

ニューラルアーキテクチャサーチの最適停止手法とその評価

坂本 航太郎 統計的機械学習研究センター 特任助教

概要

ニューラルアーキテクチャサーチ (Neural architecture search: NAS) はニューラルネットの構造設計を自動化するAutoMLの枠組みである。近年のNASは低コストでの構造探索に成功している一方で、精度と探索コストには依然としてトレードオフがあり、探索の停止のタイミングを決定することは重要な課題である。また、NASの代表の一つである Differentiable Architecture Search (DARTS) においては性能劣化が報告されており、早期停止による性能劣化を防ぐ手法が提案されている。本研究では構造パラメータに対する汎化誤差の期待値の差分を用いたNASの最適停止手法を提案する。複数のベンチマークデータセット、NAS-Bench-201, NATS-Bench, NAS-Bench-301, NAS-Bench-ASR, を用いて、本提案手法がNASを高精度に保ちつつ、探索コストを削減できることを評価する。

背景

Neural Architecture Search (NAS)

所定のデータセット・タスクに対して優れたパフォーマンスの得られるニューラルネットの構造を比較的短時間で導くことを目指す



NASの課題

- 膨大な探索コスト→エコなNAS
- 早期停止：コスト削減と性能改善

Xu, J., et al. KNAS: Green Neural Architecture Search. (2021) <https://proceedings.mlr.press/v139/xu21m.html>

Baker, B., et al. (2018) <https://openreview.net/forum?id=Hlnk3NivG>

Li & Talwalkar, (2019) <http://proceedings.mlr.press/v115/li20c.html>

NASの探索をいつどのように停止するか？

One-shot Neural Architecture Search (One-shot NAS)

ニューラルネットワークの構造とパラメータを一度の学習サイクルにおいて同時に最適化する枠組み

Differentiable Architecture Search (DARTS)

ニューラルネットワーク→有向非巡回グラフ (Directed Acyclic Graph : DAG)

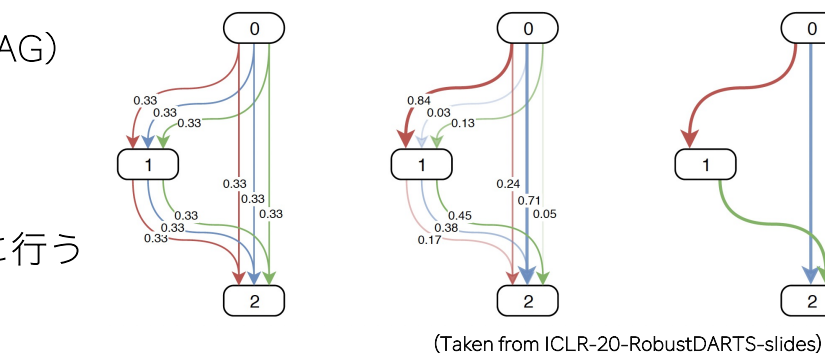
$$\text{Bi-level最適化問題} \quad \min_{\alpha} \{f(\alpha) \triangleq \mathcal{L}_{\text{valid}}(w^*(\alpha), \alpha)\}$$

$$s.t. \quad w^*(\alpha) = \arg \min_w \mathcal{L}_{\text{train}}(w, \alpha)$$

$w^*(\alpha) \approx w - \xi \nabla_w \mathcal{L}_{\text{train}}(w, \alpha)$ と近似し、以下の最適化を交互に行う

$$\nabla_w \mathcal{L}_{\text{train}}(w, \alpha)$$

$$\nabla_{\alpha} \mathcal{L}_{\text{valid}}(w - \xi \nabla_w \mathcal{L}_{\text{train}}(w, \alpha), \alpha)$$



$$x_j = \sum_{v=0}^{K-1} \frac{\exp(\alpha_v^{(i,j)})}{\sum_{v'=0}^{K-1} \exp(\alpha_{v'}^{(i,j)})} \cdot o_v(x_i) \quad o^{(i,j)} \in \arg \max_{v \in \mathcal{O}} \alpha_v^{(i,j)}$$

