

Глубинное обучение

Лекция 9: Дифференцируемое программирование

Лектор: Антон Осокин

ФКН ВШЭ, 2019



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

Объявления

- Открыт набор в группу байесовский методов

<https://bayesgroup.ru/admission/>

- Байесовские методы (Deep)
- Deep learning
- Deep learning for CV
- Deep learning for NLP
- Optimization

- Открыт набор на летнюю школу

<http://deepbayes.ru/>

План лекции


- Мотивация: зачем встраивать алгоритмы в нейросети?
- Алгоритмы в нейросетях
 - Структурный пулинг = комбинаторная оптимизация
 - Итерации алгоритмов = слои нейросетей
 - Дифференцирование по входу алгоритма
- Недифференцируемые модели
 - Что делать?

Алгоритмы в нейросетях, зачем?

- У нейросетей очень сложные структуры
- Структуры сложно создавать для новых задач
- Часто нейросеть не может «выучить всё», надо помогать
- Комбинировать существующие решения и нейросети
- Существуют очень мощные алгоритмы для сложных задач
- Deep learning meets computer science!

Задача: распознавание рукописных СИМВОЛЫ

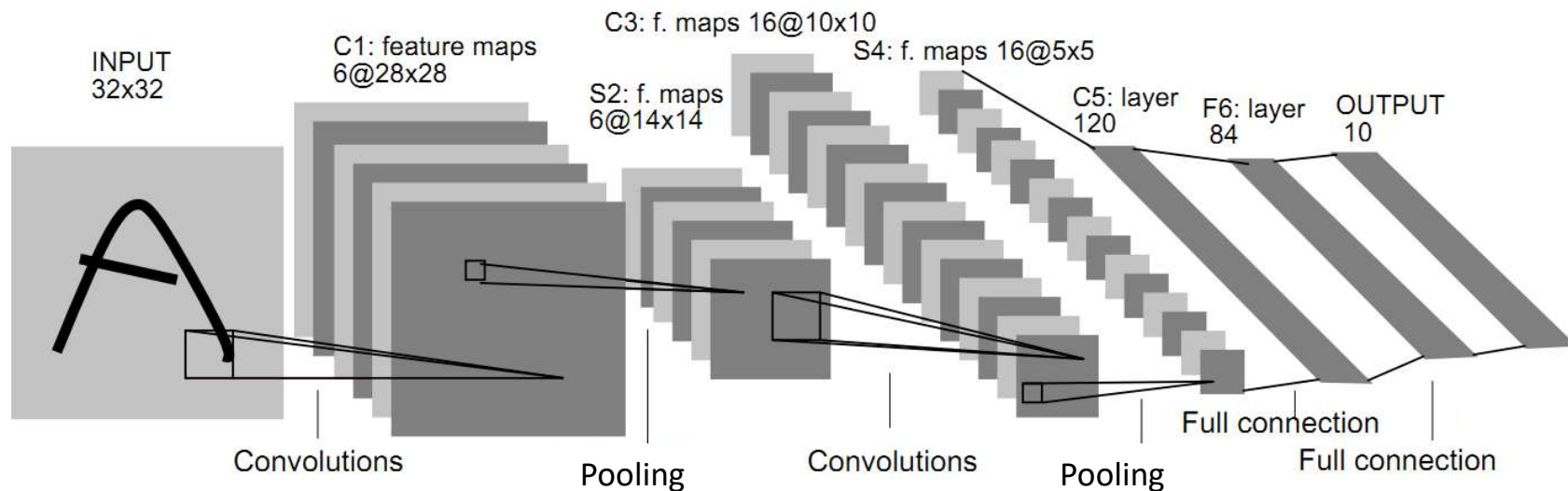
- Модельная задача – распознавание рукописных символов

 → command

- A MNIST для structured prediction
- Простая задача => очень много методов применимо
- Стандартные алгоритмы:
динамическое программирование, СЛАУ

