMの構造は、ビジネスチームが理解しやすく、アクセスしやすいデータになります。次元モデリングは、ビ
- ジネスロジックの集中的な実装と、ビジネスユーザー間での一貫した定義をサポートします。
- このデザインは、新しい主題分野の「プラグ・アンド・プレイ」をサポートしており、実際、次元を追加することでモデルの力が増していきます。

最初の反復 - ARRと顧客数の解法

2019-12月に初期イテレーションを提案し、2020-5月にARR/顧客数をサポートするモデルを展開しました。



# 便利なリンクとリソース

- dbt 最新のデータウェアハウスにおけるKimballの次元モデリングについての談話で、なぜKimballを使い続けるべきなのかという重要なアイデアを含む。
- ディメンショナル・モデリング・マニフェスト
- キンボールグループの次元モデリング技術

layout: handbook-page-toc title:" $\vec{r}$  -  $\beta$   $\vec{r}$   $\gamma$   $\beta$   $\gamma$   $\gamma$   $\gamma$   $\gamma$   $\gamma$ 

のページについて

{:.no\_toc .hidden-md .hidden-lg}。

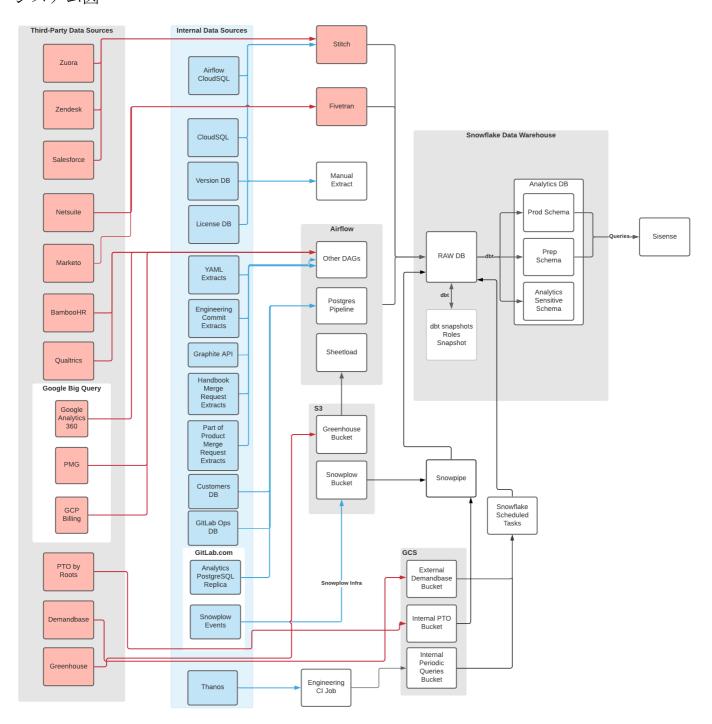
• TOC {:toc .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# クイックリンク

Airflow{:.btn .btn-purple-inv} DataImageProject{:.btn .btn-purple-inv} GitLabDataUtilsProject{:.btn .btn-purple-inv} PythonGuide{:.btn .btn-purple-inv} Meltano{:.btn .btn-purple-inv}

# システム図



### システム図の編集

システム図を編集するには、Lucidchartのドキュメントにアクセスして、必要な変更を加えます。赤枠内の変更であれば、GitLab に自動的に反映されます。表示領域を変更する必要がある場合は、Lucidchartドキュメントを使って新しいリンクを生成して公開し、画像を取り込むために上記で使用したリンクを更新するためのMRが必要です。

### エアフロー

オーケストレーションにはすべてAirflowを使用しています。

### 命名規則

DAGとタスクの名称は、エンジニアや他のユーザーが、気流UIのコンテキストの中でも外でも、DAGが何をするのかを理解できるように、できるだけ説明的であるべきです。名前には、DAGやタスクが何をどのように行っているかという情報を含める必要があります。

例として、gitlab\_com\_db\_extract DAG の DAG とタスクの名前を考えてみましょう。この DAG には gitlab-com-packages-package-files-db-incremental-pgp-extract のような名前のタスクがあります。タスクと DAG の名前の重複を避けたいところですが、タスクの名前しか見ないこともあるので(Slack や Prometheus のエラー、あるいは airflow で 'Browse' を使ったときなど)、結果的には便利です。

2021年夏以降、私たちが作成した最も一般的なDAGを、より省略して、しかし同じように説明する方法を導入しました。この規則は、それ以降に作成されたすべてのDAGで使用されるべきです。そのDAGで何が行われているかを示すために、プレフィックスを使用してください。例えば、gitlab extract DAGは抽出とロードの両方のフェーズを行うので

el\_gitlab\_com\_db\_extractという名前になりますが、 greenhouse\_extractはS3から抽出したデータをSnowflakeにロードする だけなので1\_greenhouseとなります。

### プレフィックス指標の一覧

### プリフィックスインディケータ

| eExtract   |  |
|------------|--|
| lLoad      |  |
| tTransform |  |
| pump       |  |
| mMetadata  |  |

### 製作中

すべての**DAG**はKubernetesPodOperatorを使って作成されるため、エアフローポッド自体の依存関係は最小限に抑えられ、インフラに大きな変更がない限り再起動する必要はありません。

現在のAirflowのデプロイメントでは、deployment.ymlで定義されている通り、4つのコンテナが稼働しています。

- 1. サイドカーコンテナは、レポのアクティビティフィードをチェックして、masterへのマージがないかを確認します。マージがあった場合、サイドカーはレポをリクローンし、Airflowが最も新鮮なDAGを実行するようにします。
- 2. スケジューラ「Airflow」について
- 3. Airflowウェブサーバー
- 4. Airflow が私たちの cloudsql インスタンスに接続するための cloudsql プ

ロキシ Kubernetes セットアップ

Google Cloud Platform (GCP) のgitlab-analysisプロジェクトで実行しています。Airflowはdata-opsクラスタで動作しています。

ノードプール

このクラスタ内には、highmem-pool、production-task-pool、testing-pool、sdc-1という4つのノードプールがあります。各ノードプールには、監視とリソース管理を容易にするために専用の用途があります。

- 1. highmem-pool Airflowサーバー、スケジューラー、ネットワークコンポーネントの実行に使用されます。1-2ノードからオートスケールが可能です。
- 2. production-task-pool SCDタスクを除く、ほとんどのプロダクションAirflowタスクの実行に使用されます。2~5ノードでオートスケールします。
- 3. sdc-1 SCD抽出の実行に使用されます。1~3ノードでオートスケールします。
- **4.** testing-pool 通常は稼働しているノードを持たないプールですが、エンジニアがローカルに起動したAirflowタスクを実行するために使用されます。1-2ノードからオートスケールが可能です。

highmem-poolを除くすべてのnodepoolには、どのnodepoolがどのAirflowタスクを起動するかを管理するためのラベルとテイントが

あります。タスクがノードプールでスケジュールされるためには、タスクはプールに対してnodeAffinityを持っていなければならず、また、以下の条件にマッチするtolerationを持っていなければなりません。

ティングを使用しています。postgresパイプラインジョブのSlowly-Changing Dimensionsタスクにaffinityとtolerationを追加したこのMRをご覧ください。

#### 名前空間の作成

クラスタが作成されると、デフォルトではdefaultという名前のネームスペースが1つ関連付けられます。このプロジェクトのエアフロー設定では、本番環境で2つのネームスペースを設定する必要があります。

- 1. デフォルト
- 2. テスト

名前空間:-Kubernetesでは、名前空間は、単一のクラスタ内でリソースのグループを分離するためのメカニズムを提供します。リソースの名前は、ネームスペース内では一意である必要がありますが、ネームスペース間では一意ではありません。

クラスタのプロビジョニングが完了したら、名前空間と秘密のファイルを設定します。names\_space\_testing.yamlには、名前空間の詳細に関する情報が含まれています。

airflow\_image/manifests/names\_space\_testing.yamlというフォルダからkubectl apply -f name\_space\_testing.yamlというコマンドを実行します。すると、以下のような出力が得られるはずです。 namespace/testingが作成されました。

作成/存在する名前空間のリストを確認するには、kubectl get namespaceを使用します。

kubeのシークレットの作成

名前空間の作成後、名前空間 default と testing の両方にエアフローシークレットを作成します。 default 名前空間での秘密 の作成には kube\_secret\_default.yaml を使用します。 airflow\_image/manifests ディレクトリから以下のコマンド kubectl create -f kube\_secret\_default.yaml を実行すると、以下のような出力が得られます。

シークレット/エアフロー作成

秘密が正常に作成されたかどうかを確認するには、以下のコマンドを使用します。以下のコマンドを使用します。 kubectl get secrets --namespace=default .

テスト用ネームスペースでのエアフローシークレットの作成」と同じ手順で行います。テスト中のシークレットの作成については namespace use kube\_secret\_testing.yaml.

以下のコマンドを実行してください。kubectl create -f kube\_secret\_testing.yaml airflow\_image/manifests/.

秘密ファイルのすべての値は、ドキュメントとして1パスワードデータチームセキュアボールトに存在し、そのドキュメントの1から秘密ファイルに値を追加します。ドキュメントをtestingの名前空間で編集するには、kubectl edit secret airflow -o yaml --namespace=testingというコマンドを使い、デフォルトで作成するには、kubectl edit secret airflow -o yamlというコマンドを使います。注:このコマンドは、deployment.ymlファイルを適用する前に実行して作成する必要があります。

# DNS

airflow.gitlabdata.comというURLが私たちのクラスタを指すようにするために、gcloud compute addresses create airflow-west --region=us-west1というコマンドを使って、us-west1というリージョンに固定IPをプロビジョニングしました。生成されたIPは35.233.169.210でした。これは、gcloud compute addresses listを実行することで確認できます。TLSステップが機能するためには、スタティックIPはグローバルIPではなくリージョナルIPでなければならないことに注意してください。

これは、インフラストラクチャチームによってRoute 53のドメインにマッピングされ、この問題で記録されました。

TLS

このインフラストラクチャーの問題により、airflow.gitlabdata.com ドメイン用の証明書が生成されました。この証明書は、kubectl create secret tls airflow-tls -- cert=airflow.gitlabdata.com.chained.crt -- key=airflow.gitlabdata.com.key というコマンドを実行することで、kubernetes のシークレットに保存されました - これら

は、証明書の実際のファイル名です。

証明書と鍵を連結したものです。これにより、秘密のairflow-tlsが作成されました。証明書ファイル(site、chain、chained(site+chain)、key)も1passwordのData Team Secure vaultに保存されています。

NGINX IngressControllerには、リダイレクトやTLSの強化のための優れた組み込み機能があるので、これを使うことにしました。 NGINXをクラスタにインストールするには、以下の手順に従います。

- 1. helmがシステムにインストールされていない場合は、brew install helmコマンドでインストールします。
- 2. その後、コマンド helm repo add nginx-stable https://helm.nginx.com/stable を使って、nginx-stable バージョンを helm repo に追加します。
- 3. 安定版の最新バージョンを取得するには、helm repo updateを使用します。これにより、以下のような出力が得られるはずです。

あなたのチャートリポジトリから最新の情報を取得するまでお待ちください。 ...チャートリポジトリ「nginx-stable」からのアップデートに成功しました。 ...Successfully got an update from "ingress-nginx" chart repository Update Complete.\*Happy Helming!\*

4. helm repo の状態を確認するには、コマンド helm repo list を使用します。

NAMEURL

ingress-nginx https://kubernetes.github.io/ingress-nginx

nginx-

stablehttps://helm.nginx.com/stable`.

5. NGINX Ingress Controller helm install airflownginx nginx-stable/nginx-ingress -- values nginx values.yaml 以下のような出力が得られます。

名前:エアフロージンクス

LAST DEPLOYED: Wed Oct 20 19:01:38 2021

NAMESPACE: default STATUS: deployed REVISION: 1 テストスイート: なし NOTES

NGINX Ingress Controllerのインストールが完了しました。

NGINXのバリューファイルでは、ロードバランサーのIPアドレスを何にするかを定義しています。ロードバランサーのIPには、前のセクションで生成されたアドレスが設定されます。installコマンドに渡された値は、controller-deployment.yamlファイルで展開されます。

NGINXを削除して再インストールする必要がある場合は、helm delete airflownginxで行うこ

とができます。また、ingressの定義もこれらの設定で更新されました。

- force-ssl-redirectはtrue
- tlsはairflow.gitlabdata.comに設定されており、kubernetesのsecretであるairflow-tlsを使用しています。
- ホストにマッチするすべてのトラフィックは、デフォルトのバックエンドではなく、airflow-webserverサービスに送られます。

厳密には必要ではありませんが、変更を適用する際にはingressを削除した方がすっきりします。これは、GCPのUIまたは kubectl delete ingress airflow-ingressコマンドで行うことができます。新しい設定を適用するには、kubectl apply -f ingress.yaml というコマンドを使用します。

GCPでエアフローのログを見る

Airflow のあらゆる種類のログは、GoogleCloudのLogsExplorer で確認できます。このアプリケーションは、収集されるさまざまなログに少し圧倒されることがありますので、いくつかのヒントをご紹介します。

• ログフィールドビューが開いていることを確認します。このビューは、ウェブサーバや特定のDAGの実行など、関心のあるものだけにログをフィルタリングするのに重要です。このビューを開くには、右上の「PAGE LAYOUT」をクリックし、「Log fields」がチェックされていることを確認してください。Log fieldsビューは、実際のログのすぐ左に表示されます。

- Log fields ビューを使って、特定のプロセスからのログにフィルタリングすることができます。container\_nameは、非常に有効なフィルターです。コンテナのほとんどは、エアフローの展開自体に関係しています。特定のDAGランを探している場合は、コンテナ名はすべてのDAGで共有されるベースであるため、より多くのフィルタリングが必要になります。ポッド名は各DAGで一意であるため、特定のDAGランのログを調査する場合は、ポッド名をフィルタリングすることができます。
- 閲覧中のログのタイムフレームを変更するには、現在のタイムフレームが表示された時計のアイコンがあります。このボタンをクリックすると、さまざまなオプションに基づいて新しいタイムフレームを表示することができます。

GCP IAM

### {:#gcp-iam}。

エンジニアには、GCPの以下の権限を持つことが推奨されています: Cloud

- SQL Admin
- Kubernetes Engine Admin
- Storage Admin

Airflow をローカルで開発するためには、エンジニアはサービスアカウントも必要です。これらの認証情報は、あなたのマネージャーが用意したサービスアカウントを示すものでなければなりません。アカウントは、あなたの電子メールと同じパターンでなければなりません。例えば、tmurphy@gitlab- analysis.iam.gserviceaccount.com です。推奨されるパーミッションは以下の通りです。

- \* クラウド SOL クライアント
- \* Kubernetes Engine Developer
- \* Storage Object Creator
- \* ストレージオブジェクトビューア

サービスアカウントキーは、JSONとしてダウンロードし、ユーザーのコンピュータに安全に保存するためにユーザ

一に送信する必要があります。サービスアカウントのプロビジョニング

- gitlab-analysisGCPプロジェクトのServiceAccountセクションに移動して
- Create Service Account をクリックします。
- 個人ユーザーの場合、名前はファースト・イニシャル・ラスト・ネームのパターンに
- 従う特定の人のためのサービスアカウントであることを示す
- サービスアカウントに関連する権限を付与する サ
- ービスアカウントをユーザーに付与する
- サービスアカウントのJSONキーを生成し、ユーザーに送信する

# 失敗したジョブの処理

一度に複数の失敗したDAGが表示されることはありません。gitlab.comデータベースからの抽出のようにexecution\_dateに依存するインクリメンタルジョブの場合、失敗したDAGはタスクインスタンスをクリアして、修正が適用された後に再実行できるようにする必要があります。

execution\_dateに依存しないジョブについては、修正が適用された時点でジョブを手動で再実行し、失敗したDAGrun(s)を削除する必要があります。あるDAGに対して失敗したDAGrunがある場合、それはそのDAGの現在の状態が壊れており、修正する必要があることを意味します。

これにより、AirflowのDAGリストを見て、注意が必要なものとそうでないものがすぐにわかるようになります。バックフィル

あるDAGで増分実行ができなかった場合や、テーブルにデータがない場合、バックフィルを行うには2つの方法があります。テーブルが小さく、バックフィルが比較的短時間で済む場合は、テーブルをドロップして完全同期を行う方法があります。しかし、アップストリームのエラーでDAGが停止した場合など、テーブル数が多い場合には、この方法が取れない場合があります。

後者の場合は、Airflow Schedulerのポッドコンテナでbackfillコマンドを実行するのが良いでしょう。そのコマンドは

airflow backfill gitlab\_com\_db\_extract -s 2019-10-30T00:00:00 -e 2019-11-04T12:00:00 --delay\_on\_limit 30 -- reset dagruns

これにより、指定された時間枠に既に存在するDAGrunとタスクインスタンスがクリアされる一方で、その時間枠に存在しない新しいDAGrunが生成されます。フラグが何であるかは、CLIのAirflowドキュメントで詳しく説明されています。

--task\_regex フラグを使用して、トリガーされるタスクをフィルタリングすることを検討してください。これは、gitlab.comのインクリメンタルロードの場合に、"pgp-extract"を使用して下流のdbtタスクをスキップするのに便利です。

バックフィルされるタイムフレームのDAGランがすでに存在している場合、上記のコマンドを実行すると、タスクが同時に実行されます。DAGを連続して実行する必要がある場合。

- まず、DAGのスケジューリングをオフにします。
- 次に、Airflow UIに入り、「Browse (ブラウズ)」から「DAG runs (DAGラン)」をクリックします。次に、検索バーを使って、バックフィルされるDAGにフィルタリングします。バックフィルされるDAGランをすべて選択し、「Actions --> Delete」をクリックします。
- DAGランが削除された今、エアフローはタスクを連続して実行します。しかし、タスクが削除されていないので、エアフローは実際にはタスクを実行しないでしょう。そこで、タスクをクリアするために、Browse --> Task Instances に行き、検索バーとソートを使って、再度実行する必要のあるタスクインスタンスを選択します。それらのタスクインスタンスを選択した状態で、Actions --> Clearに進み、それらのタスクインスタンスの状態をクリアします。
- この問題が解決したら、DAGに戻ってスケジューリングを再開します。
- その後、上記のbackfillコマンドを実行するが、リセットするdagrunsがないので、--reset\_dagrunsは省略する。これでbackfillが連続して実行されるはずです。

## In Merge Requests

ローカルでのAirflowの使用を容易にしつつ、KubernetesでのDAGの適切な実行をテストするために、docker-composeを使用してローカルのAirflowインスタンスをスピンアップし、KubernetesPodOperatorを使用してKubernetesでのDAGの実行を可能にしています。Dockerの項を参照して、適切な環境変数を設定してください。

Airflowでのコード変更からテストまでの流れは、以下のようになります(これは、そのタスクのDAGがすでにあることを前提としています)。

- 1. コードをコミットして、リモートブランチにプッシュします。
- 2. make init-airflow を実行して、postgres db コンテナをスピンアップし、Airflow テーブルを初期化します。Dockerが 起動していない場合は、エラーが発生します。
- 3. make airflowを実行してAirflowを立ち上げ、コンテナの1つにシェルを取り付ける
- **4. Web**ブラウザを開き、localhost:8080にナビゲートすると、あなた自身のローカル**Web**サーバが表示されます。一般的な**Admin**ユーザーは、ユーザー名とパスワードをadminに設定して、**MR** airflowのインスタンスに自動的に作成されます。
- 5. airflow シェルで、テストしたい DAG/Task をトリガーするコマンドを実行します。例えば、airflow run snowflake\_load snowflake-load 2019-01-01 (docker-compose ファイルで設定されているように、すべての kube ポッド は testing 名前空間に作成されます)。また、DAG全体(例えば、分岐ロジックをテストするためのdbt DAG)を実行したい 場合は、airflow backfill dbt -s 2019-01-01T00:00:00 -e 2019-01-01T00:00:00のようなコマンドになります。
- 6. ジョブが終了すると、DAG/Taskインスタンスにナビゲートしてログを確認することができます。

また、どのようなコマンドがあり、何をするのかを説明するmake helpコマンドもあります。いく

つかのゲッチュ。

- 最新バージョンのDockerを使用していることを確認してください。これにより、ERRORのようなエラーを防ぐことができます。./docker- compose.yml "のバージョンはサポートされていません。
- dagで新しいpythonスクリプトを呼び出す場合は、chmod +xを実行してファイルが実行可能であることを確認してください。 your python file.pyを使用します。これにより、パーミッション拒否のエラーを避けることができます。
- dag で追加された新しいシークレットが kube\_secrets.py にもあることを確認してください。これは、Airflow がどの秘密を使用するかの真実の源です。実際のシークレットの値はこのファイルには保存されず、ポインターだけが保存されます。
- ◆ イメージが古くなっている場合は、docker pull <イメージ名>というコマンドを使って、最新のイメージを強制的にプルし

ます。次のようなエラーが発生した場合「could not find an available, non-overlapping IPv4 address pool among the defaults to assign to the network」のようなエラーが出る場合は、実行中のVPNをオフにしてみてください。

ビデオウォークスルー

- エアフロー pt1
- エアフロー pt2
- Airflowのテスト環境。のAirflowとポッドオペレータ (Kubernetes) のテスト環境を動画で解説しました。 GCPです。

ローカルエアフローの設定。

すべてのDAGはKubernetesPodOperatorを使用して作成されるため、ローカルで作業する際には、ローカルでタスクを実行する際にポッドをスピンさせることができるクラスタが必要になります。これを実現するためには、クラスターに接続する必要があります。クラスタに接続するには gcloud container clusters get-credentials data-ops --zone us-west1-a -- project gitlab-analysis を使用します。cluset でローカルセットアップを動作させるためには、we name testing という名前の名前空間を作成し、testing の名前空間に airflow secret が存在することを確認する必要があります。これを作成するための手順は上記のとおりです。これらの設定が完了すると、ローカルセットアップがクラスター内のポッドをスピンできるようになります。

ローカルエアフロー設定のトラブルシューティング

No Such File or Directory: 'Users/(user)/google-cloud-sdk-bin/gcloud'

FileNotFoundErrorです。[Errno 2] No such file or directory: '/Users/(user)/google-cloud-sdk/bin/gcloud': '/Users/(user)/google-cloud-sdk/bin/gcloud'.