DARTSの課題

skip-connection (恒等写像)が多く含まれる構造が獲得されやすいなどの性能劣化 (ランダムサーチより悪化)

Yang, A., et al., Nas evaluation is frustratingly hard. (2020) <https://openreview.net/forum?id=HlygrdpVKvr>

既存手法：早期停止による解決法

- DARTS+**
単一のセル内でskip-connectionが2つ以上になったら探索を停める/構造パラメータが安定したら停める
- RobustDARTS**
損失地形のシャープさを表していると言われる検証損失のヘシアン¹の支配的な固有値をモニターし、過去5エポック分に対する最大固有値の平均値の比が0.75以下になった (固有値が大きく増加した) 時に停止する

導入

確率緩和に基づくOne-shot NAS : ASNG-NAS

- 構造をスーパーネットからサンプリングし、重み は構造間で共有
- カテゴリ変数空間で定義される確率分布族
- 確率分布のもとでの目的関数の期待値の最:

$$J(w, \theta) = \int_{c \in \mathcal{C}} \mathcal{L}(w, c) p_{\theta}(c) dc = \mathbb{E}_{p_{\theta}}[\mathcal{L}(w, c)]$$

① カテゴリカル分布から構造を選択するone-hotベクトルをサンプリング

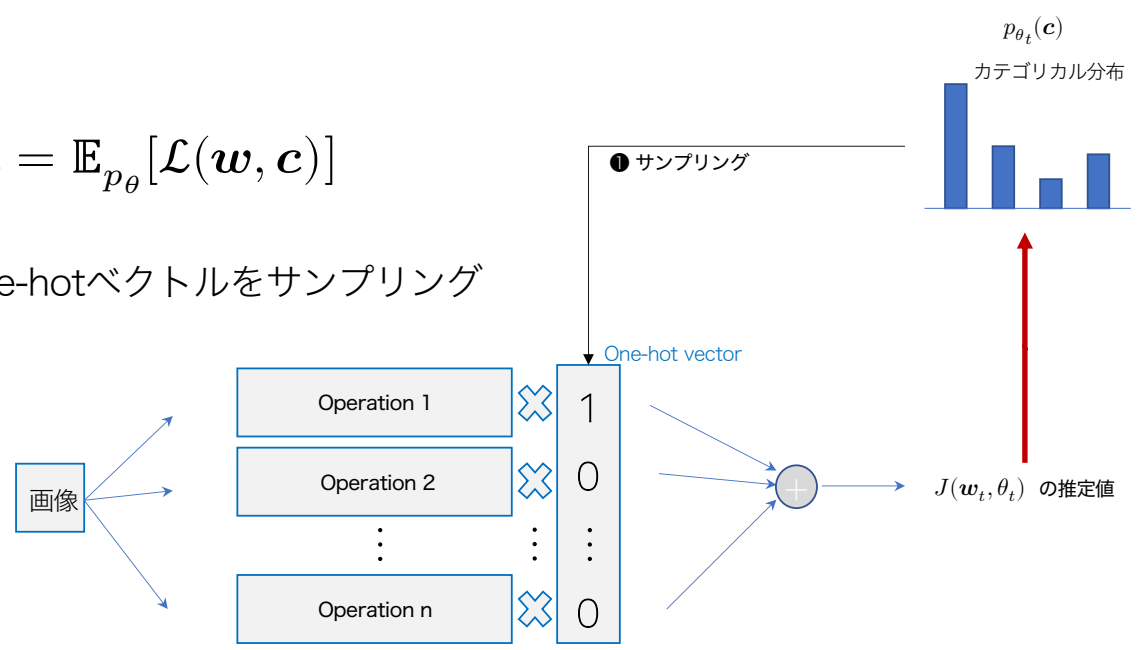
② サンプルした構造を順伝播&重み更新

$$w_{t+1} = w_t + \eta_w \nabla_w J(w_t, \theta_t)$$

③ 検証データを使って順伝播&分布更新

$$\theta_{t+1} = \theta_t + \eta_{\theta} \nabla_{\theta} J(w_{t+1}, \theta_t)$$

④ 学習率を更新



[1] Y. Akimoto, S. Shirakawa, N. Yoshinari, K. Uchida, S. Saito, and K. Nishida, "Adaptive stochastic natural gradient method for one-shot neural architecture search," ICMML 2019

提案法

停止基準の構成

- NASの目的は汎化誤差を最小化するニューラルネットの構造を獲得することなので汎化誤差の差分に注目して停止基準を構成することは自然である
- 探索ステップ前後での汎化誤差のハイパーパラメータの分布に対する期待値の差分は汎化誤差が探索時には得られないので直接計算できない
- そこでその上界を評価して閾値と比較することで停止タイミングを決定する
- 差分の上界を評価するにあたり、以下のPinskerの不等式を考える

定理 ([2], Fact 9). $p_{\theta_t}(c)$ と $p_{\theta_{t-1}}(c)$ を任意の確率分布とする。このとき任意の可測関数 $\mathcal{L}(c) \in [a, b]$ について、以下の不等式が成り立つ:

