



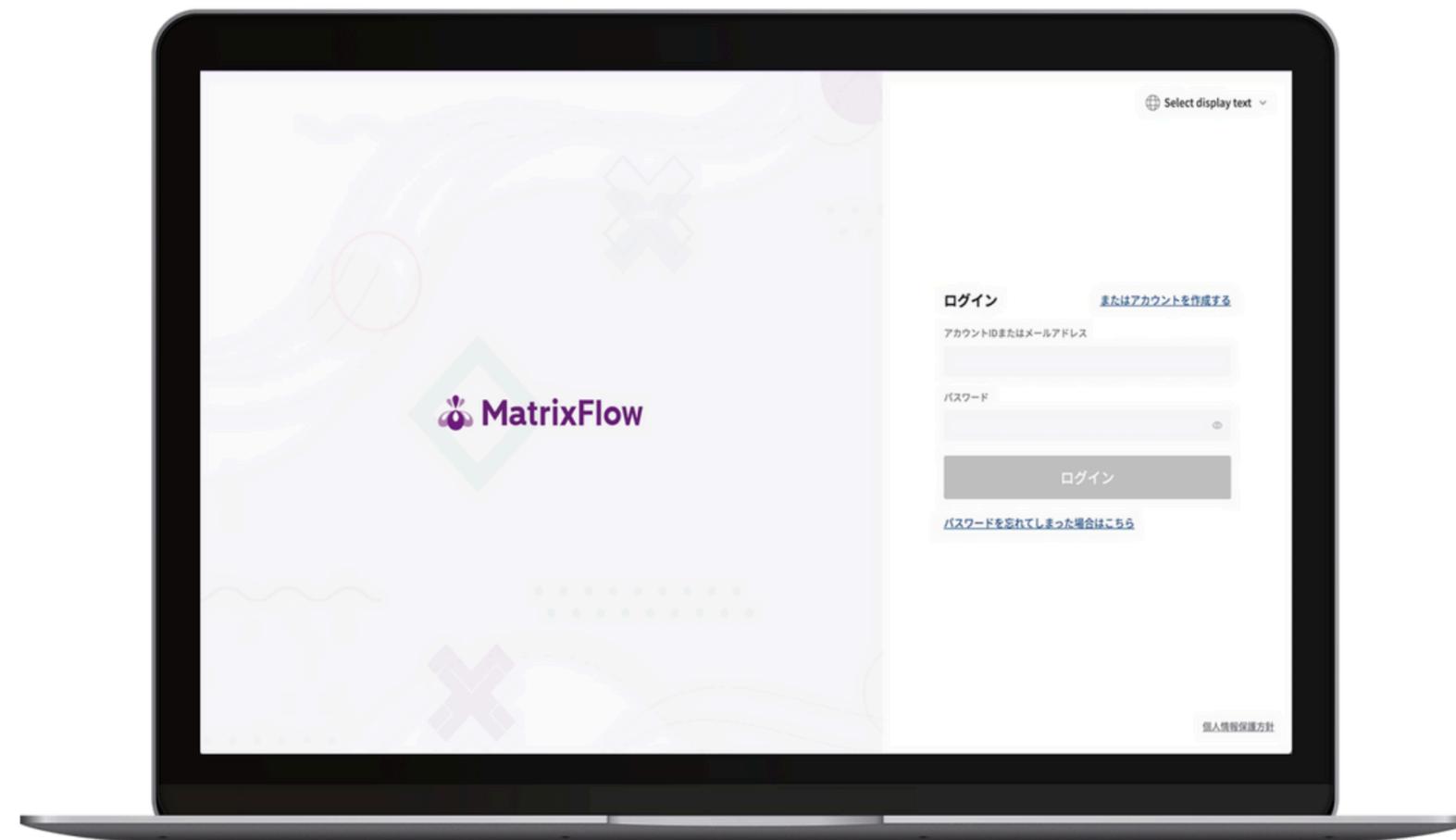
生成AIを活用した自動前処理機能

+

生成AIによる行削除機能

+

データ結合機能



## 会社概要

【社名】 株式会社MatrixFlow

【代表者】 田本 芳文

【設立】 2018年10月

【資本金】 1.71億円（資本準備金含）

【住所】

東京都台東区上野3丁目16番2号

天翔上野末広町ビル206号室

【事業内容】

『MatrixFlow』のSaaSサービス提供

AIの受託開発・研究・コンサル

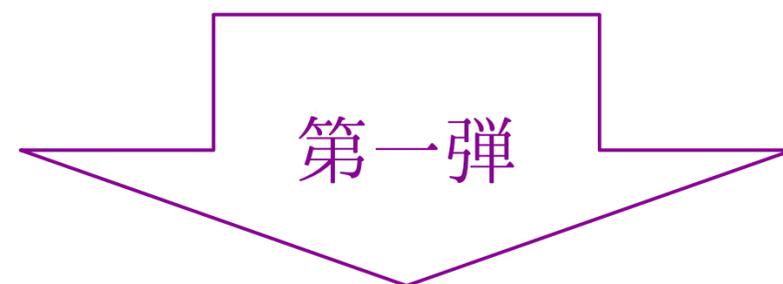
## 実績/許認可

- ・ 株式会社dipのAIアクセラレータープログラム採択
- ・ Plug and Playのアクセラレータープログラム採択
- ・ 2020年10月 1億円の資金調達を実施
- ・ セキュリティマネジメントシステム（ISMS）認証取得
- ・ プライバシーマーク（Pマーク）認証取得
- ・ 総務省後援ASPIC IoT・AI・クラウドアワード2021  
AI部門ニュービジネスモデル賞受賞
- ・ 2022年4月 資金調達を実施