Структура нейросети

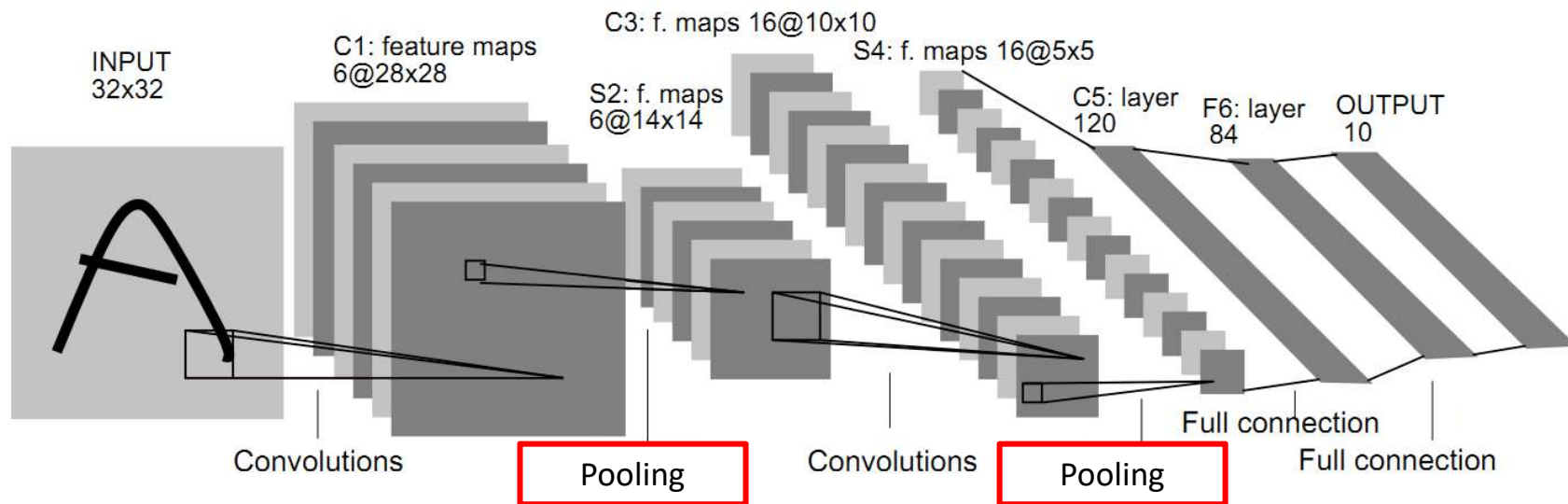
- Нейросеть для распознавания одного символа:



- Линейные операции с параметрами
 - Свёртки для изображений
- Нелинейность (sigmoid, ReLu)
- Пулинг для понижения размерности и инвариантности
- Обучение = стохастическая оптимизация

Пулинг для выбора активаций

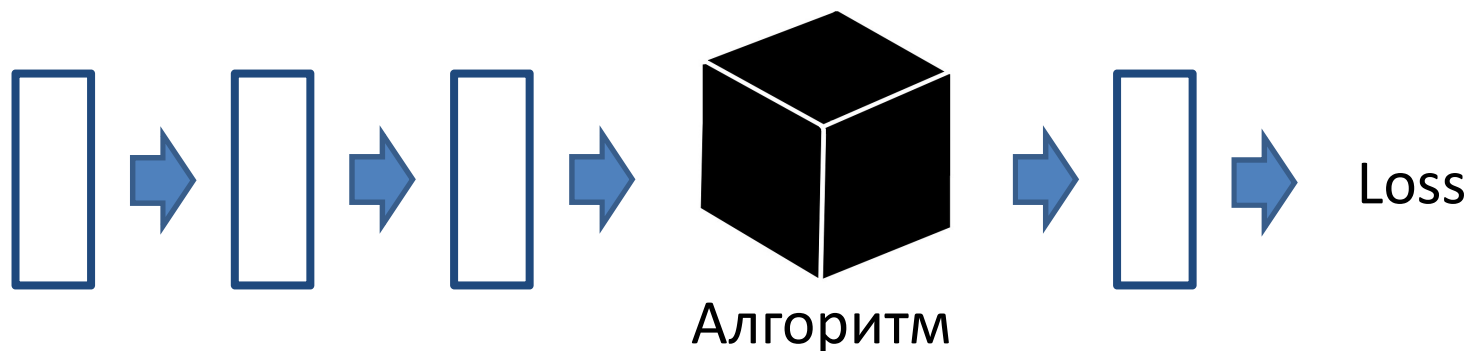
- Нейросеть для распознавания одного символа:



