

目次：

- ・ P2 : Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model__(MuZero)

(関連論文)

- ・ P3 : Deep Residual Learning for Image Recognition__(ResNet)
 - ・ P4 : Identity Mappings in Deep Residual Networks__(ResNet)
 - ・ P5,6 : A Survey of Monte Carlo Tree Search Methods__(モンテカルロ木探索)
-
- ・ P9 : Maximum_Entropy_Deep_Inverce_Reinforcement_Learning__(逆強化学習)

Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model

(2019) 著者: Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre,

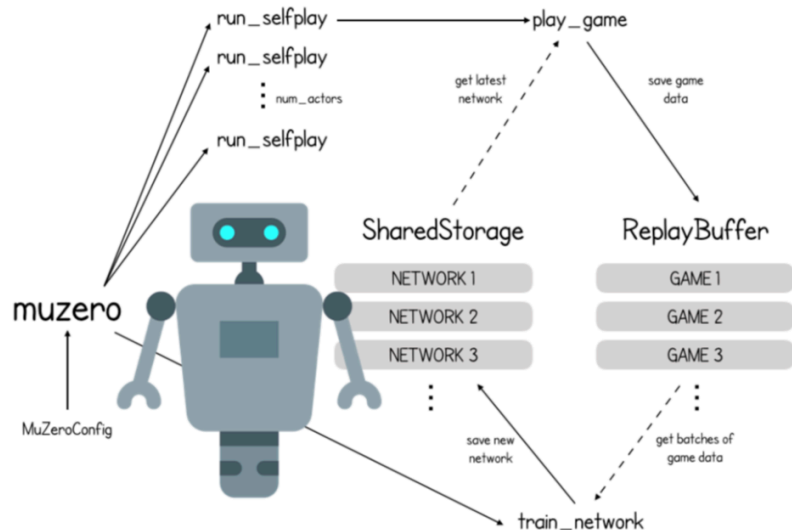
Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy Lillicrap, David Silver

概要

AIが囲碁・チェス・将棋などのボードゲームおよびゲームを学習、マスターし、ボードゲームにおいては世界チャンピオンに圧勝している。ゲームにおいても多くを人間以上のレベルに達しており、今後は実世界への応用が期待される。

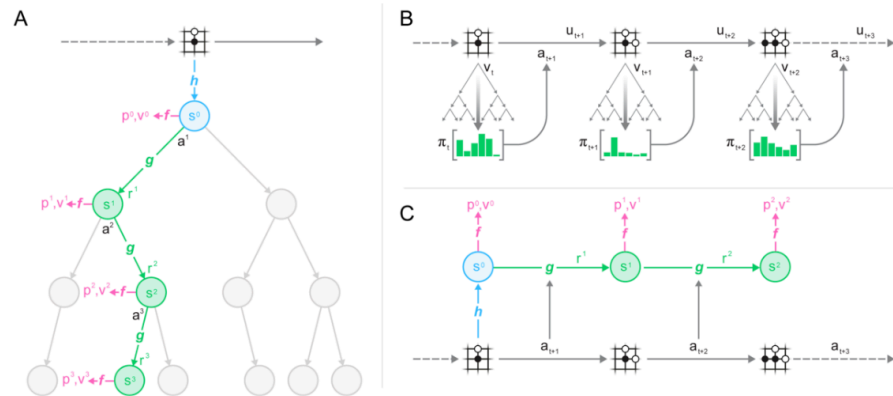
名前はMuZero。

◇MuZeroの構造



先行研究: AlphaZeroからの進化

- ・表現関数(h)と遷移関数(g)をニューラルネットワーク (ResNet) で表現。
(全部で3つのニューラルネットワークを学習)
- ・ニューラルネットワークで表現したことでAIにルールを教えなくとも学習過程で自らルールを学習。
(ただし勝敗引き分けは教える必要あり)



次に読むべき文献

- ・ Identity Mappings in Deep Residual Networks
- ・ Single-Player Monte-Carlo Tree Search
- ・ Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments
- ・ Ray- A Distributed Framework for Emerging AI Applications