掲載メディア

テクノロジーで世界をつくる



AIの民主化を推進

一部の企業や技術者がAI技術を独占するのではなく誰でもAIが活用できる社会を実現



# ビジネスのための AI活用プラットフォーム

プログラミングなしで、  
短期間でAIを構築・活用し、  
幅広い課題を解決します



## 数学・統計学，プログラミングの知識は一切不要

プログラミングすることなく（No-Code），  
データサイエンスなどの高度なスキルがなくてもAIを構築・活用できるSaaSサービスです。

MatrixFlowがない場合

数学・統計学



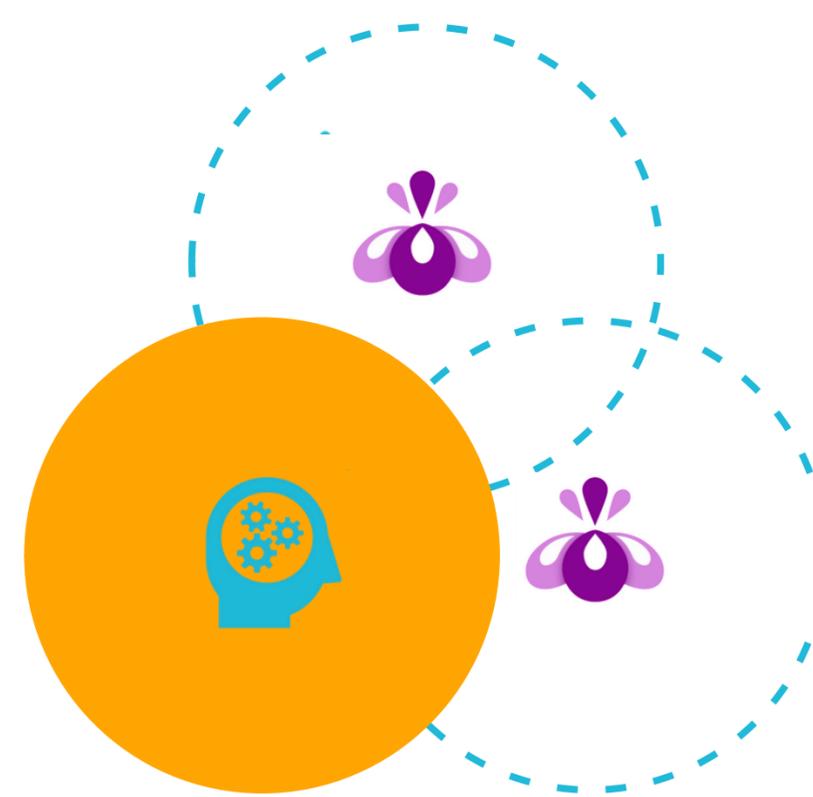
ドメイン知識 プログラミング

AI化には多くのスキルが必要



MatrixFlowがある場合

数学・統計学



ドメイン知識 プログラミング

ビジネスの知識だけあればよい

# 本日の説明会要旨

## はじめに

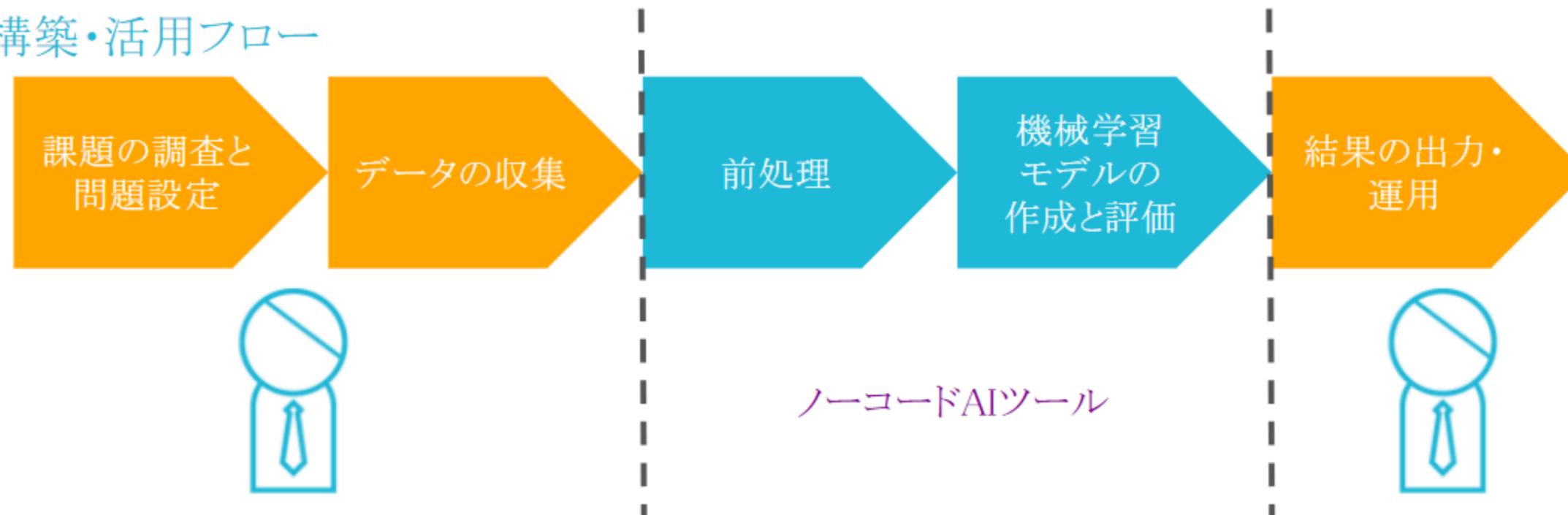
これまで、AIモデルのアルゴリズム選択やハイパーパラメータのチューニングは自動化が進んでいましたが、データの前処理は自動化が難しい領域とされてきました。

