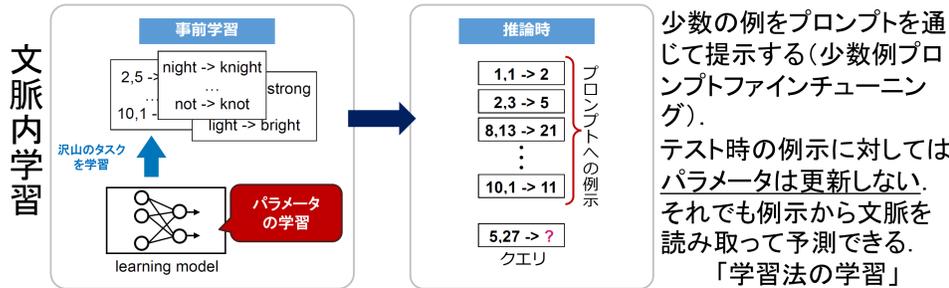


In-context learning (文脈内学習) の理論

文脈内学習の理論を展開:

なぜTransformerは文脈内学習を実現できるのか?

- 1. 統計的学習理論: Transformerは文脈内学習でミニマックス最適な誤差を達成
- 2. 最適化理論: 事前学習における勾配法により適切な特徴を学習可能→サンプル複雑度の改善



問題設定: $y_{i,t} = F_t^\circ(x_{i,t}) + \epsilon_{i,t}$ ($i = 1, \dots, n$)
 $t = 1, \dots, T$: 訓練タスクのインデックス

ミニマックス最適性

[Kim, Nakamaki, Suzuki: Transformers are Minimax Optimal Nonparametric In-Context Learners. **NeurIPS2024**]

真の関数

$$F_t^\circ(x) = \beta_t^\top f^\circ(x)$$

仮定: $\text{Cov}(\beta) = \Sigma = \text{diag}((k^{-2s-1})_{k=1}^\infty)$
 (関数空間の複雑さ)

- 事前学習 (pretraining): 特徴量 (表現) を学習 $[f^\circ]$
 > 例: Fourier, B-Spline
 > 文脈 (t) に非依存
 > データを表現する「最も効率的」な基底を学習
- 文脈内学習 (in-context): 係数を学習 $[\beta_t]$
 > 文脈 (t) に依存
 > 例示から現在の文脈 β_t を推定
 → Attentionで獲得



定理 (推定誤差): $\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim N^{-2s} + N^2 \delta_N^4 + N^{2r+1} \delta_N^2$ 特徴量の近似誤差
 $+ \frac{N}{n} \log(N) + \frac{N^{2r}}{n^2} \log^2(N)$ 文脈 (β_t) を推定する誤差
 $+ \frac{1}{T} (N^2 \log(\epsilon^{-1}) + \log(N(\frac{\epsilon}{\sqrt{N}}, \mathcal{F}_N, \|\cdot\|_\infty))) + \epsilon$ 基底(特徴量)の学習誤差

例 (B-spline基底: f_j^* がB-spline→Besov/Sobolev空間):
 $\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim n^{-\frac{2s}{2s+1}} + \frac{n^{\frac{2}{2s+1}}}{T}$

「文脈 (β_t) 」の推定誤差 「表現 (f°) 」の推定誤差
 ミニマックス最適

Tが小さい: 記憶中の状態
 Tが大きい: 記憶が完了し汎化できる状況
 → テスト時推論のスケールアップ

定理 (ミニマックス下限):
 $\inf_{\hat{\phi}} \sup_{f \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{\phi})] \gtrsim n^{-\frac{2s}{2s+1}} + \frac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT}$
 ただし, $\epsilon_{1,n} \approx \frac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT}$

直感: 十分多くの問題例を事前に見ておけば, テスト時に最適な誤差で予測できる. 事前知識が足りなければテスト時に失敗する.