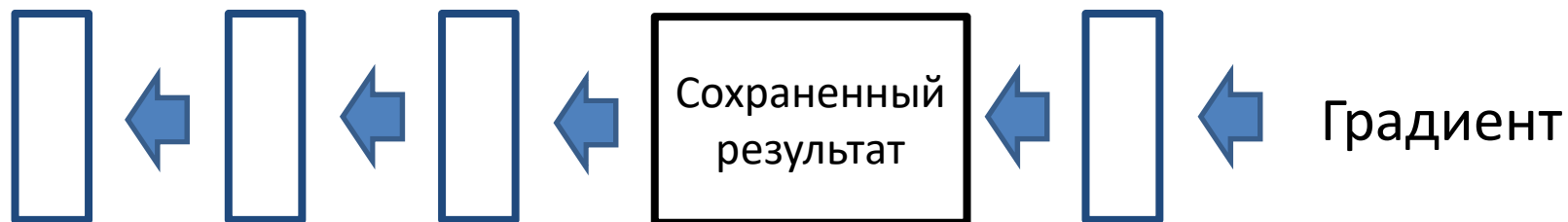
- Пулинг – процедура агрегирования активаций (max, sum)
- Алгоритмы для пулинга: квантиль, сортировка
- И это дифференцируемо?
 - «дифференцируемо» в смысле нейросетей

Структурный пулинг = комбинаторная оптимизация

Нейросеть: проход вперёд



Для прохода назад нужен только результат алгоритма



Conditional Random Field (CRF)

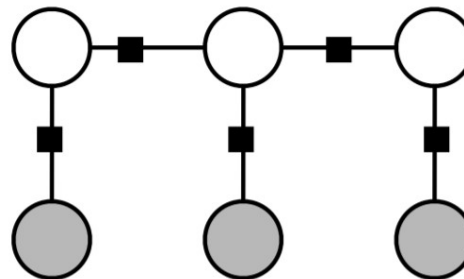


- CRF – способ учесть связи между символами
- Связи задаются функцией меток, связывающей метки

$$F(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \theta_i(y_i \mid x_i) + \sum_{i=1}^{T-1} \theta_{i,i+1}(y_i, y_{i+1})$$

- Здесь \mathbf{y} – метки символов, \mathbf{x} – изображения символов
 θ - потенциалы (унарные и парные)

- Графическая модель
(фактор-граф)



- Унарные потенциалы вычисляются нейросетью

CRF: наилучшая конфигурация

- CRF – способ учесть связи между символами
- Связи задаются функцией меток, связывающей метки
$$F(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \theta_i(y_i \mid x_i) + \sum_{i=1}^{T-1} \theta_{i,i+1}(y_i, y_{i+1})$$
- Здесь \mathbf{y} – метки символов, \mathbf{x} – изображения символов
 θ - потенциалы (унарные и парные)
- Чем F больше, тем конфигурация лучше
- **Динамическое программирование** – сведение к подзадачам
 - $V_i(y_i)$ – лучшее значение слагаемых F для слагаемых $\leq i$
 - Проход вперёд:
$$V_1(y_1) = \theta_1(y_1 \mid x_1)$$
$$V_{i+1}(y_{i+1}) = \theta_{i+1}(y_{i+1} \mid x_{i+1}) + \max_{y_i} \left(\theta_{i,i+1}(y_i, y_{i+1}) + V_i(y_i) \right)$$
 - Оптимальное значение: $F^* = \max_{y_T} V_T(y_T)$
 - Проход назад для восстановления конфигурации

Обучение CRF – структурный SVM

- Обучение по размеченной выборке $\{x_n, y_n\}_{n=1}^N$
- Обучение – задача оптимизации $\min_{\theta} \frac{1}{N} \sum_{n=1}^N \Psi(x_n, y_n | \theta)$

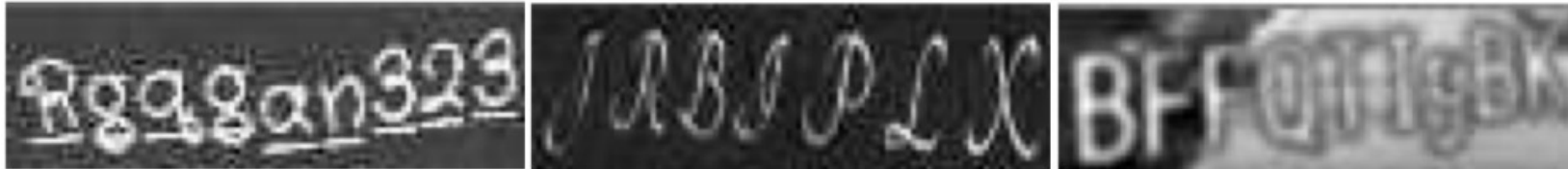
- SSVM - обобщение метода опорных векторов

$$\Psi(x_n, y_n | \theta) = \max_y \left[F(y | x_n, \theta) + \Delta(y, y_n) \right] - F(y_n | x_n, \theta)$$