$$\mathbb{E}_{p_{\theta_t}(c)} [\mathcal{L}(w, c)] - \mathbb{E}_{p_{\theta_{t-1}}(c)} [\mathcal{L}(w, c)] \leq (b - a) \sqrt{\frac{1}{2} D_{KL}[p_{\theta_t}(c) \parallel p_{\theta_{t-1}}(c)]}$$

[2] D. Russo and B. V. Roy, "An information-theoretic analysis of Thompson sampling," Journal of Machine Learning Research, vol. 17, no. 68, pp. 1–30, 2016.

- NASにおいて重みパラメータは更新されるので定理は単純適用ができない
- そこでテ일러展開によりその影響を近似する方針をとる

$J(w, \theta)$ を w について w_t まわりで n 次のテ일러展開したときの値の値を $T_n(w, \theta)$,

剰余項 $R_n(w, \theta) := J(w, \theta) - T_n(w, \theta)$ を $T_n(w, \theta)$ から 0 次の項を取り除いたものを $T'_n(w, \theta)$ とする。

$$J(w_t, \theta_t) - J(w_{t-1}, \theta_{t-1}) = T_n(w_t, \theta_t) + R_n(w_t, \theta_t) - J(w_{t-1}, \theta_{t-1})$$

$$= J(w_{t-1}, \theta_t) - J(w_{t-1}, \theta_{t-1}) + T'_n(w_t, \theta_t) + R_n(w_t, \theta_t)$$

上式の第1項、第2項に定理を適用し以下の不等式を得る：

$$J(w_t, \theta_t) - J(w_{t-1}, \theta_{t-1}) \leq (b - a) \sqrt{\frac{1}{2} D_{KL}[p_{\theta_t}(c) \parallel p_{\theta_{t-1}}(c)]} + T'_n(w_t, \theta_t) + R_n(w_t, \theta_t)$$

J を w について C^n 級の関数とし、 $w \in [w_{t-1}, w_{t-1} + d]$ かつ $\|w_t - w_{t-1}\| \leq d$ の全てにおいて $J(w, \theta) \leq M$ と仮定したとき、剰余項 $R_n(w, \theta)$ に対してテ일러の不等式を適用することで以下の不等式が成り立つ：

$$\mathbb{E}_{p_{\theta_t}(c)} [\mathcal{L}(w_t, c)] - \mathbb{E}_{p_{\theta_{t-1}}(c)} [\mathcal{L}(w_{t-1}, c)]$$

$$\leq (b - a) \sqrt{\frac{1}{2} D_{KL}[p_{\theta_t}(c) \parallel p_{\theta_{t-1}}(c)]} + T'_n(w_t, \theta_t) + \frac{M}{n!} \|w_t - w_{t-1}\|^n$$

