



AWS Summit

Tokyo





今日から始められる、機械学習！ Amazon Machine Learningのご紹介

Toshiaki Enami, Yuta Imai
Solutions Architect, Amazon Data Services Japan, K.K.

■ Gold Sponsors



Empowered by Innovation



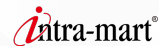
■ Global Sponsors



■ Silver Sponsors



■ Bronze Sponsors



■ Global Tech Sponsors

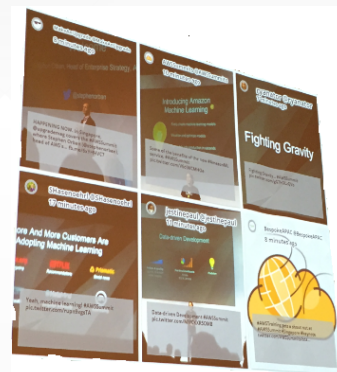


■ Logo Sponsors



ハッシュタグ **#AWSummit**

と **#DevCon**で、皆さんのツイート
が展示エリアの大画面に表示されます



公式アカウント **@awscloud_jp**
をフォローすると、ロゴ入り
コースターをプレゼント



【コースター配布場所】

メイン展示会場、メイン会場1F受付、デベロッパーカンファレンス会場

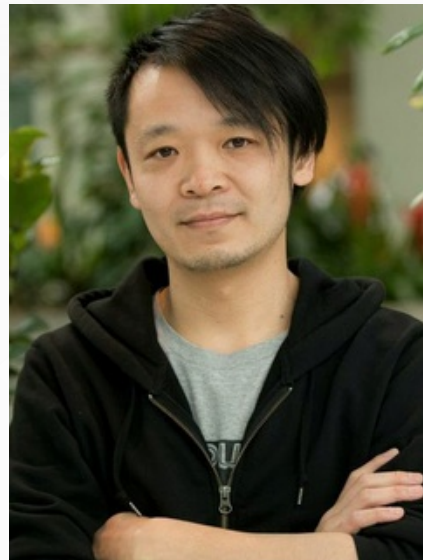


自己紹介

今井雄太 (いまいゆうた)

yuimai@amazon.co.jp

- ソリューションアーキテクト
- ゲーム、広告系のお客様を主に担当
- 好きなAWSのサービス：
 - Amazon Elastic MapReduce



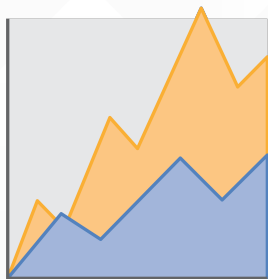
アジェンダ

1. 機械学習とは？
2. Amazon Machine Learning
3. アーキテクチャへの組み込み



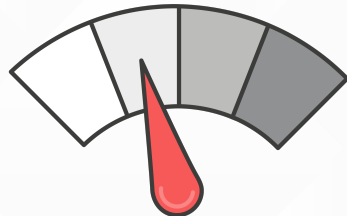
Part1: 機械学習とは？

3種類のデータ駆動型アプリケーション



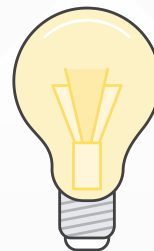
遡及的分析とレポーターティング

Amazon Redshift,
Amazon RDS
Amazon S3
Amazon EMR



即時の判断
リアルタイム処理と
ダッシュボード

Amazon Kinesis
Amazon EC2
AWS Lambda



予測
スマートアプリ
ケーションを作
成可能にする

Amazon ML

機械学習の例

- すべてAmazon Machine Learningで実現可能な例-

- このメールはスパムメールか？
 - 過去のメールアーカイブをもとにYes/Noを予測する
- この商品は本、日用品、食品のいずれなのか？
 - 多くの商品データをもとにその商品のカテゴリを予測する
- 明日の売上はどのくらいになるか？
 - 過去の売上データなどをもとに明日の売上を予測する

機械学習の例

- このメールはスパムメールか？
 - 過去のメールアーカイブをもとにYes/Noを予測する

教師データ

送信者domain	サーバーIP	送信時刻	言語	Spam?
A	123.123.123.123	...	JP	1
B	111.111.111.111	...	EN	
B	111.111.111.111	...	EN	
D	123.456.789.012	...	FR	1