- Общая схема метода
 1. Вычисление потенциалов (проход вперёд через нейросеть)
 2. Вычислении функции потерь (дин. программирование)
 3. Градиент по выходу нейросети
 4. Градиент по параметрам нейросети (backprop)
 5. Шаг оптимизации

Примеры использования

- Свободный текст (Jaderberg et al., ICLR 2015)



- Детектирование нескольких объектов (Vu et al., ICCV 2015)



- Тэггирование изображений (Chen et al., ICML 2015)



female/indoor/portrait



sky/plant life/tree



water/animals/sea



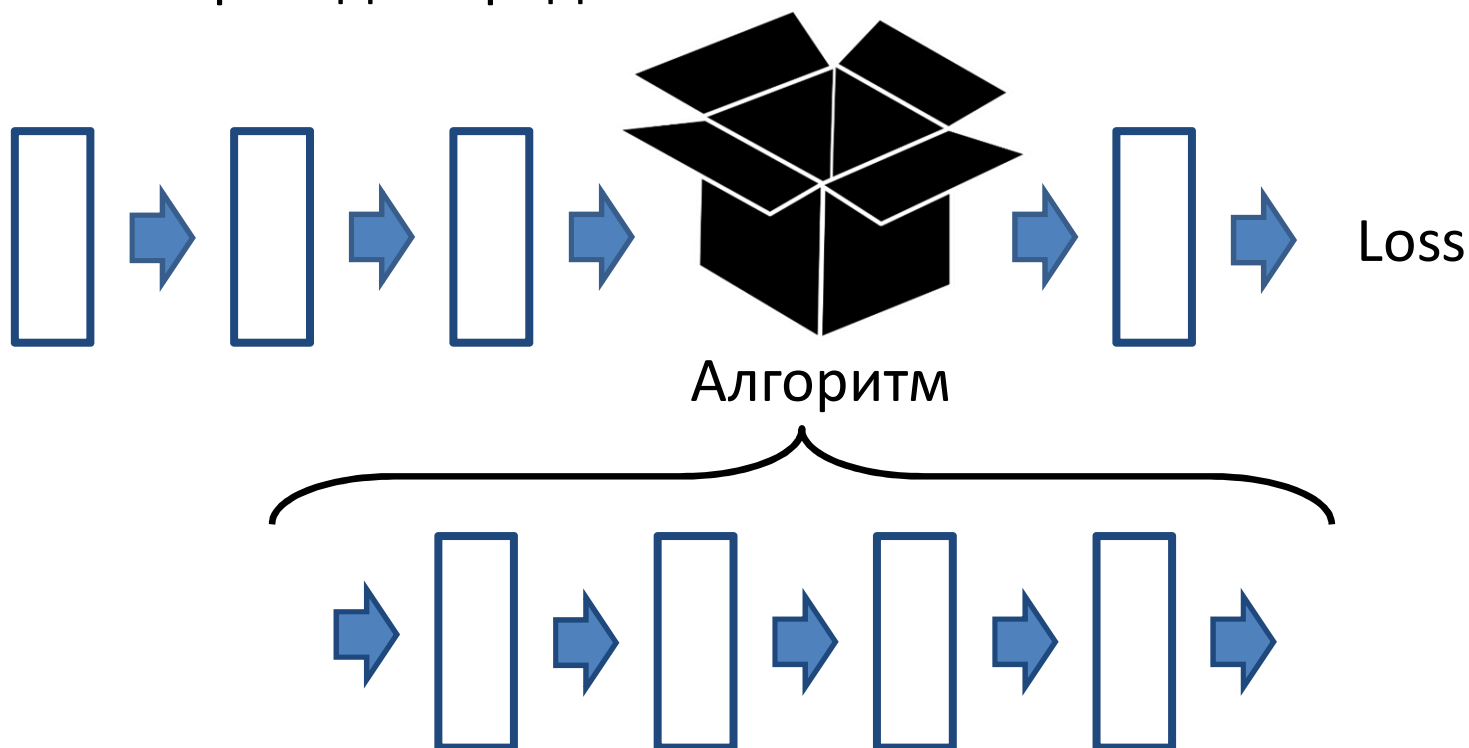
animals/dog/indoor



indoor/flower/plant life


Итерации алгоритма как слои сети

Нейросеть: проход вперёд



Проход назад – обычный back propagation

Обучение CRF – максимальное правдоподобие

 → command

- Обучение по размеченной выборке $\{x_n, y_n\}_{n=1}^N$
- Обучение – задача оптимизации $\min_{\theta} \frac{1}{N} \sum_{n=1}^N \Psi(x_n, y_n | \theta)$
- Метод максимального правдоподобия
 $\Psi(x_n, y_n | \theta) = -\log P(y | x_n, \theta), \quad P = \frac{1}{Z(x_n, \theta)} \exp(F(y | x_n, \theta))$
- Z – нормировочная константа
$$Z(x_n, \theta) = \sum_y \exp(F(y | x_n, \theta))$$
- Сумма экспоненциального числа слагаемых
- Используем Алгоритм!

Обучение CRF – максимальное правдоподобие

- Z – нормировочная константа

$$Z(\mathbf{x}_n, \boldsymbol{\theta}) = \sum_{\mathbf{y}} \exp(F(\mathbf{y} \mid \mathbf{x}_n, \boldsymbol{\theta}))$$