$J(w, \theta)$ を w について w_t まわりで n 次のテ일러展開したときの値から 0 次の項を取り除いたもの

$M \geq J(w, \theta)$

上記の不等式を停止基準として採用するためには \mathcal{L} の範囲 $[a, b]$ と M を予め決めないといけないが $M = (b - a)$ と仮定でき、0 次の項でテ일러展開したときは $(b - a)$ を消去できる。ただし、高次展開では、 $(b - a)$ は残るので、目的関数の上界を決めておく必要がある。

$$r_t^1 = (b - a) \sqrt{\frac{1}{2} D_{KL}[p_{\theta_t}(c) \parallel p_{\theta_{t-1}}(c)]} + \nabla_w J(w_{t-1}, \theta_t) \cdot (w_t - w_{t-1}) + (b - a) \|w_t - w_{t-1}\|$$

$$r_t^2 = (b - a) \sqrt{\frac{1}{2} D_{KL}[p_{\theta_t}(c) \parallel p_{\theta_{t-1}}(c)]} + \nabla_w J(w_{t-1}, \theta_t) \cdot (w_t - w_{t-1}) + \frac{1}{2} (w_t - w_{t-1})^T \nabla_w^2 J(w_{t-1}, \theta_t) (w_t - w_{t-1}) + \frac{(b - a)}{2} \|w_t - w_{t-1}\|^2$$

- 誤差関数の有界性が求められる
→ bi-tempered logistic loss を採用する
→ より一般的な交差エントロピー等も仮定することで使用可能

停止基準の自動決定

- 上界が閾値を下回ったら構造探索を停止する
- 検証データが利用可能であるので検証データを用いて停止基準を自動決定することを考えたい
- Makarov, et al. (2021)は予測モデルのハイパーパラメータのベイズ最適化における過剰適合の解決策として、simple regret の上界が汎化誤差の標準偏差を下回ったときに訓練を早期停止することを提案している
- 汎化誤差の標準偏差は k-fold 交差検証、訓練及び検証データから以下の式により推定可能である

$$\lambda_t = \sqrt{\frac{1}{k} + \frac{|S_{\text{val}}|}{|S_{\text{train}}|}} \hat{s}_t$$

- 提案法の停止基準は汎化誤差の期待値の差分をもとにしており、汎化誤差の標準偏差の推定値と比較することは自然である
- そこで汎化誤差の標準偏差の推定値を閾値として採用する

