

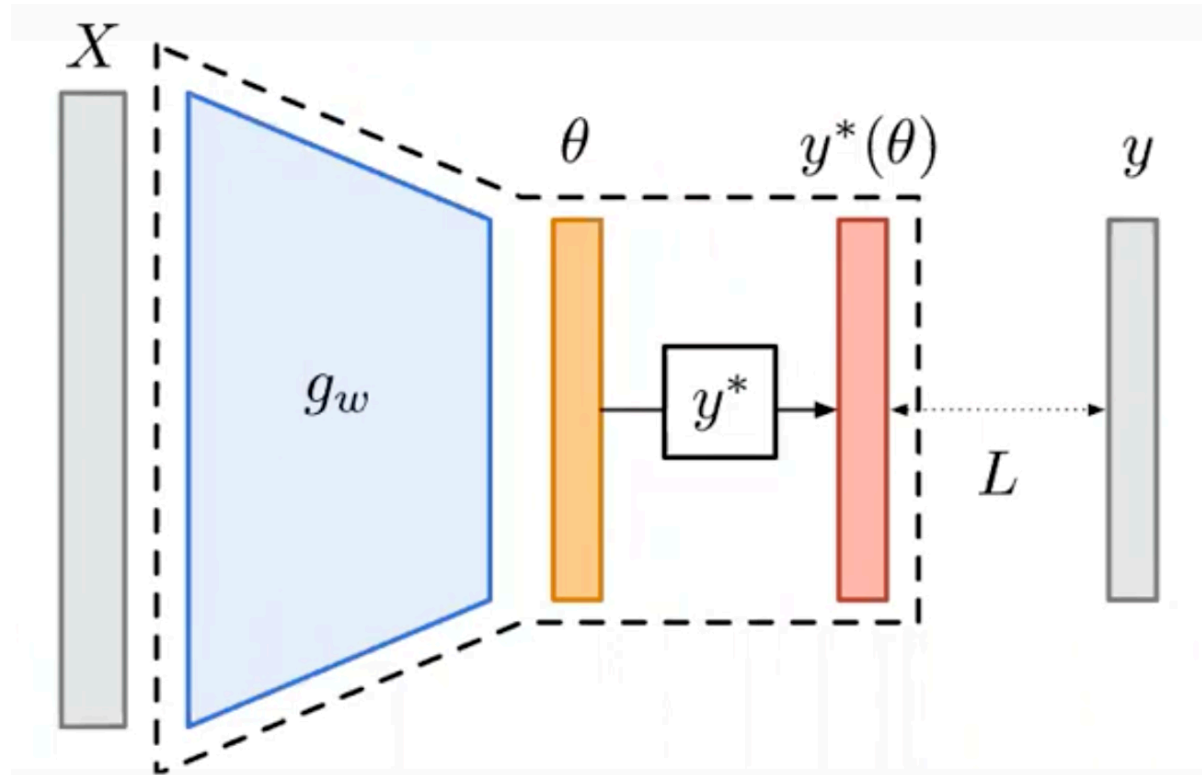
Learning with Differentiable Perturbed Optimizers

Обучение с дифференцируемыми возмущенными
оптимизаторами

Quentin Berthet, Mathieu Blondel, Olivier Teboul, Marco Cuturi, Jean-Philippe Vert, Francis Bach

Zakieva Azaliya, AMI172
HSE Research Seminar
December 2020

Проблема: дискретность



Примеры:

- Задача ранжирования (предсказания перестановок)
- Выбор ближайших соседей
- Поиск кратчайшего пути

Возмущенный максимизатор (Perturbed maximizer)

Общая задача дискретной оптимизации:

