

深層学習 DAY4 レポート

2024 年 06 月 21 日 新規作成

- 深層学習 DAY4 レポート
 - 強化学習
 - 考察
 - AlphaGo
 - 考察
 - 計量化・高速化技術
 - 考察
 - 応用技術
 - 考察
 - ResNet(転移学習)
 - 実装演習
 - 事前学習なし
 - 転移学習
 - ファインチューニング
 - 結果
 - wide ResNet
 - 転移学習
 - ファインチューニング
 - 結果
 - 考察
 - EfficientNet
 - Vision Transformer
 - 物体検知と SS 解説
 - Mask R-CNN
 - 補足
 - Faster R-CNN
 - YOLO(You Only Look Once)
 - FCOS(Fully Convolutional One-Stage Object Detection)
 - Transformer
 - BERT
 - 実装演習

- 考察
 - GPT
 - 音声認識
 - CTC
 - DCGAN
 - Conditional GAN
 - Pix2Pix
 - A3C
 - Metric-learning(距離学習)
 - MAML(メタ学習)
 - グラフ畳込み(GCN)
 - Grad-CAM, LIME, SHAP
 - 実装演習
 - Docker
-

強化学習

強化学習(Reinforcement Learning)はエージェントが環境と相互作用しながら、試行錯誤を通じて行動を学習する機械学習の一つ。エージェントに最適な行動方針（ポリシー）を学ばせるこの手法は深層学習と組み合わせて深層強化学習と呼ばれる。教師あり学習教師なし学習とは異なり、環境の中で目的として設定された報酬を最大化させるための行動を学習する。

基本概念

- エージェント
学習を行う主体。環境の状態を観察し行動を選択する。
- 環境
エージェントが相互作用する対象。エージェントの行動に応じて状態を変化させ報酬を与える。
- 状態
環境の状況を表す情報。エージェントはこの情報に基づき行動を決定する。
- 行動
エージェントが取ることができる操作や動作。
- 報酬
エージェントの行動に対して環境から与えられるフィードバック。エージェントの目標は累積報酬を最大化させること。
- ポリシー
状態に基づいてエージェントが行動を選択する戦略。

- 価値関数

ある状態における累積報酬の期待値。エージェントはこの関数を学習することで最適な行動を決定する。

主要なアルゴリズム

- Q 学習

モデルフリーの強化学習アルゴリズムで行動価値関数を学習する。

- SARSA

オンポリシーのモデルフリー強化学習アルゴリズム

- DQN

Q 学習をディープラーニングと組み合わせたアルゴリズムニューラルネットワークを使って行動価値関数を近似する。

- ポリシー勾配法

ポリシーを直接最適化する手法

強化学習の応用

- ゲーム AI（チェス、囲碁）

- ロボット制御（動作計画や制御）

- 自立走行車（経路計画や運転行動）

- 金融取引（最適な投資戦略の学習）

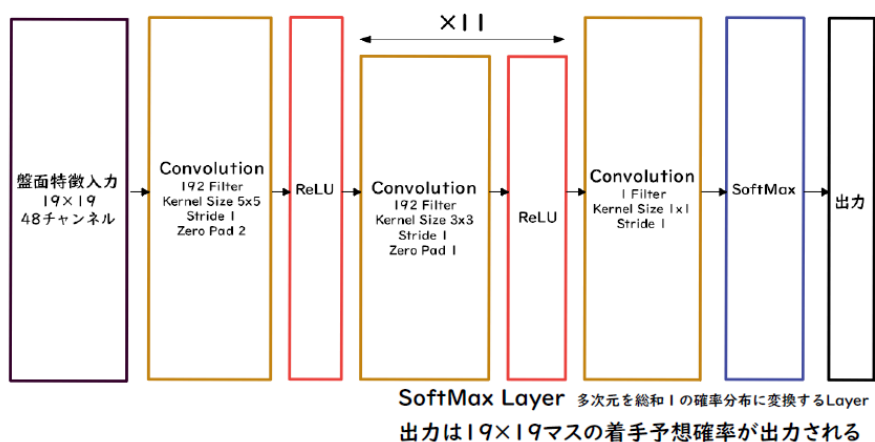
考察

深層学習や計算能力の進歩により、強化学習では単純な分類や予測タスクだけでなくポリシーに基づいて環境へ適応しながら報酬を最大化させるように学習していくことができる。テーブルゲームに代表されるような人間の知的行動を機械的に再現できるようになった。似たものに遺伝的アルゴリズムがある。以前はこれも強化学習の一種だと思っていたが、厳密には遺伝的アルゴリズムは進化計算に基づいた最適化手法であるため強化学習とは異なる。

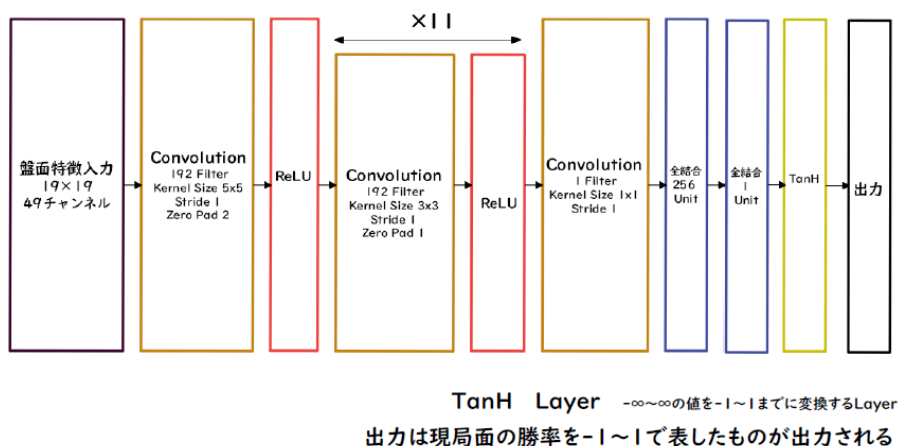
AlphaGo

Google DeepMind によって開発された囲碁 AI。2016 年にはトップ囲碁プレイヤーに勝利する。

Alpha Go (Lee)のPolicyNet



Alpha Go (Lee) のValueNet

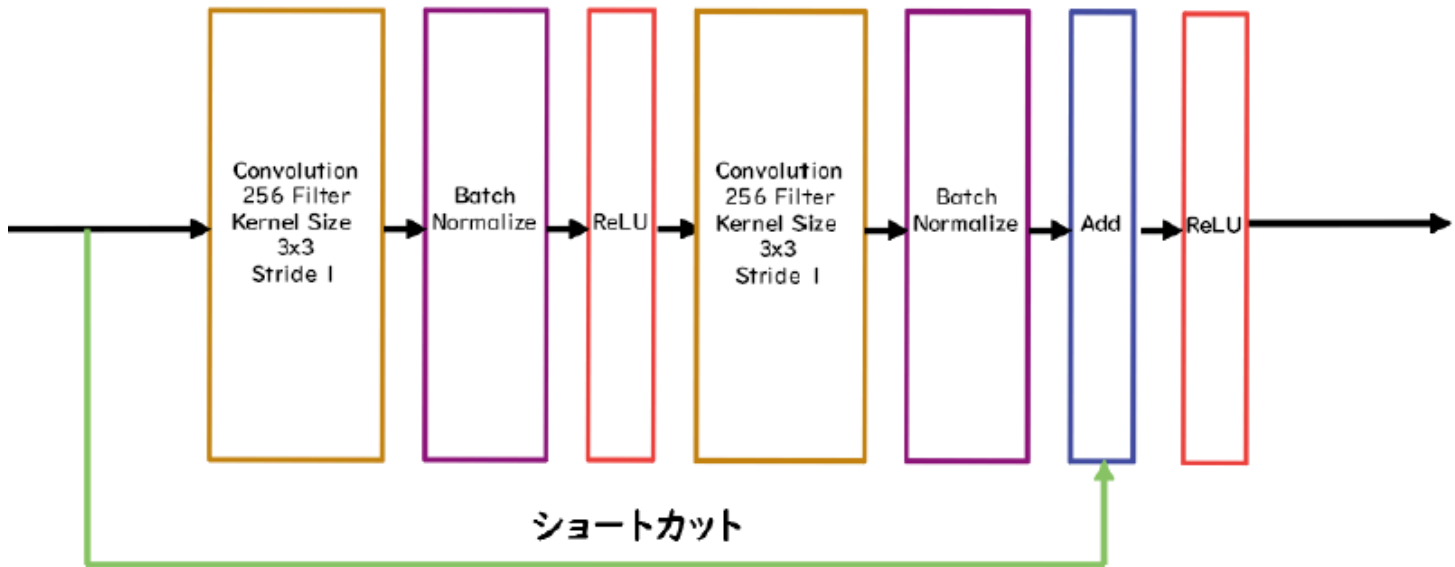


特徴

- 次の一手を予測する PolicyNet と現局面の勝敗確率を予測する ValueNet から構成される。
- モンテカルロ木探索という手法でランダムシミュレーションを行い、その時々勝敗確率を集計して最適手を探索する。
- 深い層を構築するために ResidualNetwork というショートカット機構を導入している。これにより 100 層を超えるネットワークでも勾配消失勾配爆発を抑えて安定した学習が可能となった。信号が通過するショートカットの組み合わせによりアンサンブル効果が期待できる。