データ前処理は、データの構造や意味を深く理解した上で行う必要があり、これがデータサイエンス業務全体の8割を占めるとさえ言われています。

しかし、大規模言語モデル（LLM）の登場により、データの構造や意味を理解し、人間のような判断をもとに処理を行うことが可能になりました。

この技術的進化は、前処理を自動化する新たな道を切り開き、データ分析の生産性を飛躍的に向上させる可能性を示しています。

### AI構築・活用フロー



- 生成AIによる自動前処理機能

- 読み込んだデータを解析し、次ページ以降でご説明する「MatrixFlowがもつ前処理機能」を最適に自動的に実行してくれる機能です。

- 生成AIによる行削除機能

- 前処理「行の削除」が追加されました。この機能は従来の空白の行を消す機能ではなく、生成AIに希望する「列」と「行」を伝えることで、該当箇所の行を消してくれる機能です。

- データ結合機能

- データ結合機能は、データ管理にある2つのデータを統合する機能です。対象のデータ種類はテーブルデータ（CSV）です。

# 生成AIによる自動前処理機能

- データクレンジングの自動化

- 各種エンコーディング、名寄せ、対数変換、単位の統一、損値や外れ値の検出・処理、重複データの削除などのクレンジング処理を自動的に実行し、正確で信頼性の高いデータを生成します。