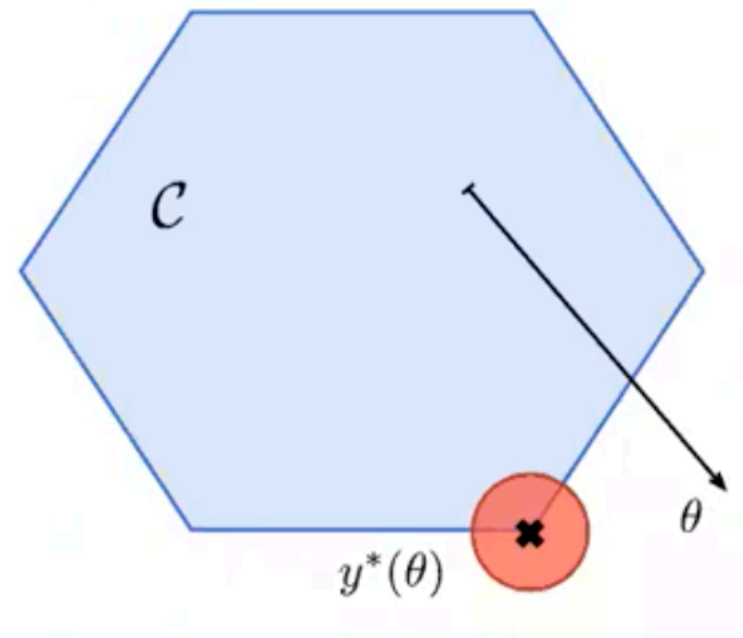
$\mathcal{Y} \subset \mathbb{R}^d$ - конечное мн-во точек

\mathcal{C} - выпуклая оболочка \mathcal{Y}

$\theta \in \mathbb{R}^d$ - входной параметр

$$F(\theta) = \max_{y \in \mathcal{C}} \langle y, \theta \rangle \quad y^*(\theta) = \arg \max_{y \in \mathcal{C}} \langle y, \theta \rangle$$

$$y^*(\theta) = \nabla_{\theta} F(\theta)$$



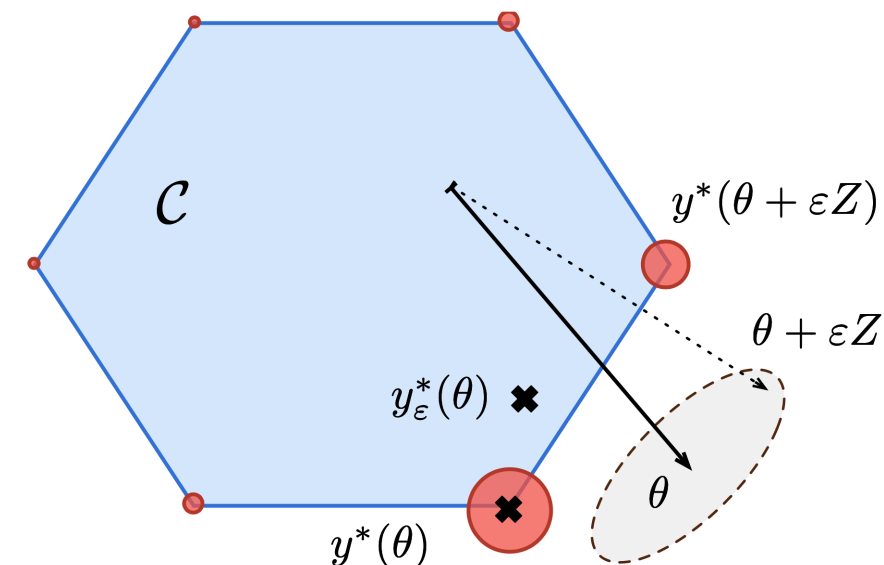
Возмущенный максимизатор (Perturbed maximizer)

Добавляем шум:

εZ - вектор шума

$\varepsilon > 0$ - температура

Z имеет положительную дифференцируемую плотность



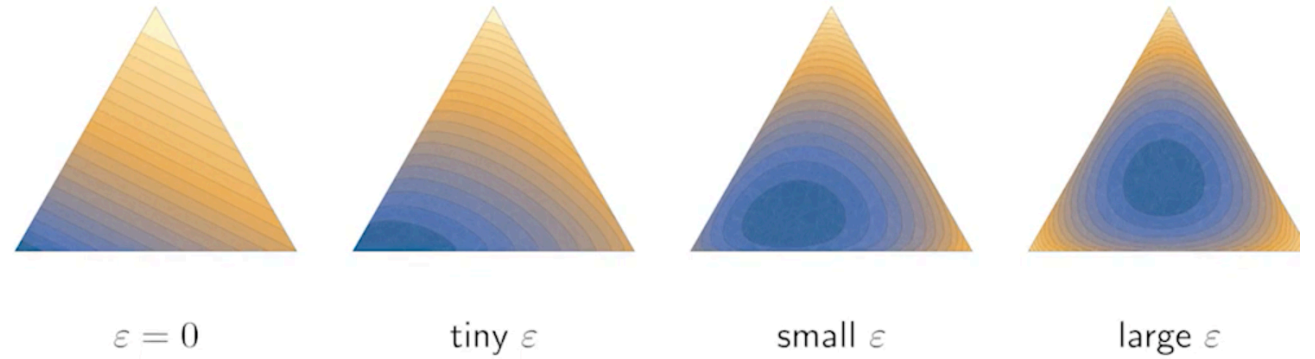
$$F_\varepsilon(\theta) = \mathbf{E}[F(\theta + \varepsilon Z)] = \mathbf{E}[\max_{y \in C} \langle y, \theta + \varepsilon Z \rangle]$$

