

# PythonとTensorFlowによるオートエンコーダ実装入門: 基本から畳み込みまで

## 1. はじめに

### オートエンコーダとは何か？ (What is an Autoencoder?)

オートエンコーダ (Autoencoder, AE) は、ニューラルネットワークの一種であり、主に教師なし学習 (より正確には自己教師あり学習) に用いられます。その基本的な目的は、入力されたデータを効率的に圧縮 (エンコード) し、その後、圧縮された表現から元のデータを可能な限り忠実に復元 (デコード) することです<sup>1</sup>。

簡単なアナロジーとして、長い文章を要点にまとめる作業 (エンコード) と、その要点から元の文章に近いものを再現する作業 (デコード) を考えてみましょう。この過程で、元の文章の細かなニュアンスが失われる可能性があるように、オートエンコーダによる圧縮・復元も通常、完全な再現ではなく、ある程度の情報損失を伴います (非可逆圧縮)<sup>8</sup>。

オートエンコーダはしばしば「教師なし学習」と呼ばれますが、技術的には「自己教師あり学習」の一形態と見なすのがより正確です<sup>1</sup>。なぜなら、外部から与えられるラベル (正解データ) を使うのではなく、入力データそのものを再構築の目標 (教師信号) として利用するためです。ネットワークは、入力データを自分自身の「教師」として学習を進めます。

### なぜオートエンコーダを学ぶのか？ (Why learn Autoencoders?)

オートエンコーダは、そのシンプルな構造にもかかわらず、多くの強力な応用可能性を持っています。主な用途をいくつか紹介します。

- **次元削減 (Dimensionality Reduction):** 画像のような高次元データを、重要な特徴を保持したまま低次元の表現に圧縮します。これは主成分分析 (PCA) と似た目的を持ちますが、ニューラルネットワークを用いることで、データ間の複雑な非線形関係も捉えることができます<sup>1</sup>。
- **特徴抽出 (Feature Extraction):** エンコーダ部分は、データから本質的な特徴を学習します。学習済みのエンコーダは、分類などの他の機械学習タスクのための特徴抽出器として利用できます<sup>1</sup>。
- **データノイズ除去 (Data Denoising):** ノイズが含まれたデータを入力とし、元のクリーンなデータを復元するように学習させることで、データのノイズを除去するモデルを構築できます<sup>4</sup>。
- **異常検知 (Anomaly Detection):** 正常なデータのみを使ってオートエンコーダを学習させると、異常なデータ (学習データに含まれないパターン) が入力された際にはうまく復元できず、再構築誤差が大きくなる傾向があります。この性質を利用して異常を検知できます<sup>4</sup>。
- **生成モデルの基礎 (Foundation for Generative Models):** 変分オートエンコーダ (VAE) のような発展形は、学習したデータ分布から新しいデータを生成することも可能です<sup>1</sup>。

## 本記事で学ぶこと (What you will learn in this article)

本記事では、深層学習の初心者の方を対象に、オートエンコーダの基本的な概念から実践的な実装までを段階的に解説します。具体的には、以下の内容を扱います。

1. オートエンコーダの基本的な仕組み(アーキテクチャ、学習プロセス、用途)
2. ニューラルネットワーク構築におけるPython、TensorFlow、Kerasの役割
3. 実践的なデータ準備(Fashion MNISTデータセットの読み込みと前処理)
4. 基本的な(全結合層を用いた)オートエンコーダのエンコーダ・デコーダの構築
5. モデルのコンパイル(オプティマイザ、損失関数の設定)
6. モデルのトレーニング(学習の実行)
7. 畳み込みオートエンコーダの構築と基本的なオートエンコーダとの比較
8. 学習済みモデルを用いた画像の再構築と結果の評価(可視化)

このチュートリアルを通して、オートエンコーダの理論を理解し、実際にコードを書いて動かすことで、その動作原理を深く理解することを目指します。

## 2. オートエンコーダの仕組み (How Autoencoders Work)

### 基本的な構造: エンコーダ、ボトルネック、デコーダ (Basic Architecture: Encoder, Bottleneck, Decoder)

オートエンコーダは、大きく分けて3つの主要なコンポーネントから構成されます<sup>2</sup>。

- **エンコーダ (Encoder):** 入力データをより低い次元の表現に圧縮する役割を担います<sup>2</sup>。高次元の入力情報を受け取り、それを凝縮します。
- **ボトルネック (Bottleneck) / 潜在空間 (Latent Space):** エンコーダによって圧縮された情報が保持される層です<sup>2</sup>。この層の次元数(潜在次元、latent\_dim)は、通常、入力データの次元数よりも大幅に小さく設定されます。この次元数の制約が、ネットワークに入力データの最も重要な特徴を学習させる駆動力となります<sup>2</sup>。この圧縮された表現は、「コード」や「潜在ベクトル」とも呼ばれます。
- **デコーダ (Decoder):** ボトルネック層の圧縮された表現を受け取り、それを元の高次元データ形式に復元(再構築)する役割を担います<sup>2</sup>。

この一連の流れを図で示すと、以下のようになります。

コード スニペット

graph LR

A[Input Data (高次元)] --> B(Encoder);

B --> C[Bottleneck / Latent Space (低次元)];

C --> D(Decoder);  
D --> E[Output Data (再構築された高次元データ)];

(図の概念は<sup>1</sup>に基づく)

## 学習プロセス: 入力の再構築と損失関数 (Learning Process: Input Reconstruction and Loss Function)

オートエンコーダの学習における主要な目標は、元の入力データ  $x$  と、デコーダによって再構築された出力データ  $x'$  との間の差(誤差)を最小限にすることです<sup>1</sup>。ネットワークは、ボトルネックによる制約の中で、入力と出力ができるだけ近くなるように、恒等関数(入力をそのまま出力する関数)を近似することを学習します<sup>3</sup>。