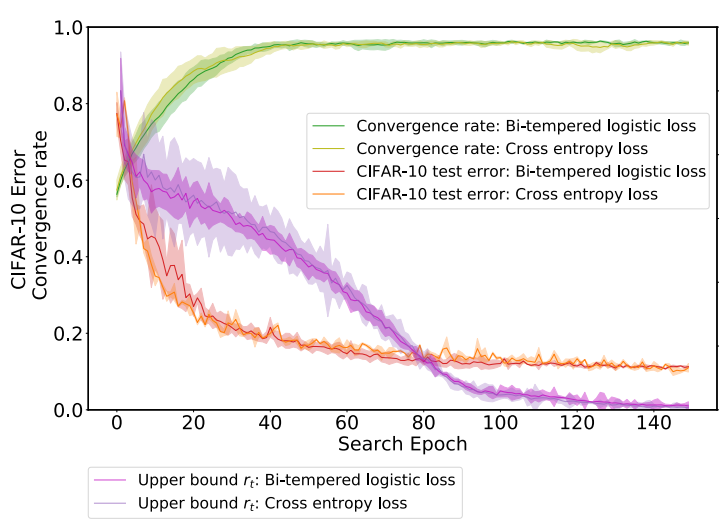
DARTSの停止基準への応用

- 提案法はASNG-NASの設定を停止基準の定式化に用いており、よりポピュラーなone-shot NASの手法であるDARTSへ直接適用はできない
- DARTSでは(1)最適化対象が汎化誤差の期待値ではない(2)構造パラメータは確率緩和をベースにしたものではない(連続緩和)
- そこでオペレーション候補の選択を担うsoftmaxを確率分布と見做して探索エポック前後の変化を評価することで停止基準として用いる方針を取る
- 停止基準に用いるはsoftmaxにより計算される事後分布である
- 閾値の自動決定についてはone-hot化の離散化を行った後に5-fold交差検証によって計算する