ResidualNetwork

Residual Block 基本形



Residual Network の派生形

- Residual Block の工夫
 - Bottleneck
1 層目で次元削減を行って 3 層目で次元を復元する
 - PreActivation
BatchNorm->ReLU->Conv->BatchNorm->ReLU->Conv->Add の並びにしたもの
- Network 構造の工夫
 - WideResNet
Conv の filter 数を k 倍にした ResNet。段階的に幅を増やしていく。
 - PyramidNet
WideResNet を改良して各層で filter 数を増やしていく ResNet

考察

囲碁には詳しくないが将棋ソフトなら Bonanza や ELMo が有名。エルモ囲いなど将棋ソフトが発祥の戦術もある。単純に人類の競争相手というだけでなくプロ棋士が研究するためのツールとしても使われている。

計量化・高速化技術

深層学習ではパラメータの最適化のために多くの計算を必要とする。学習を高速化するための手法や計算リソースを以下に述べる。

分散深層学習

- データ並列化

親モデルを各ワーカーに子モデルとしてコピーする。データを分割し各ワーカーごとに計算させる。

- 同期型

各ワーカーの計算が終わるのを待ってから全ワーカーの勾配平均を計算し親モデルのパラメータを更新する。

- 非同期型

各ワーカーはお互いの計算終了を待たない。学習が終われば都度パラメータサーバーに結果をプッシュする。新たに学習を始めるときはパラメータサーバーからポップして学習する。

処理スピードは非同期型の方が速いが最新のパラメータを利用できないので学習が不安定になりやすい。精度は同期型の方が良いことが多いので現在はこちらが主流。

- モデル並列化

親モデルを各ワーカーで分割し学習させる。学習が終了すると一つのモデルに復元する。**モデルが大きいときにはモデル並列化、データが大きいときにはデータ並列化する。**

プロセッサの違い

- CPU(Central Processing Unit)

ほとんどの計算タスクを処理できる汎用プロセッサ。シングルスレッド性能が高いが並列処理性能が低く、行列演算の速度は遅い。

- GPU(Graphics Processing Unit)

多数のコアを持ち大規模な並列計算を高速に行える。CUDA や cuDNN など深層学習向けライブラリが豊富。高性能 GPU だと高価で運用コストが大きく消費電力も多い。

- TPU(Tensor Processing Unit)

深層学習の行列演算（テンソル演算）に特化しており非常に効率が良い。TensorFlow に最適化されており、汎用性が互換性に課題がある。

深層学習モデルそのものを軽量化する手法には以下がある。

- 量子化

パラメータのデータ型を float64 から float32 などの低い精度に落とすことでメモリと演算処理を削減する。極端な量子化でなければそれほど精度は低下しないことが知られている。

- 蒸留

精度の高い大規模モデルから軽量のモデルを作成する。これはモデルを簡約化することであり、学習済みモデルの知識を小規模なモデルに継承させることで計量でありながら高い精度を発揮するモデルの構築が期待できる。

- プルーニング

モデル精度への寄与度が低いニューロンを削減し、モデルの圧縮を行う。

考察

深層学習の推論をモバイル端末やマイコンで実施するための軽量化方法として TensorFlow Lite がある。TensorFlow のモデルを TensorFlow Lite モデルに変換することもでき、よりユーザーに近い環境で深層学習モデルを利用できるようになる。

応用技術

- MobileNet

モバイルデバイス向けの画像認識モデル。畳込み処理を工夫して計算量を削減した。通常の畳込み演算とは異なり、深さ方向の畳込みと点方向の畳込みを組み合わせることで効率化している。

- 通常の畳込み計算量: $H \times W \times K \times K \times C \times M$
- MobileNet の計算量: $H \times W \times K \times K \times C + H \times W \times C \times M$
 - Depthwise Conv: $H \times W \times K \times K \times C$
 - Pointwise Conv: $H \times W \times C \times M$

- DenseNet

DenseBlock と呼ばれるモジュールを用いて層を深くしても学習が進むようにした。DenseBlock では前の層の入力を足し合わせることをしている。前方の各層から出力すべてが後方の層への入力となる。

- WaveNet

音声波形を生成する深層学習モデル。時系列のデータに対して畳込み処理を適用している。受容野を増やすため層が深くなるにつれて畳込むリンクを離す Dilated Convolution という処理をしている。

考察

時系列データだからといって必ず RNN を使うわけではなく、タスクに応じて様々な工夫が必要になる。計算量を減らせば軽量化することは理論的に間違いはないが、いろいろな工夫を試してみても結果的に精度向上できたが理論的に説明できていないことも多いように感じた。

ResNet(転移学習)

ResNet(ResidualNetwork)は残差ブロックと呼ばれる機構を用いることで層を深くしても勾配消失問題が起きにくく高い性能を示した。残差ブロックの基本的なアイデアは各層に対してその入力を出力に直接加えるショートカットを設けたことである。

実装演習

事前学習なし

```
resnet = tf.keras.applications.resnet.ResNet50(weights=None)
resnet.trainable = True
x1 = resnet.layers[-2].output # avg_poolまでのoutputを取得します。
out = tf.keras.layers.Dense(class_num, activation='softmax')(x1) # avg_poolから出力層に繋がります。
model = tf.keras.models.Model(inputs=resnet.input, outputs=out)

model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001), loss='categorical_crossentropy', metrics=
model.summary()
```

転移学習

```
resnet = tf.keras.applications.resnet.ResNet50(weights='imagenet')
resnet.trainable = False
x1 = resnet.layers[-2].output # avg_poolまでのoutputを取得します。
out = tf.keras.layers.Dense(class_num, activation='softmax')(x1) # avg_poolから出力層に繋がります。
model = tf.keras.models.Model(inputs=resnet.input, outputs=out)

model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001), loss='categorical_crossentropy', metrics=
model.summary()
```

