

ネットワーク侵入検知データセットにおける 半教師あり学習を用いた 学習データの少ないデータセットでの機械学習精度向上

サイトの利用者

アクセス



アクセスログが残る。
Ex. プロトコルの種類、
データ量、接続時間・・・



AI

正常なアクセス? or 異常なアクセス?

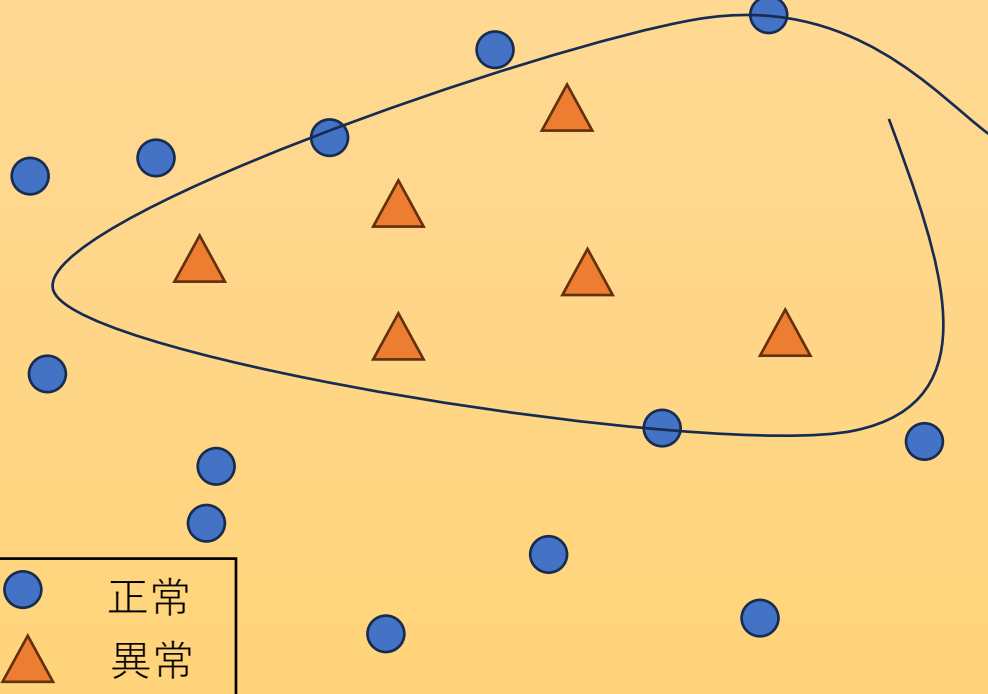


キーワード：機械学習，AI，ニューラルネットワーク，オートエンコーダ，勾配ブースティング決定木，勾配降下法，アンサンブル学習，決定木，特徴量エンジニアリング，オーバーサンプリング，不均衡データセット，TensorFlow(Keras)，LightGBM，SMOTE

機械学習とは

学習(分類問題の場合)

=機械学習モデルを作る
=仕切り線を入れること



機械学習モデルとは？

データ $x(x_1, x_2, x_3 \dots)$



機械学習モデル

$y=0$ (正常)

$y=1$ (異常)

重要!