$$y_\varepsilon^*(\theta) = \mathbf{E}_{p_\theta(y)}[Y] = \mathbf{E}[\arg \max_{y \in C} \langle y, \theta + \varepsilon Z \rangle] = \mathbf{E}[\nabla_\theta \max_{y \in C} \langle y, \theta + \varepsilon Z \rangle] = \nabla_\theta F_\varepsilon(\theta)$$

$$\downarrow$$
$$p_\theta(y) = P(y^*(\theta + \varepsilon Z) = y)$$

Свойства

Связь с регуляризацией



$\epsilon \Omega = (F_\epsilon)^*$ выпуклая функция с областью определения \mathcal{C} :

$$y_\epsilon^*(\theta) = \arg \max_{y \in \mathcal{C}} \{ \langle y, \theta \rangle - \epsilon \Omega(y) \}$$

Поведение при экстремальных температурах

При $\epsilon \rightarrow 0$: $y_\epsilon^*(\theta) \rightarrow y^*(\theta)$ для уникального \max

При $\epsilon \rightarrow \infty$: $y_\epsilon^*(\theta) \rightarrow \arg \min_{y \in \mathcal{C}} \Omega(y)$

Свойства

Дифференцируемость

Для шума Z с распределением $d\mu(z) \propto \exp(-\nu(z))dz$ и дважды дифференцируемой v , справедливо следующее утверждения:

$$F_\varepsilon(\theta) = \mathbf{E}[F(\theta + \varepsilon Z)]$$

$$y_\varepsilon^*(\theta) = \nabla_\theta F_\varepsilon(\theta) = \mathbf{E}[y^*(\theta + \varepsilon Z)] = \mathbf{E}[F(\theta + \varepsilon Z) \nabla_z \nu(Z) / \varepsilon]$$

$$J_\theta y_\varepsilon^*(\theta) = \mathbf{E}[y^*(\theta + \varepsilon Z) \nabla_z \nu(Z)^\top / \varepsilon] = \mathbf{E}[F(\theta + \varepsilon Z) (\nabla_z \nu(Z) \nabla_z \nu(Z)^\top - \nabla_z^2 \nu(Z)) / \varepsilon^2]$$

Практическая реализация

$y_\varepsilon^*(\theta) = \arg \max_{y \in \mathcal{C}} \{ \langle y, \theta \rangle - \varepsilon \Omega(y) \}$ - хотим оценку на возмущенный максимизатор и его Якобиан

Нужно: $(Z^{(1)}, \dots, Z^{(M)})$ - независимо одинаково распределённые с.в., $Z \sim \mu$

$$y^{(m)} = y^*(\theta + \varepsilon Z^{(m)}) = \arg \max_{y \in \mathcal{C}} \langle y, \theta + \varepsilon Z^{(m)} \rangle$$

$\mathbf{E}[y^{(m)}] = y_\varepsilon^*(\theta)$ для любого $m \in \{1, \dots, M\}$

Несмещенная оценка Монте-Карло для $y_\varepsilon^*(\theta)$: $\bar{y}_{\varepsilon, M}(\theta) = \frac{1}{M} \sum_{m=1}^M y^{(m)}$

Обучение: функция потерь

Функция потерь Фенхеля-Янга: $L_\varepsilon(\theta; y) = F_\varepsilon(\theta) + \varepsilon \Omega(y) - \langle \theta, y \rangle$

- Неотрицательна
- Выпукла
- Минимум равный 0 тогда и только тогда, когда θ : $y_\varepsilon^*(\theta) = y$.