この入力と出力の差を定量的に測るのが損失関数 (Loss Function) です。画像データの再構築では、平均二乗誤差 (Mean Squared Error, MSE) がよく用いられます<sup>5</sup>。MSEは、元の画像と再構築された画像の対応するピクセル値の差を二乗し、その平均を計算します。ピクセル値を(例えばシグモイド活性化関数を用いて)0から1の間の確率として扱う場合は、バイナリクロスエントロピー (Binary Cross-Entropy, BCE) も用いられることがあります<sup>8</sup>。損失関数の選択は、データの正規化方法や出力層の活性化関数に依存することがあります。

学習は、この損失関数によって計算された誤差を最小化するように、エンコーダとデコーダのネットワークの重みを同時に調整していくプロセスです。この調整には、勾配降下法 (Gradient Descent) やその派生アルゴリズム(例えば Adam) が用いられます<sup>1</sup>。

## 主な用途 (Main Applications - Revisited with more detail)

オートエンコーダの基本的な仕組みから、以下のような応用が生まれます。

- **次元削減 (Dimensionality Reduction):** エンコーダは、データをより低い次元の潜在空間に写像します。これは、データの可視化や、他の機械学習モデルへの入力として利用する際に役立ちます。ニューラルネットワークを用いるため、PCA(主成分分析)のような線形手法では捉えきれない、データの非線形な構造を保持した次元削減が可能です<sup>10</sup>。潜在空間のデータをt-SNEなどの手法で2次元にマッピングして可視化することも行われます<sup>8</sup>。
- **特徴抽出 (Feature Extraction):** 学習済みのエンコーダは、新しいデータから圧縮された、意味のある特徴量を抽出するためのツールとして独立して利用できます。これらの特徴量は、例えば画像分類器などの教師あり学習モデルの入力として使うことができます<sup>5</sup>。
- **データノイズ除去 (Data Denoising):** 意図的にノイズを加えたデータをエンコーダに入力し、デコーダには元のノイズがないデータを再構築するように学習させます<sup>4</sup>。これにより、ネットワークはデータの本質的な構造(信号)とノイズを分離することを学習し、ノイズ除去器として機能します。
- **異常検知 (Anomaly Detection):** 正常なデータのみでオートエンコーダを学習させます。学習データに含まれないパターンを持つ異常データが入力されると、エンコーダはそれをうまく圧

縮できず、デコーダもそれを正確に再構築できません。結果として再構築誤差（損失）が大きくなるため、閾値を設けることで異常データを検出できます<sup>4</sup>。

## ボトルネック制約の重要性

オートエンコーダの能力の鍵を握るのは、ボトルネック層における情報の圧縮です。もしボトルネック層の次元数が入力層と同じかそれ以上であれば、ネットワークは単に入力データをそのままコピーすることを学習してしまい、データの本質的な特徴を捉えるインセンティブが働きません<sup>2</sup>。ボトルネック層の次元数を入力よりも小さくするという制約を課すことで、エンコーダは再構築に必要な最も重要な情報だけを選び出して効率的に表現する方法を学習せざるを得なくなります。この強制的な圧縮こそが、次元削減や特徴抽出といったオートエンコーダの有用な機能を生み出す源泉なのです。これは、潜在空間の次元数 (latent\_dim) が、圧縮の度合いと再構築の忠実度との間のトレードオフを決定する重要なハイパーパラメータであることを示唆しています。

## 3. 開発環境とツール (Development Environment and Tools)

### Python, TensorFlow, Keras の役割 (Role of Python, TensorFlow, Keras)

オートエンコーダの実装には、以下のツールが一般的に使用されます。

- **Python:** 実装に使用する主要なプログラミング言語です。豊富なライブラリとコミュニティサポートがあり、機械学習開発で広く採用されています。
- **TensorFlow:** Googleによって開発された、強力な数値計算ライブラリであり、深層学習モデルの構築と実行のための基盤となります<sup>21</sup>。自動微分(勾配計算のため)や、GPU/TPUといったハードウェアアクセラレーションによる高速化、モデルのデプロイメント機能などを提供します。また、TensorBoardという可視化ツールも含まれています<sup>21</sup>。
- **Keras:** TensorFlow(または他のバックエンド)上で動作する高レベルAPIです<sup>1</sup>。ニューラルネットワークの層(Layers)、活性化関数、オプティマイザ、損失関数などをモジュール化された使いやすいコンポーネントとして提供し、特に初心者にとって、モデルの構築とトレーニングをより簡単かつ迅速に行えるように設計されています<sup>21</sup>。

### TensorFlow/Keras を使う理由 (Why use TensorFlow/Keras?)

オートエンコーダの実装にTensorFlowとKerasの組み合わせを選択する理由はいくつかあります。

- **使いやすさ (Ease of Use):** Kerasは、Sequential API、Functional API、Model Subclassing APIといった直感的なインターフェースを提供し、複雑なネットワーク構造も比較的簡単に定義できます<sup>21</sup>。これにより、プロトタイピングや実験を迅速に行うことができます<sup>21</sup>。