これは、MacでのGoogle Cloud SDKのデフォルトのインストール先が上記のようになったためですが、linuxやコンテナでは/usr/lib/google-cloud-sdk/bin/gcloudにインストールされます。この値は、コンテナに渡される/.kube/configファイルです。

このエラーを修正するには、/.kube/config を編集して、コマンドパスのパラメータをコンテナ内の場所に更新するだけです。
/usr/lib/google-cloud-sdk/bin/gcloud このファイルは、SDKをインストールしたり、次のコマンドを実行したりするたびに更
新されます: gcloud container clusters get-credentials data-ops.詳細については、関連する問題を参照してください。

プロジェクト変数

現在の実装では、以下のプロジェクト変数を使用しています。

- snowflake\_account
- snowflake\_report\_warehouse
- snowflake\_{flavor}\_user
- snowflake\_{flavor}\_password
- snowflake\_{flavor}\_database
- snowflake\_{flavor}\_role
- snowflake\_{flavor}\_warehouse

以下のようなフレーバーが定義されています。

- LOADフレーバーは、Extract & Loadプロセスで使用さ
- れます TRANSFORMフレーバーは、Transformプロセス
- で使用されます TESTフレーバーは、Snowflakeを使
- 用したテスト用です PERMISSIONフレーバーは、許可 ボット用です
- SYSADMINフレーバーは、ハウスキーピング・タスク(レビュー・インスタンスの設定など)のためのものです。このフレーバーでは SNOWFLAKE\_SYSADMIN\_DATABASE」と「SNOWFLAKE\_SYSADMIN\_WAREHOUSE」です。

以下の変数は、実行環境に応じてジョブレベルで設定されるため、プロジェクトの設定では設定しないでください。

- snowflake\_user
- snowflake\_password
- snowflake\_role
- snowflake database

 $snowflake\_warehouse$ 

エアフローのモニタリングとアラート

### モニタリング

- エアフローの指標値はThanosで探索することができます。
- Airflow は社内の Prometheus クラスタによって監視されています。 Prometheus が Airflow を監視するには、次のような GET リクエストを行います。

/admin/metricsエンドポイントを定期的に使用しています。

- メトリクスエンドポイントは、**GitLab Airflow**の**Docker**イメージにパッケージされている**airflow-prometheus-export**をカスタマイズして利用可能にしたメトリクスを公開します。
- airflow-prometheus-exporterパッケージは、airflowサーバー、タスク、DAGに関する様々なメトリクスを取得するようにあらかじめ設定されています。また、このパッケージは、XCOMと呼ばれるものからメトリクスを引き出すこともできます。タスクは、XComを使って相互にデータをやり取りすることができます。タスクから XCom が返されると、それは Airflow データベースに保存されます。prometheus exporterは、Airflowデータベースを読み込むので、XComの値をメトリクスとして使用することができます。
- KubernetesPodOperatorがXComを返すようにするには、まずタスクがこのようにXComの受け渡しを有効にしている必要があります。私たちが使っているAirflowのバージョンは、XComを有効にする方法のバージョンが混在しているため、このようにXComパッシングを有効にするためのブーリアンが2種類ある場合があります。KubernetesPodOperatorでXComを有効にした後、実行するコードでは特定のファイルにjsonオブジェクトを書き込む必要があります。これを簡単にし、ファイルの場所を抽象化するために、gitlabdata.orchestration\_utils.push\_to\_xcom\_file が作成されました。この関数は、JSON オブジェクトを受け取り、それを XCom ファイルに書き込みます。この関数は、タスクごとに一度だけ呼び出す必要があります。XCom内の値をメトリックとして使用するには、メトリックがXComファイルに書き込まれたJSONオブジェクトのファーストクラスのメンバーである必要があります。たとえば、以下のようになります。

{"record\_count":5, "other\_record\_count":6}であれば、record\_countとother\_record\_count as metrics.

• prometheus-exporter で XCom をメトリックとして使用できるようにするには、この設定にタスク名とメトリック名を追加するだけです。すべての気流タスクにインポートされるべきメトリックがある場合は、タスク名 all を使用できます。設定を変更した後、メトリックを表示させるには、「デプロイメントとポッドの再起動」で説明したように、エアフローイメージを再構築して再デプロイする必要があります。

## アラート

- 上記のメトリクスに基づくアラートは、runbooks リポジトリの rules/airflow.ymlファイルに定義されています。新しいアラートは、ルールリストに追加することで追加できます。expr要素はPromQL式で、アラートが発生した場合は1を返し、それ以外の場合は0を返す必要があります。for要素は、アラートが実際にトリガーされる前に、どのくらいの時間、式が1と評価されなければならないかを定義します。ラベルは重大度を含みます。重大度は現在GitLabシステム全体で定義されているので、待機中のGitLabインフラメンバーが注意を払わない限り、Airflowルールの重大度は低くすべきです。チームラベルは、最終的にどの slack チャンネルがアラートを受け取るかを決定します。チームラベルは以下のように設定します。data-analyticsはランブックのサービスカタログで設定されたチーム名であるため、データチームのすべてのアラートにはdata-analyticsを使用します。
- データ分析チームに割り当てられたPrometheusアラートは、data-prom-alerts Slackチャンネルに送信され、トリアージを担当するチームメンバーが調査・対処する必要があります。

AirflowとKubernetesの共通タスク

## ヒント

◆ kubectlをkbcとエイリアスすることをお勧めします。

## エアフローの問題を解決する

時には、何が起こっているのかわからないような故障が起こることがあります。ここでは、よくあるチェック項目をご紹介します。

- ◆ GitLab.comとGoogleCloudの状況確認
- 分析プロジェクトやデータイメージプロジェクトへのマージ要求が最近あったかどうかを
- 確認する コンテナログを確認する
- GKEダッシュボード 現状の把握に役立つ
- GCPエラー報告 Airflowやプロジェクト全体でどのようなエラーが発生しているかを判断するのに非常に役立ちま

す。

• **GCP**のエラーレポートは、**2020年2**月にクラスタがダウンしたときに役立ちました。エラーは、アプリケーション がストレージ不足であることを報告していました。これにより、「Managing Persistent Volume Claim」で説明されている修正が行われました。

- GCPのKubernetesセクションにある「Workloads」と「Services & Ingress」セクションを確認します。警告やエラーが出ていないか確認してください。
- kubectl get podsを実行して、airflow-deploymentというラベルが付いたものが戻ってくるかどうかを確認します。ポッドに入ってみる
- kubectl get pods -Aを実行し、名前空間のボトルネックの原因となっている古いポッドや陳腐化したポッドがないかどうかを再確認します。

## セミオーソドックスな問題

- ログに<task> had an event of type DPendingDと表示され、タスクが開始されません。
  - 。 k8sがレジストリからコンテナを取得できないことが原因の可能性があります。最近の事例については、この問題
  - を参照してください。 GKEのリソース不足が原因で発生する可能性があります。

### Kubernetes Airflow Clusterへの接続。

- 1. Kubectlのインストール
- 2. gcloud container clusters get-credentials data-ops --zone us-west1-a --project gitlab-analysis を 実行して、データチームのクラスタに接続します。
- 3. kubectl get podsを実行して、正常に返ってくることを確認します。
- 4. すべてのコマンドは本番環境に接続されるため、現在Kubernetesにはテスト環境がありません。典型的なテスト方法は、ローカルのdocker-composeセットアップを使用することです。

# Airflow Webserver UIへのアクセス

• kubectl port-forward deployment/airflow-deployment 1234:8080を実行します。これで、ブラウザでlocalhost:1234 に移動すると、ポートフォワードしたインスタンスのウェブサーバに移動します。注:安定したURLにアクセスできるようになったので、もうこの作業は必要ありません。

# エアフローの更新

- 主要なステークホルダーにアップグレードのタイミングを知らせるイシューを用意する。アップグレードを実行する前に、少なくとも1日前に告知すること
- airflow\_image/Dockerfileにあるバージョンをバンプさせます。そのラインはARG AIRFLOW\_VERSION=<version\_number>
- のようになります。この手順の具体的な方法については、以下の「デプロイメントとポッドの再起動」セクションを参照してください。
- ポッド内のいずれかのコンテナに実行し、エアフローのアップグレードを行います。

# リソースを見る

- kubectl get allを実行します。これで、すべてのポッド、デプロイメント、レプリカセットなどが表示されます。
- kubectl get podsコマンドで、現在の名前空間にあるすべてのポッドのリストを見ることができます。

### パーシステント・ボリュームの表示

• パーシステントボリュームやパーシステントボリュームクレーム (ログが保存されている場所) の一覧を見るには、それぞれ kubectl get pvとkubectl get pvcというコマンドを使います。パーシステントボリュームは名前空間に属していないため、 get persistent volumesのコマンドは名前空間に関係なくすべてのボリュームを表示します。しかし、パーシステントボリュームの主張は特定の名前空間に属しているため、現在のコンテキストの名前空間内のもののみを表示します。

# デプロイメントとポッドの再起動

**kubernetes**のリソースマニフェストは、airflow-image/manifests/に格納されています。**kubernetes**でこれらのリソースを作成または更新するには、まずkubectl delete deployment airflow-deploymentを実行し、次にkubectl apply -f <manifest-file.yaml>を実行します。

一度に1つのポッドにしか要求できない永続的なボリュームを使用しているため、通常のkubectl apply -fによる修正はできません。毎回、新しいデプロイメントを設定する必要があります。

• 例えば、Airflowのロックアップや再起動の継続、コンテナが使用しているAirflowイメージの更新など、Podを強制的に 再起動する必要がある場合は、kubectl delete deployment airflow-deploymentを実行します。これにより、すべて のポッドが消去されます (airflowによって実行されているものも含まれるので、注意が必要です)。その後、データイメージ リポジトリのルートフォルダからkubectl apply -f airflow\_image/manifests/deployment.yamlを実行すると、マニフェストがKubernetesに送り返され、ポッドが再起動します。

ポッド付きアクセスシェル

### {:#access-pod}になります。

- **kube**ポッド内に存在するシェルに入るには、**kubectl exec -ti <pod-name> -c <container-name>**というコマンドを使います。 /bin/bashです。これにより、選択したポッドやコンテナ内のシェルに入ることができます。この機能は、ウェブサー バの**UI**ではなく、シェルで直接エアフローのコマンドを実行したい場合に便利です。
  - kubectl exec -ti airflow-deployment-56658758-ssswj -c scheduler /bin/bash そのポッドとschedulerというコンテナにアクセスするコマンドの例です。コンテナ名は、airflow\_image/manifests/deployment.yaml に記載されています。この情報は、kubectl describe pod を実行した場合にも得られます。
     <pod-name>さんは、読みにくくなったと思ったようです。
    - 補足:リソース/名前の形式で引数を渡す場合、リソースタイプを別の引数として指定する必要はありません (例: 'kubectl get resource/<resource\_name>'ではなく'kubectl get resource resource\_name>'
  - また、ローカルでコンテナを実行する際には、docker composeコマンドを使って行うこともできます。
  - 起動している既存のコンテナに入るには
    - docker-compose exec airflow\_scheduler bash
  - スピンアップしてバシッと決めること。
    - docker-compose run airflow\_scheduler bash
- シェルに入ってからやることダグで特定の
  - タスクを起動する。
    - テンプレート: airflow run <dag> <task name> <execution date> -f -A
    - 具体的な例: Airflow run dbt dbt-full-refresh 05-02T15:52:00+00:00 -f -A
    - また、-A フラグは依存関係を無視します (つまり、上流にブランチされていなくても気にしません)。

## 秘密の更新

- prod環境でsecretを更新する最も簡単な方法は、kubectl edit secret airflow -o yamlというコマンドを使うことです。このコマンドを使うと、テキストエディタでsecretが開かれ、そこから編集することができます。新しいシークレットはbase64エンコードされている必要があり、最も簡単な方法はecho -n <secret> | base64 -です。シークレットファイルを編集する際に、いくつかのヌル値があります。ファイルを正常に保存するためには、ヌル値を"" に変更しなければなりません。
- 新しいシークレットを追加する際には、そのシークレットを適切な値でテスト環境に追加することも忘れてはいけません。そのためのコマンドがkubectl edit secret airflow -o yaml --namespace testingです。このコマンドは、上記の本番環境用のガイドラインと同じです。テスト環境に新しいシークレットを追加しないと、そのシークレットを使用するDAGはテスト時に実行されません。

### 稼働中のDAGを停止する

- 問題のダグのグラフビューに移動する グラフビュ
- ーでタスクを選択する
- ポップアップしたモーダルで、「Mark Failed」または「Mark Success」のいずれかを選択し、「
- **Downstream**」オプションを選択します。**Kubernetesworkloads**」タブで、該当するポッドが停止していることを確認します。必要に応じて削除します。

## DAG内で特定のタスクを実行する

• kubectl exec -ti <pod\_name> -c <webserver|scheduler> というコマンドで、コンテナ内で動作しているシェルを取得します。

/bin/bash

merged.md

2021/11/8

• 最も単純なケースでは、Airflow run <dag\_id> <task\_id> <execution\_date>で十分です。

• dbt full-refreshを実行する必要がある場合など、より複雑なケースでは、さらにいくつかのフラグが必要です。 airflow run dbt dbt-full-refresh <execution\_date> -f -A.f フラグは、すでに成功または失敗とマークされていても、タスクを強制的に実行します。A フラグは、持つべき依存関係に関係なく実行するように指示します。

## エアフロー変数の設定・変更

- 時折、Airflow 変数を使用することがあります。現在、BambooHR データ抽出時のデータ整合性チェックをオーバーライドするために、Airflow 変数を使用しています。Airflow の変数は Airflow のデータベースに格納されており、Airflow で実行されるコードから簡単にアクセスすることができます。
- 変数の設定や値の変更を行うには、Airflow UIに入ります。 "Admin" タブをクリックし、"Variables "をクリックします。新しい変数を追加するには、+ボタンを押します。必要なキーと値を入力します。値を暗号化する場合は、ボックスをチェックします。その後、保存を押します。変数を編集するには、対象となる変数の横にある編集アイコンをクリックし、必要な内容を変更します。その後、保存ボタンを押します。
- BambooHRの整合性チェックのバイパスでは、変数のキーをBAMBOOHR\_SKIP\_TESTとします。値は、一時的にテストをスキップしたいテーブルの名前をコンマで区切ったリストにします。BambooHR の抽出が成功したら、この変数を削除するか、値を空の文字列に変更してください。

パーシステントボリュームクレームとそのデータの管理

### {:#managing-pvc}になります。

Airflow の persistent volume claim は、persistent\_volume.yaml マニフェストで定義されます。これは、ログ用のデプロイメントで使用され、次のディレクトリに書き込まれます。 /usr/local/airflow/logs ディレクトリに書き込みます。 パーシステントボリュームクレームがいっぱいになった場合、2つの解決策があります。

- クレームに利用できる容量を増やすクレー
- ムからデータを削除する

### データの削除

- airflow-deployment podに入って、logs
- ディレクトリに移動する。
- ファイルの削除

### クレームサイズの拡大

- オブジェクトブラウザを開き、core > PersistentVolumeでDataOpsクラスタのクレームを見つけます。
- PersistentVolumeをクリックします pvc < GUID > である必要があります。
- そのボリュームのYAMLを見て、gcePersistentDiskの下にあるpdNameを見つけます。これが、アップデートが必要なGCEディスクです。
- GCPのComputeDisksセクションに行き、参照されているディスクを
- 見つけます。 そのディスクをクリックして、サイズを更新するために 編集します。
- PersistentVolumeのYAMLに戻り、spec:capcacityの下にあるstorageを更新します。
- これでクレームサイズが大きくなります。

または、プロジェクト内のpersistent\_volume.yaml定義を更新することもできます。ただし、これを再デプロイすると、すでにクレームにあるデータが削除される*可能性があります。*これはまだテストされていません。

# Postgres パイプライン

## 開発

postgres pipelineによって引き出される新しいテーブルやフィールドを追加するには、特定のソースデータベースのマニフェストファイルを変更する必要があります。これらのマニフェストファイルは、ここのフォルダにあります。新しいテーブルを追加するには、テーブルリストに項目を追加します。import\_db の値は、そのマニフェストの他の項目と一致する必要があります。import\_query はソース・ターゲット・データベースに対して直接実行されます。

可能であれば、次のようなWHERE句を追加して、インポートクエリをインクリメンタルにしてください。

WHERE updated\_at BETWEEN '{EXECUTION\_DATE}'::timestamp - interval '{HOURS} hours' and '{EXECUTION\_DATE}'::timestamp

export\_schema の値は、マニフェストの他の項目と一致する必要があります。export\_table値は、インポートされるテーブルの名前と一致する必要があります。export\_table\_primary\_keyは、postgresデータベース内の特定のテーブルの主キーに設定する必要があります。インポートクエリがインクリメンタルでない場合は、データがいつロードされたかを区別する別のカラムを作成するadvanced\_metadata: true を追加します。

マニフェストの変更が完了したら、add\_manifest\_tables テンプレートを使用してマージリクエストを作成します。その後、そこでの指示に従います。

技術的な実装の詳細については、こちらのREADMEをご覧ください。

エアフローでのテスト

Airflowでローカルにpostgres pipeline(pgp)をテストする場合、いくつか注意すべき点があります。

• プールは、どのpgp DAGにも存在する必要があります。プールを追加するには、[Admin] > [Pools]に進み、https://airflow.gitlabdata.com/pool/list/ に記載されているのと同じプールを追加します。

# Bashスクリプト

• データチームでは、基本的なコマンドラインの使い方と高度なコマンドラインのに関するトレーニングを実施しています。このトレーニングには、アナリティクス・レポ内での高度なbashスクリプトの使用に関する詳細も含まれています。

ドッカー

コンポーズ

私たちはDocker composeを使って、さまざまなイメージやアプリケーションを定義し、実行しています。これらはユーザーがMakefileを介して起動します。

## 環境変数

MakefileとDocker Composeを適切に使用するために、いくつかの環境変数が必要になります。

- GIT\_BRANCH 通常、2112-my-featureのような機能ブランチです。
- GOOGLE\_APPLICATION\_CREDENTIALS GCPに接続するためのCloudSQL認証情報。通常は、JSONファイルを指します。 GOOGLE\_APPLICATION\_CREDENTIALS="/Users/tmurphy/Projects/GitLab/gcloud\_service\_creds.json"
  - これらの認証情報は、GCP IAMセクションの指示に従ってプロビジョニングする必要があります。

イメージ

### データ

data\_imageディレクトリには、data-imageの構築とプッシュに必要なものがすべて入っています。バイナリをインストールする必要がある場合は、Dockerfileの中で直接行います。pythonパッケージはrequirements.txtファイルに追加し、動作確認済みのバージョンに固定します。

エアフロー

airflow\_imageディレクトリには、airflow-imageだけでなく、対応するKubernetesのデプロイメントマニフェストをビルドしてプッシュするために必要なものがすべて含まれています。新鮮なデプロイメントのために必要な唯一の手動作業は、airflowのシークレットを設定することです。必要なシークレットは、airflow\_image/manifests/secret.template.yamlに記載されています。

は

airflow\_testing db.

デフォルトのインスタンスのログは gs://gitlab-airflow/prod に、テスト用のインスタンスのログは gs://gitlab-airflow/testing に保存されます。

#### dbt

dbt\_imageディレクトリには、データイメージのビルドとプッシュに必要なものがすべて入っています。バイナリをインストールする必要がある場合は、Dockerfileで直接行います。pythonパッケージはrequirements.txtファイルに追加し、動作確認済みのバージョンに固定します。このイメージはデータアナリストが使用するものなので、イメージの中にdbt以外のものがあってはいけません。

新規画像の作成

{:#new-images}をご覧ください。

本番用のイメージは、master ブランチに git タグがプッシュされたときにのみ作成されます。現在のリリースフローは以下の通りです。

# master ブランチか feature ブランチのどちらかで実行: git tag v<sem\_ver>

# タグを押す

git push origin \$(git describe --tags --abbrev=0)

ビルドパイプラインが自動的に実行されます。これらの画像が参照されている場所で新規にMRを作成し、タグを更新します。

# パイソンのハウスキーピング

レポのpythonをきれいに維持するために、複数のmakeコマンドとCIジョブが設計されています。Makefileの中の以下のコマンドは、レポの分析に役立ちます。

- make lint は black python linter を実行し、ファイルを更新します (これは単なるチェックではありません)
- make pylintは、pylintチェッカーを実行しますが、コードのフォーマットをチェックするためにblackを使用しているため、チェックしません。重複したコード、エラー、警告などをチェックします。コードの品質を高めるための一般的なものです。DAGs dirは、一般的なコード標準に従うことが期待されていないため、無視されます。
- make radonは、関連するPythonコードのサイクロマティックな複雑さをテストし、スコアがBの関数やモジュールを表示します。 またはそれ以下です。
- make xenonは複雑性のチェックを行い、閾値を満たさない場合はゼロ以外の終了コードを返します。を無視します。 shared\_modulesやtransform reposは、削除されたり、廃止されたり、後日更新されたりするまで。

# ピアリングされたVPCへのアクセス

GitLab固有のELTの中には、ピアリングされたGCPプロジェクト内のデータベースに接続するものがあります(pingなど)。接続を許可するために、いくつかのアクションが取られています。

- 1. ランナーが実行されるKubernetesクラスターは、IPエイリアシングを使用するように設定されているので、各ポッドは GCP内で実際にルーティング可能なIPを取得します。
- 2.2つのプロジェクトとそのネットワークの間には、VPCピアリング関係が確立されています。
- **3.** ランナーとなる**Kubernetes**クラスターのポッドサブネットからのアクセスを許可するために、上流のプロジェクトでファイアウォールルールが作成されました。

# Extractorのドキュメント

- BambooHR
- コミット統計 (エンジニアリ
- ング) グラファイト (エンジ
- ニアリング) SheetLoad
- ハンドブックのYMLファイル

ランナーの更新

**CI**ジョブはgitlab-dataグループで、**GCP**プロジェクトのgitlab-analysisを介して**Kubernetes**で実行しています。すべてのリポジトリで共有するためにグループランナーを設定しています。

新しいグループのランナー・トークンを関連付ける必要がある場合や、ランナー・イメージを更新する必要がある場合などです。以上が基本的な手順です。注 - helm 3のリリース以降、これらのコマンドはすべてGCPのCloud Shellコンソールで実行することが推奨されています。ランナー用のデプロイメント(現在はgitlab-data-gitlab-runner)に移動し、kubectlのドロップダウンを使用してシェルに入ります。

クレデンシャルを取得するには

gcloud container clusters get-credentials bizops-runner --zone us-west1-a --project gitlab-analysis

ヘルスリリースを見るには

helm list --namespace <namespace> を表示します。

特定のリリースのチャート値を取得するには

helm get values --namespace <namespace> <chartname>

準備のためのコマンド

helm repo add <chart>

<url> helm repo update

ランナーを削除するには

helm delete --namespace <namespace> <chartname>

Helmチャートでランナをインストールするには

helm install --namespace <namespace> --name <chartname> -f <valuesfile> <chartrepo/name>

ランナーのバージョンやグループトークンを更新する例

```
helm list --namespace gitlab-data
helm get values --namespace gitlab-data gitlab-runner
helm get values --namespace gitlab-data gitlab-runner >
values.yaml helm repo add gitlab https://charts.gitlab.io
helm repoの更新
helm delete --namespace gitlab-data gitlab-runner
helm install --namespace gitlab-data --name gitlab-runner -f values.yaml gitlab/gitlab-runner
```

## 私たちのYAML設定は以下の通りです。

```
アフィニティです。{}
checkInterval:30
コンカレント:10
gitlabUrl: https://gitlab.com
hostAliases:[]
imagePullPolicy:IfNotPresent
metrics:
 enabled: true
nodeSelector:{}
podAnnotations:{}
podLabels:{}
rhac:
 clusterWideAccess: false
 create: false
 podSecurityPolicy:
   enabled: false
```

```
resourceNames:
   - gitlab-runner
 serviceAccountName: gitlab-data-gitlab-runner
runnerRegistrationToken: <token found in https://gitlab.com/groups/gitlab-data/-/settings/ci_cd> ラン
 のビルドを行います。{}
 キャッシュ
 を使用しま
 す。{}
 helpers:{}
 image: ubuntu:16.04
 outputLimit: 4096
 pollTimeout: 180
 privileged: false
 services:{}
 タグ:アナリティクス,ハウスキーピ
ングセキュリティコンテクスト
 fsグループです。65533
 runAsUser: 100
```

Snowflakeにデータをロードする自動化プロセス

バージョンDBとライセンスDBのロードタスク

ソースPostgres DBデータエクスポート

バージョンプロジェクトとライセンスプロジェクトで実行される1日1回のCIジョブがあります。このジョブでは、データベースエクスポートのバージョンスクリプトまたはライセンススクリプトを実行し、GCSバケットにCSVファイルをエクスポートします。これらのファイルはgitlab-version-{table\_name}-という名前です。

それぞれ{monday\_of\_week}またはgitlab-license-{table\_name}-{monday\_of\_week}となります

。スノーフレークステージ

バージョンDBのCSVファイルがエクスポートされるGCSバケットは、Snowflakeではステージとして設定されています。 raw.version\_db.version\_dump.ライセンスDBファイルがエクスポートされるGCSバケットは、ステージ raw.license\_db.license\_dumpとして設定されています。つまり、Snowflakeからは、すべてのファイルを一覧表示したり、ファイル内のデータをテーブルにコピーしたり、ファイルを削除したりすることができるのです。

カラムマッチングとテーブル生成

**CSV**ファイルは自己記述的ではありません。どの列がどの位置にあるかを知るための列へッダーがありません。このため、RAWのテーブルは、**CSV**ファイルの順番と正確に一致した列の順番を持つ必要があります。このテーブルを簡単に作成するために、このbashスクリプトを作成しました。create文を生成するために。