Градиент функции: $\nabla_\theta L_\varepsilon(\theta; y) = \nabla_\theta F_\varepsilon(\theta) - y = y_\varepsilon^*(\theta) - y$

Обучение

$$L_{\varepsilon, \text{emp}}(w) = \frac{1}{n} \sum_{i=1}^n L_{\varepsilon}(g_w(x_i); y_i)$$

$$\nabla_w L_{\varepsilon, \text{emp}}(w) = \frac{1}{n} \sum_{i=1}^n J_w g_w(x_i) \cdot (y_{\varepsilon}^*(g_w(x_i)) - y_i)$$

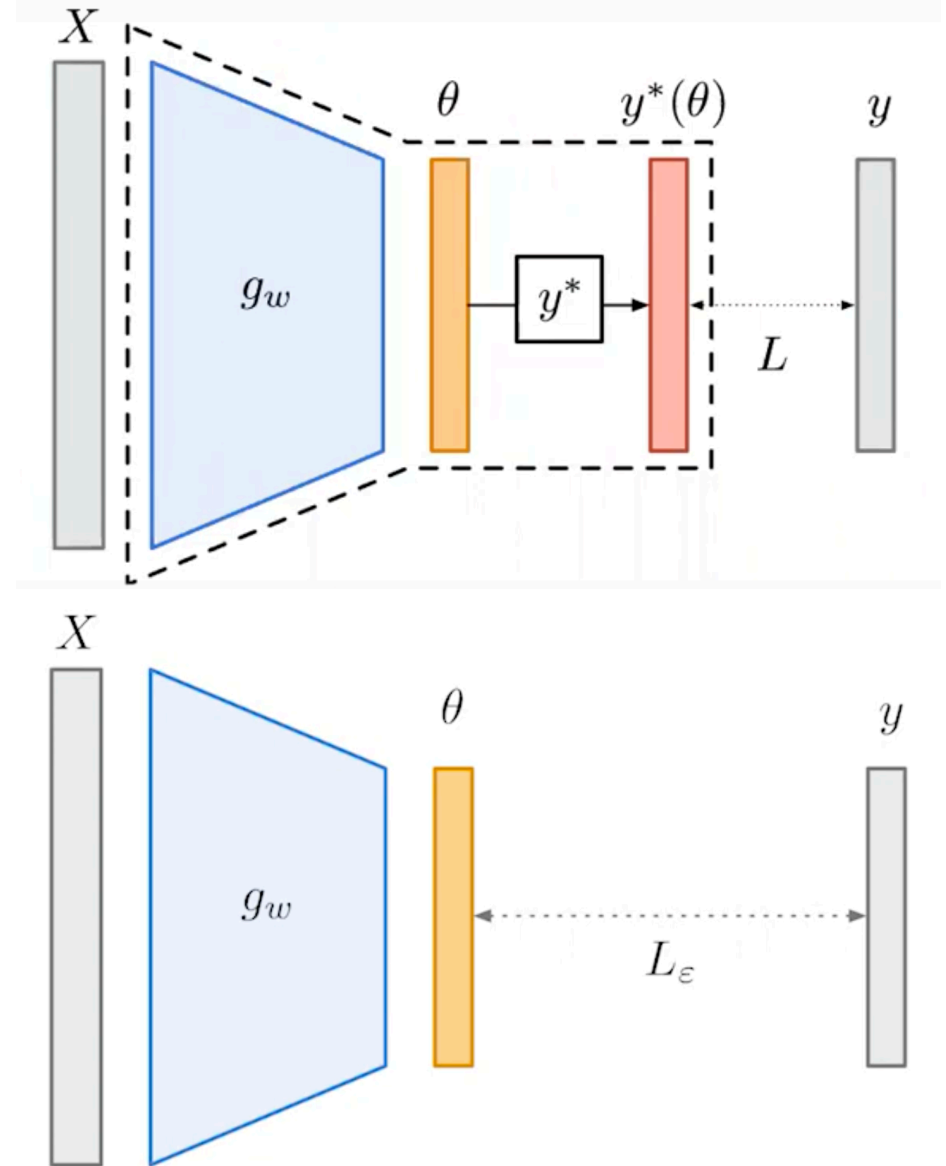
+

Дважды стохастичный градиент (doubly stochastic scheme) $\nabla_w L_{\varepsilon}(g_w(x_i); y_i)$

$$\bar{\gamma}_{i,M}(w) = J_w g_w(x_i) \left(\frac{1}{M} \sum_{m=1}^M y^*(g_w(x_i) + \varepsilon Z^{(m)}) - y_i \right)$$