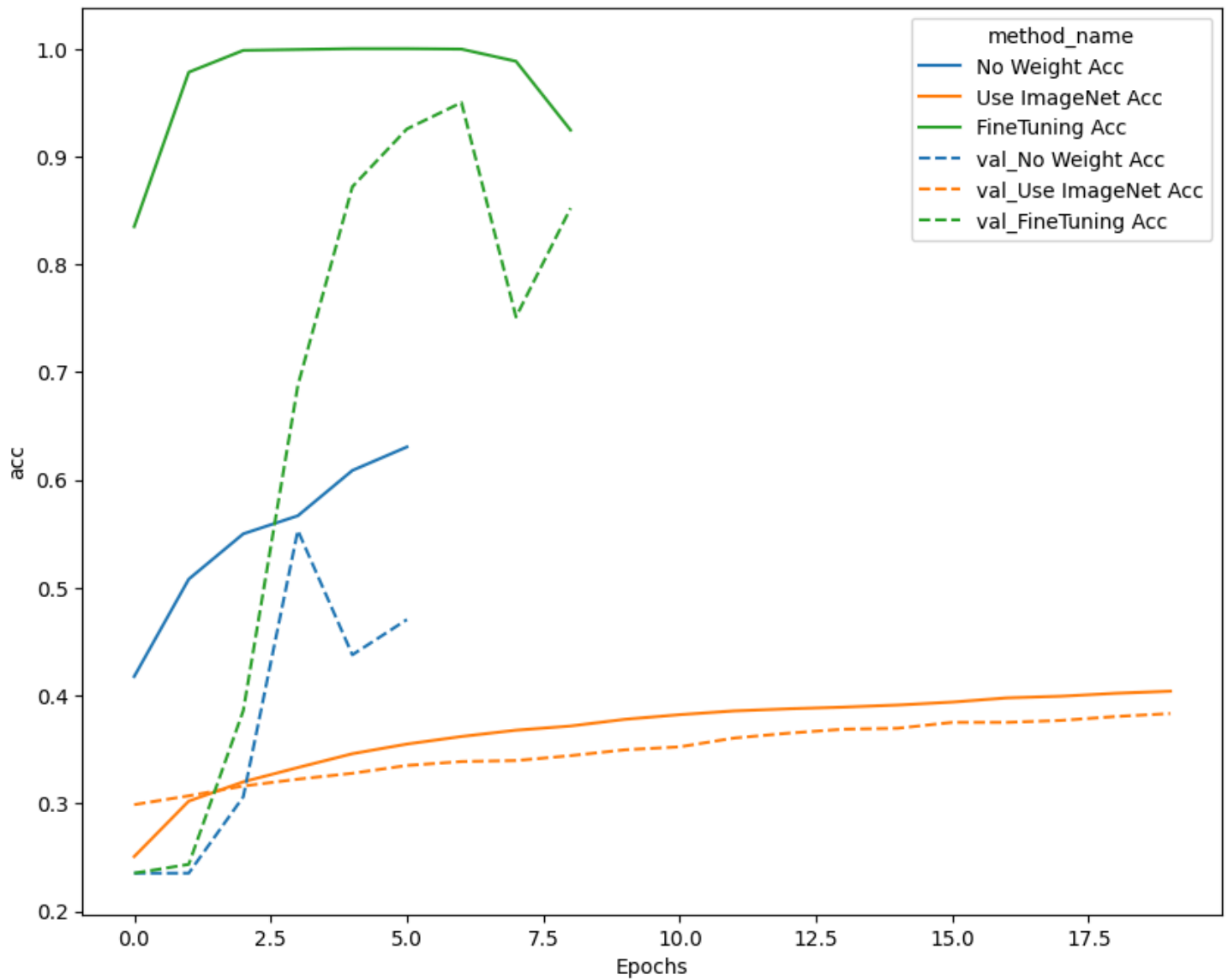
ファインチューニング

```
resnet = tf.keras.applications.resnet.ResNet50(weights='imagenet')
resnet.trainable = True
x1 = resnet.layers[-2].output
out = tf.keras.layers.Dense(class_num, activation='softmax')(x1)
model = tf.keras.models.Model(inputs=resnet.input, outputs=out)

model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001), loss='categorical_crossentropy', metrics=
model.summary()
```

結果

アルゴリズム	経過時間
事前学習なし	4min.
転移学習	4min.
ファインチューニング	4min.



wide ResNet

Residual Block 内の畳み込みに対してチャンネル数を増やしたモデル

転移学習

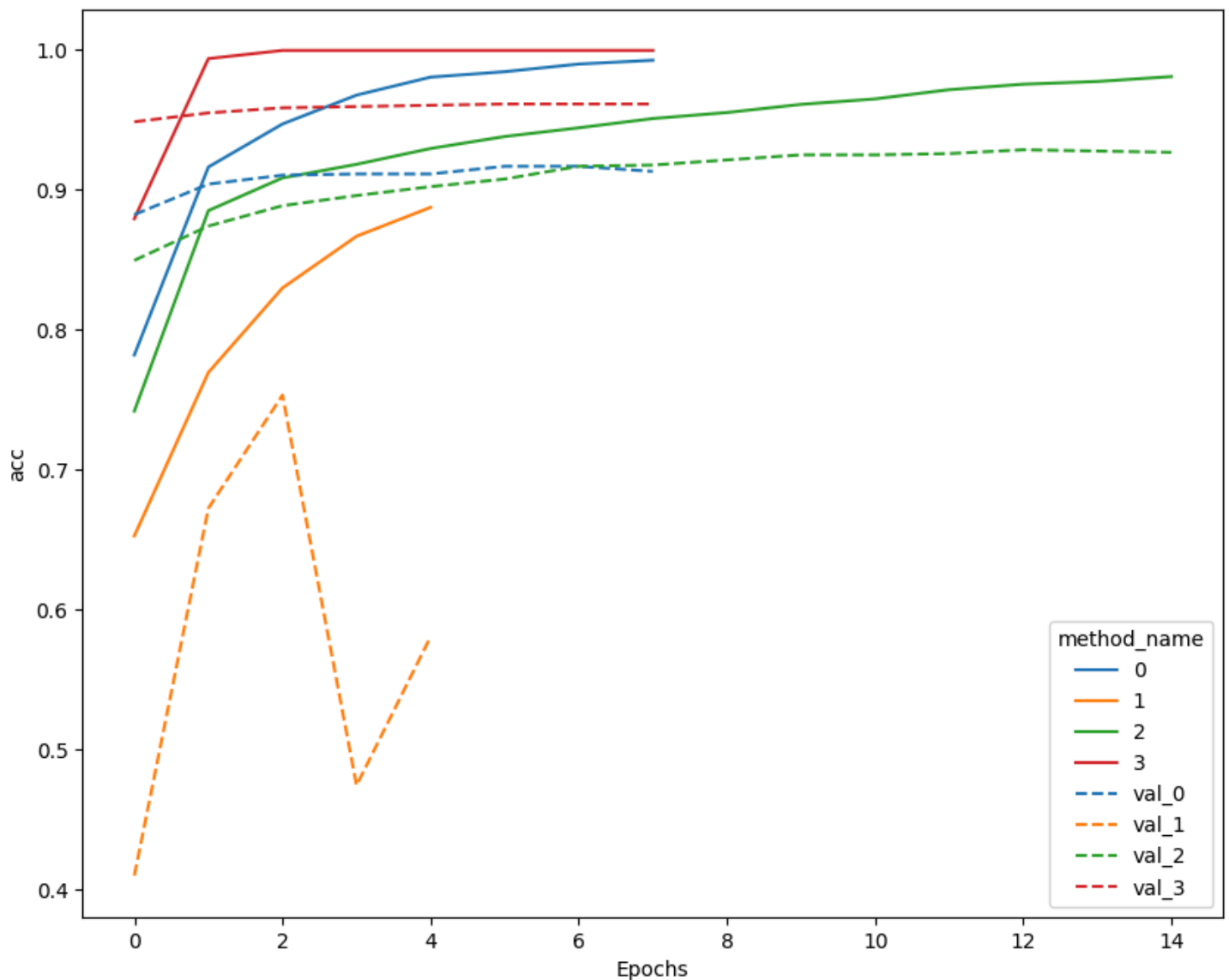
```
model = tf.keras.Sequential([
    tf.keras.layers.InputLayer(input_shape=(224, 224, 3)),
    hub.KerasLayer("https://tfhub.dev/google/bit/s-r50x3/1", trainable=False),
    tf.keras.layers.Dense(class_num, activation='softmax')
])
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001), loss='categorical_crossentropy', metrics=
model.summary())
```

ファインチューニング

```
model = tf.keras.Sequential([
    tf.keras.layers.InputLayer(input_shape=(224, 224, 3)),
    hub.KerasLayer("https://tfhub.dev/google/bit/s-r50x3/1", trainable=True),
    tf.keras.layers.Dense(class_num, activation='softmax')
])
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.00001), loss='categorical_crossentropy', metrics=
model.summary()
```

結果

No.	アルゴリズム	経過時間
0	ResNet 転移学習	-
1	ResNet ファインチューニング	-
2	Wide ResNet 転移学習	32min.
3	Wide ResNet ファインチューニング	40min.