- 1. 目的のデータベースのdeclare文のコメントを外す
- 2. 環境変数PATH\_TO\_MANIFESTSに、クローン化したバージョンのsqlソースフォルダへのパスを設定します。version dbでは、このフォルダのクローンバージョンを指すようにします。license dbの場合は、このフォルダを指定します
- 3. do\_create\_tables.shを実行すると、すべてのcreate table文が出力されるはずです。

スノーフレークのタスク

CSVファイルは、スノーフレークステージからスノーフレークタスクで毎日読み込まれています。タスクは以下のようなSQL を実行して生成されています。

```
タスク users_load_task WAREHOUSE =
LOADING を作成または置換します。
SCHEDULE = '1440
minute' AS
COPY INTO users
from @raw.version_db.version_dump/gitlab-version-users-
file_format = (TYPE = CSV
FIELD_OPTIONALLY_ENCLOSED_BY='"')
```

を LOADER ロールで実行します。これらのタスクは毎日実行され、新規または更新されたファイルのみをロードします。タスクの定義を確認するには、LOADERロールでshow tasksを実行し、コンテキストの一部としてversion\_dbまたはlicense\_dbスキーマを指定します。

タスクモニタリング

タスクは**dbt**テストで監視されており、各タスクが過去1日以内に正常に実行されたことを確認します。バージョンテストとライセンステストは、それぞれ内部のデータテストプロジェクトのこことここにあります。タスクの問題を診断するために、raw.snowflake.task\_history\_viewを照会して、失敗したタスクのerror\_message列を検査することができます。

PTO by Roots Snowpipe

本号で紹介したように、peopleグループが維持しているプロセスでは、定期的にRoots APIでPTOに問い合わせを行い、その結果をGCSバケットのgitlab-ptoにあるJSONファイルにダンプしています。このバケットはgitlab-analysisというGCPプロジェクトの中にあります。

このデータをSnowflakeに読み込むために、Snowflakeのドキュメントを参考にして、GCS用のSnowpipeを設定しました。gitlab-pto-snowpipeという名前のGCP pubsubトピックは、新しいファイルがgitlab-ptoバケットに書き込まれるたびに新しいメッセージを受け取るように設定されました。gitlab-pto-snowpipeという名前のGCP Pub/Subサブスクリプション(配信タイプはPull)を作成し、gitlab-pto-snowpipeトピックを購読しました。

そして、Snowflakeで以下のコマンドで通知統合を作成しました。

```
通知統合の作成 pto_snowpipe_integration type = queue
notification_provider = gcp_pubsub
enabled = true
gcp_pubsub_subscription_name = 'projects/gitlab-analysis/subscriptions/gitlab-pto-snowpipe';
```

そして、このコマンドでSnowpipeを作成しました。

```
CREATE OR REPLACE PIPE
raw.pto.gitlab_pto AUTO_INGEST =
true
インテグレーション = 'pto_snowpipe_integration'
AS copy into raw.pto.gitlab_pto (jsontext, uploaded_at)
from (select $1, current_timestamp() as uploaded_at from
@raw.pto.pto_load) file_format=(type='json' strip_outer_array = true);
```

Snowpipeが有効になっていたので、gitlab-ptoバケットに書き込まれた新しいファイルをうまく拾うことができました。

```
alter pipe raw.pto.gitlab_pto enable;
```

デバッグ

PTOのスノーパイプに問題がある場合は、まずスノーパイプの状態を確認します。これはSnowflakeで次のように実行することで可能です。

```
select SYSTEM$PIPE_STATUS('gitlab_pto')。
```

**Demandbase Load Tasks** 

今回の課題の一環として、Demandbaseのデータロードを実施しました。Demandbaseのデータは、Demandbaseによって毎日、datastream-hosted-gitlab-3750というDemandbaseが所有するGCSバケットにparquet形式でロードされます。GitLabのスノー 273/243

フレークGCS

統合サービスアカウントは、DemandbaseからSnowflakeにデータをロードするために、このバケット内のファイルの読み取りとリスト化の権限を与えられました。

デマンドベースのデータとのインターフェイスとして、Snowflakeステージを作成しました。

```
CREATE STAGE "RAW".demandbase.data_stream
STORAGE_INTEGRATION = GCS_INTEGRATION
URL = 'gcs://datastream-hosted-gitlab-3750/datastream-gitlab/';
```

その後、GCSからSnowflake RAWデータベースにロードするために、各demandbase関係のSnowflakeタスクを作成しました。例えば、Snowflakeタスクのロードアカウントは次のように定義された。

```
create task demandbase_account_load_task
    WAREHOUSE = LOADING
    SCHEDULE = '1440
    minute' AS
copy into raw.demandbase.account (jsontext, uploaded_at)
    from (select $1, current_timestamp() as uploaded_at
    from)
@raw.demandbase.data_stream/db1_accounts/)
    file_format=(type='parquet');
```

その後、alter task DEMANDBASE\_ACCOUNT\_LOAD\_TASK resume; を実行して、タスクを有効にしました。各タスクは毎日実行し続け、GCS からの新しいファイルを Snowflake の raw.demandbase スキーマにロードします。

Thanos ロードタスク

この問題の一部として、thanosのメトリクスを毎日取得し、GCSのバケットに書き込むプロセスが設定されて

います。GCSからSnowflakeにメトリクスを取り込むために、ステージを作成しました。

```
create stage "raw"."PROMETHEUS".periodic_queries
STORAGE_INTEGRATION = GCS_INTEGRATION URL = 'gcs://periodic-queries/';
```

その後、Snowflakeタスクが設定され、新しいデータファイルが毎日読み込まれます。

```
タスク prometheus_load_task WAREHOUSE =
LOADING を作成または置換する。
SCHEDULE = '1440
minute' AS
raw.prometheus.periodic_queries (jsontext, uploaded_at)にコピーします。
from (select $1, current_timestamp() as uploaded_at from
@raw.prometheus.periodic_queries) file_format=(type='json');
```

データの更新

スティッチ管理データ

SLO違反につながるテスト失敗のためにフルリフレッシュが必要な場合、テスト失敗の根本原因を調査する時間が必要です。フルリフレッシュを行う前にインシデントを起こし、発見事項を文書化する。また、データを別の一時的なテーブルにコピーして、フルリフレッシュが妨げられないようにすることも有効です。

Stitchで抽出されたデータソース(プラットフォームページの抽出テーブルを参照)の場合、フルリフレッシュを行う推奨方法 は以下の通りです。

• Stitchのレプリケーションジョブが現在実行されていないことを確認します。統合を一時停止して、このプロセス以外の実行を開始しないようにします。

• STITCHロールとして、リフレッシュ・プロセスが本番環境に何の影響も与えないように、ターゲット・スキーマをクローンする。これは、以下のSQLを実行し、関連するスキーマ名を記入することで行うことができる。

```
create schema clone_<DATA_SCHEMA> clone <DATA_SCHEMA>;
```

• 次のSQLクエリを実行して、すべてのテーブルを取得します(関連するスキーマに変更します)。

```
SELECT
'clone_<DATA_SCHEMA>が存在すればテーブルを切り捨てる。|| LOWER(テーブル名) ||
';' FROM "RAW".「INFORMATION_SCHEMA".テーブル
WHERE LOWER(table_schema) = '<DATA_SCHEMA>';
```

- Snowflakeですべての切り捨てクエリを実行する
- Stitchに戻り、必要なタイプの新しい統合を作成します。統合の名前を「Clone <古い統合の名前>」とします。例えば、Salesforceの場合、統合の名前を「Clone SalesForce Stitch」とします。 統合を認証し、現在のStitchジョブの抽出設定に合わせて抽出を設定します。抽出の実行
- 抽出が完了すると、すべてのデータがclone\_<DATA\_SCHEMA>スキーマに入っているはずです。SQLクエリでサニティチェックを行い、データが期待通りに表示されていることを確認します。抽出が完了した後にStitchのデータが到着する場合がありますので、Stitchからのデータがすべて到着していることを確認するための時間を確保してください。
- 各テーブルに必要なSQLを取得して、新しいデータと古いデータを入れ替えます(<DATA\_SCHEMA>の置き換え 適切に)。)

```
SELECT
'alter table <data_schema>.|| LOWER(テーブル名) || ' SWAP WITH CLONE_<DATA_SCHEMA>.'.||
LOWER(テーブル名) || ';'
FROM "RAW"."INFORMATION_SCHEMA".tables
WHERE LOWER(table_schema) = '<DATA_SCHEMA>';
```

- 上記のSQLで得られた各SQL文を実行します。
- ・ 抽出」→「設定」→「削除」で、クローン化された抽出ジョブを削除します。 古い
- Stitchレプリケーションジョブを再開し、増分データが再び流れ始めるようにします
- リフレッシュされたデータが完成し、期待通りであることが明らかになった時点で、drop schema clone\_<DATA\_SCHEMA> cascade; を実行して古いデータを削除します。

dbtモデルフルリフレッシュ

dbt\_full\_refresh DAGを使用して、dbtにインクリメンタルモデル全体を一から作り直させます。

- 1. Airflowで、変数DBT\_MODEL\_TO\_FULL\_REFRESHに、リフレッシュするモデルの名前をdbtmodelselection構文に従って設定します。例えば、バージョンモデルを更新するには、sources.versionstaging.versionとします。gitlab\_dotcomモデルを更新する場合、値はsources.gitlab\_dotcom staging.gitlab\_dotcomとなります。 ☑airflow\_variable\_setting
- 2. DAGの電源を入れ、起動させま
- す。後ろで実行されるdbtコマンドは

```
dbt run --profiles-dir profile --target prod --models DBT_MODEL_TO_FULL_REFRESH --full-refresh
```

GitLabデータユーティリティ

チーム内で使用している便利な機能を集約するためのプロジェクトです。プロジェクトは https://gitlab.com/gitlab-data/gitlab-data-utils

リリースのカット

- **1. setup.py**のマイナーバージョンをインクリメントします。マイナーバージョンは**2**番目の数字で、**0.1.0**の場合は**1**となります。現在、マイナーバージョンの更新のみサポートしています。
- 2. この変更に対してMRを作成し、masterにマージします。
- 3. チェックアウトしてローカルに master をプルし、今行ったバージョンの変更を手元に残しておきます。gitlab-data-utilsのルートフォルダでmake releaseを実行すると、保留中の変更がないことが確認され、新しいタグがGitLabリポジトリにプッシュされます。これにより、pypi に公開するためのパイプラインが実行されます。
- 4. pypiとの一貫性を保つために、GitLabのリリースを以下のようなコマンドでカットします。

```
curl --header 'Content-Type: application/json' --header "PRIVATE-TOKEN: <your-private-token>" ●.
    --data '{ "name":"GitLab Data Utilities v0.0.1", "tag_name":"v0.0.1",
"description":"Initial tagged release"}'\
    --request POST https://gitlab.com/api/v4/projects/12846518/releases
```

5. このマージリクエストで行われたように、Dockerイメージを更新します。

本番用にdbtをアップグレード

- 1. 主要なステークホルダーにアップグレードのタイミングを知らせるイシューを用意する。アップグレードを実行する前に、少なくとも1日前に通知すること
- 2. 新規イメージの作成」の手順に従って、新しいバージョンのdbt-imageを作成します。
- 3. アナリティクスプロジェクトにMRを作成し、以下の項目を更新します。例として、このMRをご覧ください。
- airflow\_utils.py dbt\_project.yml 以
  - 下の場所で dbt-image を使用し
  - o ます。
  - o snowflake-dbt-ci.yml
  - o .gitlab-ci.yml
  - o docker-compose.yml
- packages.ymlに含まれる任意 のパッケージ
- index.htmlを以下のようにして作成します。
  - make dbt-docs コマンドで docs サイトを読み込む。
  - 右クリックして「View Page Source」を選択するか、Macの場合は「CMD + u」、Linuxの場合は「CTRL + u」を押すと、ページのソースコードが表示されます。
  - <body>からファイルの一番下までを更新する
  - </head>からファイルの先頭までのコードに大きな変更がないことを確認します。
- 1. data-testsプロジェクトにMRを作成し、imageバージョンとrequire-dbt-versionを更新します。例として、このMRを参照してください。
- 2. リリースの機能アップデートに基づいて、関連するdbtジョブを実行する。
- 3. マージされたら、全員にアップグレードを通知し、ジョブの失敗やユーザーの問題を監視する

layout: handbook-page-toc title:"Jupyter Guide" 説明"JupyterLabを使ってSnowFlakeと内部でやりとりするためのガイダンス"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}

関連リポジトリを見る

### 特徴

- 一般的なpythonのDS/MLライブラリ(pandas、scikit-learn、sci-pyなど)。
- dbtの認証情報を使ってSnowflakeにネイティブ接続。ログインは必要ありません。
- Git機能: JupyterLab内でGitlabリポジトリへのプッシュ/プルが可能(ssh認証が必要)。
- コンピュータやGitlabのリポジトリにある任意のpythonファイルやノートブックを実行します。
- 使っているのに見当たらない機能が必要ですか?ぜひ#bt-data-scienceでお知らせください!追加します。

### はじめに

data-scienceプロジェクトでjupyterをセットアップする際には、2つのオプションがあります。以下のいずれかを選択

- してください。フル・インストール(推奨)。**Pipfile**で定義されたすべてのライブラリ*と*、完全な anacondaのインストールを行います。
- Lightweight install: Pipfileで定義されたライブラリ*のみを*インストールします。ローカルマシンにpython環境があり、それをベースイメージとして使用したい場合に使用します。

### インストール手順

- 1. レポをローカルマシンにクローンします git clone git@gitlab.com:gitlab-data/data-science.git
- 2. 実行 cd データサイエンス
- 3. どのバージョンをインストールするかに応じて、以下のいずれかを実行してください。
  - *完全にインストールするには*、make setup-jupyter-localを実行します。
  - 軽量インストールの場合: run make setup-jupyter-local-no-conda
- 4. make jupyter-local を実行します。
- 5. Jupyterが自動的に起動するはずです。起動しない場合は
  - 1. まず、Google Chromeがデフォルトのブラウザであることを確認します(「システム環境設定」から「一般」をクリックし、「デフォルトのウェブブラウザ」セクションのドロップダウンメニューからGoogle Chromeを選択します)。
  - 2. 次に、Chromeで、dockerイメージが作成された後にターミナルで見つけたURLとトークンをコピーペーストします。http://127.0.0.1:8888/lab?token=5c7f7da79f4a0968501f087f3c79ee4dd8bd7a63e0f088a8 のように表示されるはずです。トークンは、dockerコンテナをスピンアップするたびに変わります。

## Snowflakeへの接続

- 1. パスワードが含まれていない{User}/.dbt/profiles.ymlファイルを設定していることを確認してください。ここで提供されている例を参考にしてください。
- 2. レポにあるauth\_exampleノートブックを実行して、すべての設定が正常に行われたことを確認します。初めて実行すると、Okta経由でsnowflakeのクレデンシャルを認証するためのブラウザのリダイレクトが表示されます。その後、ノートブックを再度実行すると、Snowflakeからデータを照会できるようになるはずです。
- 3. エラーが出る場合は、Snowflakeがマシン上で正しく設定されていない可能性があります。DataOnboardingIssueのSnowflakeとdbtのセクションを参照してください。.dbt/profiles.ymlが正しく設定されていない可能性があります。

## ローカルディレクトリのマウント

デフォルトでは、ローカル・インストールでは、data-scienceフォルダをjupyterのルート・ディレクトリとして使用します。これは、すべてのコード、データ、ノートブックがコンピュータの他の場所にある場合、あまり便利ではありません。これを変更するには、jupyter notebook configファイルを作成して変更する必要があります。

- 1. jupyter-lab --generate-configを実行します。すると、/Users/{user}/.jupyter/jupyter\_lab\_config.pyというファイルが 作成されます。
- 2. ファイルの場所を参照し、エディターで開く
- 3. ファイル内で以下の行を検索します。c.ServerApp.root\_dir = ''と書き換えてください。よくわからない場合は、自分のレポのディレクトリに設定してください(例:c.ServerApp.root\_dir = '/Users/{user}/repos')。行頭の#が削除されていることを確認してください。
- **4.** パスには必ずフォワードスラッシュを使用してください。フォルダ名にスペースが含まれていても、二重引用符で囲まれていれば、バックスラッシュを使うことができます。
- 5. data-scienceディレクトリからmake jupyter-localを再実行すると、ルートディレクトリが上記で指定したものに変更され

るはずです。

### Dockerのメモリ割り当ての増加

デフォルトでは、dockerはコンテナを実行するために2GBのメモリを割り当てます。jupyterやpythonではデータがメモリ内に保持されるため、これでは十分なRAMが確保できない可能性があります。メモリ不足のエラーを避けるために、dockerのメモリ割り当てを増やすことをお勧めします。

- 1. Dockerのダッシュボードを開く。
- 2. 右上の歯車アイコンをクリックすると、設定が表示されます。
- 3. リソース」では、Dockerが使用するメモリを追加で割り当てます。8GBを推奨しますが、大規模なデータセットを扱う場合はさらに増やす必要があるかもしれません。
- 4. Dockerを再起動します。

### Jupyter Extensionsのセットアップ

- data-scienceのリポジトリには、git、変数インスペクタ、折りたたみ可能な見出し、実行時間、システムモニタなど、 多くの便利なJupyter Lab拡張機能がプリインストールされています。
- これらを最大限に活用するために(そして、コンテナを実行するたびに設定しなくて済むように)、以下のファイルを作成します。/Users/{user}/.jupyter/lab/user-settings/@jupyterlab/notebook-extension/tracker.jupyterlab-settings
- そのファイルの中に、以下を貼り付けて保存します。

```
{
    "codeCellConfig":{
        "codeFolding": true,
        "lineNumbers": true,
},
    "recordTiming": true,
}
```

興味深いライブラリーがあります。

### データ/モデル解析

- **ELI5**
- QuickDAビジュアラ

# イゼーションツール。

- Plotly
- Seaborn

## MLライブラリ

- SKlearn
- Tensorflow
- Torch
- パイアース (線形回帰、ロジスティック回
- 帰) プロフェット (時系列
- オートット(時系列
- \* XGBoost (強力なブラックボックス手法) 簡

## 単なコンカレンシー

- モディン
- Dask (セルフインストールが必要)

GPUスピードアップ

• プレーンML

layout: handbook-page-toc title: "Meltano At Gitlab"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .hidden-md .hidden-lg}<sub>o</sub>

{::options parse\_block\_html="true" /}

クイックリンク

Airflow{:.btn .btn-purple-inv} DataImageProject{:.btn .btn-purple-inv} GitLabDataUtilsProject{:.btn .btn-purple-inv} PythonGuide{:.btn .btn-purple-inv}.

メルタノ

Meltanoは、独自のKubernetesクラスターで、デフォルトの名前空間で運用しています。現在、3つのプロジェクトリポジトリが設置されています。このKubernetesクラスタはmeltano-masheyとしてGCPで稼働しています。

MeltanoのUIはインターネットに公開されていません。ログを見るには、kubernetesコンテナのログを見る必要があります。これは、「LOGS」タブ、または概要ページの「ワークロード」でmeltano-gitlabクラスタを選択することで見ることができます。

Meltanoについては、以下のリポジトリがあります。

## RepositoryDescriptionリポジトリの

- 1 Gitlab-datameltano 彼のプロジェクトはインフラ関連のコードを持っています。 gitlab-app.yamlとmeltano.ymlの設定。
- TAPのソースコードを保持するプライマリリポジトリです。現時点では、以下のソースコードがあります。

  TAP-XACTLYとTAP-ADAPTIVE
- **tap-zengrc tap-zengrc**のソースコードを保持するプロジェクトです。

新しく開発されたタップについては、tap-zengrcのように、タップごとに新しいリポジトリを作成します。

使用している抽出器を更新するには、メインプロジェクトである**Gitlab-data-meltano**のmeltano.ymlファイルを更新します。変 更がマージされたら**git**タグを追加し、**gitlab-app.yml kubernetes manifest**を更新して新しいイメージを指すようにします

Meltanoは内部でAirflowを使用しており、メタデータのデータベースにはCloud SQLを使用しています。メルターノのデータベース。

ローカルからKubernetesクラスタに接続(Google Cloud SDKがインストールされていることが前提)。このコマンドが動作しない場合は、GCPに接続し、KubernetesのClusterを選択し、connect to clusterを選択してください。 最新のコマンドが表示されます。

gcloud container clusters get-credentials meltano-mashey --zone us-west1-a --project gitlab-analysis

現在の設定では、クラスターに以下のシークレットが定義されています。

- cloud-sql
- meltano-db
- airflow
- TAP-
- SECRET
   AIRFLOW-

DB

秘密ファイルの作成

```
kubectl create secret generic cloud-sql
kubectl create secret generic meltano-db
kubectl create secret generic airflow
kubectl create secret generic tap-
secrets kubectl create secret generic
airflow-db
```

## 編集の秘密

```
kubectl edit secrets cloud-sql
kubectl edit secrets meltano-db
kubectl edit secrets airflow
kubectl edit secrets tap-
secrets kubectl edit secrets
airflow-db
```

# コンテナへの格納

app.yaml

```
kubectl exec -it gitlab-production-5dd4c79694-vwwtm -c tap-actly /bin/bash
```

注意事項ポッド名は変更される可能性がありますので、UIまたはkubectl get podsコマンドで確認してください。

```
kubectl apply -f gitlab-app.yaml

# 名前空間のデプロイを削除する kubectl

delete-f gitlab-
```

GKEでのMeltanoセットアップのウォークスルーのビデオを見る

メルトノの設定に、TAPで使用するConfig変数を追加。

# 更新されたマニフェストの適用には、名前空間は必要ありません

コンフィグをMeltanoのタップに渡す必要があるときはいつでも、以下の5つのアクティビティを行う必要があります。以下の5つのアクティビティを行う必要があります。

- 1. ユニークな変数名を決め、それらの変数値をクラスタ内のtap-secretsという名前のKubernetesのシークレットに追加します。シークレットファイルを編集するコマンドは、kubectl edit secret tap-secrets -o yamlです。
- 2. これらの変数を gitlab-app.yaml ファイルに以下のように参照して追加します。

```
# Kubernetes Secrets::zengrc
- name: ZENGRC_USERNAME # 環境で一意になるようにする
valueFrom:
secretKeyRef:
name: tap-secrets # これは、秘密を追加するための秘密の名前で、存在するどのような秘密ファイルでもよい。
キーを使用します。ZENGRC_USERNAME # シークレットファイルでユニークにしておく。
- name: ZENGRC_PASSWORD # 環境で一意になるようにする
valueFrom:
secretKeyRef:
name: tap-secrets # これは、秘密を追加するための秘密の名前で、存在するどのような秘密ファイルでもよい。
キーを使用します。ZENGRC_PASSWORD # 環境に合わせてユニークに保つ
```

3. ファイルを修正した後、新たに用意した配置ファイルをクラスタに適用する必要があります。これを行うには、既存のデプロイメントを削除して、新しいデプロイメントを適用する必要があります。以下のコマンドを使用して、この作業を行います。

- kubectl delete -f gitlab-app.yamlで既存のデプロイメントを削除します。
- kubectl apply -f gitlab-app.yaml で新しいデプロイメントを適用します。
- 4. 設定ファイル (meltano.yml) の中で、これらの定義された変数を以下のように参照します。

config:

ユーザー名: \$ZENGRC\_USERNAME パス

ワード: \$ZENGRC PASSWORD

このようにする理由は、ユーザー名とパスワードは環境内で一意のキーではなく、他のタップでも使用されているため、正しい TAPS設定を渡すために、lubeの秘密変数名への参照を渡します。5) 最新のmeltano.ymlファイルをコンテナにコピーします。 kubectl cp meltano.yml default/<\*\*pod-name\*\*>:/projects

すべてのタップのスケジュールを実行して、問題がないことを確認します。

meltano.ymlファイルへのタップの追加

現在の設定では、MRを作成してタップ情報をmeltano.ymlファイルに追加しています。必要な情報が meltano.yml ファイルに 追加されたら、以下の手順で TAP を有効にしてください。以下は、各TAPに追加する必要があるRequire情報のサンプルです。

## The Name of the TAP under plugins:-- > エクストラクタ

- 名前:tap-zengrc 名前空
間:tap\_zengrc
pip\_url: git+https://gitlab.com/gitlab-data/tap-zengrc 実行ファイル:tap-zengrc
の能力を発揮します。
- の設定を
発見しまし
た。
- 名前:ベースURL
- 名前:ユーザー名
- name:

password config:

base\_url: https://gitlab.api.zengrc.com/api/v2 ## スケジュールセクションのスケジュールの下にあるタッ プのスケジュール。

 name: zengrc-to-snowflake extractor: tap-zengrc loader: target-snowflake transform: skip

interval: '@daily'
start\_date: 2021-07-13

その後、以下の手順でmeltano.ymlを実行中のコンテナにコピーします。

- 必要なクラスタに接続します。gcloud container clusters get-credentials meltano-mashey --zone us-west1-a --project gitlab-analysis kubectl exec -it gitlab-production-5f8fd9ccb-npvxl -c tap-xactly コマンドを使
- って、現在アクティブなクラスタである meltano-mashey に接続します。 /bin/bashを使用しています。注:- ポッド名が変更されている場合がありますので、正しいポッド名を取得するにはkubectl get podsを使用してください。
- 修正したmeltano.ymlをローカルからコンテナにコピーする kubectl cp meltano.yml default/gitlab-production- 5f8fd9ccb-npvxl:/projects
- 最初にschedule meltano schedule run zengrc-to-snowflakeを実行してみると、インストールを求められます。

merged.md

2021/11/8

抽出器を使ったセッションでも、以下のようなエラーが発生しました。

メルト | Running extract & load...