データ x を入れたら、予測 y を返す
関数 $y=f(x)$ を作ること

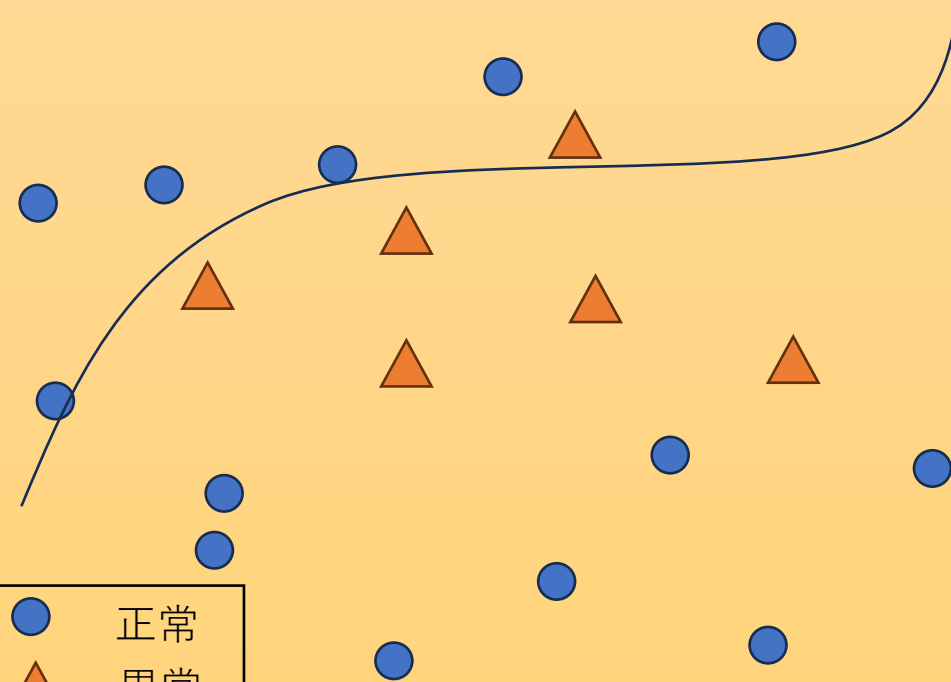
学習の流れ



境界線の関数 $y=f(x)$ のパラメータ($a_1, a_2 \dots$)を変えていくよー
Ex. $f(x) = a_1x_1 + a_2x_2 + \dots$