平均場NNによる非線形特徴学習の大域的最適性

[Kim, Suzuki: Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape. **ICML2024, oral presentation**]

非線形特徴学習 (ϕ) に2層NN(平均場NN)を考える:
 $h_\mu(x) = \int h_\theta(x) d\mu(\theta) \in \mathbb{R}^k$
 $h_\theta(x) = \alpha \sigma(\mathbf{w}^\top x)$ ($\theta = (\mathbf{a}, \mathbf{w}) \in \mathbb{R}^k \times \mathbb{R}^d$)

損失関数 $\mathcal{L}(\mu, \Gamma) := \mathbb{E}_{x_{q,r}} [\|f^\circ(x_{q,r}) - \mathbb{E}_x [f^\circ(x) h_\mu(x)^\top] \Gamma h_\mu(x_{q,r})\|^2]$

定理 (損失関数の性質):
 損失関数の局所最適解はすべて大域的最適解:
 (1) 停留点は鞍点か大域的最適解.
 (2) 鞍点には負の曲率方向があり, 指数関数的速度で鞍点から脱出可能.

Koopman作用素の理論によるDNNの汎化誤差解析

[Hashimoto, Sonoda, Ishikawa, Nitanda, Suzuki: Koopman-based generalization bound: New aspect for full-rank weights. **ICLR2024**]

DNNの一層分をKoopman作用素とみなし, Koopman作用素の作用素ノルムを用いて Rademacher複雑度を抑えた.

Rademacher complexity $\leq O\left(\frac{\|g\|_{H_L}}{\sqrt{n}} \prod_{j=1}^L \frac{G_j E_j \|W_j\|^{s_j-1}}{\det(W_j^* W_j)^{1/4}}\right)$

これまでの結果とは異なり, 中間層の行列の最小特異値が「大きい方が良い」バウンドを与えた.

勾配法による特徴学習の最適性

[Jason D. Lee, Kazusato Oko, Taiji Suzuki, Denny Wu: Neural network learns low-dimensional polynomials with SGD near the information-theoretic limit. **NeurIPS2024**]

シングルインデックスモデル:

$$f_*(x) = \sigma_*(\langle x, \theta \rangle) \quad (x \sim N(0, I_d))$$

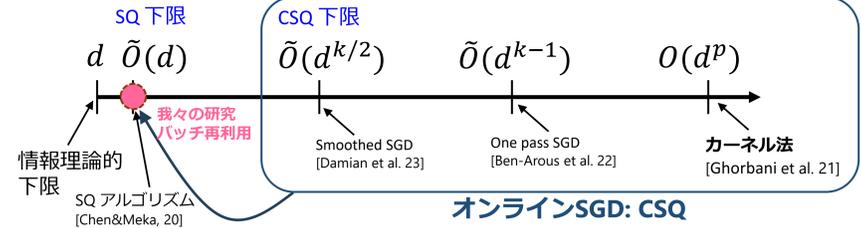
> 方向 $\theta \in \mathbb{R}^d$ とリンク関数 σ_* を推定
 → 特徴学習が重要
 $\sigma_*(z) = \sum_{i=k}^p \alpha_i^* \text{He}_i(z)$
 > σ_* は次数 p かつ情報指数 k .
 He_i : 次数 i のエルミート多項式

Q: ニューラルネットワークはどれくらいの計算/サンプル-複雑度で学習できるか?

定義 (情報指数 k [Ben-Arous et al. 2021])
 σ_* の情報指数は非ゼロ係数 α_i^* をもつ最小次数である:
 $k := \text{IE}(\sigma_*) = \min\{i \in \mathbb{N} \mid \alpha_i^* \neq 0\}$

情報指数は問題の難しさを計る指標
 $[k$ が大きい \Rightarrow 高次の情報を取り出す必要がある \Rightarrow より多くのデータが必要]

サンプル複雑度の比較



- 通常のオンライン確率的勾配降下法 (SGD) はCSQアルゴリズムに含まれる.
- SGDを用いた学習法のサンプル複雑度の下限は $\Omega(d^{k/2})$ である. 一方で, 情報理論的下限は $O(d)$ である. \Rightarrow 通常のSGDは情報理論的下限を達成できない.

Q: 通常のSGDを少し修正することでこのギャップを埋めることは可能か?
 我々の結果: ミニバッチを二回使いまわす(再利用する)だけで $O(d^{k/2})$ を $\tilde{O}(d)$ に改善できる.