/ | ELTを完了できませんでした。Cannot start extractor:実行可能な'tap-zengrc' が見づかりませんでした。エクストラクタ 'tap-zengrc' が `meltano install' を使ってまだインストールされていな以可能性があります。

抽出器 tap-zengrc`、または実行ファイル名が間違っている可能性があります。

ELTを完了できませんでした。Cannot start extractor:実行可能な 'tap-zengrc' が見つかりませんでした。 meltano install extractor tap-zengrc`で抽出器'tap-zengrc'がまだインストールされていないか、実行ファイル名が間違っている可能性があります。

• エキスパンダーをインストールした記事

root@gitlab-production-5f8fd9ccb-q6gt4:/projects# meltano install extractor tap-zengrc Installing 1 plugins...
エキスパンダー「tap-zengrc」のインストール...エ
キスパンダー「tap-zengrc」のインストール

• その後、コマンド meltano schedule run zengrc-to-snowflake を再実行します。これで snowflake にデータがプッシュされるようであれば、TAP は期待通りに動作しています。

layout: handbook-page-toc title:"Sisense For Cloud Data Teams" description:"GitLabにおけるSisense For Cloud Data Teams"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# クイックリンク

Sisense{:.btn .btn-purple-inv} dbtdocs{:.btn .btn-purple-inv}

アクセス

Sisenseは、当社のエンタープライズ標準のデータ可視化アプリケーションであり、当社のエンタープライズデータウェアハウスへの接続が承認されている唯一のアプリケーションです。

セルフサービス・ダッシュボードへのアクセス

- 1. GitLabの全員に、SisenseのView-onlyライセンスが提供されています。
- 2. 表示のみ」では、キーレビューで定期的に確認されるものを含む、機密性のないダッシュボードにアクセスできます。
- 3. OktaでSisenseにログインします。

セルフサービスのダッシュボード開発

- 1. Sisenseのレポートやダッシュボードを開発するには、エディターライセンスが必要です。エディターライセンスはお金がかかりますので、リクエストする前に必要性を確認してください。私たちの倹約の価値に留意してください。
- 2. エディターを効果的に使うには、ちょっとした技術やデータモデリングの経験が必要です。
- 3. 正しいユーザーの役割を確認するなど、アクセスリクエストを作成します。

## ライセンスの再取得

- 1. データチームは、定期的にライセンスを見直し、90日間使用されていないライセンスを回収しています。
- 2. View-onlyやEditorのライセンスが不要になった場合は、DataTeamIssueを作成していただければ、あとは私たちにお任せください。

# トレーニングリソース

セルフサービス・ダッシュボードへのアクセス

- Sisense (GitLab内部) 製品ドキュメント
- へのアクセス
- はじめに
- はじめに ビューア

セルフサービスのダッシュボード開発

- 製品ドキュメント はじめに
- GettingStarted-セルフサービス・ダッシ
- ュボード GitLab'sSisenseProject
- GitLabのSisenseEditorトレーニング (GitLab社内)
- SisenseDataCommunity
- Sisense Plotlyテンプレート

# Sisenseへのアクセス

**GitLab**の全員が**Sisense**への閲覧のみのアクセス権を持っています。**Okta**を使ってログインしてください。自分のチャートを作成するエディター権限など、高度なアクセスが必要な場合は、アクセスリクエストを作成してください。ユーザーロールのセクションもご覧ください。

# Sisenseリソース

- Sisenseプロジェクト
- Sisenseトレーニング(GitLab社内)
- Sisenseエディタトレーニング(GitLab社内)
- Sisenseデータオンボーディング。チャート/ダッシュボードの作成と分析
- Sisenseコミュニティ

セルフサービスのダッシュボード開発

このセクションの目的は、読者がGitLabのデータに関する疑問に答える独自のSisenseダッシュボードを構築できるようにすることです。最後に紹介する例は、製品組織に特化したものですが、GitLabの他のチームにも一般化できます。

Sisenseの基礎知識

独自のSisenseダッシュボードを構築する最初のステップは、正しい権限を持っているかどうかを確認することです。

Oktaを使ってSisenseにログインすると、右上にNew Chartボタンが表示されるはずです。何も表示されていない場合は、表示専用のアクセス権しかないので、上記の手順で編集者のアクセス権を取得してください。

New Chart」が表示されるようになったら、自分のダッシュボードを作成してみましょう。左側のナビバーでDashboardsを見つけ、 +アイコンをクリックしてダッシュボードを作成します。ディレクトリを整理するために、必ず名前を付けてください。

ダッシュボードの構築と名前の決定が完了したら、右上の「New Chart」をクリックして、チャートの追加を開始できます。これで、クエリを書き始める準備が整いました。

適切なデータソースを見つける

データに関する質問に答えるための次のステップは、クエリを実行するための関連テーブルを見つけることです。そのためには、Snowflakeデータウェアハウスとそこに供給されるデータソースの背景を知る必要があります。Snowflakeに保存するデータには、一般的に3つのタイプがあります。外部、内部フロントエンド、内部バックエンドです。

外部データ

外部データとは、GitLabで使用しているサードパーティのソフトウェアによって生成されたすべてのデータのことですが、 本番データを自分たちで保存しているわけではありません。外部データには、Salesforce、Zuora、Netsuite、Greenhouse、 BambooHRなどがあります。これらのデータは、APIを使用してデータウェアハウスにロードします。

内部のバックエンドデータ

**GitLab.com**は**Ruby on Rails**のアプリで、バックエンドに**Postgres**データベースを使用しています。**GitLab.com**のユーザーが新しい**MR**、課題、コメント、マイルストーンなどを作成するたびに、データベースに新しい行が作成されます。データチームはカスタム**ELT**を作成し、これらの**Postgres**テーブルをデータウェアハウスに同期させています。

自己管理型のインスタンスでは、Usage pingを使ってこれらのバックエンドデータベースの匿名化されたサマリーを毎週取得するようにしています。

内部のフロントエンドデータ

さらに、gitlab.com のフロントエンドのインタラクションを追跡するために Snowplow というツールを有効にしました。 Snowplow は、ページビューの自動追跡に加え、フォームやリンククリックの追跡も行います。 Snowplow は、ユーザーのセッションやブラウザに関する情報を含むメタデータを、すべてのイベントとともに送信します。

注: Snowplowはバックエンドのイベントを捕捉することもできますが、現時点では主にjavascript(フロントエンド)のトラッキングに使用しています。 {: .alert .alert-info}。

フロントエンドデータとバックエンドデータの違いは何ですか?バックエンドデータとは、アプリケーションにとって何らかの目的を果たすために、すでにアプリケーションのデータベースに保存されているデータのことです(MR、課題、パイプライン)。それに対して、フロントエンドのトラッキングの主な目的は分析です。

\*\*dbtドキュメント\*\* {: .panel-heading}。

私たちの[dbt Docs site](https://dbt.gitlabdata.com/)には、snowflakeでクエリを実行できるすべてのテーブルがリストアップされています。これらの多くは、テーブルとカラムの両方のレベルで文書化されており、クエリを書くための素晴らしい出発点となっています。

発見データセット

データディスカバリーはSisenseの機能で、SQLのスキルが低いユーザーでも、ドラッグ&ドロップのインターフェイスで特定のデータセットのビジュアライゼーションを作成することができます。

現在、この機能の価値を理解するためにテストを行っています。社内の誰もがアクセス可能なテストディスカバリーテストをいくつか作成し、私たちの支援なしにチャートを構築できるようにしました。

GitLab.com 利用データイベント

このデータセットは、ここに書かれているモデルに基づいて作成されています。ユーザーは、特定のnamespace\_idによって 実行されたすべてのイベントと、そのnamespaceに関する追加のメタデータを見つけることができます。

このデータセットでどのような質問に答えることができますか?

- 名前空間の作成後、最初の14日間に実行されたイベントの平均数 名前空間の作成後、最
- 初の14日間に発見されたステージの平均数 コホートビュー
- 毎日どれだけのCIパイプラインが実行されているのか?
- 毎月どれだけの有料ネームスペースが新しいマイルストーンを作っているのか?

実現可能なビジュアライゼーションの例を含むダッシュボードを作成しました。このデータセットにア

クセスするには、主に2つの選択肢があります。

- ダッシュボード上で編集者権限が与えられている場合は、右上の「新規チャート」をクリックすることができます。 また、左
- メニューのコンパスのアイコンをクリックすると、新しい探査を開始することができます。

チャートエディターが開いたら、以下のように「発見」ボタンをクリックします。

リストメニューからgitlab\_dotcom\_usage\_data\_eventsというデータセットを探します。

```
ジディスカバリー・データセット
```

例

\*\*Question 1:\*\*ページビューについて質問しているので、この質問に答えるために内部のフロントエンドデータを使用できます。以下のようにSnowplowのページビューを照会することができます。

```
アのカーリナビュースのCOLマニュカで中午ナイト、カーリの下のチャ、トガハアまが出わせるナナースともと、Cinnasは三
```

このクエリをSisenseのSQLエディタで実行すると、クエリの下のチャート部分に表が出力されます。そこから、Sisenseはデータを視覚化するための様々なオプションを提供します。チャートの構築について学ぶに、Sisenseの10分間DataOnboardingビデオを見るのが良いでしょ

これを見ると、毎月約11万人のユーザーがマージリクエストを作成していることがわかります。

質問 3: 過去 30 日間に 1 つのトライアルフォームから別のフォームに変換した「ユーザー」の数は? (Conversion Funnel) snowplow CTE を使用して、2 つのステップを別々に照会し、次にそれらを結合することができます。

```
WITH first_trial_form AS
  ( SELECT)
    user_snowplow_domain_id AS user_id,
    min_tstamp::DATEAS day_of,
                           tstampAS sent_at
 FROM analytics.snowplow_page_views_30 -- 必要に応じて、`_30`を`_all`に変更してください。カス
タムの日付範囲の場合は、derived_tstampでフィルタリングします。 `WHERE derived_tstamp BETWEEN
'2019-10-01' AND '2019-12-01'
 WHERE page_url_path = '/-trials/new'
),
second_trial_form AS
  ( SELECT
   user_snowplow_domain_id AS user_id,
   derived_tstamp::DATEAS day_of,
   derived
                          tstampAS
   sent at
  FROM analytics.snowplow_unstructured_events_30
  WHERE event_name = 'submit_form'
   AND page_url = 'https://gitlab.com/-/trials/apply'
)
SELECT
 first_trial_form.day_of,
  COUNT(DISTINCT first_trial_form.user_id) AS "View first trial form",
  COUNT(DISTINCT second_trial_form.user_id)AS "Submit last trial form",
  COUNT(DISTINCT second_trial_form.user_id) * 100 / COUNT(DISTINCT first_trial_form.user_id)
AS "Pct"
FROM first_trial_form
 LEFT JOIN second_trial_form
   ON first_trial_form.user_id = second_trial_form.user_id
    AND first_trial_form.day_of = second_trial_form.day_of
    AND first_trial_form.sent_at <= second_trial_form.sent_at
GROUP BY 1
ORDER BY 1
```

質問 4: /merge\_requests ページでは、どの「ステータス」タブ ('Open', 'Merged', etc.) がクリックされますか?以前と同様に、スノープラウデータを使ってmerge\_requestsページのクリック数とページビューを測定することができます。

| タブ名            | 総クリック<br>数 | クリック数の割合  |
|----------------|------------|-----------|
| ステートマー<br>ジド   | 200923     | 65.304300 |
| ステートオー<br>プン   | 50291      | 16.345700 |
| ステート・ク<br>ローズド | 31957      | 10.386700 |
| ステート・オール       | 24501      | 7.963400  |

この4つのタブのクリック数のうち、65%が「統合」タブへのクリックです。

これらのタブがクリックされる頻度を、マージリクエストページの*総ページビュー数に対する割合として*知りたい場合は、 このクエリにいくつかの変更を加えます。

```
WITH link_clicks AS (
  SELECT
   TRY_PARSE_JSON(unstruct_event): "data": "elementId"::VARCHAR AS element_id,
   COUNT(*)
                                                                     AS total clicks
 FROM analytics.snowplow_unnested_events_30
  WHERE event_name = 'link_click'
   AND element_id IN ('state-closed', 'state-all', 'state-opened', 'state-
   merged') AND page_urlpath LIKE '%/merge_requests'
  GROUP BY 1
),
page_views AS
 ( SELECT
   COUNT(*) AS total_views
 FROM analytics.snowplow_unnested_events_30
 WHERE event = 'page_view'
    AND page_urlpath LIKE '%/merge_requests'
)
SELECT
 link clicks.element id
                                                         AS "Tab
 link_clicks.total_clicks / page_views.total_views * 100Nge"ページビューの割合"
FROM link_clicks
 INNER JOIN
    page_views ON 1=1
```

| タブ名          | ページビューの割合 |
|--------------|-----------|
| ステートマー<br>ジド | 5.423417  |
| ステートオー       | 1.357480  |

プン ステートクロ 0.862600 ーズド

#### タブ名 ページビューのパーセ

ンテージ state-all0.661344

マージリクエストページのページビューのうち、「マージ」タブをクリックしたのは5%でした。

Sisenseのチャートパフォーマンス

Sisenseチャートがタイムアウトしたり、実行に時間がかかる場合は、Sisenseチャートを生成するためのSQLクエリをリファクタリングして、クエリを最適化する必要がある場合がほとんどです。データチームのプロジェクトで課題を作成してください。

ダッシュボードの公式バッジ

Sisenseの一部のダッシュボードには、公式バッジ(TwitterのVerified Checkmarkに似たもの)が表示されます。 Sisense 検証済みチェックマーク

これは、これらの分析がデータチームによってクエリの正確性が確認されたことを意味します。検証済みのチェックマークが付いていないダッシュボードは、必ずしも不正確なものではなく、データチームによるレビューが行われていないだけです。 公式バッジを追加・削除できるのは、Dataロールのメンバーのみです。

スペース

Sisenseのスペースは1つです。

• GitLab

この2人は、それぞれperiscopeとperiscope sensitiveという異なるユーザーでデータウェアハウスに接続します。

ほとんどの作業はGitLabスペースで行われますが、一部の*非常にセンシティブな分析*はGitLabセンシティブに限定されます。 その例としては、契約者や従業員の報酬に関わる分析や、匿名化されていないインタビューデータなどが挙げられます。

スペースはタグで整理されます。タグは機能(プロダクト、マーケティング、セールスなど)とサブ機能(クリエイト、セキュア、フィールドマーケティング、EMEA)に対応している必要があります。タグは課題のラベルと大まかに一致している必要があります(優先順位はつけません)。タグは無料です。人々が探している情報をできるだけ簡単に見つけられるようにしてください。現時点では、タグの削除や名前の変更はできません。

ダッシュボードを自動的にSlackにプッシュする

例えば、製品担当者は、製品の理由で失われた機会についての最新情報を毎週知りたいと思っているでしょう。この情報を定期的にSlackに流すことが望ましい場合は、Slackのネイティブな/remindを利用してURLを印刷することができます。ダッシュボードが自動更新されていないように見える場合は、Sisense管理者に連絡して更新スケジュールを更新してもらってください。

ユーザーの役割

{:#user-roles}

Sisenseには、admin、SQL、View Onlyの3つのユーザーロール(アクセスレベル)

があります。Sisenseのライセンスの現在の状況は、分析プロジェクトで確認できます。

- \*\*Sisenseのユーザーを更新する\*\* {: .panel-heading}。
  - コンソールにjqueryを注入する(StackOverflowより)。

```
let jq = document.createElement("script");
jq.src = "https://ajax.googleapis.com/ajax/libs/jquery/3.2.1/jquery.min.js";
jq.onload = function() {.
ロードを入力してください。
};
document.body.appendChild(jq)です。
```

• **DOM**で表示されるリストから全ユーザーのリストを取得するには、コンソールでこれを実行し、クラス名を実際のラベルに置き換えます。

```
$('div.list_class_namer_replace_this').map(function(i, el) {
  return $(el).text()})
).toArray()
```

アドミニストレーター

これらのユーザーは、新規ユーザーのプロビジョニング、パーミッションの変更、データベース接続の編集などの機能を持っています。(典型的な管理者の仕事)

リソースアドミンオンボーディング

編集者アクセス

ユーザーは、チャートやダッシュボードの基盤となるアナリティクスデータベースの analytics スキーマおよび analytics\_staging スキーマに対して SQL クエリを書くことができます。また、これを容易にするために、SQLスニペットを作成または利用することもできます。SQLアクセスライセンスの数には限りがありますので、現時点では、ディレクターを中心 としたチームごとに1つに制限することを目指しています。チーム内でSQLアクセスを持つ最適な候補者を決定するのは、ディレクターに任されています。

閲覧のみのユーザー

これらのユーザは、既存のすべてのダッシュボードを利用できます。ダッシュボードのフィルタを変更することができます。 最後に、ドリルダウン機能を利用してダッシュボードを掘り下げることができます。

部門別ユーザー

エディターロールをさらに細分化した追加のロールがあります。特定のチャートは、誰も編集できないようにします。例えば、「FinanceKPIs」ダッシュボードは、「Data」および「Finance」ロールのメンバーのみが編集できるようにします。

ユーザーをプロビジョニングする際の注意点

すべてのユーザーは、Oktaを介して表示のみのアクセス権限を持っています。

ユーザーをアップグレードするには、Sisense UIで「ロールとポリシー」のセクションに移動します。そして、ユーザーを関連するグループ(Admin/Editor)とその部門(例: Marketing、Productなど)または部署(例: UX、Securityなど)に追加します。

ユーザーは、どのグループに属していても、最高のアクセス権を継承します。これが、すべての機能がデフォルトで「表示のみ」となっている理由です。

新しいダッシュボードの作成とレビューのワークフロー

このセクションでは、Sisenseでダッシュボードを「本番」にプッシュする方法のワークフローを詳細に説明します。現在、MRファーストのワークフローを持つ機能はありません。このワークフローは、社内のすべてのダッシュボードの品質を高いレベルで確保することを目的としています。ダッシュボードのビジュアル、SQL、Python、UXがデータチームのメンバーによってピアレビューされ、ハンドブックに詳述されている基準を満たしていれば、ダッシュボードは本番に間に合います。

- 1. WIP: を名前に持つダッシュボードを作成し、WIPトピックに追加する
- 2. dbtやウェアハウスのドキュメントを活用して、クエリやチャートを構築する
- 3. ダッシュボードをレビューする準備ができたら、Sisenseダッシュボードレビューテンプレートを使用して、データチームプロジェクトにMRを作成します。何も変更する必要はなく、これは空のMRであるべきです。
- 4. テンプレートの指示に従う
- 5. データチームのメンバーにテンプレートを割り当て、レビューしてもらいます。誰にメールを送ればいいかわからない場合は、@gitlab-dataを使ってください。
- 6. すべてのフィードバックが与えられ、適用されたら、データチームのメンバーは右上のテキストタイルを更新し、ダッシュボードの作成者とレビュー者、最終更新日、関連する問題のクロスリンクを詳細に表示します(詳細は「データ分析プロセス」を参照してください)。

- 7. データチームメンバーのレビュアーは
  - ダッシュボードの名前を変更して「WIP:」ラベルを削除する

- ◇ WIPトピックからダッシュボードを削除する
- ダッシュボードに承認バッジを追加する
- 権限があればMRをマージするか、マージ権を持つ人に割り当てる
  - MRは、意味のある変更がない場合にも閉じることができます。その場合、データチームのレビュアーは、MRを閉じる前に「Approved」というラベルを追加します。

既存のダッシュボードのワークフローの更新

このセクションでは、データ・チームのピアレビュー・プロセスを経た既存のダッシュボードを更新する方法についてのワークフローを説明します。ダッシュボードが運用されると、新しいダッシュボードの作成とレビューのワークフロー全体を経ることなく、データアナリストとDRI/優先順位付けの所有者がダッシュボードへの段階的な追加を実施することができます。既存のダッシュボードを更新するには、以下の手順に従ってください。

- 1. 更新中のチャートのタイトルにWIP: を追加します。新しいチャートが追加されている場合は、タイトルにWIP: を追加します。
- 2. 新しいダッシュボードの作成とレビューのワークフローで定められているSQLスタイルガイドとSisense可視化のベストプラクティスに従って、既存のチャートに変更を加えるか、新しいチャートを追加します。
- 3. ダッシュボードをDRI/優先順位付け所有者がレビューする準備ができたら、Sisenseダッシュボードレビューテンプレートを使用して、データチームプロジェクトでMRを作成します。DRI/優先順位付け所有者に、MRの上部にある自分のセクションを記入してもらいます。データアナリストは、「既存のダッシュボードの更新」セクションのみを記入します。
- 4. ビジネス・ステークホルダー/DRIが更新をレビューして承認し、既存のダッシュボードへの更新チェックリストが完了したら、新規または更新されたチャートから「WIP:」を削除して、MRを閉じることができます。

マージリクエストによるダッシュボード、チャート、スニペット、またはビューの追加または更新

また、GitLab Data - Periscope Projectへのマージリクエスト (MR) により、チャート、スニペット、ビューを修正することができます。次の例をご覧ください。 例 または以下の手順に従ってください。

- 1. MRをGitLabDataに提出 -Periscopeプロジェクト
- 2. MRをスニペット/ダッシュボード/チャート/ビューのオーナーに割り
  - 当てる Periscopeプロジェクトでのスニペットオーナーの確認
  - 例 Sisenseでのスニペットオーナーの確認例
- 3. オーナーは "Approve "ボタンで承認する必要があります。
- 4. cc @gitlab-data
- 5. マージ権限を持つGitLab Dataチームのメンバーが、それをマー

ジします。 プロジェクトのリポジトリからの同期は双方向です。

ダッシュボードの自動更新を依頼する

データチームのプロジェクト課題を作成することで、ダッシュボードの自動更新を要求することができます。課題テンプレートを使用して、1つの特定のダッシュボードまたはダッシュボードの一括リストについて、指定された間隔での自動更新を要求できます。

Sisenseのチャートをハンドブックに埋め込む

これらのチャートをハンドブックに埋め込むのは、データチームではなく、ビジネスユニットの責任です。Sisenseには素晴らしい埋め込みドキュメントとアプリを通じたチャットサポートがあります。ハンドブックにチャートやダッシュボードを埋め込むには、主に3つの方法があります。

