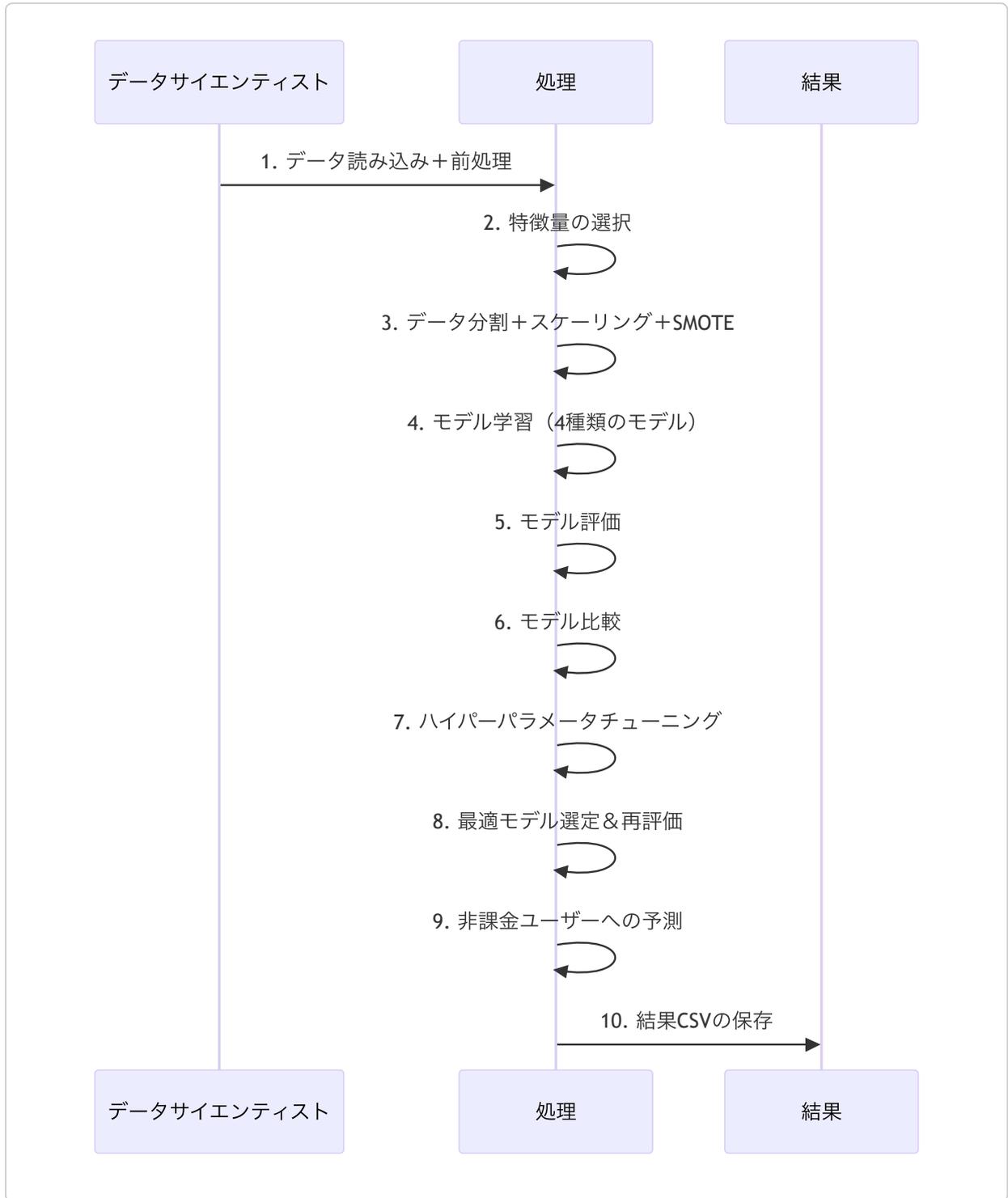


課金化予測モデル 概要

1. 処理フロー



2. 処理フローの解説

1. データ読み込み + 前処理

DynamoDB からエクスポートされた JSON ファイルを読み込み、**欠損値の補完** や **不要な列の削除** を行います。また、**status_code=2** を課金ユーザーとして扱うなど、予測に必要な前処理を実施します。