- Функция F обладает структурой:

$$\exp(F(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})) = \prod_{i=1}^T \exp(\theta_i(y_i \mid x_i)) \prod_{i=1}^{T-1} \exp(\theta_{i,i+1}(y_i, y_{i+1}))$$

- **Алгоритм динамического программирования (sum-product)**

- $V_i(y_i)$ – сумма произведений Z для слагаемых $\leq i$

- Проход вперёд:

$$V_{i+1}(y_{i+1}) = \exp(\theta_{i+1}(y_{i+1} \mid x_{i+1})) \sum_{y_i} \left(\exp(\theta_{i,i+1}(y_i, y_{i+1})) V_i(y_i) \right)$$

- Ответ: $Z = \sum_{y_T} V_T(y_T)$

Обучение CRF – вариационный вывод

- Вар. вывод – один из основных методов обучения байесовских вероятностных моделей

- Приближение P простым распределением Q : $Q(\mathbf{y}) = \prod_{i=1}^T q_i(y_i)$
$$\text{KL}(Q \parallel P) = - \sum_{\mathbf{y}} Q(\mathbf{y}) \log \frac{P(\mathbf{y})}{Q(\mathbf{y})} \rightarrow \min_Q$$

- Используем структуру модели

$$P(\mathbf{y}) \propto \exp(F(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})) = \prod_{i=1}^T \exp(\theta_i(y_i \mid x_i)) \prod_{i=1}^{T-1} \exp(\theta_{i,i+1}(y_i, y_{i+1}))$$

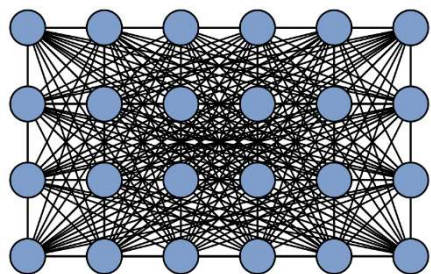
- Формулы пересчёта:

$$q_i(y_i) \propto \exp\left(\theta_i(y_i) + \sum_{y_{i-1}} \theta_{i-1,i}(y_{i-1}, y_i) q_{i-1}(y_{i-1}) + \sum_{y_{i+1}} \theta_{i,i+1}(y_i, y_{i+1}) q_{i+1}(y_{i+1})\right)$$