考察

ファインチューニングを使ったアルゴリズムのものはスコアが良い傾向にある。転移学習よりファインチューニングの方が少ないエポック数で early-stopping が効いている。また、wide ResNet モデルが一番スコアが高かったものの訓練時間が大幅に増えた。

EfficientNet

EfficientNet はモデルのスケーリングに焦点を当てており計算コストを最小化しながら高い精度を達成することを目的としている。単一の次元でモデルをスケーリングするのではなく、幅、深さ、解像度を同時にバランスよくスケーリングする。この複合スケーリングにより、パフォーマンスの向上と効率性の両立が可能となった。

- 幅のスケーリング

ネットワークの各層のユニット数を増やしてより細かい特徴表現を獲得する。しかし、深さに対

してユニット数が大きすぎると高レベルな特徴表現を獲得しにくくなる。

- 深さのスケーリング

ネットワークの総数を増やしてより深いモデルを作成する。

- 解像度のスケーリング

入力画像の解像度を高くしてより詳細な特徴を捉える。

- **複合スケーリング**

上記はある程度まで増やすと精度向上は横ばいになった。そこで次のスケーリング式を用いて幅、深さ、解像度を同時にスケーリングして効率よく精度を最大化させる。

$$\begin{aligned}\text{width} &= \alpha^k \\ \text{depth} &= \beta^k \\ \text{resolution} &= \gamma^k\end{aligned}$$

α, β, γ はスケーリング係数で k はスケールステップ。モデルのサイズと計算量が増えるごとにパフォーマンスが一貫して向上するようになる。

Vision Transformer

Vision Transformer(ViT)は自然言語分野でよく使われている Transformer を画像認識に応用したモデル。画像をパッチと呼ばれる小さな領域に分割し、そのそれぞれをトークンとして扱い Transformer の入力とする。

処理手順

1. 画像のパッチ分割

入力画像を固定サイズの小さなパッチに分割する。

2. パッチの埋め込み

各パッチを線形変換し固定次元長のベクトルに変換する。

3. 位置埋め込み

各トークンに位置情報を付与する。

4. Transformer Encoder

自己注意機構を用いて各トークン間の関係を学習する。

5. 分類

分類トークン (CLS トークン) を追加し最終的な出力として分類タスクを行う。

メリット

- 大規模データセットで特に優れた性能を発揮する。
- 自己注意機構により画像の広範な依存関係を捉えられる。
- 自然言語モデルでの大規模な事前学習を画像認識に応用できる。

デメリット

- 大量のデータが必要となる。
- 小規模データセットでは CNN より性能が劣る。
- Transformer は訓練コストが高い。

物体検知と SS 解説

物体検出タスクの種類

- 物体検出
画像内の物体を検出しそれぞれの物体に対して境界ボックスを指定する。
- セマンティックセグメンテーション
画像内のすべてのピクセルをカテゴリに分類する。同じカテゴリの物体は個々に区別しない。
- パノプティックセグメンテーション
すべての物体を個別に認識しそれぞれのピクセルにカテゴリを割り当てる。

物体検知の種類

- 単段検出器(Single-Stage Detectors)
 - YOLO(You Only Look Once)
画像をグリッドに分割しそれぞれのグリッドセルで複数のバウンディングボックスとクラス確率を予測する。リアルタイムの物体検出に向いている。
 - SSD(Single Shot Multibox Detectors)
画像全体を見て異なるスケールとアスペクト比のバウンディングボックスを予測する。
YOLO と同様に高速。
- 二段階検出器(Two-Stage Detectors)
 - R-CNN(Region-based Convolutional Neural Networks)
画像内の物体検出候補領域を提案し、その後これらの候補領域で詳細な分類と回帰を行う。
 - Fast R-CNN
R-CNN の改良版で ROI プーリングを用いて計算効率を向上させている。
 - Faster R-CNN
物体候補領域の提案部分をニューラルネットワークで行うことでさらに高速化を実現している。
- その他手法
 - Mask R-CNN
R-CNN の拡張で物体検知と同時に物体のセグメンテーションも行う。
 - RetinaNet
単段検出器の中で特にフォーカルロスを用いてクラスの不均衡問題に対処したモデル。

$$\text{IoU} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}}$$

Mask R-CNN

物体検知とセグメンテーションを同時に行うことができるモデル。物体検知に加えて各物体のピクセル単位のマスクを生成する。物体検出の結果として得られた領域についてのみセグメンテーションを行う。

Mask R-CNN の構造

- バックボーンネットワーク
画像特徴を抽出するためお畳み込みニューラルネットワークが使われる。
- 地域提案ネットワーク
Faster R-CNN と同様に地域提案ネットワークが特徴マップから物体候補領域(ROI)を提案する。
- ROI アライメント
ROI プーリングに代わり ROI アライメントが使用される。候補領域を固定サイズの特徴マップに変換する際のずれを防ぐ。
- 分類とバウンディングボックス回帰
物体のクラスとバウンディングボックスを予測する。
- マスク予測
ピクセル単位のマスクを予測する。このマスク予測は物体検知と平行して行われ各クラスに対して独立したマスクを生成する。

特徴

- ピクセルレベルで精度が高い。物体の形状や細部を詳細にとらえることができる。
- 物体検知とセグメンテーションを統合的に行うため自動運転の障害物検知、医療画像の解析など応用分野に適している。
- 基本的な Faster R-CNN に対する拡張なので既存の物体検知システムに容易に組み込むことができる。

ROI プーリング

- 入力画像内の異なるサイズの候補領域を固定サイズの特徴マップに変換する
- 提案された候補領域を固定サイズのグリッドに分割し、各グリッドセル内で Max プーリングを行い、固定サイズの特徴マップを生成する。
- 元の候補領域の情報が一部失われることがある

ROI アライメント

- 補間処理によってより多くのピクセルの享保を使うことで推定精度を上げる。
- グリッドセル内の特徴を種痘する際に双線形補間などを用いてピクセル間のサンプリングを行う。これによって誤差を最小限に抑えより正確な特徴マップを生成する。
- 計算コストが高いが ROI プーリングよりも精度の向上が期待できる。

補足

Faster R-CNN

Faster R-CNN は R-CNN(Region-based Convolutional Neural Network)シリーズの改良版であり、高速かつ精度の高い物体検知を実現している。

Faster R-CNN の構造

- バックボーンネットワーク
入力画像から特徴マップを抽出する。
- 地域提案ネットワーク(RPN)
画像全体から物体が存在する可能性が高い候補領域を提案する。特徴マップ上でアンカーと呼ばれる固定サイズのボックスをスライドさせ物体が含まれているかどうか評価する。
- ROI プーリング
候補領域を固定サイズの特徴マップに変換する。
- 分類と回帰ヘッド
各候補領域に対して物体のクラス分類とバウンディングボックスの予測をする。

YOLO(You Only Look Once)

YOLO は画像を一度だけ見ることで複数の物体をリアルタイムで検出できる。他の物体検知手法と比べて高速であるためリアルタイムアプリケーションに適している。画像を複数のグリッドに分割し、各グリッドが特定の物体を検出する。

アーキテクチャ

- 画像入力
入力画像は固定サイズにリサイズする。
- 畳込み層
特徴抽出のため一連の畳込みそうとプーリング層を使用する。
- 完全結合層
特徴マップをフラットにして最終的に完全結合層を通して固定長の出力ベクトルを得る。
- 出力層
出力ベクトルは $S \times S \times (B \times 5 + C)$ のサイズになる。ここで S はグリッドのサイズ、 B は各グリッドセルが予測するバウンディングボックスの数、 C はクラスの数を表す。各バウンディングボックスは 5 つの要素($x, y, w, h, \text{信頼度}$)を持つ。

動作

- グリッド分割
入力画像を $S \times S$ のグリッドに分割する。
- バウンディングボックスの予測
各グリッドセルは B 個のバウンディングボックスを予測する。
- クラス予測
各グリッドセルは C 個のクラススコアを予測する。
- 最終出力
各グリッドセルから予測されたバウンディングボックスの信頼度スコアとクラススコアを組み合わせ、最終的な物体検出結果を生成する。

YOLO は一度の予測で画像全体を処理しリアルタイムで物体検出できるほどの高速性がある。アーキテクチャもシンプルでエンドツーエンドのトレーニングが可能のため実装が比較的簡単となっている。ただ、小さな物体の検出が難しいことや物体が密集していたり重なり合っている場合には検出精度が低下する欠点がある。