別の方法でSpam判定済みの教師データ(過去のデータ)をもとに

予測対象データ

送信者domain	サーバーIP	送信時刻	言語	Spam?
B	123.123.123.123	...	JP	
A	123.456.789.012	...	FR	

予測対象データ(新しいデータ)のSpam判定をする

機械学習の例

- この商品は本、日用品、食品のいずれなのか？
 - 多くの商品データをもとにその商品のカテゴリを予測する

教師データ

商品名	価格	大きさ	メーカー	カテゴリ
A	123	...	AA	Book
B	456	...	BB	Book
C	100	...	CC	Food
D	500	...	DD	Grocery

大量の商品のデータをもとに

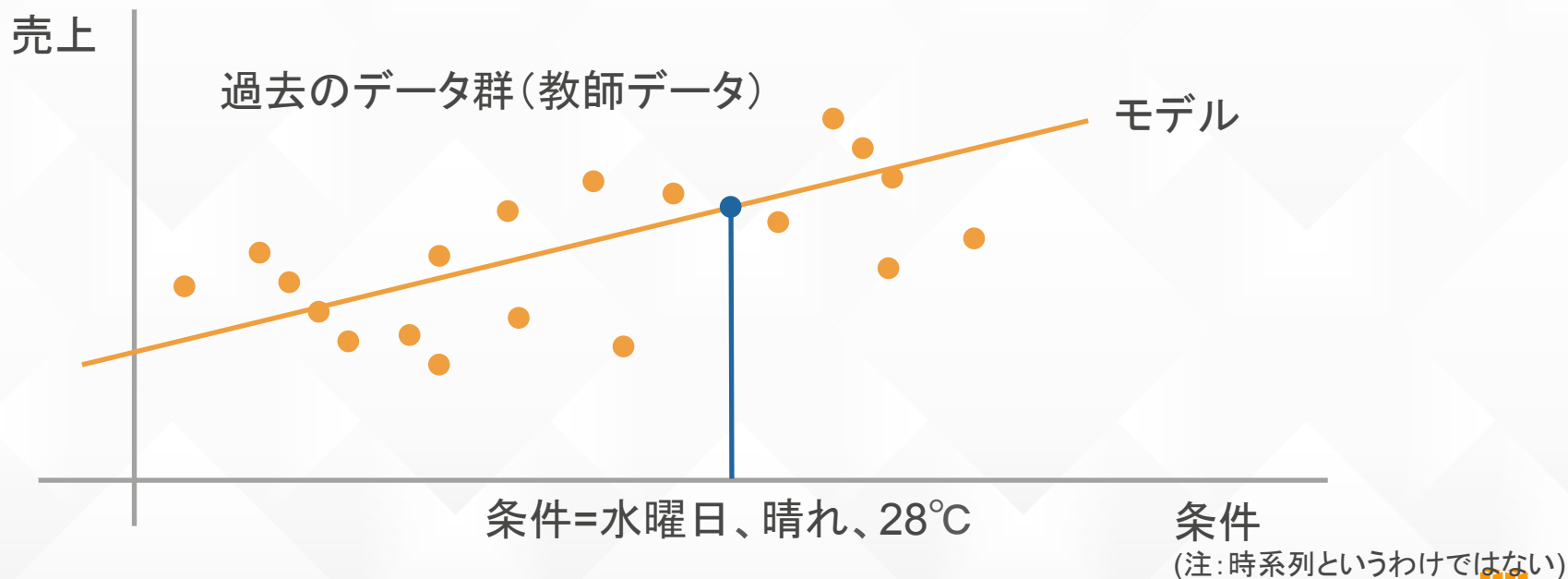
予測対象データ

商品名	価格	大きさ	メーカー	カテゴリ
E	500	...	EE	
F	600	...	FF	

未知の商品のカテゴリを予測する

機械学習の例

- 明日の売上はどのくらいになるか？
 - 過去の売上データなどをもとに明日の売上を予測する



その他の例 ...