- **柔軟性 (Flexibility):** シンプルでありながら、Keras(特にFunctional APIやModel Subclassing)は、層の共有、複数の入力や出力を持つモデルなど、複雑でカスタムなモデルの構築にも対応できます<sup>21</sup>。
- **エコシステム (Ecosystem):** TensorFlowは、モデルの監視ツールであるTensorBoard<sup>21</sup>、データセット読み込みのためのTensorFlow Datasets、事前学習済みモデルを利用できるTensorFlow Hub、モバイルデバイスへのデプロイを可能にするTensorFlow Liteなど、開発を支援する包括的なエコシステムを提供しています。
- **コミュニティとドキュメント (Community & Documentation):** 豊富な公式ドキュメント、チュートリアル、そして活発なコミュニティが存在し、学習や問題解決のためのサポートが充実しています<sup>4</sup>。

## KerasによるTensorFlowの簡略化

TensorFlow自体は、低レベルではグラフ計算やセッション管理など、初心者にはやや複雑な概念を含んでいます。Kerasは、これらの複雑さを覆い隠す抽象化レイヤーとして機能します<sup>21</sup>。Kerasが提供するLayer、Optimizer、Lossといった既製の構成要素を使うことで、開発者はTensorFlowの内部的な詳細(例えば勾配計算の具体的なプロセスなど)を深く理解していなくても、機能するオートエンコーダを構築できます。Kerasのユーザーエクスペリエンスとモジュール性への注力<sup>21</sup>が、開発サイクルの短縮<sup>22</sup>と学習の容易さを実現しており、本記事のような教育目的や、初心者が深層学習を始める際の理想的な選択肢となっています。

## 必要なライブラリのインストール (Installing necessary libraries)

実装を始める前に、必要なライブラリをインストールします。ターミナルまたはコマンドプロンプトで以下のコマンドを実行してください。

Bash

```
pip install tensorflow matplotlib numpy
```

## 4. 実践: 基本的なオートエンコーダの実装 (Hands-on: Implementing a Basic Autoencoder)

ここでは、最も基本的な形式のオートエンコーダを、全結合層(Denseレイヤー)のみを使用して実装します。データセットには、手書き数字のMNISTデータセットと似た構造を持つ、衣料品の画像データセットであるFashion MNISTを使用します<sup>18</sup>。

## データセットの準備: Fashion MNIST の読み込みと前処理 (Dataset Preparation: Loading and Preprocessing Fashion MNIST)

まず、必要なライブラリをインポートし、Fashion MNISTデータセットを読み込みます。

Python

```
import matplotlib.pyplot as plt
import numpy as np
import tensorflow as tf
from tensorflow.keras import layers, losses
from tensorflow.keras.datasets import fashion_mnist # Fashion MNIST を使用
from tensorflow.keras.models import Model

# データの読み込み (ラベルは使用しないため _ で無視)
(x_train, _), (x_test, _) = fashion_mnist.load_data()

# 正規化 (ピクセル値を 0-1 の範囲に)
x_train = x_train.astype('float32') / 255.
x_test = x_test.astype('float32') / 255.

# データの形状を確認
print(x_train.shape) # 出力例: (60000, 28, 28)
print(x_test.shape) # 出力例: (10000, 28, 28)
```

解説:

- `fashion_mnist.load_data()`: Kerasに含まれる関数で、データを訓練用(`x_train`)とテスト用(`x_test`)に分けて読み込みます。オートエンコーダの学習では通常ラベルは不要なため、`_`で受け取って無視しています<sup>18</sup>。手書き数字のMNISTデータセットを使用する場合は、`from tensorflow.keras.datasets import mnist`として`mnist.load_data()`を呼び出します<sup>12</sup>。
- `astype('float32') / 255.`: 画像のピクセル値(通常0-255の整数)をfloat32型に変換し、255で割ることで0から1の範囲に正規化します<sup>3</sup>。これは、ニューラルネットワークの学習を安定させ、効率的に進めるために重要な前処理ステップです。
- `print(x_train.shape)`: データの形状を確認します。Fashion MNISTは28x28ピクセルのグレースケール画像であり、訓練データが60,000枚、テストデータが10,000枚あることがわかります<sup>18</sup>。

## エンコーダ・デコーダの構築 (Building the Encoder and Decoder using Dense layers)

KerasのModel Subclassing APIを使用して、オートエンコーダモデルを定義します。このAPIは、モデルの構造をより柔軟に定義できる方法です<sup>18</sup>。

Python

```
# 潜在空間の次元数
latent_dim = 64
# 入力画像の形状 (高さ, 幅)
shape = x_train.shape[1:] # (28, 28)

class Autoencoder(Model):
    def __init__(self, latent_dim, shape):
        super(Autoencoder, self).__init__()
        self.latent_dim = latent_dim
        self.shape = shape
        # エンコーダの定義
        self.encoder = tf.keras.Sequential([
            layers.Flatten(), # 28x28 の画像を 784次元のベクトルに平坦化
            layers.Dense(latent_dim, activation='relu'), # 784次元 -> latent_dim (64)次元に圧縮
        ], name="encoder")
        # デコーダの定義
        self.decoder = tf.keras.Sequential(, name="decoder")

    # 順伝播の定義
    def call(self, x):
        encoded = self.encoder(x)
        decoded = self.decoder(encoded)
        return decoded

# オートエンコーダモデルのインスタンスを作成
autoencoder = Autoencoder(latent_dim, shape)
```

解説:

- latent\_dim = 64: ボトルネック層の次元数を64に設定します<sup>18</sup>。この値はハイパーパラメータであり、変更することで圧縮率と再構築性能のトレードオフを調整できます。
- shape = x\_train.shape[1:]: 入力画像の形状(ここでは(28, 28))を取得します。
- class Autoencoder(Model):: KerasのModelクラスを継承して、独自のモデルクラスを定義します<sup>18</sup>。