評価実験

ASNG-NASの最適停止

パラメータ θ の収束率との比較



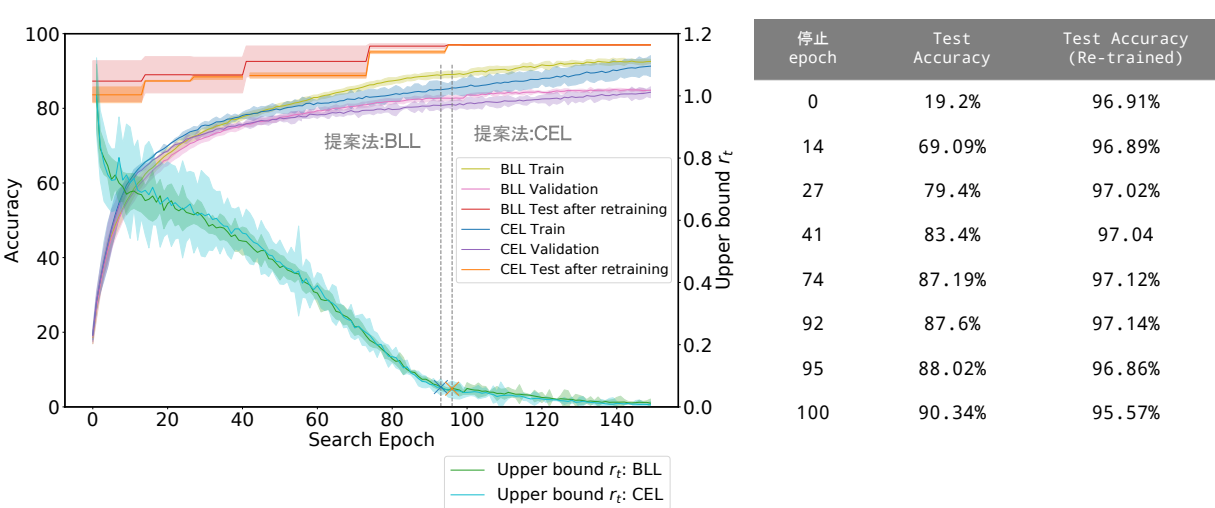
パラメータ θ の収束率

$$\frac{1}{d} \sum_{i=1}^d \max_j [\theta]_{i,j}$$

i : cell index

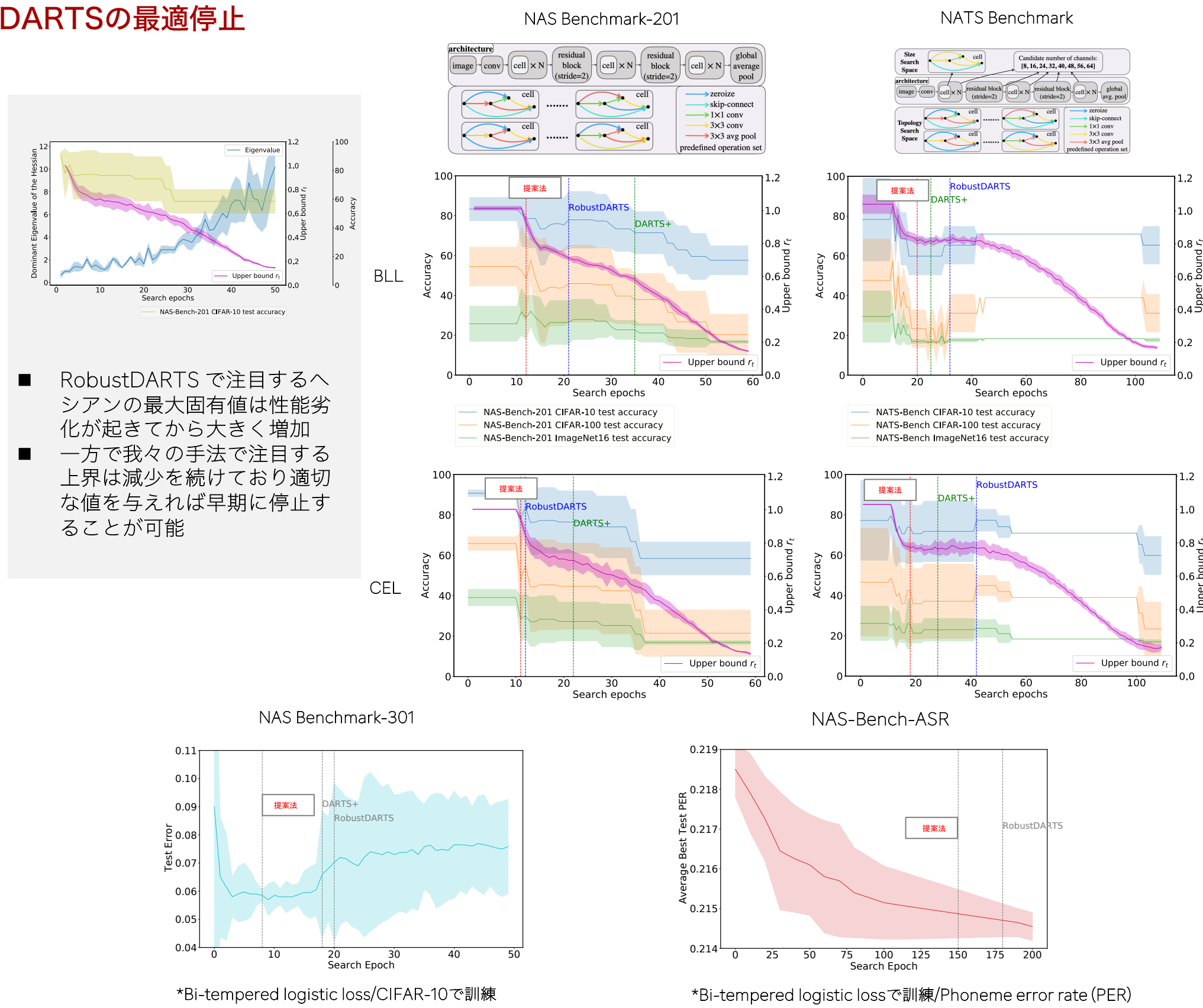
j : operation index

上界の評価値と停止タイミング



- CIFAR-10 のtrain/validation/test 精度と提案法の停止基準による停止タイミングの一つを示している
- テスト精度は停止タイミングの後も増加しているが構造探索による性能向上はこれ以上望めないと考えられる

DARTSの最適停止



- 既存の早期停止手法であるRobustDARTS及びDARTS+に比べ提案法は非常に早期に停止できている
- DARTSにおいては冗長な探索が行われていたことを示唆している

謝辞

九州工業大学石橋英朗氏、横浜国立大学白川真一氏、筑波大学秋本洋平氏、佐藤怜氏、統計数理研究所日野英逸氏との共同研究です。本研究は国立研究開発法人新エネルギー・産業技術総合開発機構 (NEDO) の委託業務 (JPNP18002) の支援を受けました。