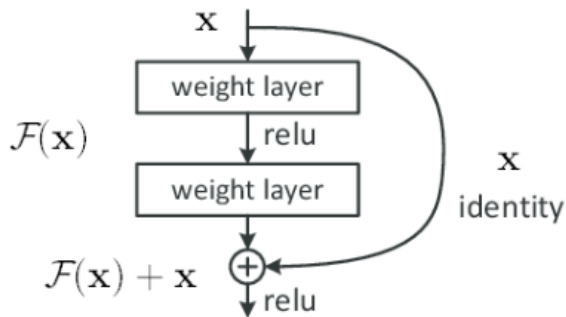
Deep Residual Learning for Image Recognition

(2015) 著者: Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun

概要

より深層なニューラルネットワーク（最大1000層以上）を構築することを可能にした論文。もともと層を深くすれば性能が悪化と言われていたが、ResNetでは「ある層で求める最適な出力を学習するのではなく、層の入力を参照した残差関数を学習する」ことで最適化しやすくなった。

◇ResNetの構造



求める写像を $H(x)$ 、入力 x との残差を $F(x)$ とすると $F(x) : H(x) - x$ で表すことができる元の画像は $F(x) + x$ で表すことができる。

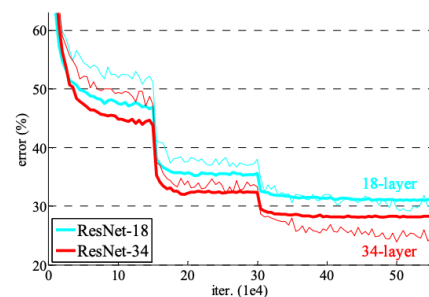
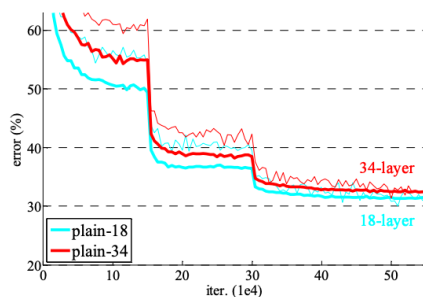
ショートカットさせているだけなのでパラメータが増えず計算が効率的に行える。

x 層を足すことで勾配消失が起こりづらい形式で学習ができる。

深いネットワークを構成する場合はブロックを3層にする。

実装による検証結果

ImageNetを使ったPlainとResNetでの比較。18層と34層で実験しており、より深い34層においてはResNetの性能が高いことがわかる。その他複数のデータで実験しているがいずれも高性能が出ている。



議論

- ・ 構造は現状が最良であるか。
- ・ 最良な最適化手法は。
- ・ 学習時間がかかりすぎる。

次に読むべき文献

- ・ Identity Mappings in Deep Residual Networks

Identity Mappings in Deep Residual Networks

(2016) 著者: Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun

概要

前述のDeep Residual Learning for Image Recognition の派生研究。

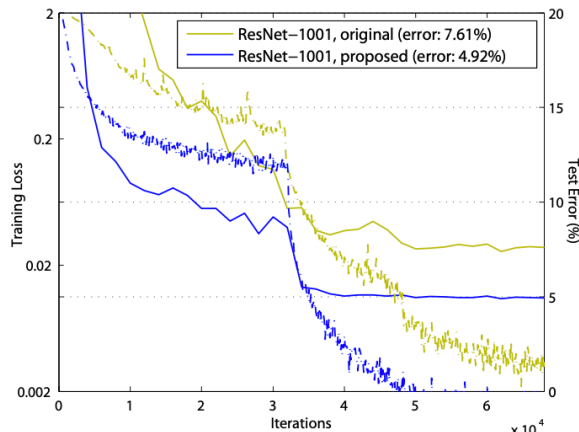
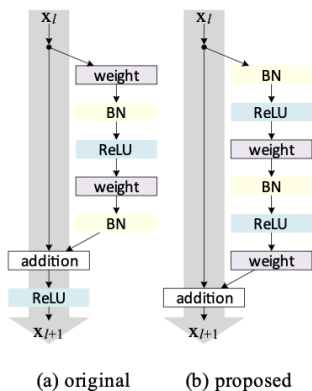
ResNetにおけるモデルの構造をどうすれば最善になるかの議論。

