

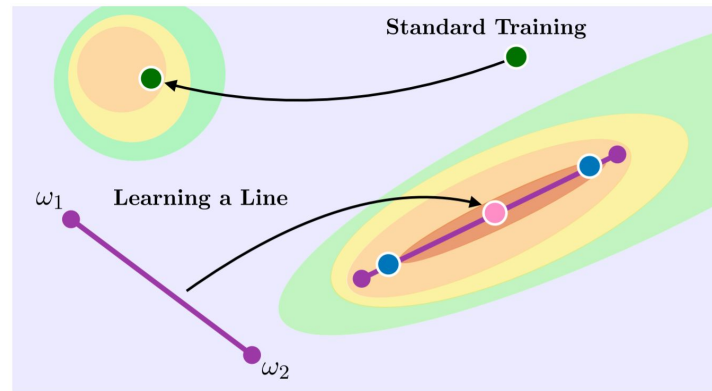
# Learning Neural Network Subspaces

Исаев Сергей  
Семерова Елена  
Ким Михаил

# Постановка задачи

При оптимизации нейросети хотим находить не точку (хороший набор весов), а подпространство весов (отрезок, кривая, симплекс), в котором выполняются полезные свойства:

- Модель из центра подпространства будет работать лучше с точки зрения метрик, калибровки и устойчивости, чем при обычном обучении.
- Из моделей, полученных из подпространства, можно сделать хороший ансамбль.



# Немного определений и наблюдений

Ансамбль моделей:  $\hat{y} = \frac{1}{2} (f(\mathbf{x}, \theta_T^1) + f(\mathbf{x}, \theta_T^2))$  — как правило улучшает метрики, калибровку и устойчивость модели.

Ансамбль весов:  $f(\mathbf{x}, \frac{1}{2}(\theta_T^1 + \theta_T^2))$  — как правило так себе идея (но мы будем искать такие веса, чтобы это правило ломалось)

Коннектор это такая непрерывная функция  $P : [0, 1] \rightarrow \mathbb{R}^n$ , возвращающая веса модели  $P(0) = \psi_1$ ,  $P(1) = \psi_2$ , для которой выполняется свойство:

$$\inf_{\alpha \in [0, 1]} \text{Acc}(P(\alpha)) \gtrsim \frac{1}{2} (\text{Acc}(\psi_1) + \text{Acc}(\psi_2))$$

Что утверждает свойство про ансамбль весов на языке коннекторов?

# Ещё немного наблюдений и определений

**Наблюдение:** линейный коннектор всё-таки можно построить, например, если уже немного немного предобученную модель начать обучать по разными траекториям. Например, за счет рандомизации семплирования примеров в батче.

**Определение:** коннектор можно обобщить и на  $m$ -мерное пространство. Введём  $m$ -мерное пространство для коннектора:

$$\Delta^{m-1} = \{\alpha \in \mathbb{R}^m : \sum_i \alpha_i = 1, \alpha_i \geq 0\}$$

А  $e_i$  — стандартный базисный вектор в  $\Delta^{m-1}$  из 0 и 1 на позиции  $i$ .