- `__init__(...)`: モデルのコンストラクタです。エンコーダとデコーダを内部で定義します。
  - `self.encoder: tf.keras.Sequential`を使ってエンコーダを定義します。
    - `layers.Flatten()`: 2次元の入力画像(28x28)を1次元のベクトル(784要素)に変換します。全結合層(`Dense`)は1次元の入力を想定しているため、この変換が必要です<sup>4</sup>。
    - `layers.Dense(latent_dim, activation='relu')`: 全結合層です。入力(784次元)を`latent_dim`(64次元)に圧縮します。`activation='relu'`は、活性化関数としてReLU (Rectified Linear Unit)を使用することを指定します。ReLUは、深層学習で広く使われる活性化関数の一つです<sup>4</sup>。
  - `self.decoder: tf.keras.Sequential`を使ってデコーダを定義します。
    - `layers.Dense(tf.math.reduce_prod(shape), activation='sigmoid')`: 全結合層です。`latent_dim`(64次元)の潜在表現を入力とし、元の画像のピクセル数( $28 * 28 = 784$ )と同じ次元数に復元します。`tf.math.reduce_prod(shape)`は`shape((28, 28))`の要素の積(784)を計算します<sup>18</sup>。`activation='sigmoid'`は、活性化関数としてシグモイド関数を使用します。シグモイド関数は出力を0から1の範囲に収めるため、正規化されたピクセル値の再構築に適しています<sup>18</sup>。
    - `layers.Reshape(shape)`: 784次元のベクトルを、元の画像の形状(28x28)に戻します<sup>18</sup>。
- `call(self, x)`: モデルが呼び出されたときの処理(順伝播)を定義します。入力`x`をエンコーダに通し、得られた潜在表現をデコーダに通して、再構築された結果を返します<sup>18</sup>。
- `autoencoder = Autoencoder(latent_dim, shape)`: 定義したAutoencoderクラスのインスタンスを作成します<sup>18</sup>。

## モデルのコンパイル (Compiling the Model - Optimizer, Loss)

次に、作成したモデルをトレーニング可能な状態にするためにコンパイルします。コンパイル時には、オプティマイザ(最適化アルゴリズム)と損失関数を指定します。

Python

```
autoencoder.compile(optimizer='adam', loss=losses.MeanSquaredError())
```

解説:

- `autoencoder.compile(...)`: モデルの学習プロセスを設定します<sup>18</sup>。
- `optimizer='adam'`: オプティマイザとしてAdamを選択します。Adamは、効率的な勾配降下法アルゴリズムであり、多くの深層学習タスクで良好な性能を示すことが知られています<sup>8</sup>。
- `loss=losses.MeanSquaredError()`: 損失関数として平均二乗誤差(MSE)を指定します。前述の通り、MSEは元のピクセル値と再構築されたピクセル値の差の二乗平均を計算し、この値を最小化するように学習が進められます<sup>18</sup>。シグモイド出力層と組み合わせる場合、



losses.BinaryCrossentropy()も選択肢となります<sup>8</sup>。

## モデルのトレーニング (Training the Model - fit, epochs, batch size)

モデルのコンパイルが完了したら、fitメソッドを使ってトレーニングを実行します。

Python

# モデルのトレーニング

```
history = autoencoder.fit(x_train, x_train, # 入力とターゲットは同じ訓練データ
                          epochs=10,      # トレーニングのエポック数
                          shuffle=True,    # 各エポックでデータをシャッフル
                          validation_data=(x_test, x_test)) # 検証データ
```

解説:

- autoencoder.fit(...): モデルのトレーニングを開始します<sup>18</sup>。
- x\_train, x\_train: 最初の引数は入力データ、2番目の引数はターゲット(教師)データです。オートエンコーダでは、入力データそのものを再構築することが目標なので、両方に同じx\_trainを指定します<sup>18</sup>。
- epochs=10: トレーニングデータセット全体を10回繰り返して学習します。エポック数は、モデルが十分に学習するまで調整する必要があります<sup>18</sup>。
- shuffle=True: 各エポックの開始時に訓練データをランダムに並び替えます。これにより、モデルがデータの順序に依存するのを防ぎ、汎化性能を高める効果が期待できます<sup>18</sup>。
- validation\_data=(x\_test, x\_test): 各エポックの終了時に、モデルの性能を評価するために使用する検証データを指定します。ここではテストデータを使用し、訓練データと同様に入力とターゲットにx\_testを指定します。検証データでの損失を見ることで、モデルが未知のデータに対してどの程度汎化できているか(過学習していないか)を監視できます<sup>18</sup>。
- history: fitメソッドは、トレーニング中の損失や(指定されていれば)評価指標の履歴を含むHistoryオブジェクトを返します。これを使って後で学習曲線をプロットできます。

## 結果の評価: 画像の再構築と可視化 (Evaluating Results: Image Reconstruction and Visualization)

トレーニングが完了したら、学習済みモデルを使ってテスト画像を再構築し、その結果を元の画像と比較して視覚的に評価します。

Python

```
# テストデータからエンコードされた表現を取得
encoded_imgs = autoencoder.encoder(x_test).numpy()
```

```
# エンコードされた表現から画像をデコード(再構築)
decoded_imgs = autoencoder.decoder(encoded_imgs).numpy()
```

```
# 結果の表示
n = 10 # 表示する画像の数
plt.figure(figsize=(20, 4))
for i in range(n):
    # オリジナル画像の表示
    ax = plt.subplot(2, n, i + 1)
    plt.imshow(x_test[i])
    plt.title("original")
    plt.gray()
    ax.get_xaxis().set_visible(False)
    ax.get_yaxis().set_visible(False)

    # 再構築された画像の表示
    ax = plt.subplot(2, n, i + 1 + n)
    plt.imshow(decoded_imgs[i])
    plt.title("reconstructed")
    plt.gray()
    ax.get_xaxis().set_visible(False)
    ax.get_yaxis().set_visible(False)
plt.show()
```