以下4つの考察を記載

考察①

ResNetブロック内の順番変更によるlossの減少

(a)のoriginalはlossが7.61%に対し、(b)proposedは4.92%と大幅に改善。



考察②

ResNetブロック内に様々なショートカットコネクートを適用。

ここでの結果は、ショートカットの乗算操作（スケーリング、ゲーティング、 1×1 畳み込み、ドロップアウト）は情報の伝播を妨げ、最適化の問題になる可能性があることが示されている。

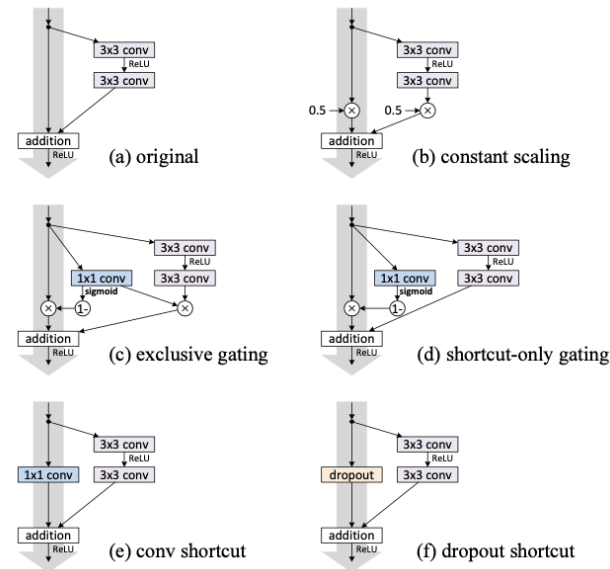


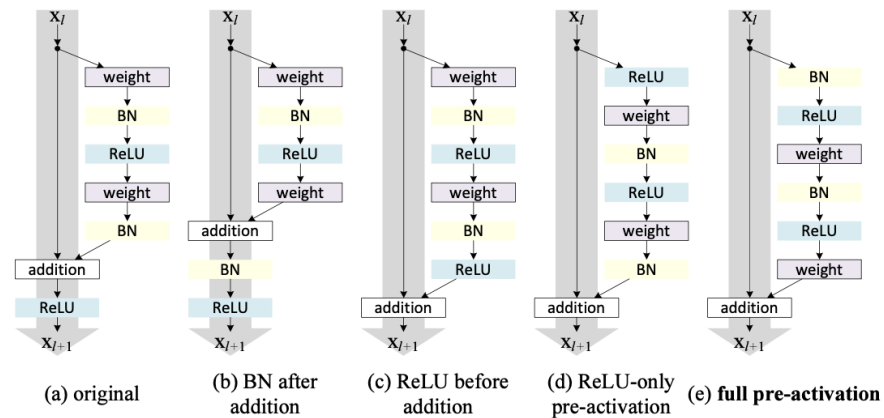
Table 1. Classification error on the CIFAR-10 test set using ResNet-110 [1], with different types of shortcut connections applied to all Residual Units. We report “fail” when the test error is higher than 20%.

case	Fig.	on shortcut	on \mathcal{F}	error (%)	remark
original [1]	Fig. 2(a)	1	1	6.61	
constant scaling	Fig. 2(b)	0	1	fail	This is a plain net
		0.5	1	fail	
		0.5	0.5	12.35	
exclusive gating	Fig. 2(c)	$1 - g(\mathbf{x})$	$g(\mathbf{x})$	fail	init $b_g = 0$ to -5
		$1 - g(\mathbf{x})$	$g(\mathbf{x})$	8.70	init $b_g = -6$
		$1 - g(\mathbf{x})$	$g(\mathbf{x})$	9.81	init $b_g = -7$
shortcut-only gating	Fig. 2(d)	$1 - g(\mathbf{x})$	1	12.86	init $b_g = 0$
		$1 - g(\mathbf{x})$	1	6.91	init $b_g = -6$
1×1 conv shortcut	Fig. 2(e)	1×1 conv	1	12.22	
dropout shortcut	Fig. 2(f)	dropout 0.5	1	fail	