ハードコードされたHTML

HTMLを受け入れられるファイルタイプであれば、いつでもHTMLをハードコードすることができます。.htmlファイルがそのわかりやすい例です。しかし、Markdown (.md) や埋め込みRuby (.erb) ファイルでは、通常のHTMLにフォールバックすることができます。

ダッシュボード

ダッシュボード全体をハンドブックに埋め込むのとても簡単ですダッシュボードを埋め込むに、まずそのダッシュボードを外部共有 ダッシュボード」にする必要があり。次にURL文字列に?embed=trueを追加して、埋め込みリンクにしますそのURLを以下のように

差し込みます。

<iframe class="dashboard-embed" src="https://app.periscopedata.com/shared/string-of-numbers-here?
embed=true" height="700"> </iframe>。

SAFEスペースからダッシュボードを使用する場合、外部共有ダッシュボードを作成すると、リンクを持っている誰もがダッシュボードを利用できるようになるので、推奨しない。ハンドブック内でダッシュボードをスクロールする必要がないようにすることを目的としているため、高さの値を適切に調整する必要があります。プログラムでこれを行う方法はありません。

#### チャート

ハンドブックの埋め込みチャートは、常にsigned\_periscope\_urlへルパー関数を使用して生成する必要があります。この関数は、データチームのメンバーが手伝わなくても、自動的に署名されたURLを生成してくれます。これは、Sisense APIにさまざまなデータオプションを渡す実験をするときに特に便利です。このヘルパー関数は、SAFEスペースからチャートを埋め込もうとするとエラーを返します。SAFEスペース内のチャートへのURLのみを共有することを推奨する。

ファイルは .erb で終わる必要があります。index.html.mdという名前のファイルで作業している場合は、.erbを追加するだけで、ファイル名がindex.html.md.erbになります。

データを引数として関数に渡すだけです。サブアレイやオブジェクトなど、Sisense APIで必要とされるあらゆるデータを取ることができます。Sisenseのドキュメントには、embedAPIで利用可能なオプションの全リストがあります。

<embed width="100%" height="400px" src="<%= signed\_periscope\_url(chart: 6114177, dashboard: 463858, embed:
'v2' ) %>">。

この方法は、プレーンなMarkdownやHTMLファイルではレンダリング時にコードが実行されないため、機能しません。

ヒント:必要なPERISCOPE\_EMBED\_API\_KEYがCI変数として設定されているため、埋め込みチャートはローカルには表示されません。グラフが正しく表示されることを確認するには、MR内でレビューアプリを起動してください。

パフォーマンス指標 YML

- 名称:MRレート

is\_key: true

data/performance\_indicators.yml は、パフォーマンス・インジケータのコンテンツを載せたハンドブック・ページを、規約に従って自動的に生成するシステムの基礎となるものです。オブジェクトに periscope\_data プロパティとサブ値を与えると、テンプレートは自動的に署名付きの URL を生成して HTML を書いてくれます.これは、上記と同じ signed\_periscope\_url ヘルパー関数を使用します。このヘルパー関数は、SAFE空間のチャートを埋め込もうとするとエラーを返します。SAFEスペース内のチャートへのURLのみを共有することをお勧めします。

の説明を参照してください。MR率は、開発エンジニアが平均してどのくらいのMRをこなしているかを毎月評価するものです。
periscope\_data:
 チャート。
 6114177
 ダッシュボード463858 エ
ンベッド: V2

CSVやGoogle SheetsでSisense*に*データを追加する

CSVを使ってSisenseにデータをアップロードする方法はいくつかあります。これらの方法はすべて文書化されており、データチームのプロジェクト課題を作成してリクエストすることができます。リンクをクリックして選択肢を確認してください。

SisenseでCSVをアップロードする場合は、New ChartまたはNew ExplorationウィンドウでSnowflakeデータベースを使用していることを確認してください。そうでない場合は、GitLabの内部データモデルにアクセスするために、常にGitLab\_(Use\_this\_one!)をデフォルトにする必要があります。Sisenseでデータベースを変更する方法については、以下の画像をご覧ください。 Sisenseのデータベース

SisenseのCSVアップロード機能を使って、個人情報や機密データをSisenseにアップロードしないようにしてください、このデータはすべてのGitLabチームメンバーが一般にアクセスできます。

Sisenseからのデータの取り出し

Sisenseからデータを取り出す必要がある場合もあるでしょう。これが単発のケースであれば、いつでもUIからCSVをダウンロードすることができます。ただし、CSVダウンロードの最大サイズは500MBで、クエリは4分以内に完了しなければならないことに注意してください。このデータを定期的にシートなどに取り込む必要がある場合は、「編集」 $\rightarrow$ 「チャートフォーマット」 $\rightarrow$ 「詳細」 $\rightarrow$ 「CSVの公開URLを公開」でCSVの公開URLを公開することができます。シートでは、

=importdata("PUBLICURLHERE")を使用できます。

エクスポートしようとしているデータに機密性がなく、サイズが**500MB**以上で**100**万行以下の場合は、データチームのプロジェクト課題を作成します。

テキストファイルを操作するツールが必要な場合は、構造化されたテキストデータを操作するためのコマンドラインツールの リストをご検討ください。 これらの素晴らしいリソース.

ヒントとコツ

正しい会計年度の情報のみを表示するダッシュボードの設置

現在の会計四半期または次の会計四半期にのみフィルタリングするダッシュボードが必要な場合があります。Sisenseの既製の日付フィルターでは、カスタムの会計年度に対応できません。

あなたの分析では、次のように追加します。(datevalue)をフィルタリングしたい日付で更新してください)

LEFT JOIN analytics.date\_details on current\_date =
date\_actual WHERE [datevalue] < last\_day\_of\_fiscal\_quarter
AND [datevalue] > first\_day\_of\_fiscal\_quarter

ダッシュボードのクエリで当月を除外する

ほとんどの場合、クエリから現在の月を除外し、完了した月のみをレポートする必要があります。当月は不完全なので、これらの数字を表示すると誤解を招く恐れがあります。ダッシュボードのクエリで以下のステートメントを使用して、当月を除外してください。

WHERE <month\_column> < date\_trunc('month', CURRENT\_DATE)</pre>

タイムゾーン

SnowflakeのタイムスタンプはすべてUTCでなければなりません。

Sisenseの表示時間はPT(太平洋時間)に設定されています。これは、通信のガイドラインに沿ったものです。

created\_date=daterange]を使用する場合、Sisenseはcurrent\_timestampを使用し、PTに変換して比較します。例えば、10月4日の13:00PT (20:00UTC) に、過去3日間のデータをリクエストした場合、Sisenseは2019-10-02 07:00:00.000から2019-10-05 07:00:00.000までのフィルターを作成します。これらの時間はUTCで、PTの2019-10-02の開始時の午前0時とPTの2019-10-04の終了時の午前0時に対応し、つまりPTの丸3日間です。データベースがUTCで値を保存している場合(私たちはそうしています)、この比較はまさにあなたが望むものです。

エンドユーザーが気にするべきことは、チャートに表示するために日付をフォーマットすることです。タイムスタンプをPTに変換するには、[created\_at:pst]という構文を使用します。また、次の構文を使って日付に変換することもできます。

:日付を[created\_at:pst:date]のように指定します。これは、Salesforceなどのソースシステムの日付とSisenseで表示される日付を比較する際に必要です。

覚えておくべき重要なことは

merged.md

2021/11/8

• 特に指定や宣言がない限り、データベースやクエリの日付はすべてUTCです。

• Sisenseのdaterangeフィルターは、太平洋時間に基づ

いています。日付範囲フィルターの使い方

ダッシュボードのフィルターとして使用したい集約された日付がある場合は、集約された期間 を日付範囲の開始値とし、集計の終了日より1日少ない値を日付範囲の終了値とします。日付範囲の開始値は、日付期間にマッピングすることができます。

# DRS

日付範囲の終了日については、集計フィルタで選択した値に基づいて終了日を自動的に計算するために、クエリに追加のカラムを作成する必要があります。仮に、sfdc\_opportunity\_xf.close\_dateを気になる日付として使っていた場合、以下のような例になります。 dateadd(day,-1,dateadd([aggregation],1,[sfdc\_opportunity\_xf.close\_date:aggregation])) as date\_period\_end 次に、日付範囲の終わりのマッピングを追加します。

# DRE

特定のモデルが照会される場所を見つける

チャートやスニペットの生成に使用されるすべてのクエリは、Sisenseプロジェクトのperiscope/masterブランチにあります。検索フィールドにモデルや興味のあるキーワードを入力すると、関連するクエリを見つけることができます。

#### 不明瞭な視覚化を避ける

円グラフは、データを視覚化する方法としては、一般的にはあまり良くないとされています。このブログでは、円グラフを使ってはいけない理由を紹介しています。

高画質な画像の輸出

Sisenseから静的なチャートをエクスポートする際には、スクリーンショットを撮るのではなく、内蔵のエクスポート機能を使用してください。エクスポートすると、背景が透明な高品質の画像が作成されます。Persicopeから画像をエクスポートするには、任意のチャートの右上にある「その他のオプション」を選択し、「画像のダウンロード」を選択します。

#### ジダウンロードイメージ

Sisense Housekeeping

Sisenseは私たちのビジネス・インテリジェンス・ツールであり、Single Source of Truth(真実の情報源)として機能しています。私たちのSSOTとして、Sisenseは非常に高いレベルの清潔さ、整頓、正確さを維持する必要があります。

また、データチームが作成・承認したダッシュボードが正確で有益なものであることも必要です。また、データチームのメンバーが定期的にメンテナンスを行うことも必要です。

## 大原則。

- 承認されたダッシュボードはすべて機能し (チャートがすべて表示され)、正確なデータを表示しなければなりません。
- データチームのメンバーが作成したWIPダッシュボードはすべてWork in Progressでなければなりません。これは、プロジェクトがまだ進行中であり、チームの誰かがこのダッシュボードのリリースをまだ計画していることを意味します。
- 空のダッシュボードは持たないタイ
- トルのないダッシュボードは持たない
- SQLエラーが発生したダッシュボードは保存しません:エラーを修正するか、ダッシュボードをアーカ

イブします。未使用のダッシュボードの自動アーカイブ

現時点では、すべてのダッシュボードについて、自動アーカイブ機能が有効になっています。つまり、ダッシュボードがしばらく表示されないと

45日以上経過すると、自動的にアーカイブされます。アーカイブされたダッシュボードは削除されず、アーカイブ

を解除することができます。月例Sisenseクリーンアップ

エントロピーは、ビジネス・インテリジェンス・ツールの自然な状態ですが、避けることができます。この傾向に対処するた 301/243 merged.md

2021/11/8

めに、データチームはSisenseスペースで定期的にメンテナンス作業を行っています。

毎月、データチームのメンバーがメンテナンスを行います。これは、毎月末の1週間前に行われるData Opsミーティングの中で積極的に主張することができます。

このメンテナンスタスクは、毎月第1週に完了しなければなりません。そのためには、データチームのプロジェクトで新しい課題を開き、Sisense Cleanup Issueテンプレートを選択する必要があります。このテンプレートには、完了すべきタスクのリストが表示されます。すべてのタスクが完了したら、その課題を閉じることができます。

Sisense APIキー

何らかの理由でAPI Keyをローテーションする必要がある場合は、以下の場所でローテーション

- する必要があります。ハンドブックプロジェクト
- KPIスライドProject
- Okta

パフォーマンス指標のページ生成コードを担当するチームメンバーも、ローカルでページを構築するために必要になります。

layout: handbook-page-toc title:"Permifrost" 説明"スノーフレークパーミッションの管理"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

# クイックリンク

パーミフロストプロジェクト{:.btn .btn-purple-inv} PyPl{:.btn .btn-purple-inv}。

パーミフロスト

Permifrostは、Snowflakeデータウェアハウスのパーミッションを管理するためのPythonツールです。このツールを使用するための主なドキュメントは、プロジェクト内とPyPIにあります。

注意点とエラー

- 設定ファイルに存在しないオブジェクトがあってもエラーにならない ロール/ウェアハウ
- ス/データベースの作成と削除がパーミフロストで管理されない
  - ファイルからロール全体を削除しても、ロールは削除されません

開発

以下の手順で、仮想環境の作成と準備を行います。

```
# virtualenvの作成
python -m venv ~/.venv/permifrost

# virtualenvの起動
ソース ~/.venv/permifrost/bin/activate

# 依存関係のインストール
pip install -r requirements.txt

# pip3 開発に必要なものをすべてインストールする
pip install -e '.[dev]'
```

変更をコミットしたら、マージリクエストを送信し、デフォルトのテンプレートを更新します。

リリースプロセス

バージョニング

Permifrostでは、バージョン番号としてsemverを使用

しています。前提条件

続ける前に、最新のmasterブランチがローカルにあることを確認してください。

git fetch origin

ワークフロー

Permifrostはタグを使って成果物を作成しています。新しいタグをリポジトリにプッシュすると、DockerイメージやPyPIパッケージとして公開されます。

- 1. 上記の開発ワークフローに従って、インストールが最新の状態になっていることを確認します。
- 2. 以下のコマンドを実行してください。
  - # `origin/master` から `release-next` ブランチを作成してチェックアウトします git checkout -B release-next origin/master
  - # view changelog (変更点とログに記録された変更点の照合) changelog view
  - # チェンジログが検証された後、リリースにタグを付ける make type=minor release
    - # パッチリリースを行う場合は、 make type=patch release を実行 してください。
    - # メジャーリリースの場合は、make type=major release を実行してください。
  - # タグの確保 タグが作成されたら、先ほどバンプしたバージョンを確認します。 `0.22.0` => `0.23.0`.

git describe --tags --abbrev=0

- # タグを上流にプッシュしてリリースパイプラインを起動する git push origin \$(git describe --tags --abbrev=0)
- # リリースブランチをプッシュして新しいバージョンをマージし、マージリクエストを作成します git push origin release-next
- 3. マスターをターゲットとしたrelease-nextからのマージリクエストの作成
- 4. 変更がマージされたら、必ずソースブランチを削除するようにしてください。
- 5. パブリッシュパイプラインが成功すると、リリースはPyPIで公開されます。

layout: handbook-page-toc title:"データパイプラ

イン"このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .hidden-md .hidden-lg}.

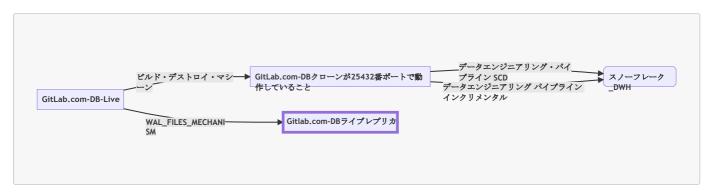
{::options parse\_block\_html="true" /}。

# 背景

データウェアハウスには、異なる抽出方法(Fivetran、Stich、Postgresパイプラインなど)を介して、異なるソースシステムからのソースデータが含まれています。このページでは、異なるデータソースと、データパイプラインを介してこれらのデータを抽出する方法について説明します。

# GitLab Postgresデータベース

Snowflakeにデータを取り込むための専用のgitlab.com read replicaデータベースインスタンスがあります。2つのレプリカが用意されており、それぞれが独自のレプリケーション頻度と動作を持っています。



- GitLab.com-DBのクローンは、10:30PM UTCに破壊され、11:15PM UTCに再構築されます。この結果、このデータベースへの問い合わせが可能な時間帯は、00:00AM UTCから11:30PM UTCとなります。この時間外にこのレプリカからデータを抽出した場合、エラーが発生しますが、データが失われることはありません。
- Gitlab.com-DB Live Replicaは、WALファイルを介して継続的にデータを取り込みます。

現在、安定したデータフィードを確保するために、インクリメンタルロードとフルロードの両方でGitLab.com-DBのクローンインスタンスを使用しています。開発やテストの段階では、Gitlab.com-DBのライブレプリカからのロードにおいて、書き込みと読み込みが同時に行われるという問題に直面しました(クエリのコンフリクト)。クエリの競合が回復するまでの時間を長くするために、max\_standby\_archive\_delayとmax\_standby\_streaming\_delayという設定があります。これはサーバー側で設定する必要があり、レプリケーション処理のラグを増やす結果になる可能性があります。これを避けるために、より静的なデータソースから読み出すようにしています。

クローンの設定に問題がないかどうかを確認するためのランブック項目。

モニタリング/アラート

Zlonkはジョブ完了モニタリングを実装しています。アラートはS4の深刻度で#alerts Slackチャンネルに流されます。以下のようなアラートが発生します。アラートが発生しました。

GitLab ジョブが失敗しました。GitLab ジョブの "clone" リソース "zlonk..." が失 敗しました。階層: db, タイプ: zlonk.postgres

私たちのパイプラインで、gitlab.comのpostgressデータベースにアクセスできない問題が発生した場合は、slackの@sre-oncallに連絡を取る必要があります。

インクリメンタルとフルロード

- 1. インクリメンタル・エキス
- これは最も便利な方法で、レコードの転送量は最小限で済みます。前提条件として、ソーステーブルレベルでデルタカラムが利用可能であること。

- 120個のテーブルを段階的に抽出 ロード時
- 間約1時間
- 6時間ごとに実行される
- 2. フル抽出 (現時点で±100テーブル&ロードタイ約4時間
- これは、ソース・テーブル・レベルでデルタ・カラムが利用できない場合や、ソースでレコードが削除された場合の代替手段です。
- 100個のテーブルを段階的に抽出していく
- ロードタイム 約4時間
- 24時間ごとに実行される

抽出方法は、マニフェストファイルによって決定されます。

PostgresからSnowflakeへのテーブルの手動バックフィル

このセクションは、gitlabのPostgresデータベースで以下の条件を満たしているテーブルに対してバックフィルを行う必要がある場合にのみ使用します。

- ◆ テーブルのエアフロータスクは最新の日付まで完了しています。
- キャッチアップがすでに完了している日付または日付範囲について、テーブルのデータが欠落している。
- テーブルをドロップしてリロードするとSLOに違反する場合。

上記のケースでは、この方法を使ってデータを取り出し、snowflakeのターゲットテーブルにロードすることができます。 この問題では、テーブルの1つに対して行われている

以下は、そのための手順です。

Step 1:-Kubernetesのポッドに接続します。

kubectl exec -ti <pod\_name> -c <webserver|scheduler> /bin/bash.

#### 現在の設定では

kubectl exec -ti airflow-deployment-7484d899c6-tfm8v -c scheduler /bin/bash

gcpクライアントが時々変更されるため、以下のようなエラーが発生することがあります。Unable to connect to the server: x509: certificate signed by unknown authority.この問題を解決するには、次のコマンドを実行します。 gcloud container clusters get-credentials data-ops --zone us-west1-a --project gitlab-analysis

Step 2:- pythonライブラリを使ってPostgresへの接続を確立します。これは、Postgresに直接アクセスできないため、Kubernetesポッドを経由する必要があります。データを抽出しなければならないテーブルは、それぞれのマニフェストファイルからクエリをピックアップします。以下は、クエリを実行して出力をcsvファイルで取得する手順です。この例では、テーブルexperiment\_subjectsに対して実行しています。

```
パイソン3
import psycopg2
conn =
psycopg2.connect(
host="",
database="",
user=""
password="",
port=,
options='-c statement_timeout=900000'
)
cur=conn.cursor()
sqlstr = "COPY (SELECT id , experiment_id, user_id, group_id, project_id, variant,
created_at, updated_at, converted_at, context FROM experiment_subjects WHERE updated_at
BETWEEN '2021-04- 05T19:14:20'::timestamp AND '2021-04-26T10:08:16'::timestamp) TO STDOUT
WITH CSV HEADER "
with open('experimental_subjects.csv','w') as f:
cur.copy_expert(sqlstr, f)
```

これにより、Kubernetesポッド自体にファイルが生成されます。注意:認証情報は安全な保管場所から受け取る必要があります。

Step 3: gsutilを使ってGCSからpostgres\_pipelineにファイルをアップロードする。

```
gsutil cp experiment_subjects.csv gs://postgres_pipeline/ です。
```

ステップ4:タイムスタンプ変数を準備します。これは、\_uploaded\_at列でファイルをリッチ化するためにデータフレームから プッシュするわけではありません。そのため、コマンドラインを使って行う必要があります。

```
import time
print(time.time())
```

ターゲットテーブルには \_uploaded\_at カラムがあり、入力する必要があるため、返された値は snowflake へのコピーの際に使用されます。ステップ 5: マニフェスト ファイルで定義された各列を以下のように読み込み、\_uploaded\_at タイムスタンプ列を以下のように追加します。これは、データが正しいかどうかを検証するためにsnowflakeで実行する必要があり、特定の監査を行うために使用することもできます。

```
SELECT $1 as id ,
experiment_idとして$2。
user_idとして$3。
$4 as group_id,
project_idとして$5。
バリアントとして6ドル。
created_atとして$7。
updated_atとして8ドル。
converted_atとして9ドル。
コンテキストとして $10, _uploaded_at として
'1621841251.8090081' FROM
`@stage`/experiment_subjects.csv;
```

ステージは、Postgresのパイプラインのステージです。

ステップ6: copy intoを使ってデータを読み込み、RAW.tap\_postgres.gitlab\_db\_experiment\_subjectsテーブルにこのデータを読み込みます。

GitLabデータベーススキーマの変更とデンジャーファイル

データエンジニアリングチームは、GitLabプロジェクトにデンジャーファイルを作成しています。ソーススキーマの変更が通知されることは、GitLab.com データベースからの抽出処理でエラーが発生しないようにするために不可欠です。なぜなら、抽出は一連の select 文を実行することで行われるからです。その日のトリアージを担当するデータエンジニアは、スキーマの変更を調査し、データチームが必要とするアクションのための課題を作成するDRIです。

トラステッドデータフレームワークテスト

このデータパイプラインでは、3種類のTrusted Data Frameworkテストが実行されます。

- 1. ソースの鮮度。
- 2. ロウカウントテスト。
- 3. データの実測値です。

dbt機能によりソースの鮮度を確認。

行数のカウントとデータの取得は、Airflowのgitlab\_com\_data\_reconciliation\_extract\_loadという追加のDAGを介して行われ、結果はRAW.TAP\_POSTGRES.GITLAB\_PGP\_EXPORTというsnowflakeテーブルに格納されます。

サービス ping

サービスpingは、GitLab Incが特定のGitLabインスタンスの使用状況データを収集するための方法です。サービスPing(旧称: Usage Ping)について、製品の観点からの詳しい情報はこちらをご覧ください。豊富なドキュメントを含む包括的なガイドは、ServicePingGuideに掲載されています。

サービスピンには大きく分けて2種類あります。

- セルフマネージドサービス
- Ping SaaSサービスPing

詳しくは「4種類のサービスPingプロセス」をご覧ください。

セルフマネージドサービスPing

セルフマネージド・サービス**Ping**は、**Versions**アプリからデータウェアハウスにロードされ、**VERSION\_DB**に格納されます。 データベースを使用しています。

SaaSサービスのPing

SaaS Service Pingは、2つの方法でデータウェアハウスにロードされます。

- Gitlab Postgres Database ReplicaのSQL文 (SQLベース) とRedisのRestFUL API
- ◆ コール (Redisベース) を使用した場合

データチームによる実装の詳細は、Readme.mdファイルに記載されています



インスタンスのSQLベースのメトリクスの読み込み

データはPostgresのSQLレプリカから読み込まれます。クエリは、この抽出物の中にある非常に大きなJSONファイル(ファイル内に数百のクエリがある)でバージョン管理されています。クエリは、インスタンスクエリと名前空間クエリの2つのカテゴリに分かれています。インスタンスクエリはGitLab.com全体のデータを生成し、名前空間クエリはGitLab.comの各名前空間のデータを生成します。データはテーブル(RAW.SAAS\_USAGE\_PINGスキーマ内)に格納されます。

- raw.saas\_usage\_ping.instance\_sql\_metrics
- RAW.SAAS\_USAGE\_PING.INSTANCE\_SQL\_ERROR このテーブルには、**SQL**メトリクスのデータ処理中にエラーがポップアップする **SQL**コマンドが含まれています。
- raw.saas\_usage\_ping.gitlab\_dotcom\_namespace