2. 特徴量の選択

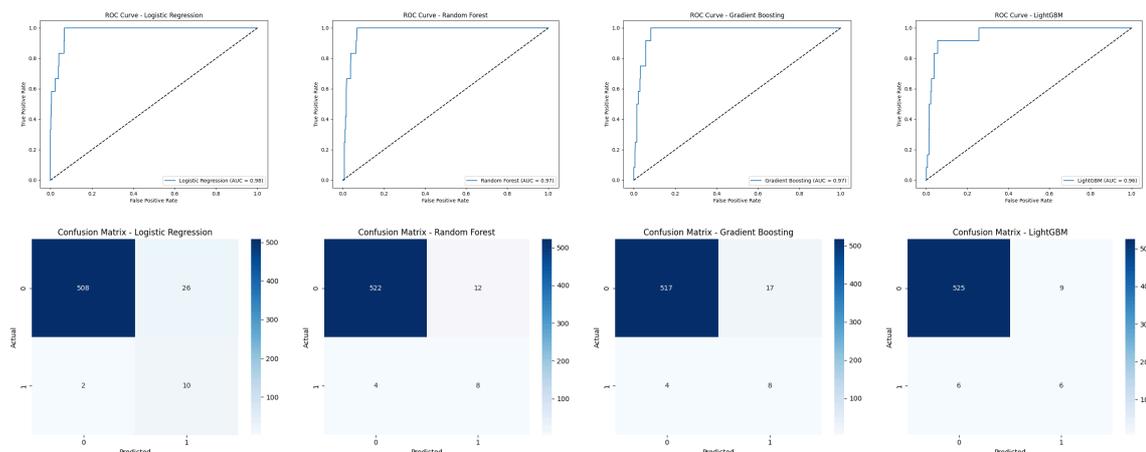
画像ラベル数 (image_label_count)、年齢 (age)、経過時間 (updated_time_diff)、**テキストから抽出した「依存度」「経済的余裕」スコア** など、モデルに取り込みたい特徴量を抽出します。

3. データ分割 + スケーリング + SMOTE

データを「学習用」と「テスト用」に分割し、数値を **標準化** します。さらに、課金ユーザーが少ない不均衡データに対して、**SMOTE** (少数クラスのオーバーサンプリング) を適用し、学習バランスを改善します。

4. モデル学習 (4種類)

ロジスティック回帰・ランダムフォレスト・勾配ブースティング・LightGBM の 4 種を同時に学習させます。

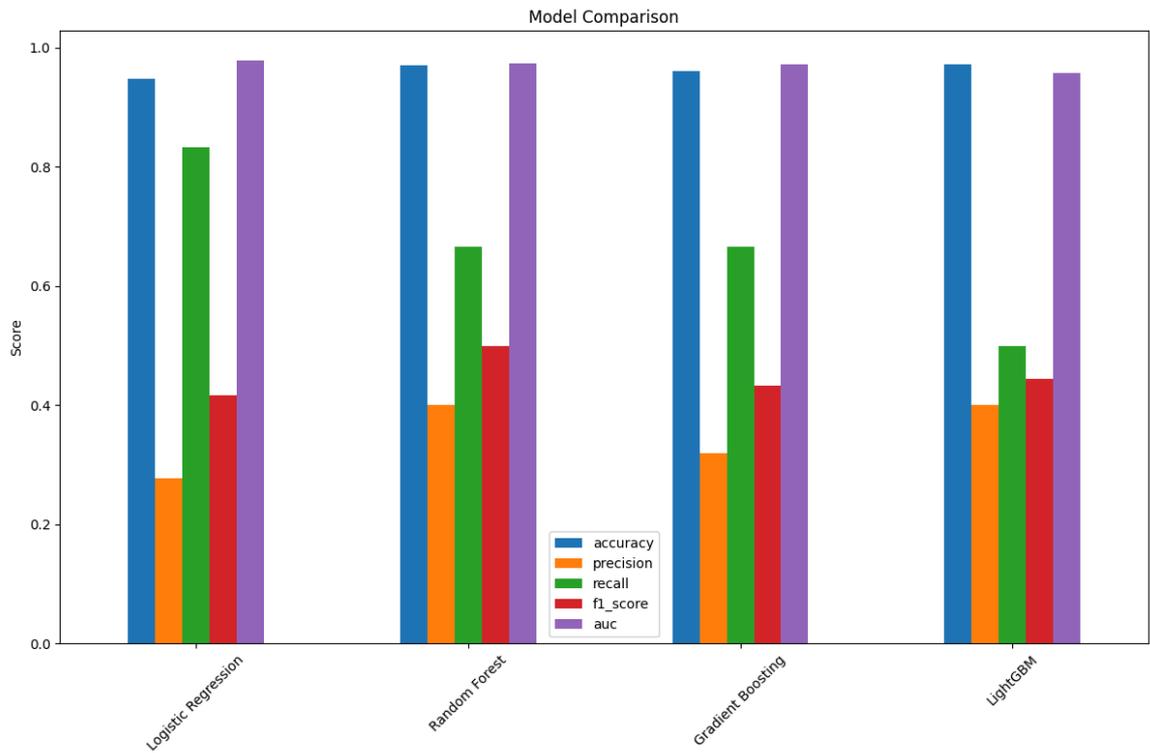


5. モデル評価

テストデータに対して、精度 (Accuracy)、適合率 (Precision)、再現率 (Recall)、F1スコア、AUCなど、多面的な指標で性能をチェックします。