- 自然言語処理による列名の意味解析

- カラム名やデータの文脈を理解し、自動的に適切な前処理方法を選択しています。

- 直感的なUIで操作を簡単に

- データサイエンスの専門知識がないユーザーでも、ワンクリックで前処理を完了できます。

- 前処理の透明性を担保

- 前処理の実行内容は明示されますので透明性を担保できます。

### 数値に変換する

AIモデルが学習に使えるように、文字列を数値に置き換える

- ・ One-Hotエンコーディング
- ・ ダミーコーディング
- ・ ラベルエンコーディング

### 欠損値を変換する（欠損値がある列のみ表示されます）

AIモデルが学習に使えるように、欠落している値を補完する

- ・ 欠損値を含む行を削除する
- ・ 欠損値を（数値で）埋める

### 列を削除する

### 外れ値を除去する

測定ミス、記録ミス、データの異常、データの偏りなどに起因した外れ値でAIモデルの学習にゆがみを生じさせるおそれがある値を除去する

### データ変換

- ・ 正規化：表記ゆれの名寄せ、重複データの削除、データの単位やスケールを共通の基準に整える
- ・ 標準化：平均が0、標準偏差が1の正規分布にスケーリングする
- ・ 対数変換：Yeo-Jonson変換（負値にも対応した対数変換）を実行する
- ・ 離散化：グループにまとめて離散化する

### 数値に変換する

- One-Hotエンコーディング

具体例:

- 「動物」列が ["犬", "猫", "鳥"] の場合、以下のように変換

動物			
犬	1	0	0
猫	0	1	0
鳥	0	0	1

理由:

- 順序性がないカテゴリ変数を適切に処理し、モデルが順序性を誤解しないようにするため
- カテゴリ変数をバイナリ特徴量に変換し、各カテゴリを独立した次元で表現する
- 効果的に機械学習アルゴリズム（特に線形モデルや距離計算に基づくモデル）で利用可能

### 数値に変換する

- ダミーコーディング

具体例:

- One-Hotエンコーディングと同様ですが、基準カテゴリを1つ削除（例: 「鳥」を基準）する

動物		
犬	1	0
猫	0	1
鳥	0	0



理由:

- 順序性がないカテゴリ変数を適切に処理し、モデルが順序性を誤解しないようにするため
- One-Hotエンコーディングに似ているが、1つのカテゴリを基準として省略する（kカテゴリに対しk-1次元に削減）
- 基準カテゴリをモデルの基準に設定することで、多重共線性を回避しやすくなる

### 数値に変換する

- ラベルエンコーディング

具体例:

- カテゴリごとに数値を割り当てる

評価	評価
低	0
中	1
高	2

理由:

- 順序があるカテゴリ（例: 評価スコア「低」「中」「高」）に適している
- 各カテゴリに一意の数値を割り当て、カテゴリ変数を数値変数に変換する
- 軽量で、カテゴリ数が多くてもメモリ消費が少ない

### 欠損値を変換する

- 欠損値を含む行を削除する

具体例:

- 「年齢」がNaNの行を削除する

名前	年齢	性別
山田	25	男性
佐藤	NaN	女性
高橋	30	男性



名前	年齢	性別
山田	25	男性
高橋	30	男性

理由:

- 欠損値の多い行がデータ全体に与える影響を排除する
- 小規模データセットではデータ不足のリスクがあるが、大規模データでは欠損値の影響を小さくできる

### 欠損値を変換する

- 欠損値を（数値で）埋める

具体例:

- 「年齢」が欠損している場合、平均値（27.5）で埋める

名前	年齢	性別
山田	25	男性
佐藤	NaN	女性
高橋	30	男性



名前	年齢	性別
山田	25	男性
佐藤	27.5	女性
高橋	30	男性

理由:

- データ量を維持し、モデル学習に利用できるデータを最大化する
- 平均値、中央値、最頻値、または他のデータから推定された値で埋めることで、統計的な一貫性を保つ
- 欠損がデータの一部の特徴を反映する場合（例: 欠損が特定の条件下で起こる）には慎重に処理が必要

### 外れ値を除去する

具体例:

- 年齢が範囲外（例: 正常範囲を0~120歳とする）の値を削除

名前	年齢	性別
山田	25	男性
佐藤	27	女性
高橋	300	男性



名前	年齢	性別
山田	25	男性
佐藤	27	女性

理由:

- 外れ値はモデルの学習を歪め、性能低下や誤解を引き起こす可能性がある
- 外れ値を削除または適切に処理することで、モデルの汎化性能（モデルが訓練データだけでなく、未知のデータに対してもどの程度適応し、正確な予測ができるかを示す指標）を向上させる
- ただし、外れ値がデータの重要な特徴を示す場合があるため、慎重に判断する