実装の詳細はsql-metrics-implementationで公開されています。 Redisベースのメト

リクスのインスタンスの読み込み

データはAPI経由でダウンロードされますので、APIの仕様をご参照ください。UsageDataNonSqlMetricsAPIを参照してください。データはJSON形式で保存され、約2k行のサイズです。通常、1回のロードで1ファイルとなります (現在は週1回のロードとなっています)。Redisからデータをロードする主な目的は、SQLクエリでは取得できない細かい粒度のメトリクスを確保することです。データはテーブル(RAW.SAAS\_USAGE\_PINGスキーマ内)に格納されます。

• raw.saas\_usage\_ping.instance\_redis\_metrics

実装についての詳細は redis-metrics-implementation で公開されています。

シートロード

SheetLoadは、GCSやS3からのGoogleシートやCSVをデータウェアハウスに取り込むためのプロセスです。SheetLoad

の使用方法に関する技術文書は、データチームプロジェクトのReadmeに記載されています。

GoogleシートやCSVをウェアハウスにインポートしたい場合は、データチームのプロジェクトで「CSVまたはGSheetsデータアップロード」の課題テンプレートを使って課題を作成してください。このテンプレートには、インポートしたいデータの種類と、そのデータを使って何をしたいかに応じて、詳細な手順が記載されています。

考察

{:#mind-about-sheetload}

SheetLoadは、主にスプレッドシートを正規のソースとするデータ(例:売上見積書)に使用されます。スプレッドシート以外のデータソースがある場合は、少なくとも新しいデータソースを作成してデータを自動的に取得するようにしてください。しかし、スプレッドシートがこのデータのSSOTであるならば、SheetLoadはデータをウェアハウスに取り込むための適切なメカニズムです。

#### スノーフレークへの読み込み

SheetLoadは、データベース内のテーブルを、読み込み元のシートの正確なコピーにするように設計されています。SheetLoad は、読み込み元のシートに変更があったことを検知すると、データベースのテーブルを削除し、更新されたスプレッドシートのイメージでテーブルを再作成します。つまり、列が追加されたり変更されたりしても、すべてデータベースに反映されるのです。変更は24時間以内に検出されます。

#### SheetLoadの準備

どうしても不可能な場合を除いて、SheetLoadシートはimportrange関数を使って元のGoogleシートから直接インポートするのがベストです。これにより、上流のシートを残したまま、SheetLoad版をプレーンテキストにフォーマットすることができます。追加のデータタイプの変換やデータのクリーンアップは、ベースとなるdbtモデルで行うことができます。(これはBoneyardには適用されません)。

読み込みのためにデータを準備する際の注意点をいくつかご紹介します。

- 数字の列はすべて、カンマやその他の記号を使わないようにフォーマットします \$1,000.46 の代
- わりに **1000.46** 可能な限りシンプルなヘッダーを使用します GitLab User Name の代わりに user\_name
- データがないことを示すために、空白のセルを使用します。空白のセルは、データベースにNULLとして保存されます。

#### モデリング

通常のSheetLoadデータがSisenseでアクセスできるようになる前に、dbtでモデル化する必要があります。ソースモデルとステージングモデルの最低2つのモデルがシートに対して作成されます。これらのモデルは、データチームのプロジェクトで課題を作成した後、データチームのメンバーによって作成されます。

## ボンヤード

ボーンヤード・スキーマは、スプレッドシートからデータをアップロードすることができ、Sisense内で直接クエリを利用することができます。ただし、これは1回限りの分析が必要な場合や、すでにウェアハウスにあるデータに結合する必要がある場合のためのものです。このデータは、アドホック/ワンオフのユースケースにのみ関連し、比較的短期間で陳腐化することを強調するためにBoneyardと呼ばれています。我々は定期的にBoneyardスキーマから古くなったデータを削除します。

## 証明書

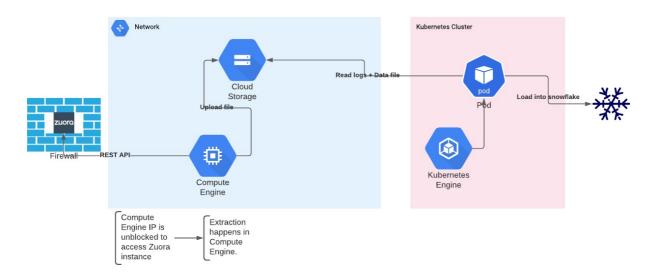
SheetLoadに証明書を追加する場合は、People Groupページの説明を参照してください。

## Zuoraの収益

Zuora Revenueは、最新の収益基準(ASC606、IFRS15)に準拠して、複雑な収益管理プロセスを自動化できるアプリケーションです。Zuora Revenue の抽出プロセスのデータパイプラインの一部として、REST Call を使用して Zuora Revenue の BI ビュー (Zuora Revenue が製品にデフォルトで作成したもので、我々はこれらのビューのみを使用し、これらのビューを作成または変更することはできませんでした) からデータを抽出しています。Zuora Revenue BIビューのデータエンティティは、主要な物理テーブルに基づいています。一部のBIビューは、ベースとなっている物理テーブルと類似しています。その他のBIビューは、物理テーブルに基づく計算から派生します。

#### Zuoraネットワークアーキテクチャ

Zuora システムはファイアウォールの後ろに置かれています。ファイアウォールを通過してアクセスするためには、GitLab のIPアドレスを許可リストに載せる必要があります。Kubernetes Engineは固定IPを持っていないので、Zuoraにアクセスするために固定IPを持ったコンピュートエンジンを追加しています。



このリストから、以下のテーブルだけが、Zuora RevenueにGitlabのデータが読み込まれているか、統合のために現時点で持っているデータです。

以下は、データがあり、snowflakeで作成されるテーブルのリストです。データを持たないテーブルは作成されません。これは、Zuora APIが提供するテーブル定義が、カラムの序列位置ではなく、アルファベット順になっているため、カラムのラベル付けが正しくないためです。

| シリ<br>ア小<br>え | データ主体<br>in Zuora        | 物理テーブル<br>での名前             | ZuoraViewでの名前          | S<br>ZuoraTableでの名前 | SFHas<br>レコード | プレゼント<br>スノーフレ<br>ーク |
|---------------|--------------------------|----------------------------|------------------------|---------------------|---------------|----------------------|
| 1             | AccountRPR(<br>タイプ       | O_BI3_ACCT_TYPE_VBI3_ACCT_ | _TYPEBI3_ACCT_TYPEYesY | es                  |               |                      |
| 2             | 会計<br>プレ<br>概要           | rpro_bi3_ri_acct_summ_v    | BI3_RI_ACCT_SUMM       | BI3_RI_ACCT_SUMM    | はい。           | はい。                  |
| 3             | 承認                       | rpro_bi3_appr_dtl_v        | BI3_APPR_DTL           | BI3_APPR_DTL        | はい。           | はい。                  |
| 4             | ビル                       | rpro_bi3_rc_bill_v         | BI3_RC_BILL            | BI3_RC_BILL         | はい。           | はい。                  |
| 5             | カレンダー                    | rpro_bi3_calendar_v        | BI3_CALENDAR           | BI3_CALENDAR        | はい。           | はい。                  |
| 6             | 削除されたス<br>ケジュール          | RPR_BI3_RC_SCHD_DEL_V      | BI3_RC_SCHD_DEL        | BI3_RC_SCHD_DEL     | はい。           | はい。                  |
| 7             | ヘッダー                     | rpro_bi3_rc_head_v         | BI3_RC_HEAD            | BI3_RC_HEAD         | はい。           | はい。                  |
| 8             | ホールズ                     | rpro_bi3_rc_hold_v         | BI3_RC_HOLD            | BI3_RC_HOLD         | はい。           | はい。                  |
| 9             | ライン                      | rpro_bi3_rc_lns_v          | BI3_RC_LNS             | BI3_RC_LNS          | はい。           | はい。                  |
| 10            | MJE                      | RPRO_BI3_MJE_V             | BI3_MJE                | BI3_MJE             | はい。           | はい。                  |
| 11            | РОВ                      | rpro_bi3_rc_pob_v          | BI3_RC_POB             | BI3_RC_POB          | はい。           | はい。                  |
| 12            | スケジュー<br>ル               | RPR_BI3_RC_SCHD_V          | BI3_RC_SCHD            | BI3_RC_SCHD         | はい。           | はい。                  |
| 13            | ウォーターフ<br>ォール(派生<br>型    | rpro_bi3_wf_summ_v         | BI3_WF_SUMM            | BI3_WF_SUMM         | はい。           | いいえ                  |
| 14            | オーガ                      | RPRO_BI3_ORG_V             | BI3_ORG                | BI3_ORG             | いいえ           | いいえ                  |
| 15            | アカウント<br>概要<br>(Derived) | rpro_bi3_In_acct_summ_v    | Bi3_In_acct_summ       | Bi3_In_acct_summ    | いいえ           | いいえ                  |
| 16            | 本                        | rpro bi3 book v            | BI3 BOOK               | BI3 BOOK            | いいえ           | いいえ                  |

| 17 | コスト      | rnro hi3 rc ln cost v | RI3 RC IN COST  | RI3 RC IN COST      | いいタ  | いいき   |
|----|----------|-----------------------|-----------------|---------------------|------|-------|
| え  |          | C 4742 Hil            |                 |                     |      | ーク    |
| 不水 | in Zuora | での名前                  | Zuoraview Coo和前 | Zuorarable (V)和制    | レコード | スノーフレ |
| シリ | データ主体    | 物理テーブル                | ZuoraViewでの名前   | S<br>ZuoraTableでの名前 | FHas | プレゼント |

Zuora収入エキス

▶ 以下は、「Zuora Revenue」の抽出周辺情報です。 抽出パイプラインは、ここをクリックすると拡大します。

Compute engineでの環境設定

自分のサービスアカウントを使ってzuora compute engineにSSH接続します。以下はGCPのサーバーの詳細です https://console.cloud.google.com/compute/instancesDetail/zones/us-west1-a/instances/zuora-revenue-extract-server? project=gitlab-analysis&rif\_reserved

ssh -o UserKnownHostsFile=/dev/null -o CheckHostIP=no -o StrictHostKeyChecking=no -i HOME/.ssh/google\_compute\_engine -A -p 22 @<external\_ip> です。

別の端末から~/repos/analytics/extract/zuora\_revenue/srcディレクトリに移動し、以下のコマンドを実行してディレクトリ全体をcompute engineにアップロードします。このディレクトリには、抽出処理のコードが含まれています。

gcloud compute scp --recurse src --zone "us-west1-a" zuora-revenue-extract-server:/home/vedprakash/zuora\_revenue 🗅

れにより、ローカルのブランチからcompute engineのブランチにsrcフォルダがアップロードされます。

接続後、ファイルをアップロードする。

Step 1:コンピュートエンジン内に仮想環境を作成 --- zuora-revenue-extract-venv

変更を最小限にするために、同じ名前にしておきま

す。 python3 -m venv zuora-revenue-extract-

venv Step 2: venvの有効化

source /home/vedprakash/zuora-revenue-extract-venv/bin/activate ス

テップ3:pipをアップグレードした後

pip install --upgrade pip

ステップ4: srcフォルダに移動し、必要なパッケージをすべてインストールします。

pip install -r requirements.txt

注意事項ステップ1からステップ4は、環境がクラッシュしていて最初から構築しなければならない場合にのみ必要となり、一般的な 運用では必要ありません。

以下のステップは、GCSバケットフォルダを誤って削除してしまった場合に必要となります。また、新しいテーブルをシステムに追加する必要がある場合にも、以下のステップを使用することができます。

ステップ5: テーブル名とロードデートの情報を保持するstart\_date\_<テーブル名>.csvファイルを作成する。例えば、テーブルBI3\_MJEの場合、ファイル名はstart\_date\_BI3\_MJE.csvとなり、ファイルの内容は以下のようになります。

table\_name,load\_date BI3\_ACCT\_TYPE 新しいテーブルのload\_dateは、ファイルを最初からダウンロードし始めるので、空白にしておきます。それ以外の場合は、エアフローログから最終ロード日を拾うことができます。現在のテーブルについて、ファイルを作成するためのコマンドリストを以下に示します。これは、ローカルからでも、コンピュートエンジンからでも可能です

313/243

echo "テーブル名,load\_date BI3\_ACCT\_TYPE," > start\_date\_BI3\_ACCT\_TYPE.csv echo "テーブル名,load\_date BI3\_APPR\_DTL," > start\_date\_BI3\_APPR\_DTL.csv echo "テーブル名,load\_date

```
BI3_CALENDAR," > start_date_BI3_CALENDAR.csv
echo "テーブル名,load_date
BI3_MJE," > start_date_BI3_MJE.csv
echo "テーブル名,load_date
BI3_RC_BILL," > start_date_BI3_RC_BILL.csv
echo "テーブル名,load_date
BI3_RC_HEAD," > start_date_BI3_RC_HEAD.csv
echo "テーブル名,load_date
BI3_RC_HOLD," > start_date_BI3_RC_HOLD.csv
echo "テーブル名,load date
BI3_RC_LNS," > start_date_BI3_RC_LNS.csv
echo "テーブル名,load_date
BI3_RC_POB," > start_date_BI3_RC_POB.csv
echo "テーブル名,load_date
BI3_RC_SCHD," > start_date_BI3_RC_SCHD.csv
echo "テーブル名,load_date
BI3_RC_SCHD_DEL," > start_date_BI3_RC_SCHD_DEL.csv
echo "テーブル名,load_date
BI3 RI ACCT SUMM," > start date BI3 RI ACCT SUMM.csv
```

このコマンドは、各テーブルのファイルを作成し、必要なカラム名と値を入力します。load\_dateにはnullを設定していますが、これは初回実行時の扱いとなるためです。注:ロード日がわかっている場合は、2016-07-26T00:00:00形式の%Y-%m-にします。

特定のテーブルの %dT%H:%M:%S。

Step6: ステージングエリアにファイルをアップロードする必要があります。以下は、各ファイルをステージング・エリアの各テーブルにアップロードするためのコマンド群です。

ステップ7: 抽出液を実行するには、サーバーの.bash\_profileファイルで以下の変数を宣言する必要があります。

```
export zuora_bucket=""
export zuora_dns=""
export zuora_dns=""
export authorization_code=""
export python_venv="source /home/vedprakash/zuora-revenue-extract-venv/bin/activate" #From step

2
エクスポート zuora_extract_log="/home/vedprakash/zuora_revenue/src/logs/"
export zuora_src="/home/vedprakash/zuora_revenue/src" #ソースコードのパスです。
```

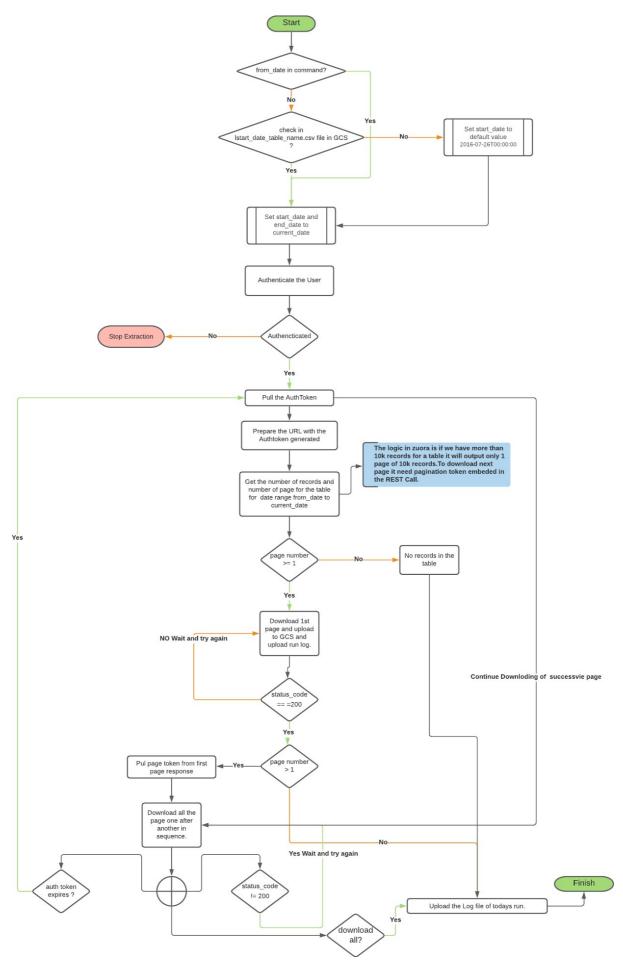
注:認証情報はzuora revenue prodの下の1つのパスワードに存在します。

ステップ8:最後のステップは、スケジュールを行うことです。crontabに以下のコマンドを追加します。編集して準備ができたら、必要なコマンドをそのマシンのcrontabに追加します。現在のスケジュールは、毎日午前2時(UTC)に実行されるように設定されています。

```
00 02 * * * . $HOME/.bash_profile;$python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3_ACCT_TYPE-バケット名 $zuora_bucket -api_dns名 $zuora_dns -
api_auth_code "$authorization_code" &>/tmp/mycommand.log
```

```
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3_APPR_DTL
                           -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3_CALENDAR
                             -バケット名 $zuora_bucket-api_dns名 $zuora_dns-api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
                            -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
-テーブル名 BI3_MJE
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3_RC_BILL
                            -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
"$authorization code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3 RC HEAD
                             -バケット名 $zuora_bucket-api_dns名 $zuora_dns-api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile;$python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3 RC HOLD
                            -バケット名 $zuora_bucket-api_dns名 $zuora_dns-api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3 RC LNS
                            -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3 RC POB
                            -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile; $python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
                             -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
-テーブル名 BI3_RC_SCHD
"$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash profile;$python venv && cd $zuora src && python3 extract zuora revenue.py
                             BI3_RC_SCHD_DEL-bucket_name $zuora_bucket -api_dns_name $zuora_dns
-table name
-api_auth_code "$authorization_code" &>/tmp/mycommand.log
00 02 * * * * . $HOME/.bash_profile;$python_venv && cd $zuora_src && python3 extract_zuora_revenue.py
-テーブル名 BI3_RI_ACCT_SUMM -バケット名 $zuora_bucket -api_dns名 $zuora_dns -api_auth_code
"$authorization_code" &>/tmp/mycommand.log
```

Zuora Extractのフローチャート



プロセスの最後には、以下が出力されます。

1. zuora\_revpro\_gitlab/RAW\_DB/staging/<table\_name>/<table\_name>\_DD-MM-YYYY.log.logという名前の成功ログファイルが存在します。たとえば、テーブル BI3\_MJE の場合、その日のログ ファイルの名前は BI3\_MJE\_21-06-2021.log で、パス gs://zuora\_revpro\_gitlab/RAW\_DB/staging/BI3\_MJE/BI3\_MJE\_21-06-2021.log にアップロードされます。

2. GCSバケットには、その日付範囲のファイルがすべて存在します。

Snowflakeに抽出用のテーブルを追加するためです。

RAW.ZUORA REVENUEスキーマにテーブルを作成します。

このリストから、いずれかのテーブルがデータを取得し、そのエントリをsnowflakeに追加する必要がある場合は、以下の手順を実行する必要があります。ステップ1:エントリーextract/zuora\_revenue/zuora\_revenue\_table\_name.ymlをzuora\_revenue load snow DAGのadd taskに追加します。ステップ2: それぞれのテーブルについて、ステップ5、ステップ6、

ステップ**8**に従います。 ステップ**3**: ダウンロードしたファイルからカラム名をピックアップし、すべてのカラムをvarcharに設定して、snowflake

注意事項テーブルの定義を作成する際に describe column API を使用しないでください。

派生したテーブルの場合

Zuora は、派生ビューのためのビュー定義を提供しています。本番環境では、派生ビューからデータを抽出することはできません。そのため、テーブルBI3\_WF\_SUMMについては、Zuoraから提供されたDDLを使ってPREP層のDBTモデルでデータを準備します。DDLの定義は、レポのextract/zuora\_revenue/README.mdにあります。

## Zoominfo

ZoomInfoは、B2BのセールスおよびマーケティングチームのためのGo-to-Market Intelligenceプラットフォームです。この統合されたクラウドベースのプラットフォームは、営業担当者やマーケティング担当者に包括的な情報を提供し、潜在的な新規顧客を見つけるのに役立ちます。このような充実したデータを得るためには、GitlabはZoominfoにアウトバウンドでデータを送信し、処理後にGitLabはZoominfoからSnowflakeデータシェアを介してインバウンドテーブルとして処理されたデータを受け取る必要があります。

Zoominfoのデータパイプラインは、Snowflakeデータ共有手法を活用した自動化された双方向のデータパイプラインです。

Snowflake Dataのシェアです。

[Snowflake data share] (https://docs.snowflake.com/en/user-guide/data-sharing-intro.html)は、1つのアカウントからsnowflake データベーステーブルの共有を可能にし、また、外部アカウントから共有されたデータへのアクセスを可能にします。これには、自分のアカウントでデータベースのアウトバウンド共有を作成し、Webインターフェース/SQLのいずれかを使用して外部アカウントに共有する必要のあるsnowflakeテーブルへのアクセスを許可する必要があります。

SQLを使ったSnowflake Data Share。

以下は、SQLを使用してsnowflakeデータシェアを介してアウトバウンド/インバウンドシェアを操作するための手順です。

SQLを使ったアウトバウンドシェア。

例えば、prod というデータベースに share というスキーマを持ち、gitlab\_user\_outbound というテーブルを消費者アカウント azitest と共有する場合を考えます。以下のSQLを実行してアウトバウンドシェアを作成します。

- ステップ1: accountadminロールを使ってShareを作成します。
  - USE ROLE accountadmin;
  - CREATE SHARE share test;
- ステップ2: 特権を与えてデータベース、スキーマ、テーブルを共有に追加する。
  - GRANT USAGE ON DATABASE prod TO SHARE share\_test;
  - GRANT USAGE ON SCHEMA prod.share TO SHARE share\_test;
  - GRANT SELECT ON TABLE prod.share.gitlab\_user\_outbound TO SHARE share\_test;
- ステップ3:消費者アカウントを共有に追加します。アカウントを共有に追加するためには、消費者はアカウントの 詳細を提供する必要があり、消費者とプロバイダーの両方のアカウントが同じスノーフレークリージョンにある必 要があります。

ALTER SHARE share\_test ADD ACCOUNTS = 'azitest';

SQLを使ったインバウンドシェア。

例えば、gitlabという名前の共有がアカウントazitestから私たちに共有されている場合、以下の**SQL**を実行して**snowflake**にデータベースを作成し、受信した共有内のテーブルとデータにアクセスします。

\* CREATE DATABASE zoominfo\_inbound FROM SHARE azitest.gitlab;

Snowflake Webインターフェースを使用したSnowflake Data Share。

以下は、Webインターフェースを使用してsnowflakeデータ共有を介してアウトバウンド/インバウンド共有を操作するための手順です。accountadminロールを使用して、snowflake Webインターフェイスの共有ページに移動し、インバウンド/アウトバウンドのデータ共有タスクを実行します。

📝イメージ-1.png

snowflakeのウェブインターフェイスを使ったアウトバウンドシェア。

- ステップ1: アウトバウンド共有を作成するには、snowflakeウェブインタフェースの共有ページで、アウトバウンドアイコンをクリックし、次に+createアイコンをクリックします。
- ステップ2:安全な共有名、データベース、テーブル/ビューの詳細を追加し、下部の「作成」ボタンをクリ

ックします。 📄

• ステップ3: お客様のアカウントを共有に追加し、リーダー/フルアカウントの作成を選択して、下部の追加ボタンをクリックします。これにより、共有が作成され、テーブルやビューが消費者に共有されます。