解説:

- `autoencoder.encoder(x_test).numpy()`: 学習済みのエンコーダ部分を使って、テストデータ `x_test` を潜在空間にエンコードします。`.numpy()` は TensorFlow のテンソル形式から NumPy 配列に変換します<sup>18</sup>。
- `autoencoder.decoder(encoded_imgs).numpy()`: エンコードされた表現 `encoded_imgs` を、学習済みのデコーダ部分を使って元の画像空間にデコード(再構築)します<sup>18</sup>。
- `plt.figure(...)`, `plt.subplot(...)`, `plt.imshow(...)`: Matplotlib ライブラリを使って、最初の `n` 枚のオリジナルテスト画像と、それに対応する再構築画像を並べて表示します<sup>18</sup>。
- `plt.title(...)`, `plt.gray()`, `ax.get_xaxis().set_visible(False)`: グラフのタイトル設定、グレースケール表示、軸の非表示など、表示を整えるためのコードです<sup>18</sup>。
- `plt.show()`: グラフを表示します。

この可視化により、オートエンコーダがどの程度元の画像を復元できているか、また、どのような情報が失われているか(例: 画像が少しぼやけるなど)を直感的に理解することができます。

## 基本的なオートエンコーダのハイパーパラメータ

この実装で使用した主要なハイパーパラメータをまとめます。

ハイパーパラメータ (Hyperparameter)	値 (Value)	説明 (Explanation)
latent_dim	64	潜在空間(ボトルネック)の次元数
epochs	10	トレーニングのエポック数
batch_size	(Keras Default)	1回の重み更新で使用するサンプル数 (fitで指定しない場合、デフォルトは32)
optimizer	'adam'	最適化アルゴリズム
loss	'mse'	損失関数(平均二乗誤差)
Input Activation	(None/Linear)	(入力層には通常なし)
Encoder Activation	'relu'	エンコーダ隠れ層の活性化関数
Decoder Activation	'sigmoid'	デコーダ出力層の活性化関数 (ピクセル値を0-1に)

(表の内容は <sup>18</sup> に基づく)

これらのハイパーパラメータは、モデルの性能に影響を与えるため、試行錯誤を通じて最適な値を見つけることが重要です。

## 5. 応用: 畳み込みオートエンコーダ (Application: Convolutional Autoencoder)

### 画像データと畳み込みニューラルネットワーク (Image Data and CNNs)

基本的なオートエンコーダでは、入力画像を1次元ベクトルに平坦化してから処理しました。しかし、この方法では、ピクセル間の空間的な位置関係(例えば、隣接するピクセルは関連性が高い、など)の情報が失われてしまいます。画像データのように、空間的な構造が重要な意味を持つデータに対しては、畳み込みニューラルネットワーク (**Convolutional Neural Network, CNN**) を利用するのが一般的です <sup>6</sup>。

CNNは、畳み込み層 (**Convolutional Layer**) とプーリング層 (**Pooling Layer**) という特殊な層を用いることで、画像の空間的な階層構造を効率的に学習できます。

- 畳み込み層: 「フィルタ(カーネル)」と呼ばれる小さなウィンドウを画像上でスライドさせながら、局所的な特徴(エッジ、コーナー、テクスチャなど)を抽出します。フィルタの重みは画像全体で共有されるため、全結合層に比べてパラメータ数が少なく、効率的に学習できます <sup>17</sup>。
- プーリング層: 畳み込み層で抽出された特徴マップの空間的な次元(解像度)を削減します。これにより、計算負荷を軽減し、位置のずれに対して頑健な特徴表現を獲得する効果があります。

す<sup>8</sup>。

これらの特性により、CNNは画像認識タスクで非常に高い性能を発揮します。オートエンコーダにCNNのアーキテクチャを組み込むことで、画像の再構築品質を向上させることが期待できます。これを畳み込みオートエンコーダ (**Convolutional Autoencoder, ConvAE**) と呼びます。

## 畳み込みエンコーダの構築 (Building the Convolutional Encoder)

畳み込みエンコーダは、入力画像を低次元の潜在表現(特徴マップ)に圧縮する役割を持ちます。通常、複数の畳み込み層とプーリング層(またはストライド付き畳み込み層)を交互に重ねて構成されます。

- **入力層 (InputLayer):** 入力画像の形状(高さ、幅、チャンネル数)を指定します。MNISTやFashion MNISTのようなグレースケール画像の場合、形状は (28, 28, 1) となります<sup>8</sup>。
- **畳み込み層 (Conv2D):** 画像から特徴を抽出します。主なパラメータは以下の通りです<sup>6</sup>。
  - filters: 適用するフィルタ(特徴マップ)の数。
  - kernel\_size: フィルタのサイズ(例: (3, 3))。
  - strides: フィルタをスライドさせるステップ数(例: (1, 1))。ストライドを2以上にすると、プーリング層と同様にダウンサンプリング効果が得られます。
  - padding: 'valid'(パディングなし)または 'same'(出力サイズが入力サイズと同じになるようにパディング)。
  - activation: 活性化関数(例: 'relu')。
- **プーリング層 (MaxPooling2D):** 特徴マップのサイズを縮小します。pool\_size(例: (2, 2))で縮小する領域のサイズを指定します<sup>8</sup>。