FCOS(Fully Convolutional One-Stage Object Detection)

FCOS はアンカーボックスを使用せずに物体検出を行う。完全畳み込みネットワークをベースにしており、各ピクセルが物体の存在を予測する責任を負っている。アンカー方式に関連する複雑な設計や調整が不要になる。

アーキテクチャ

- 特徴抽出ネットワーク
ResNet や ResNeXt などのバックボーンネットワークが使用され、入力画像を足そうの特徴マップに抽出する。
- ヘッドネットワーク
 - 分類ヘッド
各ピクセルのクラス確率を予測する、
 - 回帰ヘッド
各ピクセルから物体のバウンディングボックスのオフセットを予測する。
 - センターネスヘッド
各ピクセルのセンターネススコアを予測し物体の中心に近いほど高いスコアを与える。

動作

- 特徴抽出
入力画像をバックボーンネットワークに通し、複数の解像度の特徴マップを生成する。

- ピクセルごとの予測
各特長マップの各ピクセルに対して分類、回帰、センターネススコアを予測する。
- バウンディングボックスの生成
回帰ヘッドが予測したオフセットをもとに各ピクセルからバウンディングボックスを生成する。
- スコアの統合
各バウンディングボックスに対して分類スコアとセンターネススコアを掛け合わせて最終的なスコアを算出する。
- NMS(Non-Maximum Suppression)
一番スコアが高いバウンディングボックスとそれ以外のバウンディングボックスとで IoU を計算しその値に応じて重複するバウンディングボックスを排除し最終的な検出結果を出力する。

多様な物体サイズに対応でき、センターネスヘッドを導入したことで誤検出が少なく精度が向上した。一方で完全畳み込みネットワークとピクセルごとの予測により計算コストが増加している。

※完全畳み込みネットワーク：全結合層を持たずネットワークが畳み込み層のみで構成される

Transformer

Transformer モデルは RNN や CNN を使用せずに完全に Attention 機構のみで構築されているモデルのこと。Seq2seq と同様に Encoder と Decoder の二つの部分から構成されている。RNN や CNN の限界を克服し、非常に高い性能を発揮する。特に長いシーケンスの依存関係を効果的に扱うことができる。機械翻訳、文書要約、質問応答や自然言語生成などに応用される。

- Attention 機構
 - 入力シーケンスの異なる部分に異なる重みを与えることで、重要な情報に焦点を当てる方法。
 - Scaled Dot-Product Attention
 - Query, Key, Value という 3 つの行列を使う。
 - Attention スコアは Query と Key のドット積を計算しその結果を Key の次元の平方根で割って正規化しソフトマックス関数を適用する。
 - このスコアに Value を掛け合わせて重み付けされた出力を得る。

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

- Self-Attention
 - Self-Attention(自己注意機構)は入力シーケンス内の各要素がシーケンス内の他のすべての要素に対してどれだけ重要かを評価する。これにより各要素が他のすべての要素から情報を取得できる。

- 入力シーケンスを Query, Key, Value に変換し、これらを用いて Attention スコアを計算する。
- Multi-Head Attention
 - 単一の Attention head を使用するのではなく複数の Attention head を平行して使用することにより異なる部分空間での Attention を同時に学習できる。
- Positional Encoding
 - Transformer ではシーケンスの順序情報を持たないため入力に位置情報を追加する必要がある。代表的な手法として三角関数を使う方法がある。

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

- pos は位置インデックス
- i は埋め込み次元のインデックス
- d はモデルの埋め込み次元数

BERT

BERT(Bidirectional Encoder Representations from Transformers)は自然言語処理において広く使われている事前学習モデルの一つ。Google が 2018 年に発表した。

特徴

- 双方向性
BERT は入力文の前後の文脈を同時に考慮する双方向性を持っている。これにより文の全体的な意味をより正確に捉えることができる。
- 事前学習とファインチューニング
大規模なコーパスを使用して一般的な言語表現を事前に学習している。ファインチューニングすることで特定のタスクに対して微調整を行うことができ、少量のタスク固有データを用いて事前学習済みモデルをタスクに適応させる。

アーキテクチャ

- Transformer の Encoder 部分をベースにしている。
- 入力はトークンの埋め込み、セグメントの埋め込み、および位置の埋め込みの 3 つから構成される。
- 複数の Encoder 層から構成され、それぞれの層は Multi-Head Attention と Feed Forward NN から成る。

事前学習

- Masked Language Model
 - 入力分の一部のトークンをマスクし、そのマスクされたトークンを予測する。これによりモデルは文脈を理解しマスクされた単語を正しく予測する能力を学習する。
- Next Sentence Prediction
 - 二つの文を入力として二つ目の文が一つ目の文に続くかどうかを予測する。これにより文間関係を理解する能力を強化する。

高い汎化性能があり文の意味を正確に捉えることができる反面、計算資源を多く消費することや推論速度が遅くなることもあることが課題となっている。

実装演習

文章生成

```
BATCH_SIZE = 64
BUFFER_SIZE = 10000
dataset = dataset.shuffle(BUFFER_SIZE).batch(BATCH_SIZE, drop_remainder=True)
input_ids = tf.keras.layers.Input(shape=(None, ), dtype='int32', name='input_ids')
inputs = [input_ids]
bert.trainable = False
x = bert(inputs)
out = x[0]
Y = tf.keras.layers.Dense(len(vocab))(out)
checkpoint_dir = './training_checkpoints'
checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch}")
checkpoint_callback=tf.keras.callbacks.ModelCheckpoint(
    filepath=checkpoint_prefix,
    save_weights_only=True)
model = tf.keras.Model(inputs=inputs, outputs=Y)
def loss(labels, logits):
    return tf.keras.losses.sparse_categorical_crossentropy(labels, logits, from_logits=True)

model.compile(loss=loss,
              optimizer=tf.keras.optimizers.Adam(1e-7))
model.fit(dataset, epochs=5, callbacks=[checkpoint_callback])

text = '私は'
mecab = MeCab.Tagger("-Owakati")
text = mecab.parse(text).split()
generate_text(model, text)
```

結果

私はいん全集当隠さ不味横町湧多写し永く切るむじゃき抽象時期殺丸
裸本願縞約束嘗ちぢみ詫まら離れ離れ知る碑天井渡ついきが嫌い漕

文章分類

```
BUFFER_SIZE = 10000
dataset = train_dataset.shuffle(BUFFER_SIZE)
input_ids = tf.keras.layers.Input(shape=(None, ), dtype='int32', name='input_ids')
inputs = [input_ids]
bert.trainable = False
x = bert(inputs)
out = x[1]
fully_connected = tf.keras.layers.Dense(256, activation='relu')(out)
Y = tf.keras.layers.Dense(3, activation='softmax')(fully_connected)
checkpoint_dir = './training_checkpoints'
checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt_{epoch}")
checkpoint_callback=tf.keras.callbacks.ModelCheckpoint(
    filepath=checkpoint_prefix,
    save_weights_only=True)
model = tf.keras.Model(inputs=inputs, outputs=Y)
def loss(labels, logits):
    return tf.keras.losses.categorical_crossentropy(labels, logits)

model.compile(loss=loss,
              optimizer=tf.keras.optimizers.Adam(1e-7))
model.fit(dataset, epochs=5, callbacks=[checkpoint_callback])

text = '私は私の過去を善悪ともに他ひとの参考に供するつもりです'
encoded = tokenizer.encode_plus(
    text,
    text,
    add_special_tokens=True,
    max_length=128,
    pad_to_max_length=True,
    return_attention_mask=True
)
inputs = tf.expand_dims(encoded["input_ids"], 0)
res = model.predict_on_batch(inputs)
res
```

```
array([[0.30957258, 0.30389112, 0.3865363 ]], dtype=float32)
```

考察

文章生成タスクでは意味が通る文章ではないものの単語レベルでは抽象や離れなど生成できている。分類タスクはスコアに大きな差は出なかった。今回はエポック数が少なすぎたが十分モデルを訓練させればより良い性能が期待できる。

GPT

GPT(Generative Pre-trained Transformer)は OpenAI が開発した自然言語処理モデルの総称で、生成タスクに特化している。その高い性能と広範な応用範囲から注目されている。ある単語の次に来る単語を予測し自動的に文章を生成する。