考察③

いろいろな順番を試行

結果：BN→Relu→weight→BN→Relu→weight→ additionの順番が最善。



考察④

Residual Unitsの先頭にactivationを置くのが効率が良い

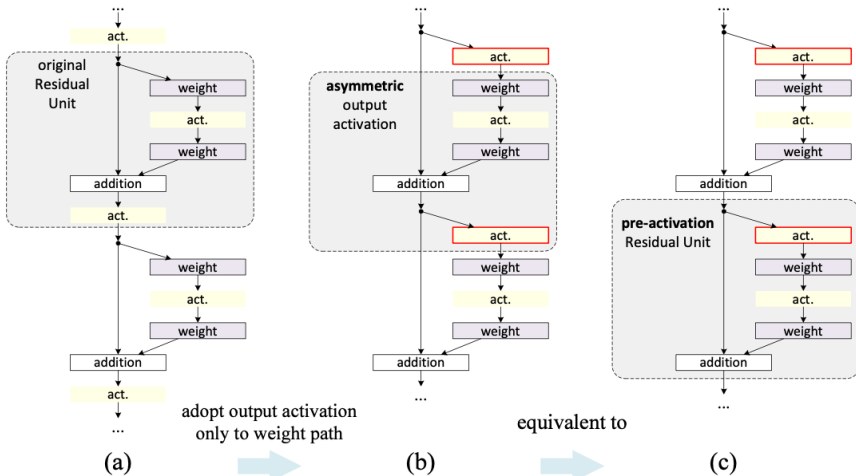


Table 2. Classification error (%) on the CIFAR-10 test set using different activation functions.

case	Fig.	ResNet-110	ResNet-164
original Residual Unit [1]	Fig. 4(a)	6.61	5.93
BN after addition	Fig. 4(b)	8.17	6.50
ReLU before addition	Fig. 4(c)	7.84	6.14
ReLU-only pre-activation	Fig. 4(d)	6.71	5.91
full pre-activation	Fig. 4(e)	6.37	5.46

Table 3. Classification error (%) on the CIFAR-10/100 test set using the original Residual Units and our pre-activation Residual Units.

dataset	network	baseline unit	pre-activation unit
CIFAR-10	ResNet-110 (1layer skip)	9.90	<u>8.91</u>
	ResNet-110	6.61	<u>6.37</u>
	ResNet-164	5.93	<u>5.46</u>
	ResNet-1001	7.61	<u>4.92</u>
CIFAR-100	ResNet-164	25.16	<u>24.33</u>
	ResNet-1001	27.82	<u>22.71</u>

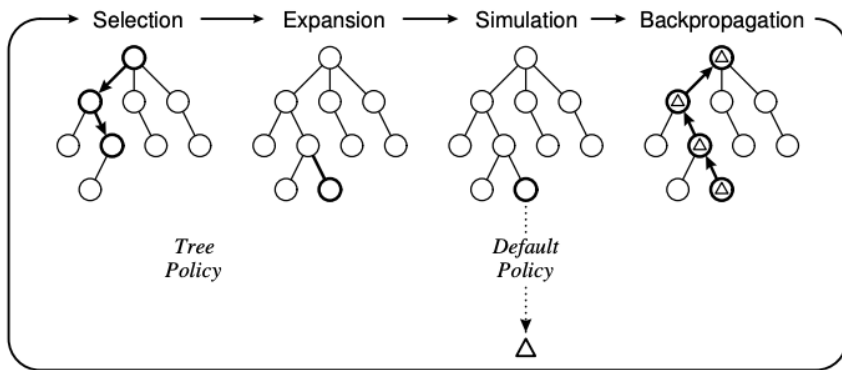
A Survey of Monte Carlo Tree Search Methods

(2012) 著者: Cameron Browne, Member, IEEE, Edward Powley, Member, IEEE, Daniel Whitehouse, Member, IEEE, Simon Lucas, Senior Member, IEEE, Peter I. Cowling, Member, IEEE, Philipp Rohlfschagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis and Simon Colton