6. モデル比較

上記の評価指標を一覧化し、どのモデルが最も優れているか可視化します。



7. ハイパーパラメータチューニング

各モデルのパラメータ（例：木の深さ、学習率等）を細かく試し、より高い精度を出せる設定を自動的に探索します。

8. 最適モデル選定 & 再評価

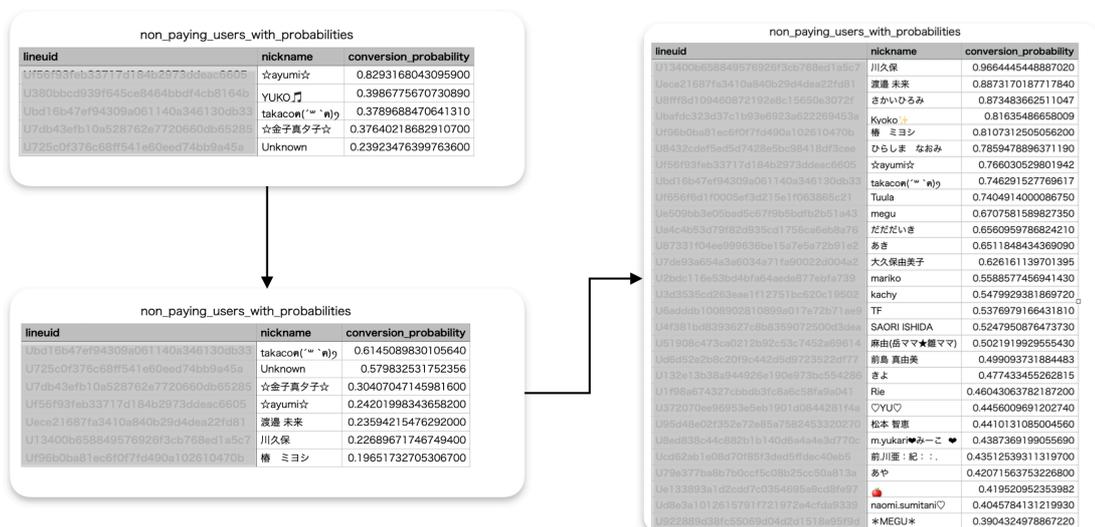
チューニングで最良の結果を出したモデル（例：LightGBM）を最終採用し、再度評価して最終的なスコアを確定します。

9. 非課金ユーザーへの予測

「現時点で課金していないユーザー」の課金化率（課金する確率）を計算し、優先的にアプローチすべき顧客を選別できます。

10. 結果CSVの保存

最終的に、課金化率が閾値以上のユーザーをCSVファイルなどで一覧出力し、マーケティング施策などに活用します。



3. ログの概要説明

ログには、実際の実行結果として以下が含まれています。

- 欠損値の確認結果、カテゴリカル列の情報、最終的に使用する特徴量のリスト
- 4種類のモデル（Logistic Regression, Random Forest, Gradient Boosting, LightGBM）の評価結果（Accuracy、Precision、Recall、F1、AUCなど）
- モデル同士の比較結果
- ハイパーパラメータチューニングの最適パラメータとスコア
- 最終的に選ばれたモデル（例：LightGBM）での評価メトリクスと混同行列
- 非課金ユーザーへの予測確率を算出し、CSVへ保存する流れ

4. ビジネス上のメリット

この仕組みにより、「課金する可能性の高いユーザー」を効率よく抽出できるため、**営業・マーケティング施策**を集中的に行いやすくなります。たとえば、**優先ユーザーには追加クーポンや特別なプランの案内**を行うなど、リソースを最適配分することで収益拡大が期待できます。

5. まとめ

- 欠損補完やSMOTEなどを活用し、不均衡データでも精度が出やすい仕組みを整備。
- *introduction*などのテキストも活用して、依存度や経済的余裕といった潜在的要素を数値化。
- 4種類のモデルを比較し、最適モデルをチューニング後に運用。
- 予測結果を基に「課金しそうなユーザー」に優先アプローチを実施。

今後も実際の施策結果をフィードバックしながら、継続的にモデル精度の向上が可能です。