基本構造

- Transformer の Decoder 部分のみを使用している。
- テキストデータをトークンに分割しそれぞれのトークンを固定長のベクトルに変換したものに、入力シーケンスの順序情報を保持するための位置情報を加えている。
- Multi-Head Self-Attention と Feed forward NN から構成される複数の Decoder レイヤを持つ。

特徴

- GPT は与えられたプロンプトから自然で一貫したテキストを生成する能力に優れている。
- モデルサイズを大きくすることで性能が向上することが示されている。大型モデルは数百億のパラメータを持つ。GPT3 では 1750 億個。
- 事前学習とファインチューニングの組み合わせにより少量のタスク固有データでも高いパフォーマンスを発揮できる。

推論

- Zero-shot
 - 特定のタスクに対して一切の訓練データを使用せずにモデルが予測を行う。
- One-shot
 - 特定のタスクに対して一つの訓練例を使用してモデルを訓練し実行させる。
- Few-shot
 - 特定のタスクに対して複数の訓練例を使用してモデルを訓練し実行させる。

音声認識

音声データはスマートスピーカーなどの自然言語処理や鳥の分類タスクなどに使うことができる。音は空気振動であるため振動強度の時間波形とも言える。機械学習では音声信号をモデルが理解しやすい別の形式に変換してから扱うことが一般的であり、kaggle BirdCLEF ではメルスペクトログラムに変換し CNN で分類したモデルが高い性能を示した。

波形->波形変換

- 標本化(Sampling)
標本化はアナログ音声信号をデジタル信号に変換すること。連続的な音声波形を一定間隔でサン

プリングする。

- サンプリング周波数

1 秒間にサンプリングする回数。元の周波数より大きい値である必要がある。

- 量子化

連続値であるサンプル値を離散値に変換すること。サンプル値を有限のビット数で表現する。

- ビット深度

各サンプルを何ビットで表現するかを決定する。

- フーリエ変換

時間領域の信号を周波数領域の信号に変換すること。

- 離散フーリエ変換(DFT)

- 高速フーリエ変換(FFT)

波形->画像変換

- スペクトログラム

時間と周波数の両方の情報を視覚化したもの。短時間フーリエ変換(STFT)を用いて信号を小さな時間窓に分割しそれぞれの窓にフーリエ変換を適用することで得られる。

- メルスペクトログラム

スペクトログラムをメル尺度に変換したもの。メル尺度は人間の聴覚感度に基づいたスケールで低周波数と高周波数の感度の違いを反映している。

CTC

CTC(Connectionist Temporal Classification)は入力信号を音声特徴ベクトルに変換、その音声特徴ベクトルの系列から対応する単語列を推定する音声認識タスクに用いられる。入力シーケンスと出力シーケンスの長さが異なる場合やラベルが時系列データ内のどこに出現するかが不明な場合に有効となる。CTC では入力シーケンスに対して出力ラベルの可能なすべてのアライメントを考慮しそれらの確率を合計して損失を計算する。

仕組み

- 入力シーケンスと出力シーケンス

入力シーケンス:時系列データ

出力シーケンス:ラベルのシーケンス

- ブランクラベルの導入

ブランクラベルは入力フレームがどのラベルにも対応しない場合を示す。

- ラベルのアライメント

該当するラベルがないことを意味するブランクの出力を許可しラベルを無理やり割り当てることがないようにしている。

- 損失計算

ターゲットラベルシーケンスに対応するすべてのアライメントの確率を合計して最終的な損失とする。

損失計算

- 前向きアルゴリズム
可能なアライメントの確率を累積的に計算する。
- 後向きアルゴリズム
逆方向に同様の計算を行う。
- 合計
前向きと後向きの確率を組み合わせて最終的な損失を計算する。

DCGAN

GAN(Generative Adversarial Nets)とは生成モデルと識別モデルの二つのニューラルネットワークを対立的に訓練することによって新しいデータを生成する手法。DCGAN(Deep Convolutional GAN)はGANの一種であり、生成モデルと識別モデルの両方に畳み込みニューラルネットワークを使用することで画像生成に特化したアーキテクチャとなる。

GAN の基本構造

- 生成モデル(Generator)
ランダムなノイズを入力として新しいデータを生成する。本物データと見分けがつかないほどリアルなデータを生成することを目標とする。
- 識別モデル(Discriminator)
本物のデータと生成モデルが生成した偽物のデータを正確に分類することを目標とする。
- 訓練プロセス
これらのモデルは交互に更新される。まず Discriminator を固定して Generator を更新し、その後 Generator を固定して Discriminator を更新する。この対立的な訓練により Generator と Discriminator は互いに進化し高品質なデータ生成が可能になる。

DCGAN の特徴

- 畳み込み層の使用
Generator と Discriminator の両方に畳み込み層を使用している。
- バッチ正規化
各層の出力を正規化する。
- スライド付き畳み込み
ダウンサンプリングやアップサンプリングにスライド付き畳み込みまたは転置畳み込み層を使

用する。プーリング層を使用せずに特徴量を圧縮または拡張する。

- ReLU および Leaky ReLU

Generator では ReLU 活性化関数、Discriminator では Leaky ReLU 活性化関数を使用する。

- Generator と Discriminator の出力の確率分布がどれくらい近いのか測る指標として JS ダイバージェンスがある。

Conditional GAN

CGAN(Conditional GAN)は GAN の一種で生成するデータに対して追加の条件を指定することで特定の属性や特徴を持つデータを生成するための手法。条件付生成を可能にすることでより制御されたデータ生成が可能となる。CGAN は基本的に GAN と同様に生成モデルと識別モデルを持つが、これらのモデルに条件が追加される。この条件情報は生成されるデータが満たすべき属性や特徴を指定する。

- 生成モデル(Generator)
 - 入力: ランダムノイズベクトル z と条件ベクトル y
 - 出力: 条件 y に従った生成データ $G(z|y)$
- 識別モデル(Discriminator)
 - 入力: データサンプルと条件ベクトル y
 - 出力: 本物か偽物かの判定 $D(x|y)$

動作原理

1. 条件付生成

生成モデルはノイズベクトル z と条件ベクトル y を入力として受け取り、条件 y に従ったデータを生成する。

2. 条件付判定

識別モデルはデータサンプルとその条件ベクトル y を入力として受け取りそのデータが本物か偽物かを判定する。

3. 損失関数

生成モデルと識別モデルは以下の損失関数を用いて交互に最適化される。

$$\begin{aligned} E_D &= -\mathbb{E}_{x \sim p_{data}(x)} [\log D(x|y)] - \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z|y)|y))] \\ E_G &= -\mathbb{E}_{z \sim p_z(z)} [\log D(G(z|y)|y)] \end{aligned}$$

応用例

- 画像生成
MNIST データセットでは特定の数字の画像を生成できる。
- テキスト生成
特定の属性を持つテキストの生成。

- データ拡張

Pix2Pix

Pix2Pix は画像から画像への変換を行うための条件付き生成モデル。これは特定のタスクに対してペアの画像を学習データとして使用し、入力画像からターゲット画像を生成することを目的としている。Pix2Pix は U-Net をベースとした生成モデルと PatchGAN をベースとした識別モデルを組み合わせて訓練される。

- 生成モデル(Generator)

U-Net アーキテクチャを使用しており、入力画像を受け取りターゲット画像を生成する。U-Net はエンコーダとデコーダの間にスキップ接続を持つ構造でエンコーダが入力画像の特徴を抽出し、デコーダがそれを元のターゲット画像に再構築する。

- 識別モデル(Discriminator)

PatchGAN アーキテクチャを使用しており、生成された画像と本物の画像のローカルパッチが本物か偽物かを判定する。PatchGAN は画像全体ではなく画像の小さなパッチごとに判定を行うため全体の詳細な特徴を捉えやすい。また、正確な高周波成分の強調により視覚的な一致性が向上している。