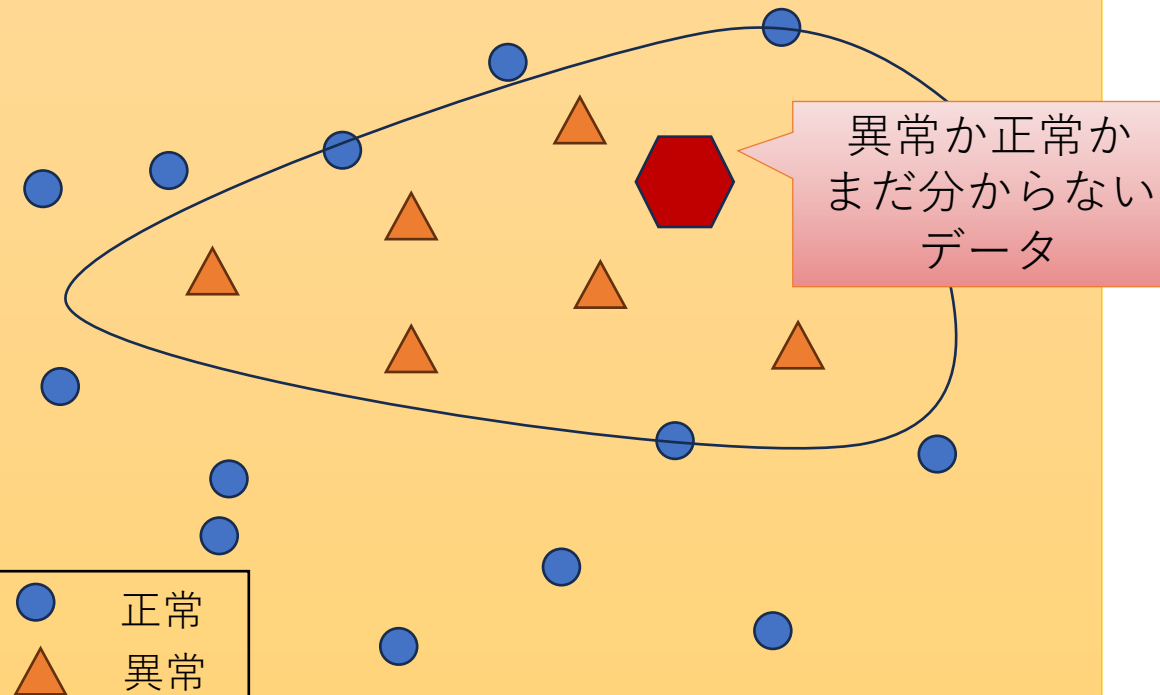
学習前

正常と異常を分類
できていない



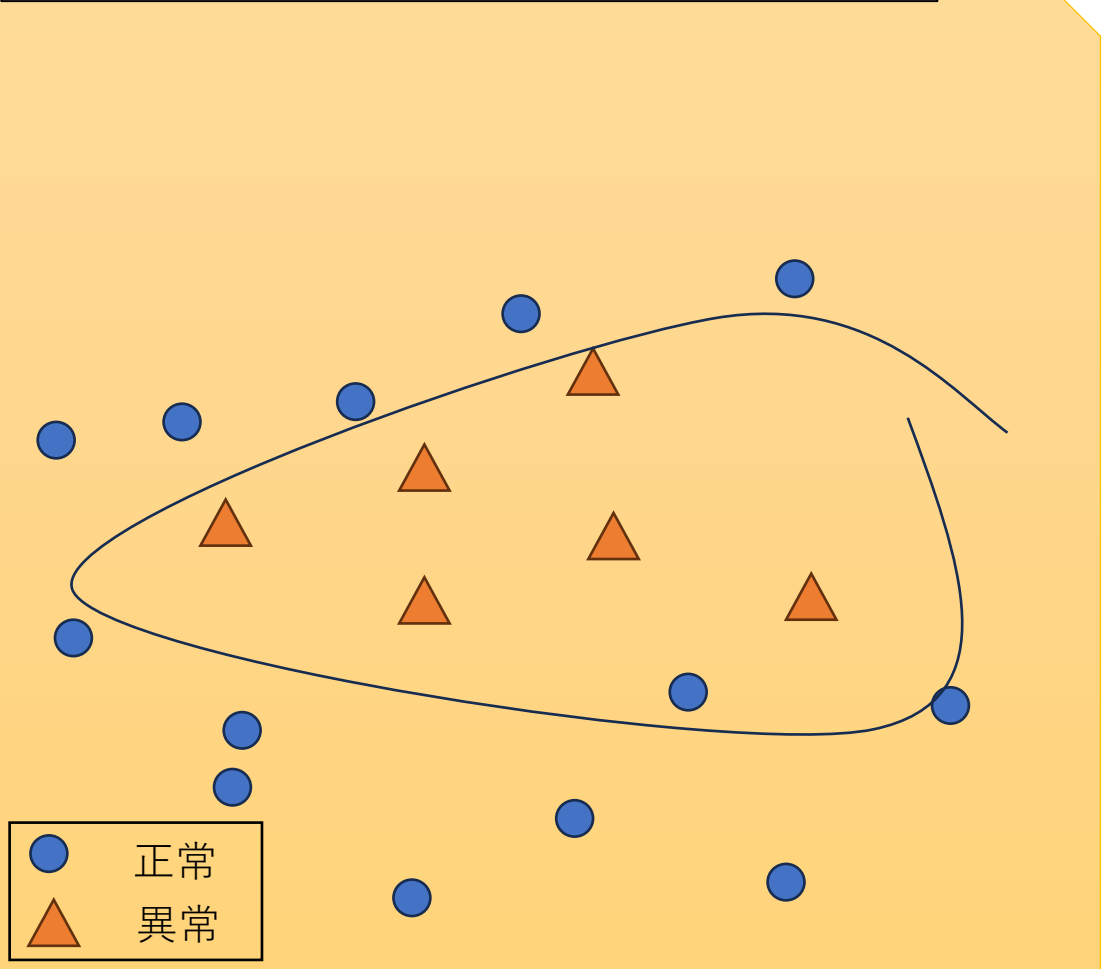
学習後

ほぼ分類できている



教師データが少ないとどうなるか

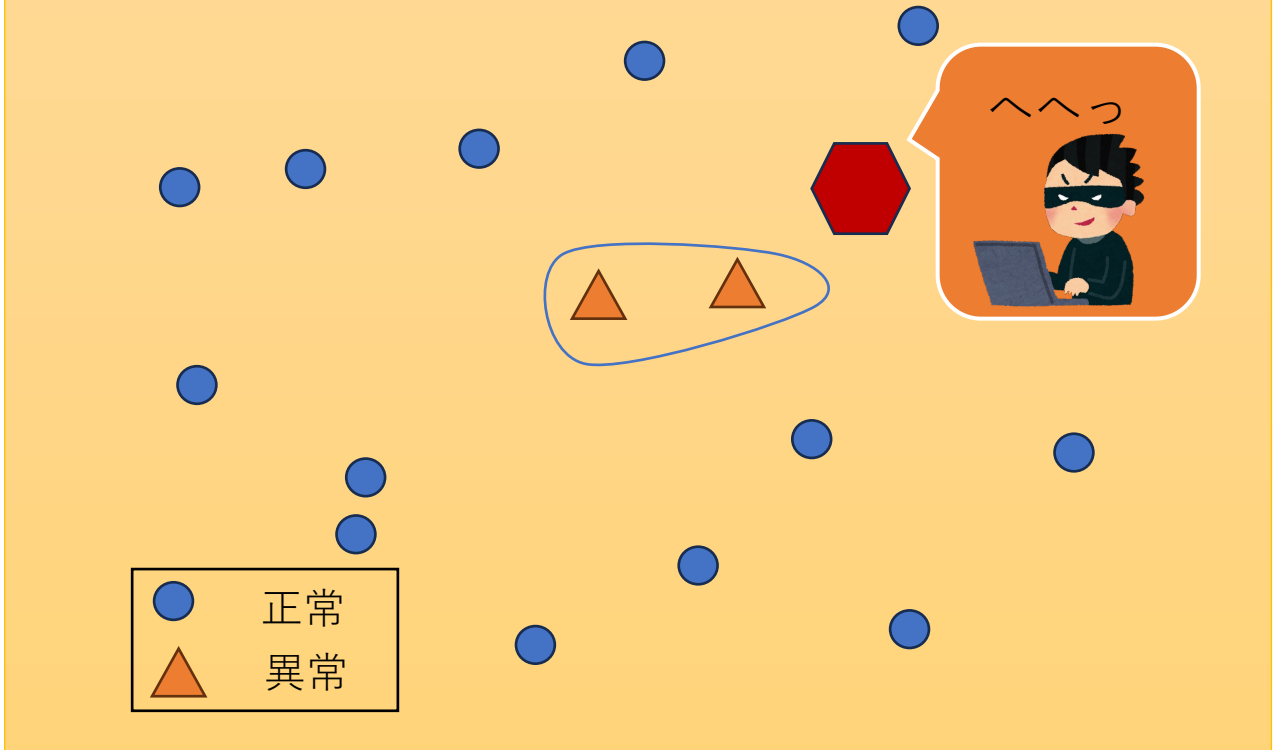
学習データが多い場合



学習データが少ない場合



どこに境界線引いたら良いか
分からないよ～

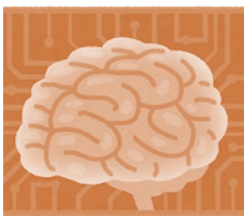


データが少なくても精度を上げるために

方法1: 不正なデータの特徴を掴んだデータを作る

※例

データx ($x_1, x_2, x_3 \dots, x_{20}$)