☑イメージ-3.png

snowflakeのウェブインターフェイスを使ったインバウンドシェア。

インバウンドシェアは、snowflakeウェブインタフェースのシェアページの「インバウンド」タブで確認できます。インバウンド共有のテーブルやデータにアクセスするためには、共有データベースを作成する必要があります。共有データベースを作成するには、「Create database from secure share」アイコンをクリックし、データベース名を入力し、アクセス権を付与して「create database」ボタンをクリックします。このプロセスにより、snowflake内に zoominfo\_inbound データベースが作成されます。受信テーブルとデータは、snowflakeのこの共有データベースの下でアクセスできます。共有データベースからのデータはプレップに取り込まれる。

Image-4.png

Image-5.png

建築シイメー

ジ.png

アウトバウンドテーブル。

• "PROD"...SHARE"...GITLAB\_USER\_OUTBOUND - Outbound tableは、First name, Last name, email address, company nameなどのGitlabのユーザー情報を持ちます。アウトバウンドテーブルは、Snowflakeのデータ共有を介してZoominfo に一度だけ共有されます。テーブルはdbtで作成されるため、時間の経過とともに変更されます。新しく更新されたデータをこのテーブルに取り込むのは、Zoominfoの責任です。

Inboundテーブルの読み込み。

Zoominfoは、インバウンドファイルをSnowflakeデータシェア経由でGitlabに送信します。共有データベース ZOOMINFO\_INBOUNDは、WebインターフェイスまたはSQLを使ってインバウンドシェアから作成されます。受信テーブルは、データをRaw層に取り込み、dbtがデータを変換するという標準的なアーキテクチャには従いません。余分なプロセスを作らず、パイプラインをより効率的にするために、共有データベースを生のデータベースとして扱います。のインバウンドテーブ

ルからのデータは

**ZOOMINFO\_INBOUND**はSnowflake Data Exchange loaderを使用してSnowflake PREPに取り込まれます。以下のインバウンド・テーブルのリストは、'zoominfo'スキーマの下でPREPデータベースに作成されます。

- "ZI\_COMP\_WITH\_LINKAGES\_GLOBAL" 国際テーブルは、情報を持っているすべての企業のリストを持っています。これは1 回の共有です。
- "ZI\_REFERENCE\_TECHS" データベースに登録されている企業が使用しているテクノロジーのリストを持つテクノグラフテーブル。
- "GITLAB\_CONTACT\_ENHANCE" User table company マッチしたテーブルで、Gitlab が zoominfo に送信するユーザーリスト に会社情報を追加します。GitlabがZoominfoに送信するのは一度だけですが、付加されたデータは四半期ごとに更新することができます。

layout: handbook-page-toc title: "Pythonガイド" 説明"私たちが選んだリンターがすべてを捕らえているわけではないので、この Python スタイルガイドを施行するのは私たちの集団の責任です。"

このページについて

{:.no toc.hidden-md.hidden-lg}

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

# パイソンガイド

私たちが選んだリンターがすべてを捕らえるわけではないので、このスタイルガイドを実施するのは*私たちの責任*です。

価値観

Campsite rule - このガイドラインはまだ完成していないため、現在スタイルガイドに準拠していないコードを使用している場合は、見かけたら更新してください。

リンティング

リンターにはBlackを使用しています。初期設定のままで使用しています。

レビュー段階では、レポ全体をlintし、ファイルをフォーマットする必要がある場合はゼロ以外の終了コードを返す手動CIジョブがあります。MRがマージされる前にこのジョブが通過するかどうかは、MRの作者とレビュアーの両方にかかっています。レポ全体をlintするには、レポの先頭からblack。を実行します。

スペーシング

PEP8では、コードの論理的な部分を空行で囲むことを推奨しています。forループやif/elseの開始時には ブロックでは、セクションの上に改行を追加して、コードに余裕を持たせます。改行のコストは安いですが、頭を使う時間は

高くつきます。タイプのヒント

すべての関数のシグネチャには、たとえNoneであっても、戻り値の型を含めた型ヒントを含めるべきです。これは良いドキュメントであり、型チェックやエラーチェックのために mypy で使用することもできます。

例

```
def foo(x: int, y: int) ->
    int: """
    2つの数字を足して返します。"""
    x + yを返す

def bar(some_str: str) ->
    None: """
```

```
文字列を印刷します。"""
print(some_str)
return
```

## インポートオーダー

輸入品はPEP8のルールに従うべきであり、さらに、前に出てくる輸入品...ステートメントと一緒に注文する必要があります。 from .... import ...

例

```
import logging
import sys
from os import environ

import pandas as pd
from requests import
get

import some_local_module
from another_local_module import something
```

#### **Docstrings**

ドックストリングはすべての関数で使用する必要があります。関数のシグネチャで型のヒントを使用しているので、各パラメータを記述する必要はありません。ドックストリングは3重のダブルクォートを使用し、句読点を含む完全な文章でなければなりません。

例

```
def foo(x: int, y: int) ->
    int: """
    2つの数字を足して、その結果を返す。"""
    x + yを返す

def bar(some_str: str) ->
    None: """
    文字列を印刷します。

これもちゃんとした文章です。"""
    print(some_str)
    return
```

## 環境変数の統合方法

関数をできるだけ再利用できるようにするために、(*よほど*の理由がない限り)関数の中で環境変数を直接使うことは非常にお勧めできません(以下にその例を示します)。その代わりに、使いたい変数を指定して渡すか、すべての環境変数を辞書として渡すのがベストな方法です。これにより、任意の辞書を渡して互換性を持たせることができるとともに、環境レベルで変数を定義する必要がなくなります。

例

```
import os
from typing import Dict
## Don't do this!
def foo(x: int) ->
   int: """
   2つの数字を足して返します。"""
   return x +
os.environ["y"] foo(1)
## Do this!
env_vars = os.environ.copy() # copyメソッドはenv varsの通常のdictを返します。 def
bar(some_str: str, another_string: str) -> None:
   2つの文字列を連結して表示します。"""
   print(some_str +
   another_string) return
bar("foo", env_vars["bar"])
## Or do this!
def bar(some_str: str, env_vars: Dict[str, str]) ->
  None: """
   2つの文字列を連結して表示します。"""
   print(some_str + env_vars["another_string"])
   return
bar("foo", env_vars)
```

#### パッケージのエイリアス

一般的なサードパーティ製パッケージには、いくつかの標準的なエイリアスを使用しています。それは以下の通りです。

- import pandas as
- pd import numpy as
   np

# 変数名の命名規則

名前に型を加えることは、自己文書化に適したコードです。可能であれば、変数には常に、特にデータ型に関して説明的なネーミングを使用してください。以下はその例です。

- data\_dfは、データフレーム
- params\_dictは、辞書
- retries\_intは、整数
- bash\_command\_strは、文字列

定数を関数に渡す場合は、渡される各変数に名前を付けて、それぞれが何であるかが明確になるようにします。最後に、

変数名が重複しないようにしましょう。

例

```
def bar(some_str: str, another_string: str) ->
None: """
2つの文字列を連結して表示します。"""
print(some_str + another_string)
```

```
## Do this!
bar(some_str="foo", another_string="bar")

## Or do this!
some_str = "foo"
another_string =
"bar"
バー(some_str, another_string)

## But don't do this!
bar(some_str=some_str, another_string=another_string)
```

### スクリプトを実行可能な状態にする

スクリプトを作成したばかりなのに実行できない場合は、実行可能な状態にする必要があります。以下を実行してみてください。

```
chmod 755 yourscript.py
```

chmod 755の説明は、このaskubuntuのページをご覧ください。

変更可能なデフォルト関数の引数

関数のデフォルト引数に可変型データ構造を使用すると、コードにバグが発生することがあります。これは、新しい ミュータブルなデータ構造は、関数の定義時に一度だけ作成され、その後の呼び出しのたびにデータ構造が使用されます。

例

```
def append_to(element,
to=[]):
to.append(element)
に戻る

my_list =
append_to(12)
print(my_list)

my_other_list =
append_to(42)
print(my_other_list)
```

出力します。

```
[12]
[12, 42]
```

参考: https://docs.python-guide.org/writing/gotchas/ 新規

抽出のためのフォルダ構成

- クライアント固有のロジックはすべて/extractに格納し、再利用可能なAPIクライアントはすべて オーケストレーション
- パイプライン固有の操作は、/extractに格納する必要があります。
- extract のフォルダ構造には、スクリプトが複数のデータセットを抽出する場合、 extract\_qualtrics\_mailingsends や extract\_qualtrics のように、 extract\_{source}\_{dataset\_name} というファイルを含める必要があります。このスクリプトは、 extractの主な機能と考えることができ、 extract DAGの出発点として実行されるファイルです。

Pythonを使用しない場合

このスタイルガイドはデータチーム全体を対象としていますので、Pythonを使うには時と場所があり、それは通常データモデリングの段階以外であることを覚えておくことが重要です。データ操作には可能な限りSQLを使用してください。

ユニットテスト

PytestはAnalyticsプロジェクトでユニットテストを実行するために使用されます。テストはプロジェクトのルートディレクトリからpython\_pytest Clパイプラインジョブで実行されます。このジョブは、テスト結果のJUnitレポートを生成し、GitLabで処理されてマージリクエストに表示されます。

新しいテストの作成

新しいテストファイルの名前は、test\_\*.py というパターンに従うべきです。そうすれば、pytest が見つけてくれますし、リポジトリの中でも簡単に認識できます。新しいテストファイルは、test という名前のディレクトリに置く必要があります。 test ディレクトリは、テストされるファイルと同じ親ディレクトリを共有する必要があります。

テストファイルは、1つ以上のテストから構成されます。個々のテストは、1つまたは複数のPythonのアサート文を持つ関数を定義することで作成されます。アサートがすべて真であれば、テストはパスします。一つでも偽のアサートがあれば、テストは失敗します。

インポートを書く際には、テストはルートディレクトリから実行されることを覚えておくことが重要です。将来的には、必要に応じてテストを容易にするために、PythonPathに追加のディレクトリを追加することができます。

例外処理

APIからデータを抽出するPythonクラスを書いた場合、そのクラスはAPIプロセスのエラーをハイライトする責任があります。 データモデリング、ソースの鮮度やフォーマットの問題は、dbtテストを使ってハイライトされるべきです。

一般的なtry/exceptブロックの使用は避けてください。

```
# Don't do
this! try:
  print("Do something")
except:
  print("あらゆるタイプの例外を捕捉")
# Do this
while maximum_backoff_sec > (2 **
   n): try:
       print("Do something")
   except APIError as gspread_error:
       if gspread_error.response.status_code in (429, 500, 502,
           503): self.wait_exponential_backoff(n)
       1 else:
           レイズ
の他にもあります。
   error(f "Max retries exceeded, giving up on {file_name}")
```

layout: handbook-page-toc title:"SiSense Style Guide" description:"SiSense Style Guide" このページでは

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

全体的なルック&フィール

KPIチャート

## 建設中

フィルタリング可能なレポート

### 建設中

適切なチャートの選択

### 建設中

カラー、ラベル、データ値

標準的なカラーパレット

SiSenseのチャートは、GitLabの標準的なカラーパレットを使用しています。

丸め

丸めるべきか、丸めないべきか?聴衆を知る。対象者や分析のユースケースに適した数値を提示する。

- グラフによっては、合計が100%になるために、パーセンテージの有効数字1桁または2桁が必要な場合があ
- ります。通貨を除いて、一般的に数字は完全な形で表示する必要があります。

# 通貨

- すべての通貨は米ドルで表示されています。
- 10,000ドル以上の場合は、各000がkに置き換えられ、例えば10,000ドルの代わりに10kとなります。
- 1,000,000ドル以上では、各000,000がmに置き換えられ、例えば10,000,000ドルの代わりに10mとなります。

# データと時間

会計日は、DIM\_DATE ディメンション・テーブルから抽出する必要があり

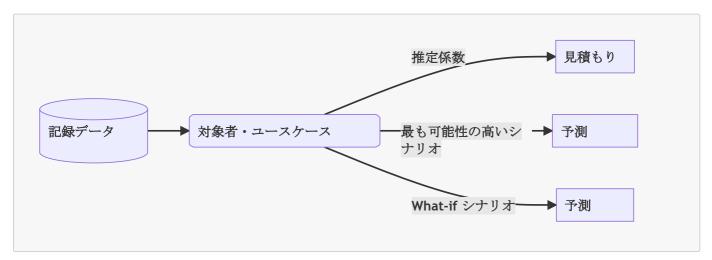
ます。日付の書式は、GitLabWritingStyleGuidelinesに準拠する必要が

## あります。

- 日付はyyyy-mm-ddとなります。
- 時刻は、UTCを使用した24時間表示です。 四半
- 期はQn、例:Q1
- 会計年度はFYyy、例: FY21
- 年度と四半期はFYyy-Qn、例: FY21-Q2

記録されたデータと計算されたデータ

データは、特定の対象者やユースケースのために、プレゼンテーションや報告のために準備されます。記録されたデータは、 すべての計算されたデータの基礎となります。



• 記録されたデータ - 検証可能な情報源と観察可能なイベントに由来する「事実」のデータ。検証を助けるために、記録データには、データ作成者/ソースの名前、データ取得日時、イベントが発生した場所などの監査メタデータが一般的に含まれている。

- 推定データ 記録されたデータに推定係数を加えて計算したデータ。推定係数は、通常、意味のある適切な期間の過去のデータトレンドに基づいています。エスティメーションは、あらゆる業界や分野で広く使用されています。
- 予測データ 過去の記録データに加えて、「最も可能性の高い」将来のシナリオに基づいて評価した計算データ。予 測はファイナンシャル・プランニングによく用いられる。
- 予測データ 過去の記録データに基づいて計算されたデータに、「what-if」の将来シナリオに基づいた評価を加えたもの。

プレゼンテーション

記録されたデータには特別な表示は必要ありませんが、計算されたデータには特別な表示が必要です。計算さ

- れたデータは、グラフのタイトルや凡例などで、*常に*明確に表示する必要があります(例: "Seats "の代わりに"Estimated Seats "を使用)。
- チャートの要素を明確に識別すること(例:同じチャートに計算タイプが混在している場合、異なるラインスタイルを使用すること

カスタムPythonモジュール

Pythonモジュールは、再利用可能で高品質でクリーンなコードを一箇所で維持するのに役立ちます。

カスタムPythonモジュールを開始/変更するには、Sisenseプロジェクトのperiscope/masterブランチでMRを開始します。Sisenseで利用可能なカスタムモジュールは、チャート作成ページの左サイドメニューでいつでも確認することができます。

シカスタムモジュール

カスタムモジュールについての公式ドキュメントはこちらをご覧ください。

layout: handbook-page-toc title:スノープラウの説明除雪機のインフラ管理

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg},

スノープラウの概要

Snowplowは、オープンソースのイベント分析プラットフォームです。SnowplowのSaaSを運営する事業体があり、彼らはオープンソース製品のコードも保守しています。Snowplowの一般的なアーキテクチャの概要はGitHubに掲載されており、どのように動作し、どのように設定されているのかという基本的な部分が詳しく説明されている。

2019年6月、Snowplowイベントの送信をサードパーティからGitLabが管理するインフラへの送信に切り替え、このページに記録しました。データチームの視点では、サードパーティの実装とあまり変わりませんでした。イベントはコレクターとエンリッチャーを経て、S3にダンプされます。

### GitLabの実装

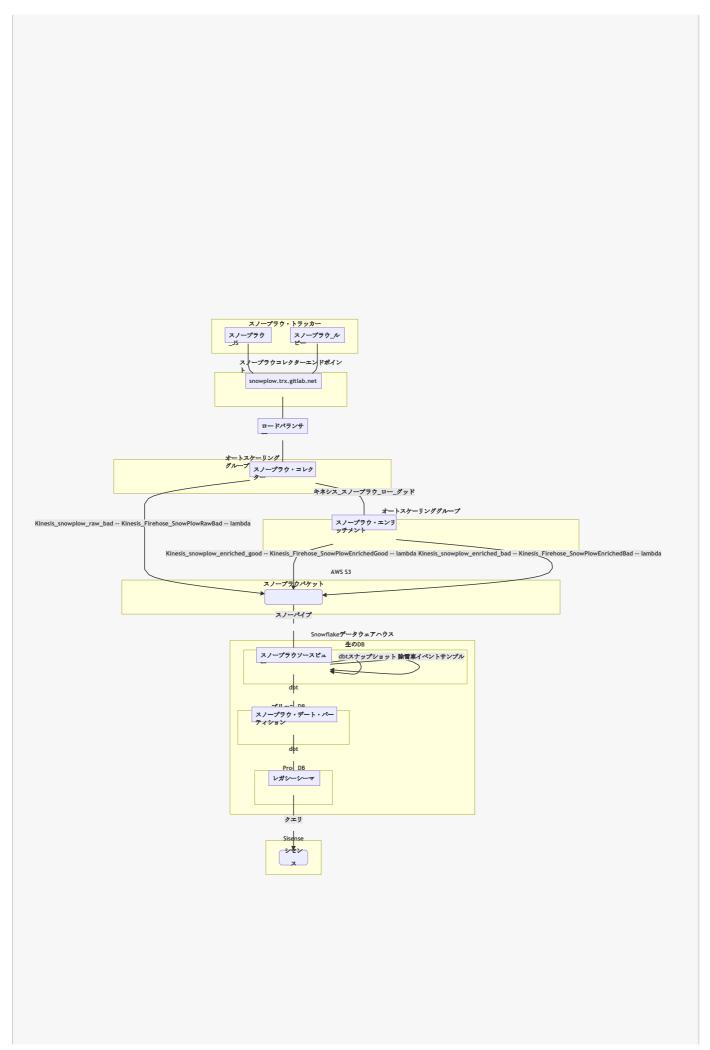
Snowplowのインフラをサードパーティのホスティングサービスから1st-partに移行するためのオリジナル設計書は、インフラ設計ライブラリに記載されています。これは構築を開始する前に書かれたもので、前提条件や設計上の決定事項の多くが含まれています。

SnowplowはTerraform on AWSで構築されており、gitlab-com-infrastructureプロジェクトでドキュ

メント化されています。セットアップの詳細については、GitLabUnfilteredの内部ビデオをご覧く

ださい。

| merged.md | 2021/11/8 |
|-----------|-----------|
|           |           |
|           |           |



**S**3

EnrichedイベントはTSV形式でs3://gitlab-com-snowplow-events/output/というバケツに格納されます。Badイベントは s3://gitlab-com-snowplow-events/enriched-Bad/ に JSON 形式で保存されます。どちらのバケツにも、/YYYY/MM/DD/HH/<data>という日付形式に沿ったパスがあります。

データウェアハウス・スノ

ーパイプ

**S3**でイベントが利用できるようになったら、**Snowpipe**を使ってデータウェアハウスにインジェストします。これは、当社の **Snowflake**データウェアハウスの機能です。良いイベントパスと悪いイベントパスのために、**AmazonSQS**のイベントキュー が設定されました。

Snowpipeが正常に動作するためには、Snowflakeの中に「ステージ」が必要で、書き込み先のテーブルも必要です。良いS3パスと悪いS3パスには、それぞれSnowflake内にステージがあります。これらはそれぞれgitlab\_eventsとgitlab\_bad\_eventsと名付けられています。これらはLOADERロールが所有しています。

グッドイベントとバッドイベントのcreate tableステートメントは以下の通りです。

```
グッド・イベント
CREATE OR REPLACE TABLE
snowplow.gitlab_events (
   app_
                             idVARCHAR,
   platformVARCHAR<sub>o</sub>
                             tstampVARCHAR,
   etl_
   collector_tstamp
                             VARCHAR,
                             VARCHAR,
   dvce_created_tstamp
   eventVARCHAR,
                             idVARCHAR,
    event
    txn_
                             idVARCHAR,
                             trackerVARCHAR,
    name_
                             trackerVARCHAR,
    v_
                             collectorVARCHAR,
    v_
                             et1VARCHARです。
    V_
                             idVARCHAR,
    user
    user_
                             ipaddressVARCHAR,
    user_fingerprint VARCHAR,
                             useridVARCHAR,
    domain_
    domain_sessionidx VARCHAR,
    network_
                             useridVARCHAR,
                             countryVARCHAR,
    geo_
                             regionVARCHAR,
    geo_
                             cityVARCHAR,
    geo
                             zipcodeVARCHAR,
    geo_
                             latitudeVARCHAR,
    geo_
    geo_
                             longitudeVARCHAR,
    geo_region_
                             nameVARCHAR,
                             ispVARCHARです。
    ip_
    ip_
                             organizationVARCHAR,
```

```
domainVARCHAR,
ip_
                          netspeedVARCHAR,
ip_
page_
                          urlVARCHAR,
                          titleVARCHAR,
page_
                          referrerVARCHAR,
page_
                          urlschemeVARCHAR,
page_
                          urlhostVARCHAR,
page_
                          urlportVARCHAR,
page_
                          urlpathVARCHAR,
page_
                          urlqueryVARCHAR,
page_
page_urlfragment VARCHAR,
                          urlschemeVARCHAR,
refr_
                          urlhostVARCHAR,
refr_
refr
                          urlportVARCHAR,
refr_
                          urlpathVARCHAR,
refr_
                          urlqueryVARCHAR,
refr_urlfragment VARCHAR,
refr_
                          mediumVARCHAR,
                          sourceVARCHAR,
refr_
refr_
                          termVARCHAR<sub>o</sub>
mkt_
                          mediumVARCHAR,
mkt_
                          sourceVARCHAR,
mkt_
                          termVARCHAR,
                          contentVARCHAR,
mkt_
mkt_
                          campaignVARCHAR,
contextsVARCHAR,
                          categoryVARCHAR,
se_
                          actionVARCHAR,
se_
se_
                          labelVARCHAR,
                         propertyVARCHARです。
se_
                          valueVARCHAR,
se_
unstruct_
                          eventVARCHAR,
                          orderidVARCHAR,
tr_
tr_
                          affiliationVARCHAR,
                          totalVARCHAR,
tr_
tr_
                          taxVARCHARです。
                          shippingVARCHAR,
tr_
                         cityVARCHAR,
tr_
tr_
                          stateVARCHAR,
                          countryVARCHAR,
tr_
ti_
                          orderidVARCHAR,
                          skuVARCHAR,
ti_
                          nameVARCHAR,
ti_
ti_
                          categoryVARCHAR,
ti_
                          priceVARCHAR,
ti_
                          quantityVARCHAR,
pp_xoffset_
                         minVARCHAR,
pp_xoffset_
                         maxVARCHAR,
pp yoffset
                         minVARCHAR,
pp_yoffset_
                          maxVARCHAR,
useragentVARCHAR,
br_
                          nameVARCHAR,
                          familyVARCHAR,
br_
br_
                          versionVARCHAR,
br_
                          typeVARCHAR,
                          renderengineVARCHAR,
br_
                         langVARCHAR,
br_features_
                          pdfVARCHAR,
br_features_flash VARCHAR,
br_features_java
                    VARCHAR
br_features_director VARCHAR \,
br_features_quicktime VARCHAR \,
br_features_realplayer VARCHAR \
br_features_windowsmedia VARCHAR
```

br\_features\_gears VARCHAR。

```
br_features_silverlight VARCHAR,
                            cookiesV
    ARCHAR,
                            colordepthVARCHAR,
    br_
                            viewwidthVARCHAR,
    br_
   br_
                            viewheightVARCHAR,
                            nameVARCHAR,
    os
    os_
                            familyVARCHAR,
                            manufacturerVARCHAR,
    os_
                            timezoneVARCHAR,
    os
   dvce_
                            typeVARCHAR,
                            ismobileVARCHAR,
    dvce_
    dvce_
                            screenwidthVARCHAR,
    dvce_
                            screenheightVARCHAR,
                            charsetVARCHARです。
    doc_
    doc_
                            widthVARCHAR,
   doc
                            heightVARCHAR,
                            currencyVARCHAR,
   tr_
   tr_total_
                            baseVARCHAR,
    tr_tax_
                            baseVARCHAR,
   tr_shipping_
                            baseVARCHAR,
                            currencyVARCHAR,
   ti_
   ti_price_
                           baseVARCHAR,
   base_
                          currencyVARCHAR,
                          timezoneVARCHAR,
    geo_
                           clickidVARCHAR,
    mkt
                            networkVARCHAR,
    mkt
    etl_
                            tagsVARCHAR,
   dvce_sent_
                            tstampVARCHAR,
   refr_domain_
                          useridVARCHAR,
   refr_dvce_
                          tstampVARCHAR,
                          contextsVARCHAR,
   derived_
   domain_
                          sessionidVARCHAR,
   derived_
                          tstampVARCHAR,
                           vendorVARCHAR,
   event_
                            nameVARCHAR,
   event_
   event_
                            formatVARCHAR,
                            versionVARCHAR,
   event_
   event_
                            fingerprintVARCHAR,
    true_
                            tstampVARCHAR,
    upload_
                            attimestamp_ntz(9) Default Cast(current_timestamp() as timestamp_ntz(9))
)
バッドイベント表
CREATE OR REPLACE TABLE
snowplow.gitlab_bad_events (
                jsontextVARIANTです
    upload_at TIMESTAMP_NTZ(9) DEFAULT CAST(CURRENT_TIMESTAMP() AS TIMESTAMP_NTZ(9))
);
```