動作原理

- 条件付生成

生成モデルは入力画像を受け取りターゲット画像を生成する。

- 条件付判定

識別モデルは生成された画像と対応する入力画像のペアが本物じゃ偽物かを判定する。

- 損失関数

以下の損失関数を用いて訓練される。

$$\begin{aligned} L_{CGAN}(G, D) &= \mathbb{E}_{x,y}[\log D(x, y)] + \mathbb{E}_{x,z}[\log (1 - D(x, G(x, z)))] \\ L_{L1}(G) &= \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1] \\ \therefore L(G, D) &= L_{CGAN}(G, D) + \lambda L_{L1}(G) \end{aligned}$$

応用例

- 画像変換

- 白黒写真をカラー写真に変換する。
- 手書きスケッチを写真に変換する。

- 画像補間

- 部分的に欠損した画像を補完する。

- セグメンテーション

- 航空写真から地図を生成する。

A3C

A3C(Asynchronous Advantage Actor-Critic)は強化学習の手法の一つであり DeepMind によって提案された。並行して動作する複数のエージェントを使用してより効率的に学習を行うことを目的としている。

基本構成

- Actor-Critic
 - Actor:行動ポリシーを決定し環境から報酬を受け取る。Actor は環境の状態を入力として各坑道の確率分布を出力する。
 - Critic:価値関数を評価し Actor の選択した行動を評価する。Critic は状態価値関数 (V) または行動価値関数 (Q) を予測する。
- Advantage 関数
 - 特定の行動が平均よりどれだけ良いかを評価する。

$$A(s, a) = Q(s, a) - V(s)$$

ここで $Q(s, a)$ は状態 s における行動 a の価値、 $V(s)$ は状態 s の価値

- 非同期更新
 - A3C は複数のワーカーが並列して動作し環境を探索して経験を集め、非同期にパラメータを更新する。これにより収束を早め局所解を回避することができる。

アルゴリズム

1. 環境の初期化

各ワーカーは独自の環境インスタンスを持ちそれぞれ独立して動作する。

2. ローカルネットワークの更新

- 各ワーカーはローカルなポリシーネットワークと価値ネットワークを持つ。
- 各ワーカーは環境から状態 s_t を観測し行動 a_t を選択する。
- 行動に基づいて環境が次の状態 s_{t+1} と報酬 r_t を返す。
- ローカルな経験を使用して Advantage 関数 $A(s, a)$ を計算し、ローカルネットワークのパラメータを更新する。

3. グローバルネットワークの更新

- 一定のステップごとに各ワーカーはローカルネットワークの更新をグローバルネットワークに反映させる。
- グローバルネットワークのパラメータはすべてのワーカーからの勾配を集約して更新される。

Metric-learning(距離学習)

Metric-learning(距離学習)はデータポイント間の類似度や距離を学習するためのアプローチのこと。距離学習では同じクラスに属するデータポイント同士が近く、異なるクラスに属するデータポイント同士が遠くなるように距離関数や類似度関数を学習する。この方法は分類、クラスタリング、検索、レコメンデーションシステムなどに用いられる。

Siamese Network

二つの入力データポイントを同じニューラルネットワークに通し、それぞれの埋め込み表現を取得する。その後これらの特徴ベクトル間の距離を計算する。目的は同じクラスに属するペアの特徴ベクトルの距離を小さく、異なるクラスに属するペアの距離を大きくすること。

- 共有ウェイト
二つの同じニューラルネットワーク（共有ウェイト）に二つの入力を与える。
- 特徴ベクトルの取得
それぞれのネットワークから出力される特徴ベクトルを取得する。
- 距離の計算
ユークリッド距離やコサイン類似度などを使用して二つの特徴ベクトル間の距離を計算する。
- 損失関数
contrastive loss などの損失関数を使用して距離が近くなるようにネットワークをトレーニングする。

$$L = \frac{1}{2} \left(yD^2 + (1 - y) \max(m - D, 0)^2 \right)$$
$$D = \|f(x_1) - f(x_2)\|_2$$
$$m : \text{マージン} \quad y : \text{識別ラベル}$$

Triplet Network

Triplet Network は Siamese Network を拡張したもの。このネットワークではアンカー、ポジティブ、ネガティブの三つの入力データポイントを使用する。アンカーとポジティブは同じクラスに属し、アンカーネガティブは異なるクラスに属する。Siamese Network で問題となった同じクラスのペアと異なるクラスのペアの間に生じる不均衡が解消される。

- 共有ウェイト
三つの同じニューラルネットワーク（共有ウェイト）にアンカー、ポジティブ、ネガティブを与える。
- 特徴ベクトルの取得
それぞれのネットワークから出力される特徴ベクトルを取得する。
- 距離の計算
アンカーとポジティブ間の距離、アンカーとネガティブ間の距離を計算する。

- トリプレットロス

以下のトリプレットロス関数を使用してアンカーとポジティブ間の距離がアンカーとネガティブ間の距離よりも一定のマージンだけ小さくなるようにネットワークをトレーニングする。

$$L = \max(0, D(a, p) - D(a, n) + m)$$

ここで $D(a, p)$ はアンカーとポジティブの距離、 $D(a, n)$ はアンカーとネガティブの距離、 m はマージン

問題点

- 学習がすぐに停滞してしまう。
- クラス内距離がクラス間距離より小さくなることを保証しない。
->Quadruplet loss 関数を導入する。

MAML(メタ学習)

MAML(Model-Agnostic Meta-Learning)はメタ学習のアプローチの一つでモデルが新しいタスクに迅速に適応できるように訓練する手法。メタ学習とはモデルが学習の方法を学習することを指し、少数のデータポイントから新しいタスクを学習する能力を高めることを目的としている。MAML はモデルが少数の勾配ステップで新しいタスクに適応できるようにパラメータを最適化する。

- Model-Agnostic
微分可能である以外モデルや損失関数の具体的な形式を仮定しない。
- Task-Agnostic
回帰、分類、強化学習など様々なタスクに適用できる。

訓練過程

1. タスクのサンプリング

トレーニングタスクの集合からいくつかのタスクをランダムにサンプリングする。

2. タスクごとのトレーニング

各サンプリングされたタスクに対してメタパラメータ（共有パラメータ）を用いてタスク固有のデータを使ってモデルを少数のステップで更新する。この更新されたパラメータを用いてタスク固有の損失を計算する。

$$\theta'_1 = \theta - \alpha \nabla_{\theta} L_{T_1}(f_{\theta})$$

3. メタ更新

すべてのタスクから損失を用いてメタパラメータを更新する。各タスクの損失に基づいてメタパラメータを最適化する。

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i L_{T_1}(f_{\theta'_i})$$

MAML は二階の勾配計算が必要なため計算コストが非常に高い。実用的には inner loop のステップ数を大きくできない。

->

- First-order MAML
2 次以上の勾配を無視する。高次の微分が不要となる。
- Reptile
inner loop の逆伝播を諦め学習前後のパラメータの差のみを取る。

MAML は few-shot learning で有効であり以下のタスクに応用される。

- 画像分類
新しいカテゴリの画像を少数の例から学習する。
- 強化学習
新しい環境やゲームに迅速に適応する。
- ロボティクス
新しいタスクや操作に適応するロボットの制御。

グラフ畳込み(GCN)

グラフ畳込みネットワーク(Graph Convolutional Network)はグラフ構造を持つデータに対して畳み込み操作を行う深層学習のアプローチ。GCN はノードの特徴をその周辺ノードの特徴とともに集約することでグラフ上でのノードの表現を学習する。グラフとはノード（頂点）とエッジ（辺）から構成される。ノードは特徴ベクトルを持ち、エッジはノード間の関係性を表している。GCN は各ノードの特徴をその周辺ノードの特徴として集約し更新する。これによりノードの新しい表現が得られる。

- 空間アプローチ Spatial GCN
グラフの空間構造を直接操作して畳み込みを実行する。これにより隣接するノードの情報を集約する。
 - メリット；直観的で異なるサイズのグラフやダイナミックグラフに対して柔軟。
 - デメリット；理論的な解析が難しい。
- スペクトルアプローチ Spectral GCN
グラフ信号処理に基づいておりグラフラプラシアンを利用して畳み込み操作を行う。このアプローチはグラフの周波数領域で層さを行うことからスペクトルと呼ばれる。
 - メリット：グラフ周波領域での解釈が可能。