Выводы

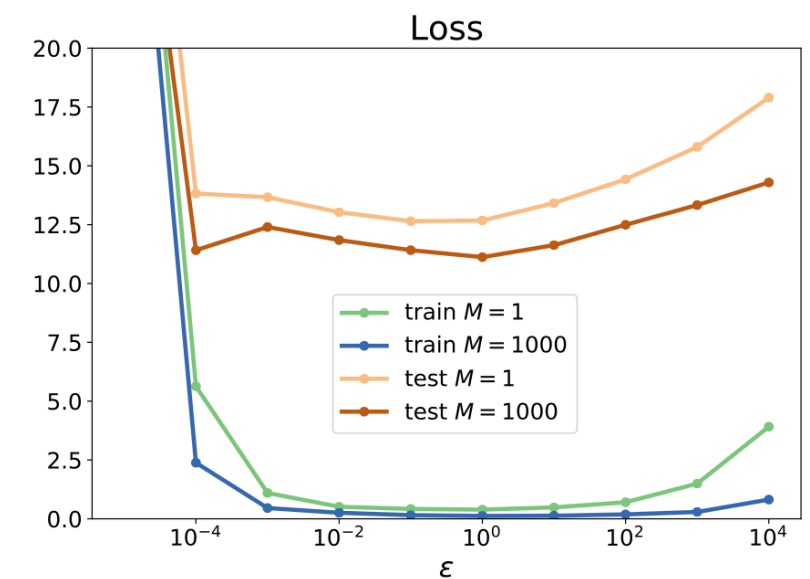
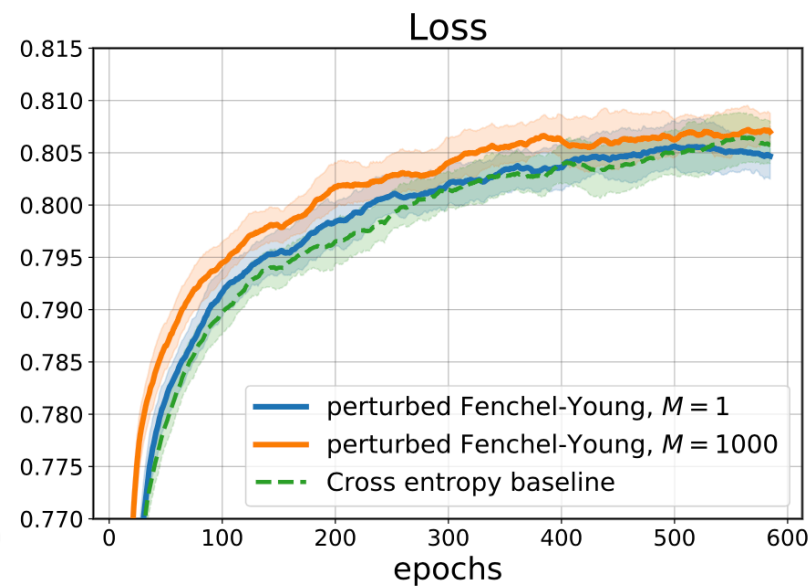
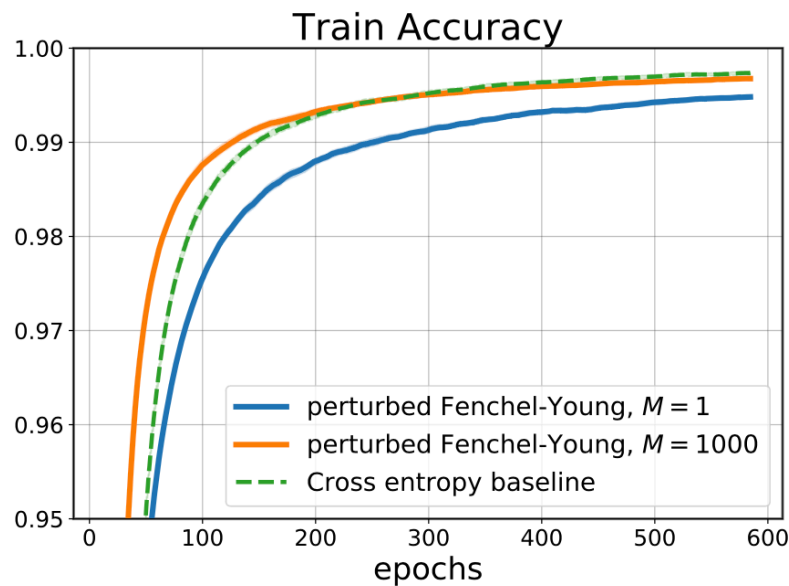
- Предложен новый универсальный метод преобразования дискретных оптимизаторов, который применим к любому солверу black-box модели.
- Метод позволяет дифференцировать argmax через возмущенный максимизатор.
- Полученные производные выражаются в виде простых мат ожиданий, которые легко аппроксимируются методами Монте-Карло.
- Произведена интеграция с функцией потерь Фенхеля-Янга, для удобного применения