詐欺の検知

不正取引の検知、不正クリック検知、スパムeメールのフィルタリング、疑わしいレビューのマーキング ...

パーソナライゼーション

コンテンツのレコメンデーション、予測的なコンテンツロード、ユーザエクスペリエンスの改善 ...

ターゲットマーケティング

オファーとお客様のマッチング、マーケティングキャンペーンの選択、クロスセリングやアップセリング ...

コンテンツ分類

ドキュメントのカテゴリ分類、履歴書と採用マネージャのマッチング ...

変動予測

サービスを使うのを止めそうなお客様の検知、無料ユーザからアップグレードのオファー ...

カスタマーサポート

お客様からのメールの適切な転送先推測、ソーシャルメディアリスニング ...

スマートアプリケーション(機械学習を有効活用したアプリケーション)が続々と登場しないのはなぜか?

1. 機械学習の専門家が少ない
2. 機械学習の仕組みを作り、スケールさせることは**技術的に困難**
3. モデルとアプリケーションのギャップを縮めるには、**長い時間と高い費用**が必要になる

スマートアプリケーションを作るには

専門家

データサイエンティスト
の人数は限られる

外注するのは高くつく

技術

多くの選択肢があるが
決定的なものが無い

使いこなし、スケールさ
せることが困難

カスタムソリューションを
作成するために多くの細
かい作業が毎回必要に
なる

使いやすさ

複雑で間違いを起こし
やすいワークフロー

特殊なプラットフォーム
とAPI

モデルライフサイクル管
理を無駄に再発明

スマートアプリケーションを作るには

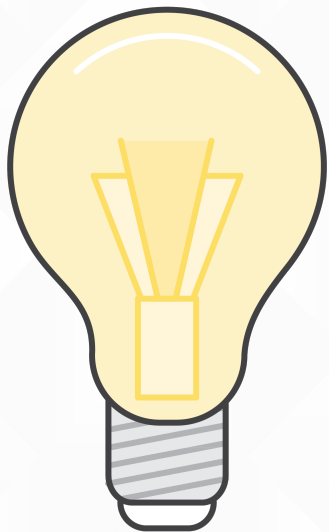
- 機械学習に強くて
- RやPython, 場合によってはHadoopやSparkに
 明るくて
- 特定のビジネス分野の経験が深い

こういう人を採用するのは難しい！



そこで . . .