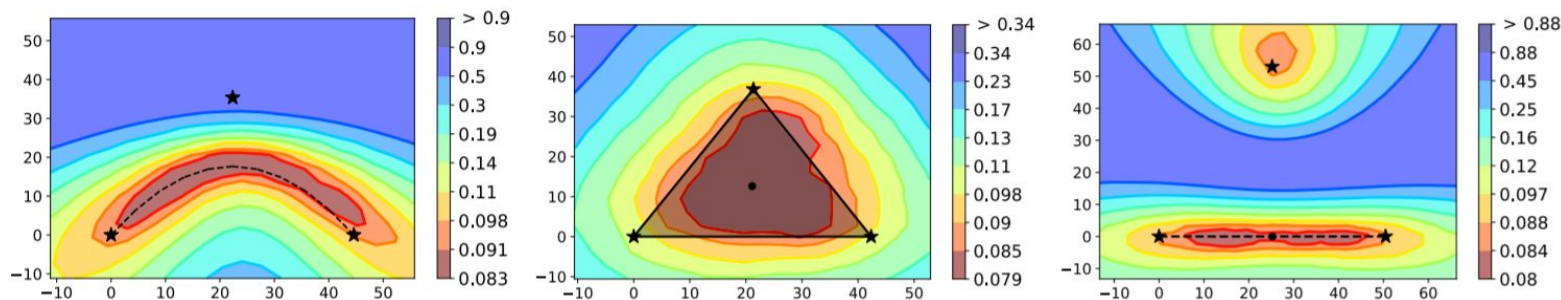
# Последний слайд с определениями, честно

Тогда  $m$ -мерный коннектор на весах модели  $\psi_1, \dots, \psi_m \in \mathbb{R}^n$  — это функция  $P : \Delta^{m-1} \rightarrow \mathbb{R}^n$ , что  $P(\mathbf{e}_i) = \psi_i$  и

$$\inf_{\alpha \in \Delta^{m-1}} \text{Acc}(P(\alpha)) \gtrsim \frac{1}{m} \sum_{i=1}^m \text{Acc}(\psi_i)$$

В статье описывается эффективный алгоритм поиска коннекторов вида

$P(\alpha) = \sum_i \alpha_i \psi_i$  (симплексы) и одномерных кривых Безье.



## А теперь формальная постановка задачи

Пусть определён “коннектор”  $P(\cdot, \{\omega_i\}_{i=1}^m) : \Lambda \rightarrow \mathbb{R}^n$

Координаты для коннектора семплируются из распределения  $\mathcal{U}(\Lambda)$ .

Примеры для модели семплируются из распределения  $\mathcal{D}$ .

Тогда мы пытаемся минимизировать такое матожидание:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[ \mathbb{E}_{\boldsymbol{\alpha} \sim \mathcal{U}(\Lambda)} [\ell(f(\mathbf{x}, P(\boldsymbol{\alpha}, \{\omega_i\}_{i=1}^m)), \mathbf{y}))] \right]$$

---

**Algorithm 1** TrainSubspace

---

**Input:**  $P$  with domain  $\Lambda$  and parameters  $\{\omega_i\}_{i=1}^m$ , network  $f$ , train set  $\mathcal{S}$ , loss  $\ell$  (e.g. a line has

$\Lambda = [0, 1]$  and  $P(\alpha; \omega_1, \omega_2) = (1 - \alpha)\omega_1 + \alpha\omega_2$ ).

Initialize each  $\omega_i$  independently.

**for** batch  $(\mathbf{x}, \mathbf{y}) \subseteq \mathcal{S}$  **do**

    Sample  $\alpha$  uniformly from  $\Lambda$ .

$\theta \leftarrow P(\alpha; \{\omega_i\}_{i=1}^m)$

$\hat{\mathbf{y}} \leftarrow f(\mathbf{x}, \theta)$

$\mathcal{L} \leftarrow \ell(\hat{\mathbf{y}}, \mathbf{y})$

    Backprop  $\mathcal{L}$  to each  $\omega_i$  and update with SGD & mo-

    mentum using estimate  $\frac{\partial \mathcal{L}}{\partial \omega_i} = \frac{\partial \ell}{\partial \theta} \frac{\partial P}{\partial \omega_i}$

**end for**

---

# Алгоритм плохой. Нужна регуляризация!

К сожалению, при таком подходе обучаемое подпространство может схлопнуться в одну точку или его точки могут быть недостаточно ортогональными для хорошего ансамбля.