Эксперименты

Классификация

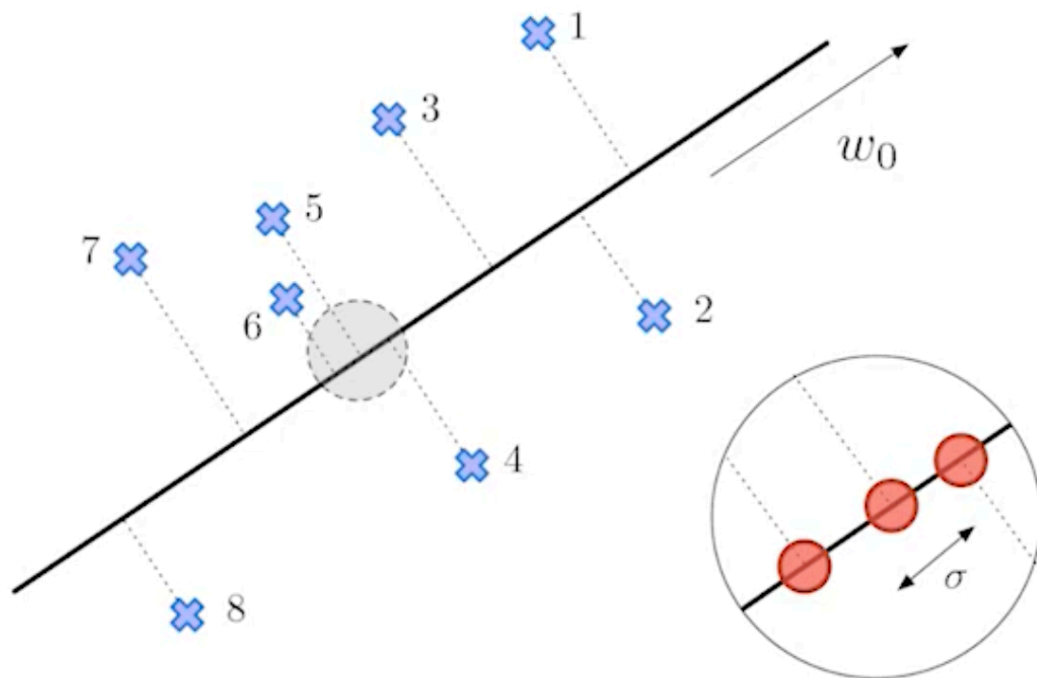
- возмущенный максимизатор с гауссовским шумом
- CIFAR-10
- vanilla-CNN (4 сверточных слоя and 2 полносвязных слоя), 600 эпох, размер батча = 32