### 列を削除する

具体例:

- 「あきらかに無意味な列」「相関が高過ぎる列」や「重要度が低い列」を削除

名前	年齢	性別	身長	体重	BMI
山田	25	男性	171	64	21.89
佐藤	27	女性	155	52	21.64
高橋	30	男性	175	105	34.29



名前	年齢	性別	身長	BMI
山田	25	男性	171	21.89
佐藤	27	女性	155	21.64
高橋	30	男性	175	34.29

理由:

- 高い欠損率を持つ列や、予測タスクに寄与しない列を削除することでデータのノイズを減らし、モデルの精度を向上させる
- 相関が高すぎる特徴量を削除することで、多重共線性のリスク（重回帰モデルにおいて、説明変数の中に、相関係数が高い組み合わせがあることをいう（例: 体重とBMI））を軽減する
- 特徴量が多すぎる場合に、次元削減によって計算効率を向上させる

## データ変換

- 正規化

具体例:

- 表記ゆれやスケールを整える：表記ゆれ修正：「東京都」「東京」を「東京」に統一
- 重複データ削除：同じレコードを1つだけ残す
- スケール統一：単位が異なる場合（例：距離がkmとm）を統一する

名前	出身地	社員番号	通勤時間
山田	東京都	E2211	1h
佐藤	東京	E2031	80min
高橋	愛知	A1701	30min
山田	東京都	E2211	1h



名前	出身地	社員番号	通勤時間
山田	東京	E2211	60min
佐藤	東京	E2031	80min
高橋	愛知	A1701	30min

理由:

- 異なる単位やスケールを持つデータを統一基準に合わせることで、機械学習モデルの計算を安定化させる
- 距離計算に基づくモデル（例：k近傍法、SVM）で、スケール差による影響を抑える
- 重複データの削除や表記ゆれの修正は、データの一貫性と品質を向上させる

### データ変換

- 標準化

具体例:

- 「収入」列: [300, 500, 700] を標準化。 平均: 500、標準偏差 (標準的な平均値との差): 200

名前	収入
山田	300
佐藤	500
高橋	700



名前	収入
山田	-1.0
佐藤	0.0
高橋	1.0

理由:

- データを平均0、標準偏差1にスケールすることで、統計的性質を標準化する
- 勾配降下法を用いるモデル (例: ロジスティック回帰、ニューラルネットワーク) の収束を安定化・高速化する

### データ変換

- 対数変換

具体例:

- 「売上高」列：[10, 100, 1000] を対数変換する

店舗名	売上高
東京店	10
名古屋店	100
大阪店	1000



店舗名	売上高
東京店	1.0
名古屋店	2.0
大阪店	3.0

理由:

- データの分布を正規分布に近づけることで、モデルの性能を向上させる
- 大きな値の影響を抑え、データを扱いやすくする
- Yeo-Johnson変換では負の値にも対応しており、幅広いデータセットに適用可能

## データ変換

- 離散化

具体例:

- 年齢データを以下のグループに分類

名前	年齢
山田	25
佐藤	27
高橋	30
鈴木	19
田中	60



10-20: 若年層  
21-40: 中年層  
41-60: 壮年層



名前	年齢
山田	中年層
佐藤	中年層
高橋	中年層
鈴木	若年層
田中	壮年層

理由:

- 数値データをカテゴリに分けることで、特定の区間に対する解釈性を高める
- 複雑な数値変数を簡素化し、モデルの学習を容易にする
- 例: 年齢を「10代」「20代」などのグループに分割する

# MatrixFlowでの「生成AI自動前処理」の実行はカンタン

← ダッシュボードに戻る  
データ前処理・結合・…

AI作成の進行状況

- プロジェクト作成  
使用するテンプレート
- データセット設定  
使用するデータセット  
退職リスク予測\_学習…
- 予測する列の選択  
選択した予測する列  
退職
- 4 前処理**  
前処理を設定
- 5 レシピの選択
- 6 学習