Amazon Machine Learningの登場

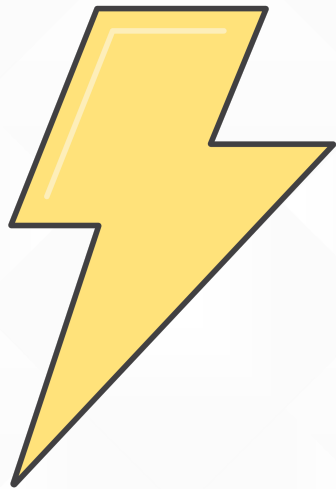


使いやすく、マネージド。開発者のために作られた機械学習サービス



Part2: Amazon Machine Learning

Machine Learning as a Service



Amazonが提供するアルゴリズム

- 利用者は自分でアルゴリズムの実装や詳細なチューニングを行う必要がない

パッケージサービスとしての提供

- 必要なワークフローが予め提供されている

スケーラビリティ

- 利用者はシステムの拡張やその運用についても考える必要がない



Amazon Machine Learningでできること

取り扱える予測モデルとアルゴリズム

二項分類

ロジスティック回帰

このメールはスパム？ Yes？ No？

多クラス分類

多項式ロジスティック回帰

これは本？ 車？ 食べ物？

回帰分析

線形回帰

あしたは水曜日。在庫はいくつくらい必要？

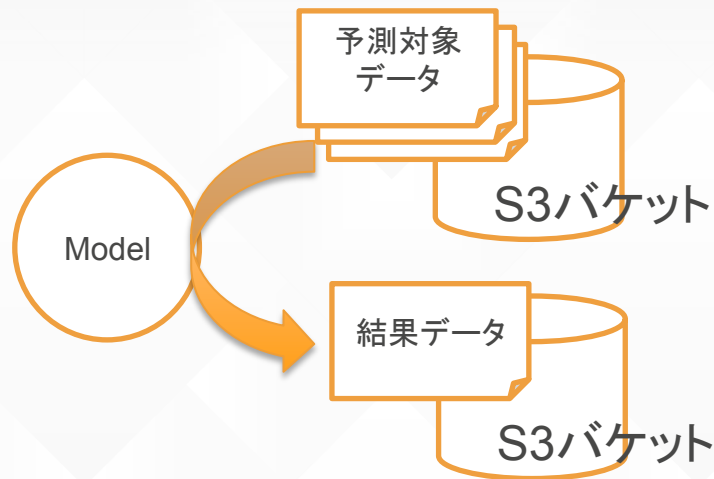


これらのモデルを使って . . .

予測手法

バッチ予測

S3等にアップロードされた予測対象データに対してまとめて予測を実施



リアルタイム予測

データ 1 件ずつAPIを使って予測を実施する

```
ml      = Aws::MachineLearning::Client.new
record = { attr_A: 'foo', attr_B: 'bar', ... }
result = ml.predict(
  ml_model_id: MODELID,
  record: record,
  predict_endpoint: ENDPOINT
)
```



Amazon Machine Learningの使い方

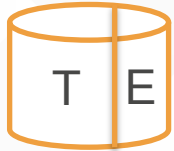
4つのステップ

1. 教師用/評価用データを準備
2. モデルを作成(学習、トレーニング)
3. モデルの品質評価
4. 実際の予測の実施

1. 教師用/評価用データを準備

Data Sourceの作成

s3://SOURCEDATA



S3、Redshift、RDSが
利用可能

S3、 Amazon Redshift、 RDS上の
MySQLに格納されたデータを指定し、
教師データ/評価用データとして利用す
る

デフォルトの設定を使うと、自動的に7
割を教師データ(T)、 3割を評価用デー
タ(E)に分割して管理してくれる

2. 教師データからモデルを作成

教師データを元にAmazon Machine Learningが自動的にモデルを選択してくれる。例えば予測対象のカラムが二値型であれば二項分類が自動的に選択される。

二項分類

ロジスティック回帰

多クラス分類

多項式ロジスティック回帰

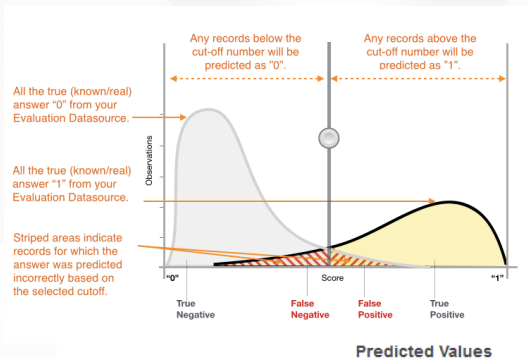
回帰分析

線形回帰

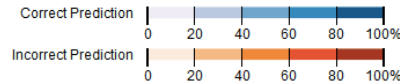
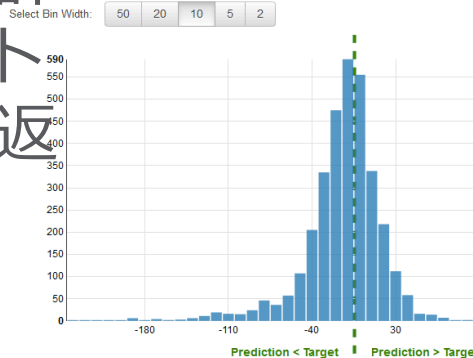
3. モデルの品質評価