Эксперименты

Ранжирование

$$y_i = \arg \max_y \langle x_i^\top w_0 + \sigma Z_i, y \rangle$$



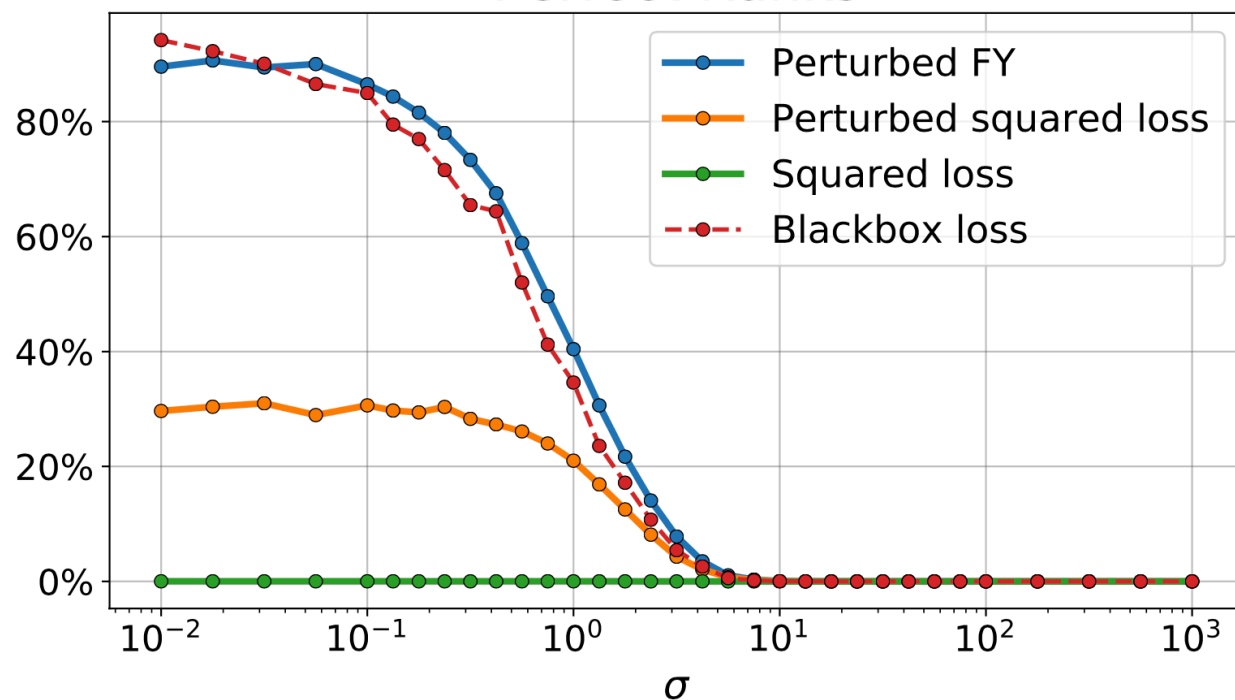
5
4
3
8
1
7
2
6

Эксперименты

Ранжирование

$$y_i = \arg \max_y \langle x_i^\top w_0 + \sigma Z_i, y \rangle$$

Perfect Ranks



Perturbed Fenchel-Young (proposed)

Perturbed + Squared loss (proposed): $\frac{1}{2} \|y_i - y_\varepsilon^*(g_w(x_i))\|^2$