生成AIで前処理 前処理を完了する

データセットを表示 統計情報を表示 行を削除

外れ値や統計情報を表示

前処理中のデータセット  
退職リスク予測\_学習用\_デモ用

予測する列  
退職

データ数  
11 詳細を開く →

保存した前処理を使用する  
前処理を選択して下さい

年齢	退職	出張	部署	通勤距離	学歴	専攻
40	退職済	まれ	営業	25km	4	マーケティング
31	在籍中		研究・エンジニア…	23km	3	薬学
44	在籍中	なし	研究・エンジニア…	200m	3	薬学
58	退職済	まれ	研究・エンジニア…	2km	4	生命科学
39	在籍中	まれ	営業	4km	4	生命科学
40	在籍中	頻繁	研究・エンジニア…	1km	4	生命科学
50	退職済	頻繁	営業			
37	在籍中	まれ	営業			

実行した前処理

ワンクリックで前処理を完了

前処理の透明性を担保

← ダッシュボードに戻る  
データ前処理・結合・…

AI作成の進行状況

- プロジェクト作成  
使用するテンプレート
- データセット設定  
使用するデータセット  
退職リスク予測\_学習…
- 予測する列の選択  
選択した予測する列  
退職
- 4 前処理**  
前処理を設定
- 5 レシピの選択
- 6 学習

生成AIで前処理 前処理を完了する

データセットを表示 統計情報を表示 行を削除

外れ値や統計情報を表示

前処理中のデータセット  
退職リスク予測\_学習用\_デモ用

予測する列  
退職

データ数  
11 詳細を開く →

保存した前処理を使用する  
前処理を選択して下さい

前処理を実行する

年齢	退職	出張	部署	学歴	月収	残業
0	1	1	1	4	33.518	0
0.414	0	0	2	3	22.077	0
1.861	1	1	2	4	41.027	0
-0.103	0	1	1	4	29.77	0
0	0	2	2	4	37.307	0
1.034	1	2	1	2	30.95	0
-0.31	0	1	1	3	31.715	0
-1.447	0	1	0	1	24.503	0

← 1-7 8-12 →

実行した前処理

前処理を保存 1つ戻る 1つ進む 加工回数 22/22 前処理をリセット

- 置換(AI: 単位統一)  
対象の列: 通勤距離
- 置換(AI: 正規化)  
対象の列: 部署
- 置換(AI: 正規化)  
対象の列: 通勤距離
- 欠損値を削除  
対象の列: 出張
- 置換(AI: 正規化)  
対象の列: 出張

# 生成AIによる行削除機能

# テキストから指定して行削除が可能

← デッシュボードに戻る

## ← 前処理

生成AIで前処理    前処理を完了する

データセットを表示    統計情報を表示    行を削除    外れ値や統計情報を表示

2021/12/31の行は不要です。

テキストから生成

テキストからの自動生成機能はAIを使用しており、間違えることもあるので、条件はからなず自分で確認してください。

結合方法  
かつ

条件  
Date    等しい    2021/12/31    条件を削除

+ 条件を追加

実行した前処理

AI作成の進行状況

- プロジェクト作成
- データセット設定
- 予測する列の選択
- 4 前処理**
- 5 レシピの選択
- 6 学習

← デッシュボードに戻る

## 前処理したデータセットを保存する

保存する前処理済みデータセット名  
Paint\_Products\_Sales\_Data\_2021231\_20241130 (前処理済み)

保存する前処理済みデータセットの説明  
データセット説明文が入ります

前処理したデータセットを保存する

実行した前処理の保存設定  
 実行した前処理も保存する

実行した前処理の内容  
入力

+ 実行した前処理の詳細設定

Date	Product A (Gallons)	Product B (Liters)	Product C (Units)
2022/01/01	749	451	280
2022/01/02	820	486	330
2022/01/03	394	302	179
2022/01/04	525	286	207
2022/01/05	533	323	240
2022/01/06	601	366	250
2022/01/07	665	395	266
2022/01/08	804	453	304
2022/01/09	828	501	334
2022/01/10	471	257	199
2022/01/11	486	310	196
2022/01/12	558	341	227
2022/01/13	619	390	248

← 1-4 →

AI作成の進行状況

- プロジェクト作成
- データセット設定
- 予測する列の選択
- 4 前処理**
- 5 レシピの選択
- 6 学習

# データ結合機能

# データ管理にある2つのデータを統合も簡単

← プロジェクト一覧

← Paint\_Products\_Sales\_Data\_2021231\_20241130 (前処理済み)