作成したモデルに対して評価
(評価用のデータを流してみ
て予測の精度を測ること)を
実施する。

精度に満足できない場合、
教師データのETLや量を精査し、
トレーニングと品質評価を
繰り返す。



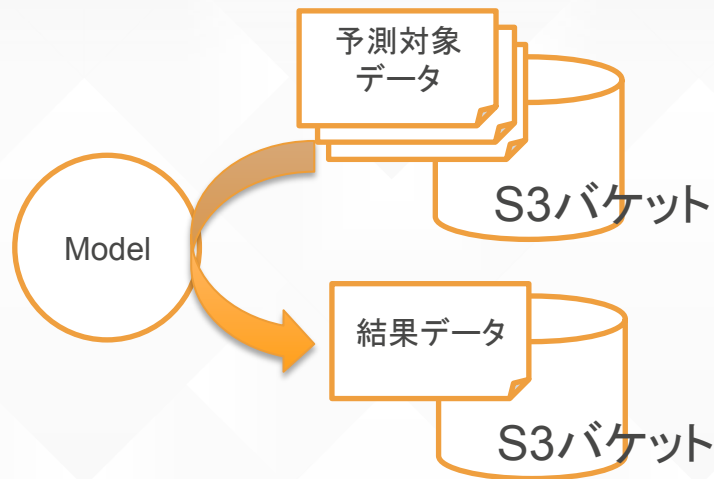
	Romance	Thriller	Adventure	Total	F1
Romance				57.92% (49.1k)	0.78
Thriller				21.23% (18.0k)	0.33
Adventure				20.85% (17.7k)	0.32
Total	77.56% (65.8k)	9.33% (7910)	13.12% (11.1k)	100.00% (84.8k)	0.47



4. 実際の予測の実施

バッチ予測

S3(S3、Redshift、RDS) 等にアップロードされた予測対象データに対してまとめて予測を実施



リアルタイム予測

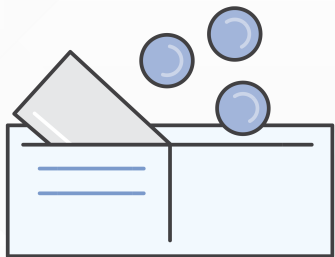
データ 1 件ずつAPIを使って予測を実施する

```
ml      = Aws::MachineLearning::Client.new
record = { attr_A: 'foo', attr_B: 'bar', ... }
result = ml.predict(
  ml_model_id: MODELID,
  record: record,
  predict_endpoint: ENDPOINT
)
```



料金

使った分だけ、安価な支払い



データ分析、モデルトレーニング、評価:

\$0.42/インスタンス時

バッチ予測: **\$0.10/1000**

リアルタイム予測: **\$0.10/1000**

+ 1時間毎のキャパシティリザベーションチャージ(モデルサイズ10MBあたり\$0.001)



リージョン

リージョン

- 現在のところus-east-1のみ
- ただし他のリージョンのS3も利用可能



ユースケースを考えてみる

機械学習を活用するために

- 機械学習をうまく活用するためには、自分たちが解決したい問題を予測モデルに対してうまく落としこんでやることが非常に重要
- ここから何枚かのスライドでは、その落とし込みの例について考えてみる

広告の不正クリック検出

- 教師データ
 - 実際のクリックログ。不正かそうでないかを別の手法で判断してフラグがついたもの。
- 問題の分類
 - 二項分類
- 出力
 - ログ一行ごとに不正クリックかそうでないか

デモグラ推定

- 教師データ
 - デモグラがわかっているユーザーの行動ログ
- 問題の分類
 - 多クラス分類
- 出力
 - ユーザーのカテゴライズ（たとえばF1,M1のような）されたデモグラ

デモグラに基づいたレコメンデーション

- 教師データ
 - 購入履歴にデモグラ情報がマッピングされたもの
- 問題の分類
 - 多クラス分類
- 出力
 - このユーザーはF1なので商品カテゴリ●●、みたいな出力