Squared loss: $\frac{1}{2} \|y_i - g_w(x_i)\|^2$

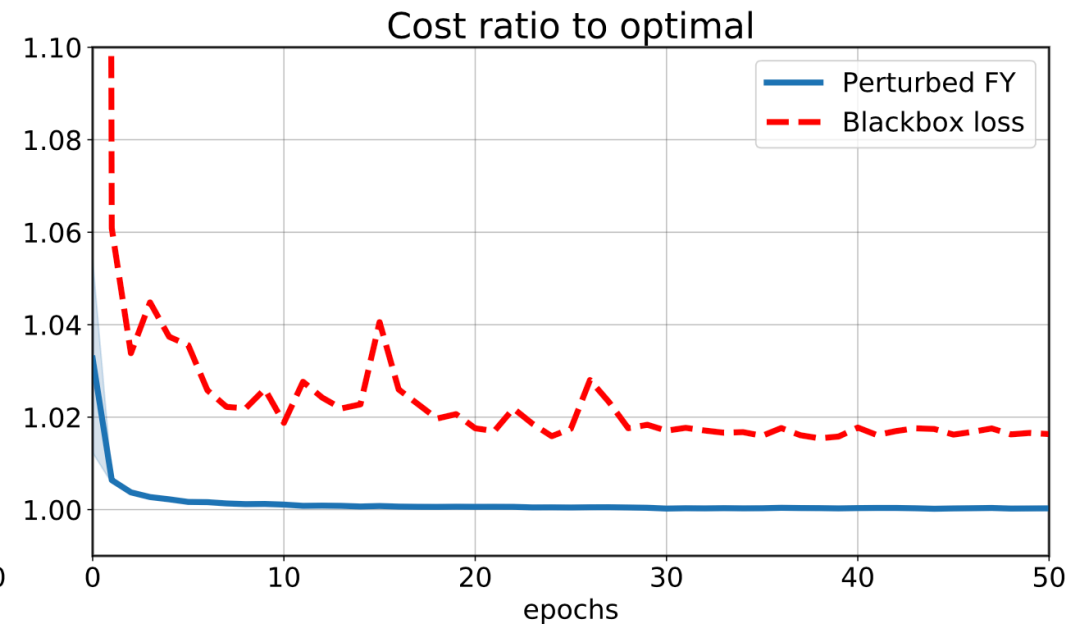
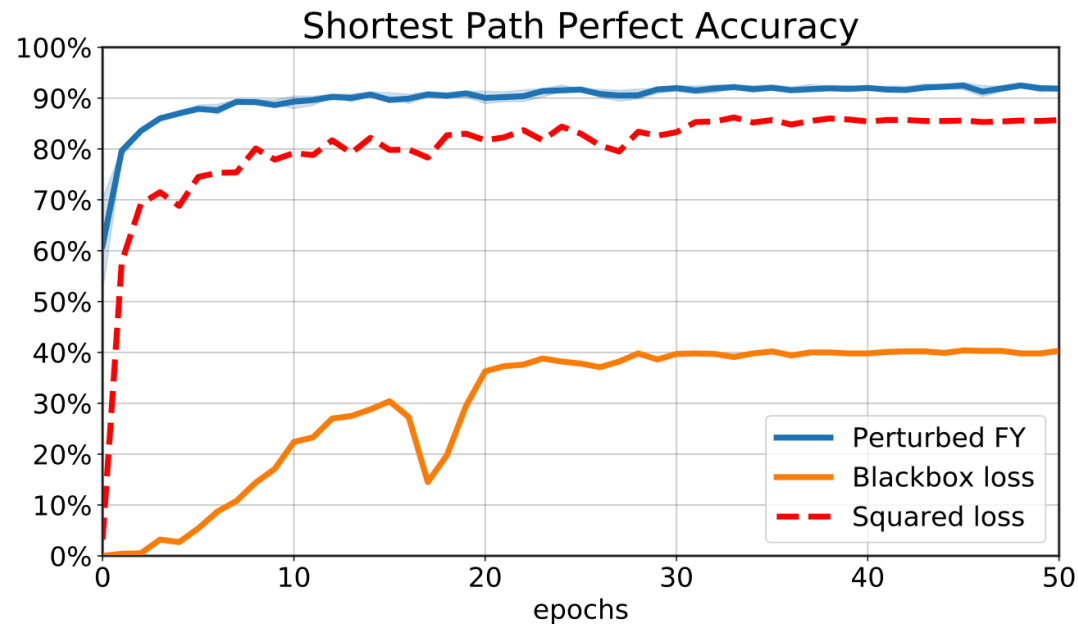
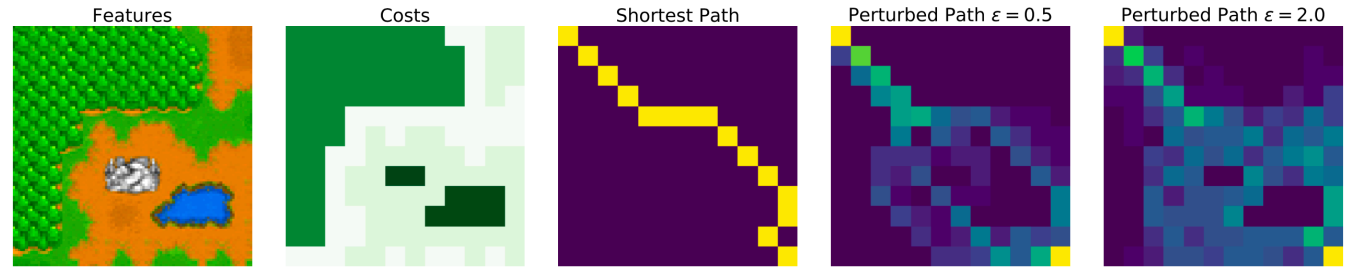
Blackbox loss: $\frac{1}{2} \|y_i - y^*(g_w(x_i))\|^2$

Используется приближение градиента из статьи Vlastelica M.
Differentiation of blackbox combinatorial solvers, 2019.

Эксперименты

Поиск кратчайшего пути

- Первые 5 слоев ResNet18
- 50 эпох, размер батча = 70
- $\epsilon = 1$, $M = 1$



Вопросы

1. Зачем авторы статьи используют возмущенный максимизатор, какие задачи решаются при использовании?
2. Выпишите формулу функции потерь Фенхеля-Янга и объясните все ее обозначения.
3. Опишите любой из 3 экспериментов и выводы, к которым они приводят.

ИСТОЧНИКИ

- Quentin Berthet, Mathieu Blondel, Olivier Teboul, Marco Cuturi, Jean-Philippe Vert, Francis Bach. Learning with Differentiable Perturbed Optimizers, 2020.