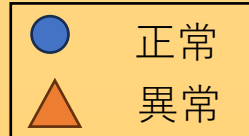
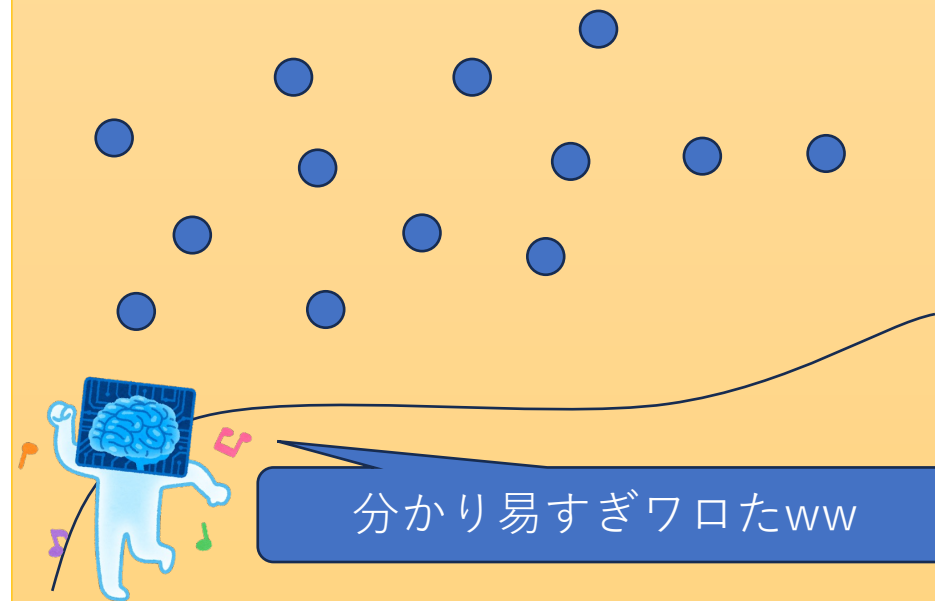
不正なデータは、 x_{14} と、 x_{18}
の値が大きい傾向がある。

新しい特徴量 x_{21} を作成

$$x_{21} = 100 \times x_{14} + 100 \times x_{18}$$

キーワード：特徴量エンジニアリング，
オートエンコーダによる特徴量抽出

x_{21} を使った分布



データが少なくても精度を上げるために

方法2: 不正なデータもどきを沢山生成する.