TSVはCSVのようにわかりやすいものではないので、次のような記述でカスタムファイル形式を作成しました。

```
CREATE OR REPLACE FILE FORMAT snowplow_tsv TYPE = CSV FIELD_DELIMITER = '\'
```

## 実際にイベント用のパイプを作ったのは

```
CREATE OR REPLACE PIPE raw.snowplow.gitlab_good_event_pipe auto_ingest= TRUE

AS COPY INTO raw.snowplow.gitlab_events
    from (select $1, $2, $3, $4, $5, $6, $7, $8,

$9,$10,$11,$12,$13,$14,$15,$16,$17,$18,$19,$20,$21,$22,$23,$24,$25,$26,$27,$28,$29,$30,$31,$32,$33,$

34,$35,$36,$37,$38,$39,$40,$41,$42,$43,$44,$45,$46,$47,$48,$49,$50,$51,$52,$53,$54,$55,$56,$57,$58,$
```

これは、JSONがより良いフォーマットであることを強調しています。これはインフラの将来のイテレーションになるかもしれません。バッドイベントパイプは以下のように作成されます。

```
CREATE OR REPLACE PIPE raw.snowplow.gitlab_bad_event_pipe auto_ingest=TRUE
AS COPY INTO raw.snowplow.gitlab_bad_events (jsontext)
FROM @raw.snowplow.gitlab_bad_events
FILE_FORMAT = (TYPE = 'JSON');
```

### パイプを表示するには

```
raw.snowplowでパイプを表示します。
```

## パイプの説明に

```
DESCRIBE PIPE raw.snowplow.gitlab_good_event_pipe;
```

走っているパイプを一時停止すること。

```
ALTER PIPE raw.snowplow.gitlab_good_event_pipe SET PIPE_EXECUTION_PAUSED = TRUE;
```

パイプを強制的に再開すること。

```
SELECT system$pipe_force_resume('raw.snowplow.gitlab_good_event_pipe')。
```

# パイプの状態を確認するには

```
SELECT system$pipe_status('raw.snowplow.gitlab_good_event_pipe')です。
```

ステージを強制的に更新し、snowpipeが古いイベントを拾うようにします。

```
ALTER PIPE gitlab_good_event_pipe refresh;
```

### dbt

RAWデータベースからPRODにデータを実体化して問い合わせを行うために、dbt内にパーティショニング戦略を実装しました。デフォルトでは、除雪機モデルとFishtown除雪機パッケージは、PREPデータベースの現在の月にスコープされたスキーマに書き込みます。2019年7月の場合、スキーマはsnowplow 2019\_07となります。

各月のパーティション内では、基本モデルとパッケージで生成されたモデルのすべてが、パーティションの日付に一致する派生タイムスタンプを持つすべてのイベントに対して書き込まれます。実行時に**dbt**に変数を渡すことで、異なる月別パーティションを生成することができます。

```
--vars '{"year":"2019", "month":"01", "part":"2019_01"}'
```

バックフィルはAirflowで行います。dbt\_snowplow\_backfillDAGは、2018年7月から当月までの各月のタスクを生成します。

### **Do Not Track**

当社の除雪機のトラッキング設定および特定の実装では、ユーザーのブラウザにDNT(Do Not Track) ヘッダーが存在する場合は常にそれを尊重します。

layout: handbook-page-toc title: "SQLスタイルガイド" 説明"私たちにはリンターがいないので、この SQL スタイルガイドを施行するのは私たちの集団の責任です。"

このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .toc-list-icons .hidden-md .hidden-lg}。

# SQLスタイルガイド

私たちにはリンターがいないので、このスタイルガイドを実行するのは*私たちの責任*です。

フィールドの名称と参照方法

- フィールド名はすべて小文字にしてください。
- ID、名前、またはtypeのような一般的に曖昧な値は、常にそれが何を識別または命名しているかを前置する必要があります。

```
グッド
セレク
ト
idAS account_id
、name AS
account_name、type AS
account_type。
・・・
バッド
セレク
ト
id、
name
、
type
```

• 異なるソースからのデータに結合する場合、フィールドの前にデータソースを付ける必要があります。 曖昧さを避けるために、sfdc\_account\_id

```
グッド
セレク
ト
sfdc_account.account_ idAS
sfdc_account_id, zuora_account.account_id AS
zuora_account_id
FROM sfdc_account
LEFT JOIN zuora_account ON ...
```

-- 悪い

```
SELECT

sfdc_account.account_id,

zuora_account.account_id AS zuora_id

FROM sfdc_account

LEFT JOIN zuora_account ON ...
```

• テーブルを結合して両方のカラムを参照する場合、エイリアスではなく完全なテーブル名を参照することを強く推奨します。テーブル名が長い(~20)場合は、可能であればCTEの名前を変更し、最後に何かわかりやすい名前にエイリアスすることを検討します。

```
すべてのフィールド名は蛇足である
グッドセグッド
```

```
セレク
  dvcereatedtstamp AS device created timestamp date_details.fiscal_quarter_name,
 cost_category.cost_category_level_1,
 cost_category.cost_category_level_2 があ
FROM budget_forecast_cogs_opex
LEFT JOIN date_details
ON date_details.first_day_of_month = budget_forecast_cogs_opex.accounting_period LEFT
JOIN cost_category
ON budget_forecast_cogs_opex.unique_account_name = cost_category.unique_account_name
-- OKだが、好ましくない。エイリアシングの代わりにCTEの名称変更を検討する
SELECT.
 bfcopex.account_id,
  -- さらに15個のカラム
 date_details.fiscal_year,
 date_details.fiscal_quarter,
 date_details.fiscal_quarter_name,
 cost_category.cost_category_level_1,
 cost_category.cost_category_level_2
FROM budget_forecast_cogs_opex bfcopex LEFT
JOIN date_details
ON date_details.first_day_of_month = bfcopex.accounting_period LEFT
JOIN cost_category
ON bfcopex.unique_account_name = cost_category.unique_account_name
バッド
セレク
F
 a.*,
 b.fiscal_year,
 b.fiscal_quarter,
 b.fiscal_quarter_name,
 c.cost_category_level_1,
 c.cost_category_level_2 の15列
FROM budget_forecast_cogs_opex a
LEFT JOIN date details b
ON b.first_day_of_month = a.accounting_period LEFT
```

JOIN cost\_category c
ON b.unique\_account\_name = c.unique\_account\_name

```
FROMテーブル
バッド
セレク
ト
dvcecreatedtstamp AS DeviceCreatedTimestamp
FROM table
```

• ブール型のフィールド名は、has\_is\_does\_のいずれかで始まります。

```
グッドセ
レクト
deleted AS is_deleted,
slaAS has_sla
FROMテーブル
バッド
セレク
ト
削除、スラ
FROMテーブル
```

• ソースデータを変換する際、大文字小文字を区別する列や、「\$」や「\_」とは異なる特殊文字を含む列を識別するために二重引用符を使用します。大文字のフィールド名に二重引用符必要ありませこれはSnowflake識別子が内部的に処理ため

```
--良い
SELECT "First_Name_&_" AS first_name,
-- 悪い
SELECT "FIRST_NAME" AS first_name,
```

# 日付

- タイムスタンプは\_atで終わります。例: deal\_closed\_at、常にUTCでなければなり
- ません。 日付は\_dateで終わります。例: deal\_closed\_date
- 月はそのように表示し、常に日付形式に切り詰める必要があります(例: deal\_closed\_month)。
- 列名に「日付」や「月」などのキーワードを使わないようにする。
- DAYOFWEEK(created\_at) >」のように、date\_partよりも明示的な日付関数を優先し、extractよりもdate\_partを優先する。
  DATE\_PART(dayofweek, 'created\_at') > EXTRACT(dow FROM created\_at)
  - date\_trunc('month', created\_at)ではカレンダーの月('2019-01-01'なら'2019-01-25') が生成され、SELECT date\_part('month', '2019-01-25'::date)では数字の1が生成されることに注意してください。

DATEDIFFを使うと、直感的でない結果になることが多いので注意が必要です。

• の例えば、SELECT DATEDIFF('days', '2001-12-01 23:59:59.999', '2001-12-02 00:00:00.000')では、次のように返されます。

タイムスタンプが1ミリ秒違っていても、1になります。

- 同様に、SELECT DATEDIFF('days', '2001-12-01 00:00:00.001', '2001-12-01 23:59:59.999') return 0 タイムスタンプがほぼ丸一日離れているにもかかわらず。
- DATEDIFF関数で適切な間隔を使用することで、正しい結果を得ることができます。例えば、DATEDIFF('days', '2001-12-01 23:59:59.999', '2001-12-02 00:00:00.000')では、1日分の間隔で、DATEDIFF('ms', '2001-12-01 23:59:59.999', '2001-12-02 00:00:00.000')では1ミリ秒の間隔で

サブクエリではなくCTE(Common Table Expressions)の使用

• CTEは、SQLをより読みやすくし、パフォーマンスを向上させます。

- 他のテーブルを参照するには、CTE を使用します。CTEはクエリの先頭
- に配置する必要があります(インポートステートメントと考えてください)。
- 性能が許す限り、CTEは単一の論理的な作業単位を実行すること CTEの名称は、明
- 確であると同時にできる限り簡潔であること
  - ∘ replace\_sfdc\_account\_id\_with\_master\_record\_id のような長い名前は避け、CTE にコメントを入れた短い名前を使うようにしましょう。これにより、結合時のテーブルのエイリアスを避けることができます。
- 紛らわしいロジックや注目すべきロジックを持つCTEは、ファイル内でコメントし、dbtドキュメ
- ントで文書化すべき モデル間で重複しているCTEは、それぞれのモデルに引き抜くべき
- クエリステートメントの上下に空の行を残す CTEは以下の
- ようにフォーマットしてください。

```
WITH events AS ( -- これらのselect文をimport文と考えてください。
...
), filtered_events AS ( -- CTEのコメントはここへ)
...
)

SELECT * -- 最後のモデルについては、常に「select * from final」を目指すべきです
FROM filtered_events
```

# 一般

• CTE内では、SQL文全体をスペース4個分インデントする必要があります。

```
--良い
WITH my_data AS (

SELECT *
FROM prod.my_data
WHERE filter = 'my_filter'
)

-- 悪い
WITH my_data AS (

SELECT *
FROM prod.my_data
WHERE filter = 'my_filter'
)
```

• クエリ内のインデント (例:列、JOIN句、複数行のGROUP BYなど) は、スペース2個分のインデントが必要です。

```
グッド
セレク
ト
column_name1,
column_name2,
```

```
カラム名3 FROM
table_1 JOIN
table_2
 ON table 1.id =
table_2.id WHERE clouds =
true
 AND gem =
true GROUP BY
1,2,3
HAVING カラムネーム1 >
 0 AND カラムネーム2 >
バッド
セレク
1
   カラム名1、カ
   ラム名2、カラ
   ム名3
FROM
table_1
JOIN
table_2
   ON table_1.id =
table_2.id WHERE clouds =
true
   AND gem =
```

AND カラムネーム2 > 0

4 タブは使用みず フペーフのひを使用してくがさい 領集者は タブをフペーフル

- タブは使用せず、スペースのみを使用してください。編集者は、タブをスペースに変換するように設定する必要があります。
- SQLの行は80文字以内にしてください。

HAVING カラムネーム1 > 0

true GROUP BY

1,2,3

• カンマは、WHEREの一時的なフィルターを除いて、行末(EOL)に右のカンマを入れてください。 節で特定の値を指定します。

```
--良い
SELECT
          AS is_deleted, -- EOL右コンマ
 削除され
            AS account_id
 た
 accountId
NAND account_id NOT IN (
                    '232'
                    234」 -- 左コンマ
                    , '425'
-- 悪い
SELECT
 削除され
           AS is_deleted, -- EOL右コンマ
            AS account_id
 た
 accountId
WHERE is deleted = false
NAND account_id NOT IN ('232', '234',
  '425')
```

• SELECTの際には、常に各列を独立した行にします。ただし、SELECT \*の場合は1つの行になります。

- DISTINCTは、SELECTと同じ行に含まれている必要があります。
- フィールド名やテーブル名を投影する際には、ASキーワードを使用する必要があります。
- ASを使用してエイリアシングを行う場合は、元のカラム名を1本の縦線に、ASキーワードを別の縦線に配置するようにしてください。

- フィールドは集計/窓関数の前に記載すること 数値による順序付けと
- グループ化 (例: GROUP BY 1, 2) が望ましい
  - o dbtモデルで3つ以上のカラムでグループ化する場合、dbt-utils group\_byマクロを使用してく
- ださい。WHEREとHAVINGのどちらかで十分な場合はPreferしてください。
- 例えば、data\_by\_row['id']::bigint as id\_valueのように、ブラケット構文を使用してJSONにアクセスすることができます。
- 結合でUSINGを使うとSnowflakeで不正確な結果が出るので絶対に使わないでください。このトピックに関するフォーラムのディスカッションを表示するには、アカウントを作成してください。
- UNIONよりもUNION ALLを優先してください。これは、UNIONが上流のデータインテグリティの問題を示している可能性があり、他の 場所で解決した方が良いからです。
- <> よりも!= を優先してください。これは、他のプログラミング言語では!= の方が一般的であり、「not equal」のように読めるため、私たちはこのような言い方をすることが多いからです。
- ◆ パフォーマンスを考慮する。LIKEとILIKE、ISと=、NOTと!と<>の違いを理解する。適切に使用する
- LOWER(column) LIKE '%match%' よりも、column ILIKE '%Match%' を優先します。これにより、誤った大文字が予期 せぬ結果を引き起こす可能性が低くなります。
- DRYプリンシパルに慣れ親しむ。dbtではCTE、jinja、マクロを、Sisenseではスニペットを活用する。同じ行を2回入力した場合、2箇所で管理する必要がある
- は、コードの行数を減らすために最適化することはありません。改行は安い。ブレインタイムは高価

## データタイプ

- デフォルトのデータ型を使用し、エイリアスは使用しないでください。詳細は、Snowflake のデータタイプの概要をご覧ください。デフォルトは
  - DECIMAL、NUMERIC、INTEGER、BIGINTなどの代わりにNUMBER。
  - ◇ DOUBLE、REALなどの代わりにFLOAT。
  - STRING、TEXTなどの代わりに
  - VARCHARDATETIMEではなくTIMESTAMP

ただし、タイムスタンプの場合は例外です。TIMEよりもTIMESTAMPを優先します。なお、TIMESTAMPのデフォルトはTIMESTAMP\_NTZです。のように、タイムゾーンが含まれていません。

## 機能

- 関数名とキーワードはすべて大文字にする IFNULL TO
- NVLを優先する
- **1**行のCASE文よりもIFFを優先する
- is\_less\_than\_tenとしてbooleanステートメント(amount < 10)を選択するよりもIFFを優先する
- 繰り返しの多いCASE文を可能な限り簡略化してください。

```
OK

CASE

WHEN Field_id = 1 THEN 'date'

WHEN Field_id = 2 THEN 'integer'

WHEN Field_id = 3 THEN

'currency' WHEN Field_id = 4

THEN 'boolean' WHEN Field_id = 5

THEN 'variant' WHEN Field_id = 6

THEN 'text'

END AS field_type
```

```
より良いケース
field_id
WHEN 1 THEN 'date'
WHEN 2 THEN 'integer'
WHEN 3 THEN
'currency' WHEN 4
THEN 'boolean' WHEN 5
THEN 'variant' WHEN 6
THEN 'text'
END AS field_type
```

# **JOINs**

- 例えば、JOINの代わりにLEFT JOINを使うなど、結合の際には明示的にしてください。(デ
- フォルトの結合はINNERです) 結合時にテーブル名をカラムの前に付け、そうでなければ省略します。
- FROMテーブルを先に、JOINテーブルを後にして結合の順番を指定します。

```
良いFROMソ
ース
LEFT JOIN other_source
ON source.id =
other_source.id WHERE ...

-- 悪い
FROMソース
LEFT JOIN other_source
ON other_source
UN other_source.id =
source.id WHERE ...
```

## コード例

• まとめてみました。

```
WITH my_data AS
    ( SELECT * )
    FROM prod.my_data
   WHERE filter = 'my_filter'
), some_cte AS (
   SELECT DISTINCT
     id,
     other field 1,
     other_field_2
    FROM prod.my_other_data
), final AS
    ( SELECT
      data_by_row['id']::NUMBER AS id_field,
     field_1
                                AS detailed_field_1です。
     field_2
                                 AS detailed_field_2,
     detailed_field_3,
       WHEN cancellation_date IS NULL AND expiration_date IS NOT
         NULL THEN expiration_date
       WHEN cancellation_date IS
         NULL THEN start_date + 7
```

```
ELSE cancellation_date
      ENDAS cancellation_date, LAG(detailed_field_3)
      OVER (
       パーティション・バイ
         id_field,
         detailed_field_1
       ORDER BY cancellation_date
                                AS previous_detailed_field_3,
      SUM(field 4)
                                AS field_4_sum,
     MAX(field_5)
    field_5_max FROM my_data
    LEFT JOIN some_cte
     ON my_data.id =
    some_cte.id WHERE field_1 =
    'abc'
     AND (field_2 = 'def' OR field_2 =
    'ghi') GROUP BY 1, 2, 3, 4, 5
    count(*) > 1 order
    by 4 desc
SELECT *
FROM final
```

## コメント

- モデル内で一行コメントを作成する場合は -- 構文を使用します。 モデル
- 内で複数行コメントを作成する場合は /\* \*/ 構文を使用します。
- コメントを書くときは、1行の制限を守ってください。コメントが長すぎる場合は、新しい行に移動するか、モデル のドキュメントに移動してください。
- dbtモデルのコメントは、モデルのドキュメントに記載されるべきです。
- **SQL**で計算する場合は、何が起こっているのかを簡単に説明し、その指標を定義したハンドブック(およびその計算方法)へのリンクを表示する必要があります。
- TODOコメントを残すのではなく、改善のために新しい課題を作る

その他のSOLスタイルガイド

- Brooklyn Data Co
- FishtownAnalytics
- Matt Mazur
- Kickstarter

layout: handbook-page-toc title:"Data & Analytics Team Operating Principles" description:"GitLab データ&アナリティクスチーム運営原則ハンドブック"

このページについて

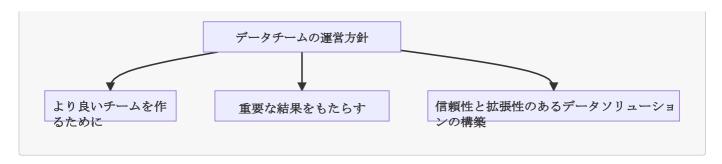
{:.no\_toc .hidden-md .hidden-lg}。

{::options parse\_block\_html="true" /}。

データ&アナリティクスの動作原理

私たちは、世界で最も革新的な企業のデータ&アナリティクスのセンター・オブ・エクセレンスであることを誇りに思うと同時に、謙虚に受け止めています。

データチームには、3つの運営方針があります。



### より良いチームを作るために

- CREDIT
- 私たちは、すべての人を尊重しています。
- 私たちは、GitLabの成功に向けて最善の意思を持
- っています。 私たちは、自分自身とお互いに正直 です。
- 私たちは意見を異にすることを許可されていますが、コミットメ
- ントに向けて行動しなければなりません 難しい会話から逃げない
- 私たちはチームとして一緒にいる
- チームでミスに取り組み、改善を目指す

### 重要な結果をもたらす

- 第一にデータドリブン、第二にデータ・インフォームド、第三に原則と
- ストーリー 私たちは、Built to Last, not built too fastを信条としています。
- 私たちは成長思考で行動し、常に水準を高めていきます 私たちは活動と
- 進歩を混同しません
- 自分の仕事がどのように戦略的イニシアティブと連携し、ビジネスパートナーがインパクトのある結果を出すために どのように役立つかを知っています。
- 私たちは、文脈、事実、洞察を提供するために余分な努力をします。
- 前もって期待値を定義します。範囲、時間、成功基準
- 重要な取り組みを定期的に進めるために、ビジネスパートナーと積極的に関わっています。

### 信頼性と拡張性のあるデータソリューションの構築

- データセキュリティとデータ品質に対する責任があります。 十
- 分にテストされた正確なデータソリューションでリードします
- 私たちが作成したプロダクションソリューションを運用・保守す
- る 私たちは、GitLab +2年のためにデザインし、今日のGitLabの ために提供する
- ベンダー・パートナーの能力とサービスを最大限に活用する
- 私たちは、複雑な結果を調整するための最良の方法は、それを描くことだと信じています。
- 自動化されたテストは最高のテストであり、私たちはデータ提供プロセスのすべてのステップでテストを実施します。

# データチームが持つその他の信念は以下の通りです。

- すべてがコードで定義され、バージョン管理されるべきである。
- データの実装には、DevOpsのベストプラクティスをワークフローに統合する必要があります。
- ◆ 高品質で保守可能なコードベースを持ちながら、各ビジネス機能と提携します。 データのセキュリ
- ティとデータの機密性を考慮しながら、コードをオープンソース化します。
- 企業内のすべての分析的な質問に対する真実のソースは単一であり、最も徹底したデータの洞察を提供するために、 複数の視点を表現することも重要である
- 糊付け作業はチームの健全性を保つために重要であり、その価値に応じて個別に評価されます。私たちは全員が貢献することを期待しています。
- 限られたリソースを、データが最大の効果を発揮する場所に集中させる
- すべてのビジネスユーザーは、簡単な統計の解釈と計算方法を学ぶことができます。

layout: handbook-page-toc title:"Data Team Programs" 説明"データプログラム"

# このページについて

{:.no\_toc .hidden-md .hidden-lg}。

• TOC {:toc .hidden-md .hidden-lg}.

データプログラムページへようこそ

データプログラムのページでは、オンボーディングから日々のオペレーションに至るまで、私たちがサポートする様々なプログラムに関する情報が掲載されています。

データプログラム

|プログラム名 | タイプ | 目的 | :--- | :--- | | 財務向けデータ | 運用 | 財務アナリストを支援するための情報 | |プロダクトマネージャー向けデータ | 運用 | プロダクトマネージャーを支援するための情報 | |マーケティングアナリスト向けデータ | マーケティングアナリストを支援するための情報 | |セールスアナリスト向けデータ | セールスアナリストを支援するための情報 | |ドータオンボーディング | 1-time/ad-hoc | 中央データチームの内外を問わず、データエンジニア、アナリスト、または開発者のトレーニング、プロビジョニング、およびイネーブルメントを行います。|データSlackチャンネル | 運用 |全社のチームメンバーの質問に対応するための#dataチャンネルの監視。|データトリアージ | 運用 | データプラットフォームが分析のために利用可能であることを保証するための日々のプロセス。|

データオンボーディング

**GitLab**にオンボーディングし、データプログラムでエンジニア、アナリスト、デベロッパーとして働く場合は、以下の手順に 従ってください。

- 1. GitLabDataAnalyticsでData Onboardingテンプレートを使って新しい課題を開きます。
- 2. 課題にわかりやすい名前をつける: Your Name Data Onboarding
- 3. マネージャーに課題を割り当て、関連するコンテンツを追加/削除する