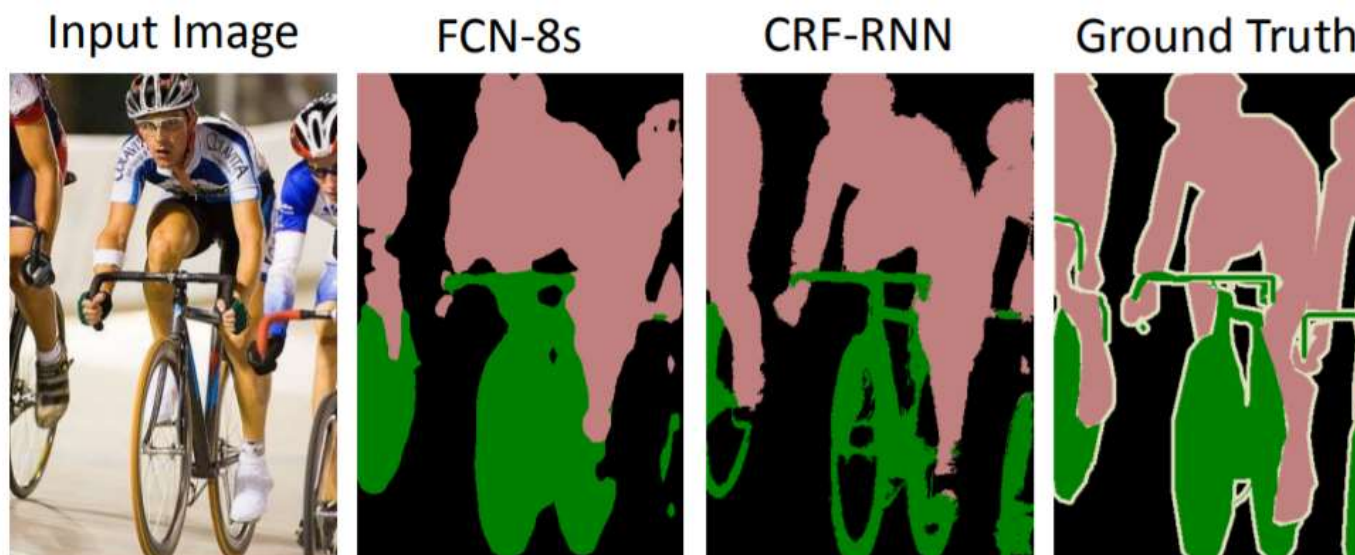
Пример: сегментация изображений

[Zheng et al., ICCV 2015]

- Вариационных вывод над полно-связной CRF



- Все операции с использованием свёрток



Пример: подавление шума

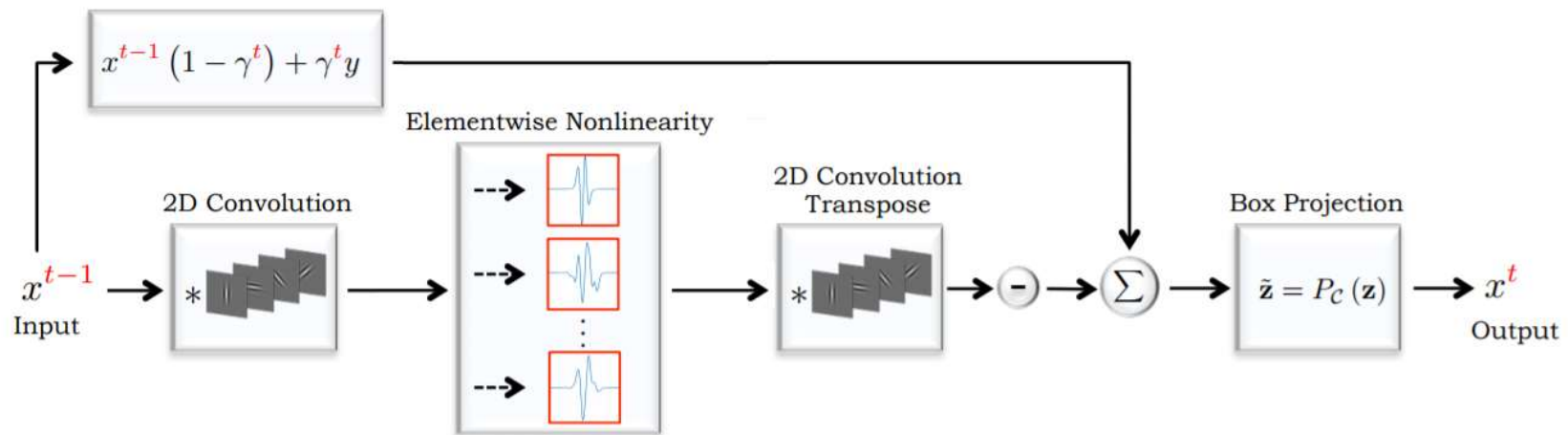
[Lefkimmiatis, CVPR 2017]