一般的な構成では、層が深くなるにつれてフィルタ数を増やし(より複雑な特徴を捉えるため)、空間的な次元を減らしていきます<sup>6</sup>。エンコーダの最終的な出力が、圧縮された潜在表現となります(これは通常、多次元の特徴マップの形をしています)。

## 畳み込みデコーダの構築 (Building the Convolutional Decoder)

畳み込みデコーダは、エンコーダによって生成された低次元の潜在表現から、元の高次元画像を再構築します。これは、エンコーダの処理を逆に行うような構造になります。

- **アップサンプリング層:** 空間的な次元を拡大するために、転置畳み込み層 (**Conv2DTranspose**) または アップサンプリング層 (**UpSampling2D**) と畳み込み層 (Conv2D) の組み合わせが用いられます<sup>6</sup>。
  - Conv2DTranspose: 畳み込みの逆演算のようなもので、学習可能なパラメータを持ちながらアップサンプリングを行います。
  - UpSampling2D: 単純に各ピクセルを繰り返すなどして次元を拡大し、その後にConv2Dで特徴を学習させます。
- **畳み込み層 (Conv2D):** アップサンプリングされた特徴マップから、より詳細な情報を再構築

するために使用されます。

- 最終層: 最後の畳み込み層(または転置畳み込み層)は、出力が元の画像の形状(例: (28, 28, 1))になるように設定し、活性化関数には入力データの正規化範囲に合わせたもの(例: 0-1の範囲なら'sigmoid')を使用します<sup>17</sup>。

デコーダの構造は、エンコーダの構造を逆にしたような形になることが多く、層が浅くなる(出力に近づく)につれて空間的な次元を増やし、フィルタ数を減らしていきます<sup>6</sup>。

## 実装例と結果の比較 (Implementation Example and Comparison with Basic AE)

以下に、Kerasを用いた畳み込みオートエンコーダの実装例を示します。基本的なオートエンコーダと同様に、Fashion MNISTデータセットを使用します。

Python

```
import matplotlib.pyplot as plt
import numpy as np
import tensorflow as tf
from tensorflow.keras import layers, losses, Input
from tensorflow.keras.datasets import fashion_mnist
from tensorflow.keras.models import Model

# データの読み込みと前処理 (ConvAE用に形状を変更)
(x_train, _), (x_test, _) = fashion_mnist.load_data()
x_train = x_train.astype('float32') / 255.
x_test = x_test.astype('float32') / 255.

# チャンネル次元を追加 -> (サンプル数, 28, 28, 1)
x_train = np.expand_dims(x_train, axis=-1)
x_test = np.expand_dims(x_test, axis=-1)

print(x_train.shape) # 出力例: (60000, 28, 28, 1)
print(x_test.shape) # 出力例: (10000, 28, 28, 1)

# --- 畳み込みオートエンコーダの定義 (Functional APIを使用) ---
input_img = Input(shape=(28, 28, 1))

# Encoder
x = layers.Conv2D(32, (3, 3), activation='relu', padding='same', strides=2)(input_img) # 28x28
```

```

-> 14x14
# x = layers.MaxPooling2D((2, 2), padding='same')(x) # MaxPoolingを使う場合
x = layers.Conv2D(64, (3, 3), activation='relu', padding='same', strides=2)(x) # 14x14 -> 7x7
# x = layers.MaxPooling2D((2, 2), padding='same')(x) # MaxPoolingを使う場合
encoded = x # 潜在表現 (7, 7, 64)

# Decoder
x = layers.Conv2DTranspose(64, (3, 3), activation='relu', padding='same', strides=2)(encoded)
# 7x7 -> 14x14
x = layers.Conv2DTranspose(32, (3, 3), activation='relu', padding='same', strides=2)(x) # 14x14
-> 28x28
# x = layers.UpSampling2D((2, 2))(x) # UpSamplingを使う場合
# x = layers.Conv2D(..., activation='relu', padding='same')(x) # UpSamplingの後にConv2D
decoded = layers.Conv2D(1, (3, 3), activation='sigmoid', padding='same')(x) # 28x28x1 に復元

# Autoencoder Model
conv_autoencoder = Model(input_img, decoded)

# --- モデルのコンパイル ---
conv_autoencoder.compile(optimizer='adam', loss='binary_crossentropy') # または loss='mse'

# モデル構造の表示
conv_autoencoder.summary()

# --- モデルのトレーニング ---
history_conv = conv_autoencoder.fit(x_train, x_train,
                                   epochs=20, # エポック数を少し増やす
                                   batch_size=128,
                                   shuffle=True,
                                   validation_data=(x_test, x_test))

# --- 結果の評価: 画像の再構築と可視化 ---
decoded_imgs_conv = conv_autoencoder.predict(x_test)

# 結果の表示 (基本的なAEと同様のコードで表示)
n = 10
plt.figure(figsize=(20, 4))
for i in range(n):
    # オリジナル画像の表示
    ax = plt.subplot(2, n, i + 1)
    plt.imshow(x_test[i].reshape(28, 28)) # reshapeでチャンネル次元を削除
    plt.title("original")
    plt.gray()

```

```

ax.get_xaxis().set_visible(False)
ax.get_yaxis().set_visible(False)

# 再構築された画像の表示 (ConvAE)
ax = plt.subplot(2, n, i + 1 + n)
plt.imshow(decoded_imgs_conv[i].reshape(28, 28)) # reshapeでチャンネル次元を削除
plt.title("reconstructed (ConvAE)")
plt.gray()
ax.get_xaxis().set_visible(False)
ax.get_yaxis().set_visible(False)
plt.show()

```

(コード例は<sup>6</sup>の概念を参考に構成)

結果の比較:

基本的なオートエンコーダ(全結合層)の結果と比較すると、畳み込みオートエンコーダで再構築された画像の方が、一般的に元の画像の詳細をより良く保持し、より鮮明になる傾向があります<sup>6</sup>。これは、畳み込み層が画像の空間的な特徴を効果的に捉え、デコーダがそれを基に再構築するためです。全結合層を用いた基本的なオートエンコーダでは失われがちだった、画像の局所的なパターンや構造が、ConvAEではより保持されやすくなります。

## 基本的なAEと畳み込みAEの比較

特徴 (Feature)	基本的なオートエンコーダ (Dense)	畳み込みオートエンコーダ (ConvAE)
入力処理	画像を平坦化(空間情報喪失)	画像を2D/3Dグリッドとして処理(空間情報保持)
エンコーダ層	Flatten, Dense	Conv2D, MaxPooling2D (または Strided Conv2D)
デコーダ層	Dense, Reshape	Conv2DTranspose または (UpSampling2D + Conv2D)
パラメータ数	大きな画像では多くなりがち	重み共有により画像に対して効率的
主な用途	シンプルなデータ、ベースライン	画像データ、空間的特徴学習

参考コード例	<sup>18</sup>	<sup>6</sup>
--------	---------------	--------------

(表の内容は <sup>6</sup> に基づく)

## アーキテクチャとデータタイプの適合性

基本的なオートエンコーダから畳み込みオートエンコーダへの移行は、深層学習における重要な原則を示唆しています。それは、ネットワークアーキテクチャはデータの構造に合わせて選択すべきであるということです。全結合層 (Dense) は入力の特徴量を個別に扱い、それらの間の空間的な関係性を無視するため、画像データには最適ではありません。一方、畳み込み層は、フィルタを用いて局所的なパターンを検出し、プーリング層が空間的な階層構造を扱うように、明示的にグリッド状データ (画像など) を処理するために設計されています<sup>8</sup>。畳み込みオートエンコーダが画像タスクでしばしばより良い結果を示すのは、このアーキテクチャとデータタイプの適合性が理由です。これは、データ (シーケンス、画像、表形式データなど) の性質を理解し、それに適した層のタイプ (RNN/LSTM<sup>11</sup>、CNN、Dense など) を選択することが、効果的なモデルを構築するための鍵であることを示しています。

## 6. まとめと次のステップ (Conclusion and Next Steps)

### 学習内容の要約 (Summary of learned concepts and skills)

本記事では、オートエンコーダの基本概念から実践的な実装までを解説しました。主な学習内容は以下の通りです。

- オートエンコーダの定義と構造: 入力を圧縮 (エンコード) し、それを復元 (デコード) するニューラルネットワークであり、エンコーダ、ボトルネック (潜在空間)、デコーダから構成されること。
- 主な用途: 次元削減、特徴抽出、データノイズ除去、異常検知など、多様な応用可能性があること。
- TensorFlow/Keras の役割: Keras が高レベル API としてモデル構築を容易にし、TensorFlow がその背後で計算エンジンとして動作すること。
- 実装手順: Fashion MNIST データセットの準備、基本的な (全結合) オートエンコーダと畳み込みオートエンコーダの構築 (エンコーダ・デコーダ定義)、モデルのコンパイル (オプティマイザ、損失関数)、トレーニング (fit メソッド)、そして結果の評価 (再構築画像の可視化) という一連の流れ。
- 畳み込みオートエンコーダの利点: 画像データに対して、空間的特徴を保持できるため、基本的なオートエンコーダよりも高品質な再構築が期待できること。

### さらなる学習のために (Resources for further learning)

オートエンコーダは非常に奥が深い分野であり、本記事で紹介したのはその入り口に過ぎません。さ



らに理解を深めるためには、以下のような発展的なトピックを探求することをお勧めします。

- 様々なオートエンコーダの派生形:
  - ノイズ除去オートエンコーダ (**Denoising Autoencoders, DAE**): 入力にノイズを加えて学習させることで、ノイズ除去能力を獲得するモデル<sup>4</sup>。
  - スパースオートエンコーダ (**Sparse Autoencoders**): ボトルネック層の活性化に対してスパース性(少数のニューロンのみが同時に活性化する)を制約として加えるモデル。これにより、より解釈しやすい特徴表現が得られることがあります<sup>1</sup>。
  - 変分オートエンコーダ (**Variational Autoencoders, VAE**): 確率的なアプローチを取り入れ、潜在空間を連続的かつ構造化されたものにするモデル。単なる再構築だけでなく、新しいデータを生成する能力を持ちます<sup>5</sup>。
  - 系列対系列オートエンコーダ (**Sequence-to-Sequence Autoencoders**): テキストや時系列データのような系列データを扱うために、エンコーダとデコーダにRNN(Recurrent Neural Network)やLSTM(Long Short-Term Memory)を用いるモデル<sup>8</sup>。
- 実験:
  - 異なるアーキテクチャ(層の数、フィルタ数、ユニット数など)を試す。
  - ハイパーパラメータ(潜在空間の次元数、学習率、バッチサイズなど)を調整する。
  - 異なるデータセット(例: CIFAR-10など)で試す。
  - 異なる損失関数(例: Binary Cross-Entropy)を試す。
- 公式ドキュメント:
  - TensorFlowとKerasの公式ウェブサイトには、詳細なドキュメント、チュートリアル、APIリファレンスが豊富に用意されています<sup>4</sup>。

オートエンコーダは、深層学習の基本的な構成要素の一つであり、その概念を理解することは、より高度なモデルや技術を学ぶ上での重要なステップとなります。ぜひ、本記事を足がかりとして、さらなる探求を進めてみてください。