Сделаем регуляризацию, которая будет давать штраф за недостаточно ортогональные точки.

$$\beta \cdot \mathbb{E}_{j \neq k} [\cos^2(\omega_j, \omega_k)] = \beta \cdot \mathbb{E}_{j \neq k} \left[ \frac{\langle \omega_j, \omega_k \rangle^2}{\|\omega_j\|_2^2 \|\omega_k\|_2^2} \right]$$



---

**Algorithm 1** TrainSubspace

---

**Input:**  $P$  with domain  $\Lambda$  and parameters  $\{\omega_i\}_{i=1}^m$ , network  $f$ , train set  $\mathcal{S}$ , loss  $\ell$ , and scalar  $\beta$  (e.g. a line has  $\Lambda = [0, 1]$  and  $P(\alpha; \omega_1, \omega_2) = (1 - \alpha)\omega_1 + \alpha\omega_2$ ).

Initialize each  $\omega_i$  independently.

**for** batch  $(\mathbf{x}, \mathbf{y}) \subseteq \mathcal{S}$  **do**

    Sample  $\alpha$  uniformly from  $\Lambda$ .

$\theta \leftarrow P(\alpha; \{\omega_i\}_{i=1}^m)$

$\hat{\mathbf{y}} \leftarrow f(\mathbf{x}, \theta)$

    Sample  $j, k$  from  $\{1, \dots, m\}$  without replacement.

$\mathcal{L} \leftarrow \ell(\hat{\mathbf{y}}, \mathbf{y}) + \beta \cos^2(\omega_j, \omega_k)$

    Backprop  $\mathcal{L}$  to each  $\omega_i$  and update with SGD & momentum using estimate  $\frac{\partial \mathcal{L}}{\partial \omega_i} = \frac{\partial \ell}{\partial \theta} \frac{\partial P}{\partial \omega_i} + \beta \frac{\partial \cos^2(\omega_j, \omega_k)}{\partial \omega_i}$ .

**end for**

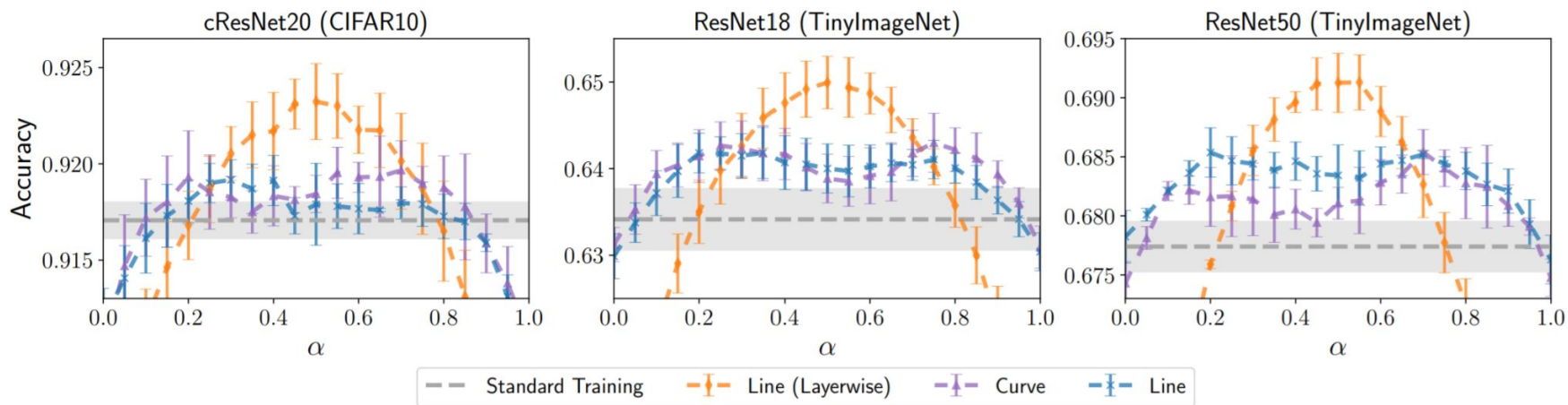
---

# Развитие подхода

Ранее мы считали веса модели просто столбцами.