顔写真から特定の人物かどうかを判定する

- 教師データ
 - 顔写真をグレースケールにしたビットマップ
- 問題の分類
 - 二項分類
- 出力
 - この写真は●●さんである/違う

問題をどうモデルに落としこむかが非常に大事

- 教師データの量
- 教師データの形式、入力値
- どのモデルを使うのか
- 何の値を予測/判定するのか

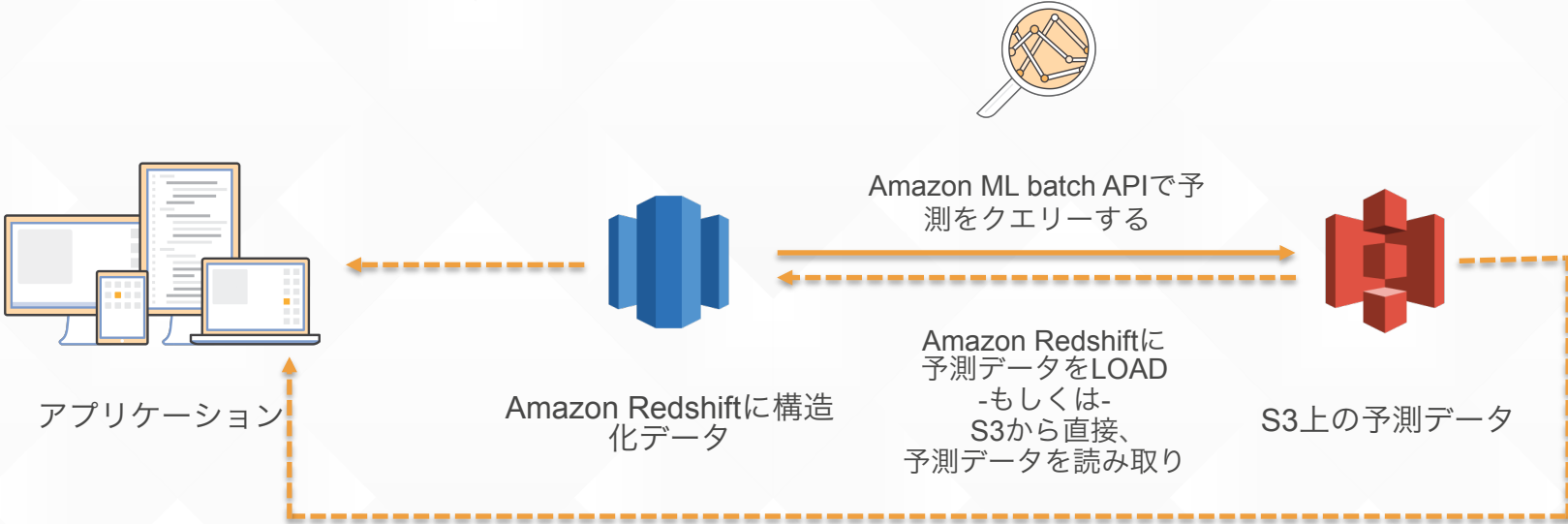


Part3: アーキテクチャへの組み込み

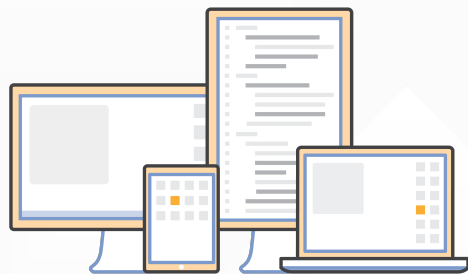
EMRを使用したバッチ予測



Amazon Redshiftを使ったバッチ予測



インタラクティブアプリケーション用のリアルタイム予測



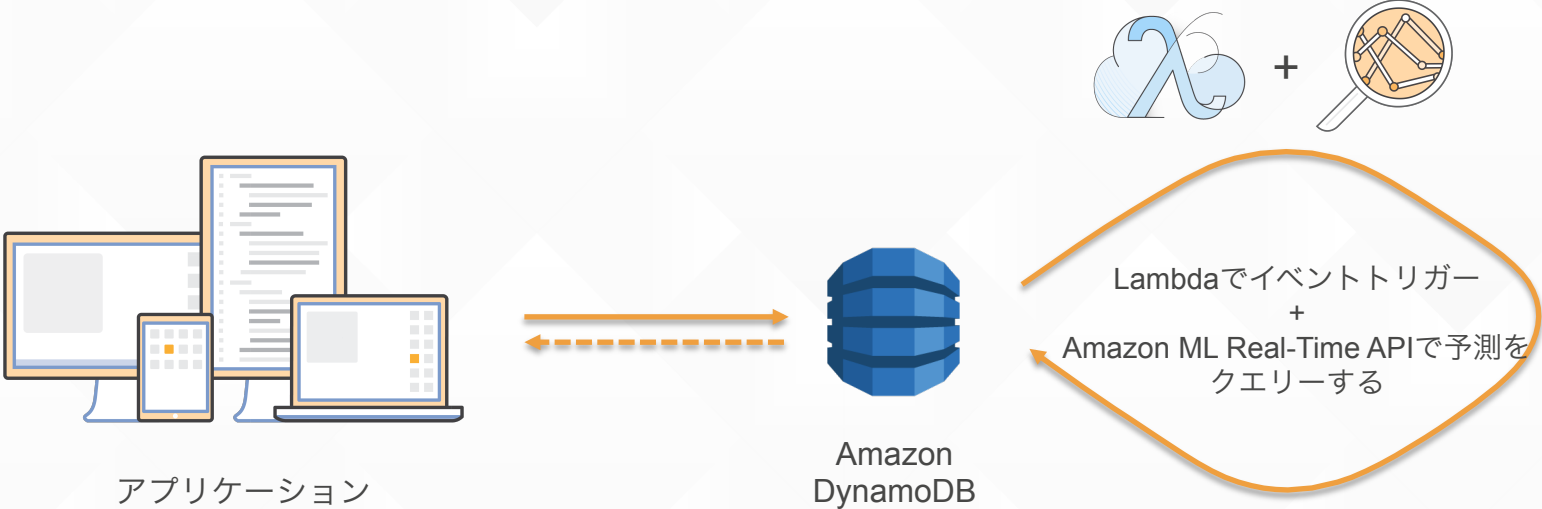
アプリケーション



Amazon ML Real-Time API
で予測をクエリーする



既存データフローに予測を追加する





さいごに

Amazon Machine Learningを使うと・・・

- 機械学習の導入を非常に容易にしてくれる。機械学習やそのためのシステムについての専門家がいなくても可能。
- S3やRedshiftにデータがあればすぐにでも使い始めることができる。

溜めていたデータの評価や、その先のビジネスへの活用を簡単に始めることができる。



Thank You



より深い情報が知りたければ

DOCUMENTATION HOME

[Amazon Machine Learning
Documentation](#) >

RELATED LINKS

[AWS Glossary](#)[Getting Started with AWS](#)[SDKs & Tools](#)[AWS Documentation on Kindle](#)[AWS General Reference](#)[AWS Training](#)[AWS Case Studies](#)[AWS Whitepapers](#)

Amazon Machine Learning Documentation

Amazon Machine Learning is a service that makes it easy for developers to build smart applications, including applications for fraud detection, demand forecasting, targeted marketing, and click prediction. The powerful algorithms of Amazon Machine Learning create machine learning (ML) models by finding patterns in your existing data. The service uses these models to process new data and generate predictions for your application.

Developer Guide

Provides a conceptual overview of Amazon Machine Learning and includes detailed instructions for using the service.

[HTML](#) | [PDF](#)

Machine Learning Concepts

Provides an overview of the basic concepts in the field of machine learning.

[HTML](#) | [PDF](#)

API Reference

Describes all the API operations for Amazon Machine Learning in detail. Also provides sample requests and responses for supported web service protocols.

[HTML](#) | [PDF](#)

Learning Algorithm

The learning algorithm's task is to learn the weights for the model. A learning algorithm consists of a loss function and an optimization technique. The loss is the penalty that is incurred when the estimate of the target provided by the ML model does not equal the target exactly. A loss function is a function that quantifies this penalty as a single value. An optimization technique seeks to minimize the loss. **In Amazon Machine Learning, we have three loss functions, one for each of the three types of prediction problems. The optimization technique used in Amazon Machine Learning is online Stochastic Gradient Descent (SGD).** SGD makes sequential passes over the training data, and during each pass, updates feature weights one example at a time with the aim of approaching the optimal weights that minimize the loss.

Amazon Machine Learning uses the following learning algorithms:

- For **binary classification**, Amazon Machine Learning uses **logistic regression (logistic loss function + SGD)**.
- For **multiclass classification**, Amazon Machine Learning uses **multinomial logistic regression (multinomial logistic loss + SGD)**.
- For **regression**, Amazon Machine Learning uses **linear regression (squared loss function + SGD)**.