## 引用文献

1. Kerasで学ぶAutoencoder - Elix Tech Blog, 4月 21, 2025にアクセス、  
<https://elix-tech.github.io/ja/2016/07/17/autoencoder.html>
2. オートエンコーダ: 抽象的な特徴を自己学習するディープラーニングの人気者 - DeepAge, 4月 21, 2025にアクセス、  
[https://deepage.net/deep\\_learning/2016/10/09/deeplearning\\_autoencoder.html](https://deepage.net/deep_learning/2016/10/09/deeplearning_autoencoder.html)
3. Autoencoders | Machine Learning Tutorial, 4月 21, 2025にアクセス、  
[https://mallahyari.github.io/ml\\_tutorial/autoencoder/](https://mallahyari.github.io/ml_tutorial/autoencoder/)
4. Intro to Autoencoders | TensorFlow Core, 4月 21, 2025にアクセス、  
<https://www.tensorflow.org/tutorials/generative/autoencoder>
5. Autoencoder in TensorFlow 2: Beginner's Guide - LearnOpenCV, 4月 21, 2025にアクセス、  
<https://learnopencv.com/autoencoder-in-tensorflow-2-beginners-guide/>
6. Autoencoders with Keras, TensorFlow, and Deep Learning ..., 4月 21, 2025にアクセス、

- <https://pyimagesearch.com/2020/02/17/autoencoders-with-keras-tensorflow-and-deep-learning/>
7. Unleashing the Power of Autoencoders: Applications and Use Cases - Analytics Vidhya, 4月 21, 2025にアクセス、  
<https://www.analyticsvidhya.com/blog/2023/05/unleashing-the-power-of-autoencoders-applications-and-use-cases/>
  8. Building Autoencoders in Keras - The Keras Blog, 4月 21, 2025にアクセス、  
<https://blog.keras.io/building-autoencoders-in-keras.html>
  9. G検定 5-4 24. オートエンコーダ | 大原かほ - note, 4月 21, 2025にアクセス、  
[https://note.com/ohara\\_designer/n/n5d407eb29802](https://note.com/ohara_designer/n/n5d407eb29802)
  10. オートエンコーダ: 定義、タイプ、アプリケーション - Ultralytics, 4月 21, 2025にアクセス、  
<https://www.ultralytics.com/ja/glossary/autoencoder>
  11. A Practical guide to Autoencoders using Keras - reckoning.dev, 4月 21, 2025にアクセス、  
<https://reckoning.dev/blog/autoencoders/>
  12. Autoencoders in Machine Learning | GeeksforGeeks, 4月 21, 2025にアクセス、  
<https://www.geeksforgeeks.org/auto-encoders/>
  13. Understanding Auto-encoders - Kaggle, 4月 21, 2025にアクセス、  
<https://www.kaggle.com/code/rockystats/understanding-auto-encoders>
  14. Autoencoder Feature Extraction for Classification - MachineLearningMastery.com, 4月 21, 2025にアクセス、  
<https://machinelearningmastery.com/autoencoder-for-classification/>
  15. TensorFlow.keras上にオートエンコーダを実装してMNIST画像内に ..., 4月 21, 2025にアクセス、  
<https://qiita.com/xtrizeShino/items/331dc0db1843071a6cd8>
  16. Autoencoders In Deep Learning | Tensorflow Training | Edureka - YouTube, 4月 21, 2025にアクセス、  
<https://www.youtube.com/watch?v=8W5Ykvg5EX0>
  17. Understanding Autoencoders With Tensorflow:Denoising ..., 4月 21, 2025にアクセス、  
<https://learnopencv.com/understanding-autoencoders-using-tensorflow-python/>
  18. オートエンコーダの基礎 | TensorFlow Core, 4月 21, 2025にアクセス、  
<https://www.tensorflow.org/tutorials/generative/autoencoder?hl=ja>
  19. ゼロからはじめる変分オートエンコーダ: 理論からTensorFlow実装、可視化、応用まで完全理解, 4月 21, 2025にアクセス、  
<https://note.com/dalhi/n/n67737adeb629>
  20. Convolutional Variational Autoencoder | TensorFlow Core, 4月 21, 2025にアクセス、  
<https://www.tensorflow.org/tutorials/generative/cvae>
  21. Keras vs Tensorflow: A Comprehensive Guide - Kanerika, 4月 21, 2025にアクセス、  
<https://kanerika.com/blogs/keras-vs-tensorflow/>
  22. Pytorch Vs Tensorflow Vs Keras: The Differences You Should Know - Simplilearn.com, 4月 21, 2025にアクセス、  
<https://www.simplilearn.com/keras-vs-tensorflow-vs-pytorch-article>
  23. The Functional API | TensorFlow Core, 4月 21, 2025にアクセス、  
[https://www.tensorflow.org/guide/keras/functional\\_api](https://www.tensorflow.org/guide/keras/functional_api)
  24. Keras 2 : examples : コンピュータビジョン - 画像のノイズ除去のための畳込みオートエンコーダ, 4月 21, 2025にアクセス、  
<https://tensorflow.classcat.com/2021/11/02/keras-2-examples-vision-autoencoder/>

25. Building a Simple Autoencoder with TensorFlow | Easy Guide - YouTube, 4月 21, 2025にアクセス、<https://www.youtube.com/watch?v=hbH2PMIMOJs>
26. Keras 2 : examples : 生成深層学習 - 変分オートエンコーダ - ClassCat® AI Research, 4月 21, 2025にアクセス、<https://tensorflow.classcat.com/2022/06/28/keras-2-examples-generative-vae/>