定理 (バッチ再利用SGDのサンプル複雑度)
 バッチ再利用型のSGDは
 $n = \tilde{O}(d\epsilon^{-2}), N = \tilde{O}(\epsilon^{-1})$
 だけのサンプル複雑度 (更新回数) で平均二乗誤差を ϵ 以下にできる:
 $\mathbb{E}_x [\|f_*(x) - f_{\hat{\theta}}(x)\|^2] \leq \epsilon$

バッチ再利用だけで情報理論的下限を達成できる.
 非線形活性化関数を使うことが証明の鍵. 非線形活性化関数は最適化の効率化にも寄与している.

Correlation statistical query (CSQ): Algorithm has an access to "noisy" information about correlation between ϕ and y . (τ is the noise level. Noise can be adversarial)
 $|\hat{q} - \mathbb{E}_{x,y}[\phi(x)y]| \leq \tau \sim n^{-1/2}$

Statistical query (SQ): General query
 $|\hat{q} - \mathbb{E}_{x,y}[\phi(x,y)]| \leq \tau \sim n^{-1/2}$ (with $|\phi| \leq 1$)

※オンラインSGDはCSQアルゴリズム:
 $\nabla_w \mathbb{E}_{x,y}[(y - f_w(x))^2] \propto -\mathbb{E}_{x,y}[y \nabla_w f_w(x)] + \mathbb{E}_x[f_w(x) \nabla_w f_w(x)]$
 ※データを二回使いまわすことで, CSQから外れる.

勾配法による特徴学習: 事前学習による学習精度改善

[Oko, Song, Suzuki, Wu: Transformer efficiently learns low-dimensional functions in context. **NeurIPS2024**]

$$F_t^\circ(x) = \sigma_*^t(\langle x, \beta_t \rangle): \text{タスク } t \text{ の関数}$$

リンク関数 σ_*^t と方向 β_t がランダムに生成される設定 (文脈から推定する)

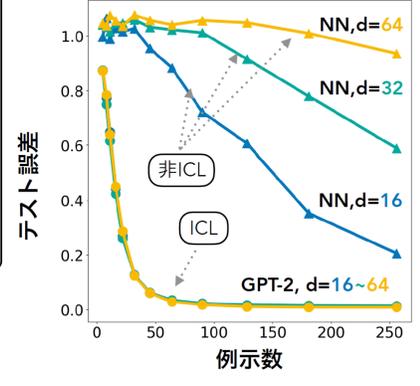
- β_t は d 次元空間のうち r 次元部分空間 S に分布 ($r \ll d$):
 $\beta_t \sim \text{Unif}(\text{Unit}(S)) \quad \dim(S) = r \ll d$
- $\sigma_*^t(z) = \sum_{i=k}^p c_i^t \text{He}_i(z)$
 c_i^t はタスクごとにランダムに生成されているとする.

$\Rightarrow \beta_t$ の部分空間を事前学習で取得

定理 (SGDによる学習の複雑度)

- 事前学習で用いるデータサイズ:
 $T = \Theta(d^{2k})$ and $n = \tilde{\Omega}(d^k)$.
 - 一層目のNNの横幅: m
 - テスト時の例示数: n^*
- SGD (勾配法) によって多項式時間で学習可能で, 以下の予測誤差を達成:
 $\mathcal{L}(\hat{W}, \hat{b}, \hat{\Gamma}) \lesssim \frac{r^{3P/2}}{\sqrt{m}} + \sqrt{r^{5P}} \sqrt{\frac{r^2}{n} + \frac{r^2}{n^*}}$
- 近似誤差 汎化誤差

必要なテスト時例示数が次元 d に依存しない.
 > 勾配法により不要な特徴量が削れている効果
 > 事前学習なしでは達成不可能



分布のミニマックス最適化問題

[Kim, Yamamoto, Oko, Yang, Suzuki: Symmetric Mean-field Langevin Dynamics for Distributional Minimax Problems, **ICLR2024**]

$$\min_{\mu} \max_{\nu} \mathcal{L}(\mu, \nu) + \lambda \text{Ent}(\mu) - \lambda \text{Ent}(\nu) \quad \mu, \nu: \text{確率分布}$$

応用:
 • 敵対的学習
 • ロバスト学習
 • 強化学習

貢献: 確率分布に対する双対平均加法を提案し, 提案手法が $O(1/K)$ で収束することを示した.
 [離散時間/分布の有限粒子近似誤差も導出]