データ前処理・結合・…

データセットを表示 統計情報を表示

AIを作成する

AIで予測する

プロジェクト管理

ダッシュボード

データセット

前処理

レシピ

学習済みAI

サービス

Date	Product A (Gallons)	Product B (Liters)	Product C (Units)
2022/01/01	749	451	280
2022/01/02	820	486	330
2022/01/03	394	302	179
2022/01/04	525	286	207
2022/01/05	533	323	240
2022/01/06			
2022/01/07			
2022/01/08			
2022/01/09			
2022/01/10			
2022/01/11			
2022/01/12			
2022/01/13			

前処理する

結合する

学習する

削除する

所有者  
hhatamoto1MF

データサイズ  
24 KB

行数  
1065

列数  
4

← プロジェクト一覧

データ前処理・結合・…

AIを作成する

AIで予測する

プロジェクト管理

ダッシュボード

データセット

前処理

レシピ

学習済みAI

サービス

データセット1

Paint\_Products\_Sales\_Data\_2021231\_20241130 (前処理済み)

データセット2

Paint\_Products\_Sales\_Data\_2021231\_20241130

処理方法

行を増やす

実行

	A	B	C	D	E
1058	2024/11/22	613	370	252	
1059	2024/11/23	726	402	276	
1060	2024/11/24	789	460	303	
1061	2024/11/25	394	236	148	
1062	2024/11/26	470	274	179	
1063	2024/11/27	521	319	230	
1064	2024/11/28	543	361	239	
1065	2024/11/29	608	377	245	
1066	2024/11/30	729	457	296	
1067	2024/12/01	786	459	295	
1068	2024/12/02	422	274	155	
1069	2024/12/03	456	288	181	
1070	2024/12/04	512	310	197	
1071	2024/12/05	559	358	218	
1072	2024/12/06	597	368	259	
1073	2024/12/07	726	431	292	
1074	2024/12/08	746	476	316	
1075	2024/12/09	424	274	161	
1076	2024/12/10	505	292	189	
1077	2024/12/11	533	318	225	
1078	2024/12/12	563	339	244	
1079	2024/12/13	623	394	271	
1080	2024/12/14	728	452	283	
1081	2024/12/15	802	475	319	
1082	2024/12/16	425	237	180	

# データ管理にある2つのデータを統合も簡単

← プロジェクト一覧

データ前処理・結合・…

データセットを表示 統計情報を表示

Date	Product A (Gallons)	Product B (Liters)	Product C (Units)
2022/01/01	749	451	280
2022/01/02	820	486	330
2022/01/03	394	302	179
2022/01/04	525	286	207

前処理する  
結合する  
学習する  
削除する

所有者  
hhatamoto1MF

データセット

データ前処理・結合・…

データセット1  
Paint\_Products\_Sales\_Data\_2021231\_20241130 (前処理済み)\_Paint\_Products\_Sales\_Data\_20241201\_20241231

データセット2  
calender\_Master

実績データとカレンダーマスタを結合

処理方法  
列を増やす

左キー  
Date

右キー  
Date

結合方式  
左優先結合 (LEFT JOIN)

列名重複時の処理  
データセット1を残す

実行

← プロジェクト一覧

データ前処理・結合・…

データセットを表示 統計情報を表示

曜日と祝日フラグのカラムが追加

Date	Product A (Gallons)	Product B (Liters)	Product C (Units)	day_of_week	holiday_FLG
2022/01/01	749	451	280	土	1
2022/01/02	820	486	330	日	0
2022/01/03	394	302	179	月	0
2022/01/04	525	286	207	火	0
2022/01/05	533	323	240	水	0
2022/01/06	601	366	250	木	0
2022/01/07	665	395	266	金	0
2022/01/08	804	453	304	土	0
2022/01/09	828	501	334	日	0
2022/01/10	471	257	199	月	1
2022/01/11	486	310	196	火	0
2022/01/12	558	341	227	水	0
2022/01/13	619	390	248	木	0

← 1/6 →

## データ収集・整理からAI構築までサポート

AI活用の壁として、「データの収集・整理」があります。

MatrixFlowではデータの収集・前処理から、AIアルゴリズム構築まで「現場で使える&高い精度のAI」を実現。

また、プロのデータサイエンティストが、現場のAI活用をサポートいたします。

### MatrixFlow AI構築サポート

課題設定/計画立案サポート

データ収集計画/データ収集サポート

データクレンジング/前処理支援サポート

AIモデル作成/推論サポート

AIモデル評価/実運用サポート



2週間~4週間に1回Web面談でのサポートを提供しております。  
ベーシックプラン月額料金に含まれています。