概要

モンテカルロ木探索(以後MCTS)のまとめ論文。MCTSが考えられてから当時までの5年間の技術を記載している。

MCTSはツリー検索の精度とランダムサンプリングの一般性を組み合わせた検索方法で、AlphaGo〜MuZeroでも採用されている。



$$UCT = \bar{X}_j + 2C_p \sqrt{\frac{2 \ln n}{n_j}}$$

モンテカルロ木探索の手順説明

- ①選択: ルートノードから開始し、UCB1に従ったポリシーが再帰的に適用されリーフに達するまでツリーを降下。
- ②展開: 利用可能なアクションに従い、ツリーを展開するために一つ（もしくは複数）の子ノードが追加される。
- ③シミュレーション: 結果を生成するためにデフォルトポリシーに従い新しいノードからシミュレーションが実行される。
- ④バックプロパゲーション: シミュレーション結果は、選択したノードを介してバックプロパゲーションされて、統計量が更新される。

X_j → 平均報酬, n = これまでのプレイ総数, $n_j = j$ が選択された回数

MuZeroのアクション選択

$$a^k = \arg \max_a \left[Q(s, a) + P(s, a) \cdot \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)} \left(c_1 + \log \left(\frac{\sum_b N(s, b) + c_2 + 1}{c_2} \right) \right) \right]$$

モンテカルロ木探索により、過去は教師付きモデルを踏襲し過去の対戦データを学習したのちにプログラムが試行を開始、アクションを指定していたのに対し、人の関与を必要としない、かつ人間の対戦に夜教師データをも使用しないself_playを確立した手法。

Maximum_Entropy_Deep_Inverce_Reinforcement_Learning

(2016) 著者: Mobile Robotics Group, Department of Engineering Science, University of Oxford

概要

マルコフ決定過程における逆強化学習。

強化学習は報酬を決め、それに伴った生涯報酬を最大にする方策を選択するが、逆強化学習は目的までの方策（エキスパート）を示すことで最適な報酬関数を導き出すもの。人間や動物の行動を調べる際には、経験に基づいて報酬関数を考慮すべきとの考え。

→ エキスパートの報酬関数を模倣する。



$$P(\varsigma|r) \propto \exp\{\sum_{s,a \in \varsigma} r_{s,a}\}.$$

$$\mathcal{L}(\theta) = \log P(\mathcal{D}, \theta | r) = \underbrace{\log P(\mathcal{D} | r)}_{\mathcal{L}_{\mathcal{D}}} + \underbrace{\log P(\theta)}_{\mathcal{L}_{\theta}}.$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}_{\mathcal{D}}}{\partial \theta} + \frac{\partial \mathcal{L}_{\theta}}{\partial \theta}.$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{\mathcal{D}}}{\partial \theta} &= \frac{\partial \mathcal{L}_{\mathcal{D}}}{\partial r} \cdot \frac{\partial r}{\partial \theta} \\ &= (\mu_{\mathcal{D}} - \mathbb{E}[\mu]) \cdot \frac{\partial}{\partial \theta} g(f, \theta), \end{aligned}$$

- ・ニューラルネットワークを使った複雑な非線形報酬関数の近似

<アルゴリズム>

①初期化したNNにより報酬を推定→②経路分布を計算→③熟練者と②の誤差を計算→④誤差逆伝搬法によってパラメータを更新

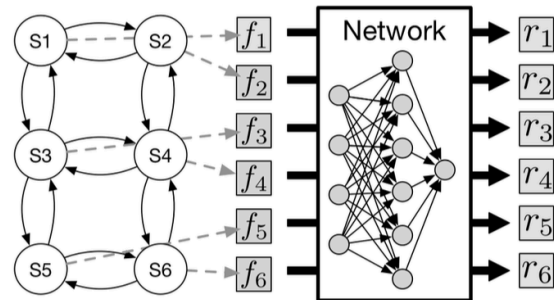


Figure 2: Schema for Neural Network based reward function approximation based on the feature representation of MDP states

- ・報酬関数は人が設定した特徴量の線形結合に限られていた。

→複雑な非線形な報酬関数でも推定可。

→高速化