- Формулировка задачи (\mathbf{x} – ответ, \mathbf{y} – шумное изображение):

$$\mathbf{x}^* = \arg \min_{a \leq x_n \leq b} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \lambda \sum_{r=1}^R \phi(\mathbf{L}_r \mathbf{x})$$

- Алгоритм решения – Proximal Gradient

$$\mathbf{x}^t = P_C \left(\mathbf{x}^{t-1} (1 - \gamma^t) + \gamma^t \mathbf{y} - \alpha^t \sum_{r=1}^R \mathbf{L}_r^T \psi(\mathbf{L}_r \mathbf{x}^{t-1}) \right)$$



Пример: подавление шума

[Lefkimmatis, CVPR 2017]

- Формулировка задачи (\mathbf{x} – ответ, \mathbf{y} – шумное изображение):

$$\mathbf{x}^* = \arg \min_{a \leq x_n \leq b} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \lambda \sum_{r=1}^R \phi(\mathbf{L}_r \mathbf{x})$$

- Результаты:



Оригинал



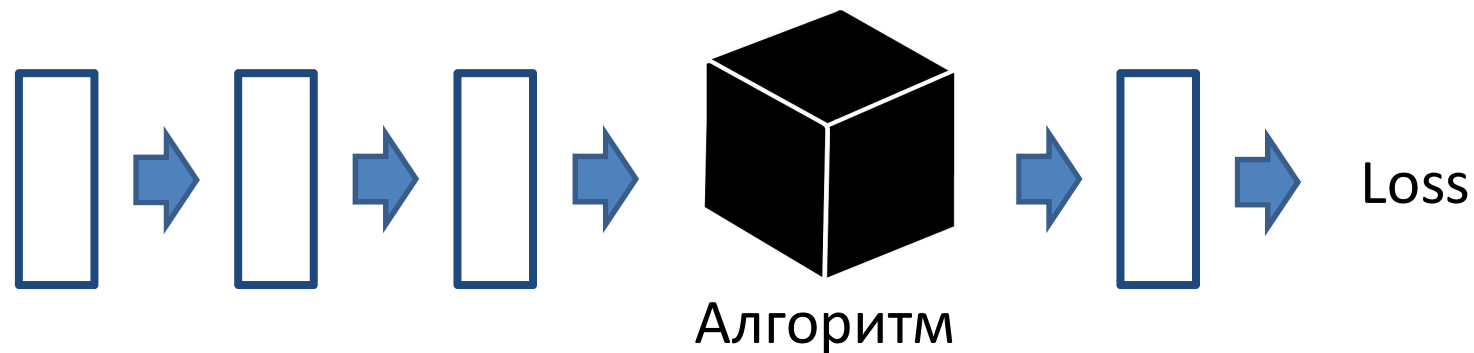
Оригинал + шум



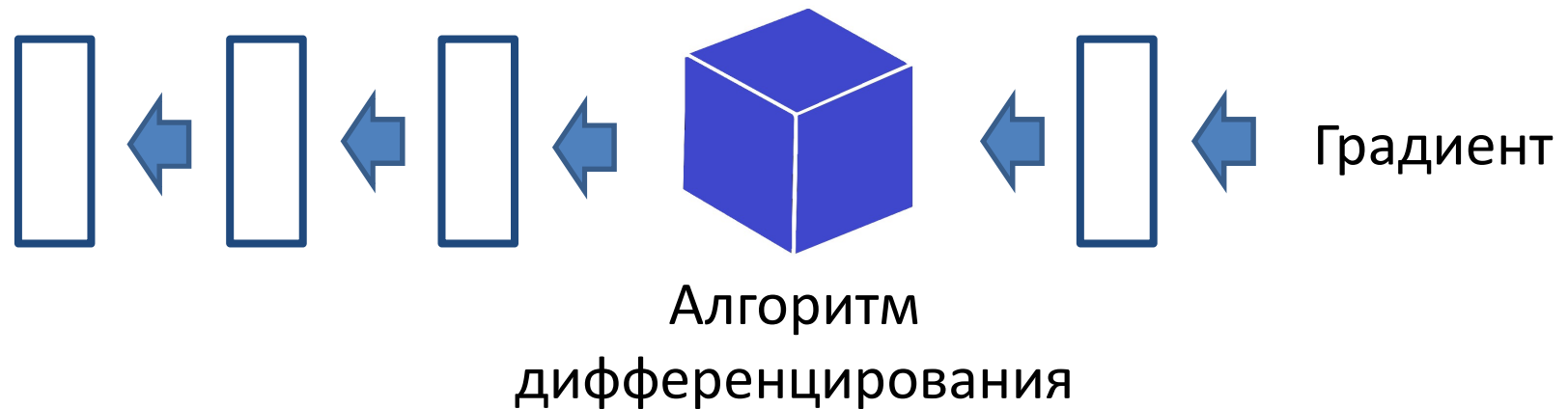
Результат

Дифференцирование по входу

Нейросеть: проход вперёд



Для прохода назад нужен другой алгоритм



Пример: Gaussian MRF

[Chandra&Kokkinos, ECCV 2016]

- Непрерывный аналог CRF

$$F(\mathbf{y}) = -\frac{1}{2}\mathbf{y}^T(A + \lambda I)\mathbf{y} + \mathbf{b}^T\mathbf{y} \rightarrow \max_{\mathbf{y}}$$

- \mathbf{y} – переменные, A , \mathbf{b} – параметры (выходы нейросети)

- Алгоритм предсказания – СЛАУ: $\mathbf{y} = (A + \lambda I)^{-1}\mathbf{b}$

- Задача дифференцирования

- Вход: параметры A , \mathbf{b} , решение \mathbf{y} , градиент $\frac{d\Psi}{d\mathbf{y}}$
- Найти: градиенты $\frac{d\Psi}{d\mathbf{b}}$ и $\frac{d\Psi}{dA}$

- Дифференцируем линейный слой (СЛАУ)

$$\frac{d\Psi}{d\mathbf{b}} = (A + \lambda I)^{-T} \frac{d\Psi}{d\mathbf{y}}$$

Пример: Gaussian MRF

[Chandra&Kokkinos, ECCV 2016]

- Непрерывный аналог CRF