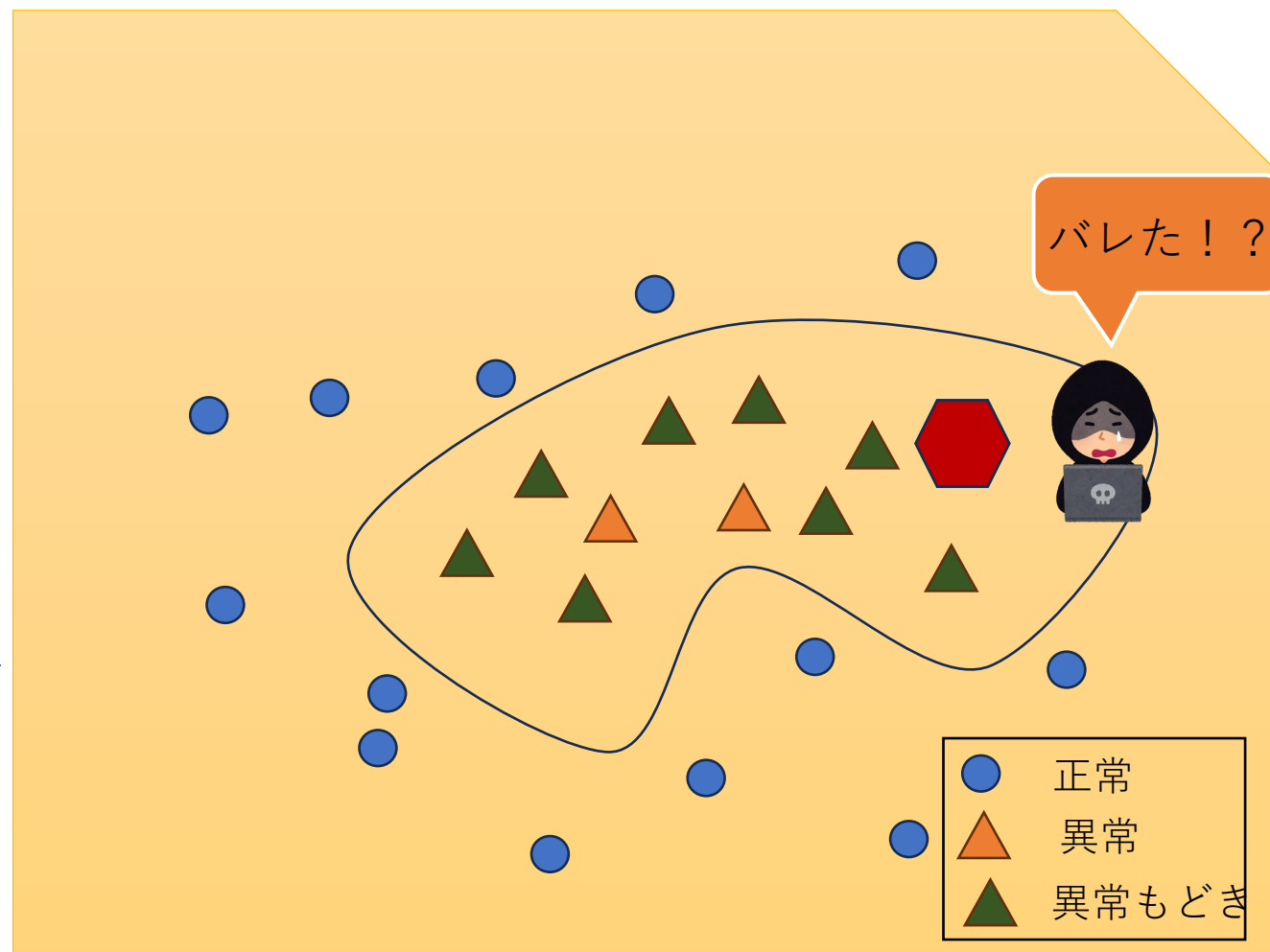
異常データに近いデータ沢山作るよ



ここら辺に線を引けばいいかな～

精度向上

キーワード：オーバーサンプリング, SMOTE



半教師あり学習とは？

教師あり学習

正常と異常の境界線を引くこと

問題

学習データが少ないと精度が低い
(上手く境界線が引けない)

教師なし学習

データの特徴ごとにまとめること

問題

正常と異常を分類できない



半教師あり学習

ポイント

正常と異常の特徴を捉えて、学習データが少なくても
精度良く境界線を引くことができる！かも？

実際、研究で書くプログラムの例

```
5 from keras.layers import Dense
```

新規 *

```
6 def main():
```

```
7     model = keras.Sequential([
```

```
8         Dense(units=19, activation='relu', input_dim=38, name='encoder1'),
```

```
9         Dense(units=10, activation='relu', name='encoder2'),
```

```
10        Dense(units=19, activation='relu'),
```

```
11        Dense(units=38, activation='relu'),
```

```
12    ])
```

```
13    model.compile(optimizer='adam', loss='mean_squared_error', metrics=['accuracy'])
```

```
14    model.fit(x_train, x_train,
```

```
15            epochs=1, # データセットを使って学習する回数
```

```
16            batch_size=32,
```

```
17            validation_data=(x_train, x_train), # 評価用データ（検証データ）の指定
```

```
18            )
```

```
19    x_pred = model.predict(x_test) # モデルを使って実際に、予測
```

```
22 if __name__ == '__main__':
```

```
23     main()
```

モデルの設定を定義

学習

未知のデータに対して予測