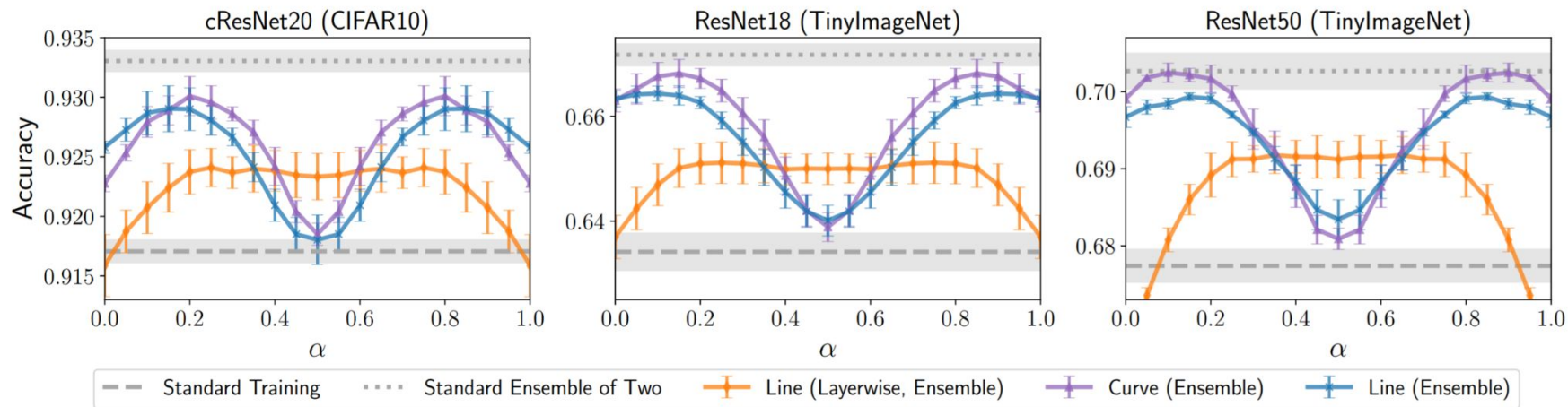
Авторы утверждают, что такой подход для обучения можно использовать отдельно для каждого слоя.

# Эксперименты



Исследование качества “средней” модели

# Эксперименты



Исследование построения ансамбля.

Тут первая модель берётся в точке  $P(\alpha)$ , а вторая в  $P(1-\alpha)$ .

# Рецензия

Содержание: предложен и эмпирически изучен метод обучения подпространств в пространстве моделей.

Достоинства:

- Хороший обзор релевантной литературы
- Понятное изложение идеи, информативные графики
- Работы по изучению loss surface важны для всей области
- Доступен понятный код

Недостатки:

- Больше proof-of-concept, чем реальный метод
- Одна задача для тестирования
- Эмпирическая работа, нет теоретических обоснований

Оценка: 8. Уверенность: 4.

# Исследование

- Статья была выпущена в 2021 году
- Первая версия на arxiv.org – февраль 2021 года
- Первый коммит на GitHub – март 2021 года
- Представлена на ICML 2021 в формате Poster



# ICML

International Conference  
On Machine Learning

Tue Jul 20 09:00 PM -- 11:00 PM (PDT)

Poster

## ***Learning Neural Network Subspaces***

Mitchell Wortsman · Maxwell Horton · Carlos Guestrin · Ali Farhadi · Mohammad Rastegari

In Poster Session 2

[\[ Visit Poster at Spot C1 in Virtual World \]](#)

 Mitchell Wortsman »

 Maxwell Horton »