- デメリット：グラフサイズに依存し固有値分解の計算コストが高い。異なるグラフに対して一般化が難しい。

応用

- ノード分類
ソーシャルネットワークにおけるユーザーの分類
- リンク予測
ソーシャルネットワークにおける友人推薦、化学分子内の結合予測など
- グラフ分類
化学分子の性質予測
- レコメンデーションシステム
ユーザーとの関係性を考慮した推薦システム

Grad-CAM, LIME, SHAP

モデル解釈性はブラックボックスとしての機械学習モデルがどのようにして予測を行っているか理解することを助ける。

- Grad-CAM(Gradient-weighted Class Activation Mapping)
Grad-CAM は特に CNN の解釈に用いられる。モデルがどの部分に注目して予測を行ったかを視覚的に示す。
 - 対象のクラスに対する出力の勾配を計算しその勾配を最後の畳み込み層の出力に対して逆伝播する。
 - 最後の畳み込み層の出力に対して勾配を集約し各フィルタの重要度を計算する。
 - 各フィルタの重要度を元に重み付けを行い重み付けされた特徴マップを合計してクラスアクティベーションマップを生成する。
 - クラスアクティベーションマップを元の画像に重ね合わせて視覚化する。
- LIME(Local Interpretable Model-agnostic Explanations)
LIME は任意の機械学習モデルに対して局所的な解釈を提供するための手法。モデルの予測に影響を与える入力特徴量の重要性を理解できる。
 - LIME は元の入力データの周辺で新しいデータポイントを生成し、そのデータポイントに対してモデルの予測を取得する。
 - 生成されたデータポイントとモデルの予測を用いて単純で解釈可能なモデルを適用する。
 - この単純なモデルの係数を用いて元のモデルの予測に対する各特徴の影響を解釈する。
- SHAP(Shapley Additive Explanations)
SHAP はゲーム理論に基づく手法で各特徴の予測に対する寄与度を計算する。SHAP は特徴の重要性を一貫して解釈できる。

- シャープレイ値を用いて各特徴の寄与度を計算する。シャープレイ値はすべての特徴の組み合わせにおける予測の変化を考慮して各特徴の寄与を公平に評価する。
- SHAP 値は特定のデータポイントに対するモデルの予測を基準にして各特徴がその予測にどれだけ貢献しているかを示す。
- 比較
 - Grad-CAM
 - 特に CNN に適している。画像データに対して視覚的な解釈を提供する。
 - LIME
 - モデル非依存で任意の機械学習モデルの適用可能。局所的な解釈を提供するため特定のデータポイントに対する予測の理由を理解しやすい。
 - SHAP
 - ゲーム理論に基づき特徴の寄与度を公平に評価する。グローバルおよび局所的な解釈が可能。

実装演習

```
import numpy as np
import cv2
import tensorflow as tf

from tensorflow.keras.applications import VGG16

class GradCam:
    def __init__(self, model):
        self.model = model
        # 畳み込み最終層の名前を確認するため
        print([layer.name for layer in self.model.layers])

    def gradcam_func(self, x, layer_name):
        # 一枚の画像だと、バッチの次元がないので足す
        X = x[np.newaxis, ...]
        # 正規化
        X = X / 255.

        # 畳み込み層の最後の層の出力を受け取る
        conv_feature = self.model.get_layer(layer_name).output
        model = tf.keras.Model([self.model.inputs], [conv_feature, self.model.output])

        # 勾配を記録するために tf.GradientTape() を使う
        with tf.GradientTape() as tape:
            # numpy配列を勾配を計算するためにtfの型に変換する
            X = tf.cast(X, tf.float32)
            conv_feature, outputs = model(X)

            # どのクラスを予測したか
            predicted_class = tf.math.argmax(outputs[0])
            # 予測したクラスの出力を取得する
            class_outputs = outputs[:, predicted_class]
            # 勾配を計算する
            grads = tape.gradient(class_outputs, conv_feature)

            print('予測クラス', predicted_class.numpy())

            # 平均を取る(GAP)
            weights = tf.math.reduce_mean(grads, axis=(1, 2))
            cam = conv_feature @ weights[..., tf.newaxis]
            cam = tf.squeeze(cam)

            # reluに通す
            cam = tf.nn.relu(cam)
            cam = cam / tf.math.reduce_max(cam)
            # 正規化を戻す
            cam = 255. * cam
```



```

        # numpy配列にする
        cam = cam.numpy()
        cam = cam.astype('uint8')

        # カラーマップを作る
        jetcam = cv2.applyColorMap(cam, cv2.COLORMAP_JET)
        # BGRからRGBに変換
        jetcam = cv2.cvtColor(jetcam, cv2.COLOR_BGR2RGB)
        jetcam = cv2.resize(jetcam, (224, 224))
        jetcam = jetcam + x / 2

    return jetcam

model = VGG16(weights='imagenet')
gradcam = GradCam(model)
image = cv2.imread('interpretability_example_input.png')
image = cv2.resize(image, (224, 224))
cam = gradcam.gradcam_func(image, 'block5_conv3')

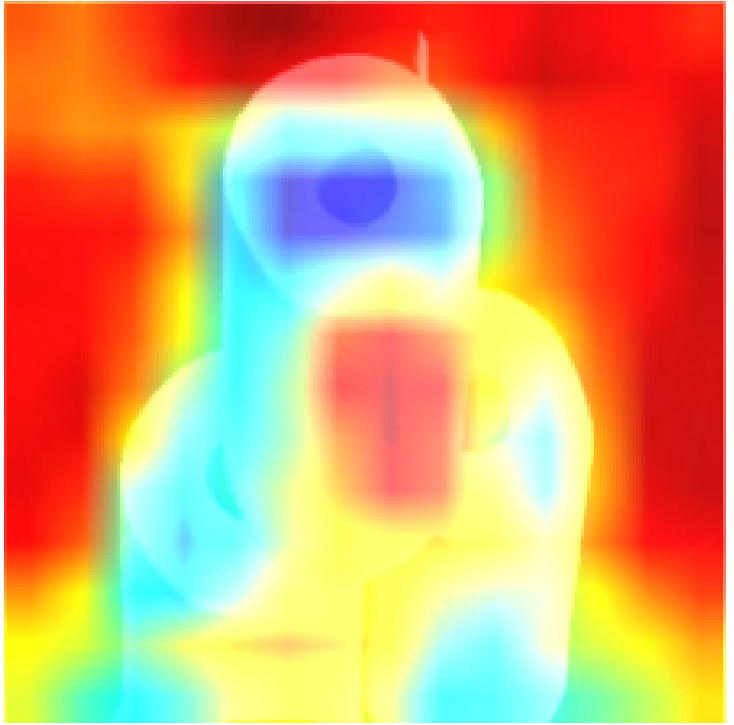
if ENV_COLAB:
    from google.colab.patches import cv2_imshow
    cv2_imshow(image)
    cv2_imshow(cam)
else: # Jupyter
    from matplotlib import pyplot as plt
    cv2.imshow('', image)
    cv2.waitKey(0)
    cv2.destroyAllWindows()
    cv2.imwrite('./data/interpretability_example_output.png', cam)
    cam_read = cv2.imread('./data/interpretability_example_output.png')
    cv2.imshow('', cam_read)
    cv2.waitKey(0)
    cv2.destroyAllWindows()

```

結果

オリジナル

Grad-CAM



Docker

Docker はアプリケーションのデプロイと管理を簡素化するためのコンテナ技術ツール。深層学習のプロジェクトは特定のライブラリやツールのバージョンに依存することが多く、異なる環境間での移植が難しいが、Docker を使用することで開発環境、テスト環境、本番環境で一貫した環境を構築することができる。

- 環境設定の管理
Dockerfile を使用して必要なライブラリや依存関係を明確に定義できる。
- 一貫性
Docker コンテナ内での動作はどの環境でも同じになる。
- 隔離された環境
Docker コンテナはホスト環境と隔離されているため、依存関係の競合を回避できる。
- バージョン管理
特定のライブラリバージョンを Dockerfile で指定することで依存関係のバージョン管理が容易になる。
- 再現性
Docker イメージを DockerHub などのリポジトリに保存共有することで誰でも同じ環境を再現できる。
- オーケストレーションツール
Kubernetes などのオーケストレーションツールと組み合わせることで複数のコンテナを効率的に管理スケールアップできる。