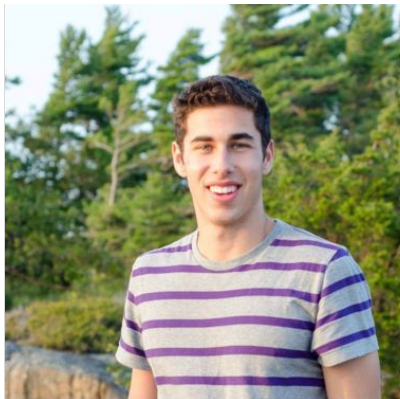
 Carlos Guestrin »

 Ali Farhadi »

 Mohammad Rastegari »

Recent observations have advanced our understanding of the neural network optimization landscape, revealing the existence of (1) paths of high accuracy containing diverse solutions and (2) wider minima offering improved performance. Previous methods observing diverse paths require multiple training runs. In contrast we aim to leverage both property (1) and (2) with a single method and in a single training run. With a similar computational cost as training one model, we learn lines, curves, and simplexes of high-accuracy neural networks. These neural network subspaces contain diverse solutions that can be ensembled, approaching the ensemble performance of independently trained networks without the training cost. Moreover, using the subspace midpoint boosts accuracy, calibration, and robustness to label noise, outperforming Stochastic Weight Averaging.

# Исследование



**Mitchell Wortsman**

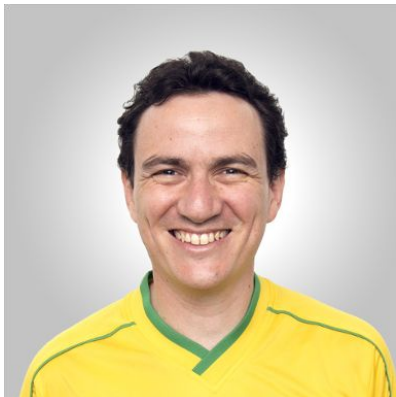
- BA Прикладная математика и компьютерные науки в Brown University in Providence, RI
- PhD студент в University of Washington
- Статья во время стажировки в Apple



**Maxwell Horton**

- BA Компьютерные науки и физика в The California Institute of Technology
- PhD в University of Washington
- Является сотрудником Apple
- Основная область -- Computer Vision

# Исследование



**Carlos Guestrin**

- Senior Director of AI and Machine Learning в Apple
- Профессор в University of Washington и Stanford University
- Один из соавторов XGBoost



**Ali Farhadi**

- Науч.рук. Митчелла и Максвелла
- Исследователь AI и ML в Apple
- Профессор в University of Washington
- Один из соавторов YOLO



**Mohammad Rastegari**

- Senior Technival Leader of AI and Machine Learning в Apple
- Профессор в University of Washington
- Ментор Митчелла в AI2



# Исследование



Frankle et. al, 2020 [1], [2]; Garipov et al., 2018 [3]; Draxler et al., 2018 [4], Fort et al., 2019/20 [5], [6], [7]; Izmailov et al., 2018 [8].

# Исследование

## Цитирования:

- **[1] LCS: Learning Compressible Subspaces for Adaptive Network Compression at Inference Time**  
Elvis Nunez, Maxwell Horton, Anish Prabhu, Anurag Ranjan, Ali Farhadi, Mohammad Rastegari; 2021.
- **[2] Subspace Learning for Personalized Federated Optimization**  
Seok-Ju Hahn, Minwoo Jeong, Junghye Lee; 2021.
- **[3] Exploring the Power of Lightweight YOLOv4**  
Chien-Yao Wang, Hong-Yuan Mark Liao, I-Hau Yeh, Yung-Yu Chuang, Youn-Long Lin; ICCV, 2021.
- **[4] Towards Better Plasticity-Stability Trade-off in Incremental Learning: A simple Linear Connector**  
Guoliang Lin, Hanglu Chu, Hanjiang Lai; 2021.

# Исследование

## **Дальнейшее развитие**

- Расширить область покрытия экспериментов
- Исследовать разные способы регуляризации
- Применение к другим задачам и областям DL