$$F(\mathbf{y}) = -\frac{1}{2}\mathbf{y}^T (A + \lambda I)\mathbf{y} + \mathbf{b}^T \mathbf{y} \rightarrow \max_{\mathbf{y}}$$

- \mathbf{y} – переменные, A , \mathbf{b} – параметры (выходы нейросети)

- Алгоритм предсказания – СЛАУ: $\mathbf{y} = (A + \lambda I)^{-1}\mathbf{b}$

- Задача дифференцирования

- Вход: параметры A , \mathbf{b} , решение \mathbf{y} , градиент $\frac{d\Psi}{d\mathbf{y}}$

- Найти: градиенты $\frac{d\Psi}{d\mathbf{b}}$ и $\frac{d\Psi}{dA}$

- Производная по A : $\frac{d\Psi}{dA} = -\frac{d\Psi}{d\mathbf{b}}\mathbf{y}^T$

Примеры дифференцирования

- SVD разложение $X = U\Sigma V^T$ [Ionescu et al., ICCV 2015]

$$\frac{\partial L \circ f}{\partial X} = DV^T + U \left(\frac{\partial L}{\partial \Sigma} - U^T D \right)_{diag} V^T + 2U\Sigma \left(K^T \circ \left(V^T \left(\frac{\partial L}{\partial V} - VD^T U\Sigma \right) \right) \right)_{sym} V^T$$

$$K_{ij} = \begin{cases} \frac{1}{\sigma_i^2 - \sigma_j^2}, & i \neq j \\ 0, & i = j \end{cases} \quad D = \left(\frac{\partial L}{\partial U} \right)_1 \Sigma_n^{-1} - U_2 \left(\frac{\partial L}{\partial U} \right)_2^T U_1 \Sigma_n^{-1}$$

- Дискретная! оптимизация [Djolonga&Krause, NIPS 2017]
 - Используется эквивалентность дискретной минимизации субмодулярных функций непрерывным релаксациям
- Решение дифференциальных уравнений [Chen et al., NIPS 2018]
 - Вычисление градиента – решение другого уравнения

Недифференцируемые модели?

- Часто недифференцируемые функции – backpropable (“нейро-дифференцируемые”)
 - Примеры: max, ReLu, медиана
- Совсем-недифференцируемые функции
 - Кусочно-постоянные функции
 - argmax
 - $f(x) = 0$ if $x < 0$ else 1
 - Сложные индексы
 - Позиция прямоугольник на изображении
 - Ответы внешних систем:
 - Программа, среда, человек

Что делать?

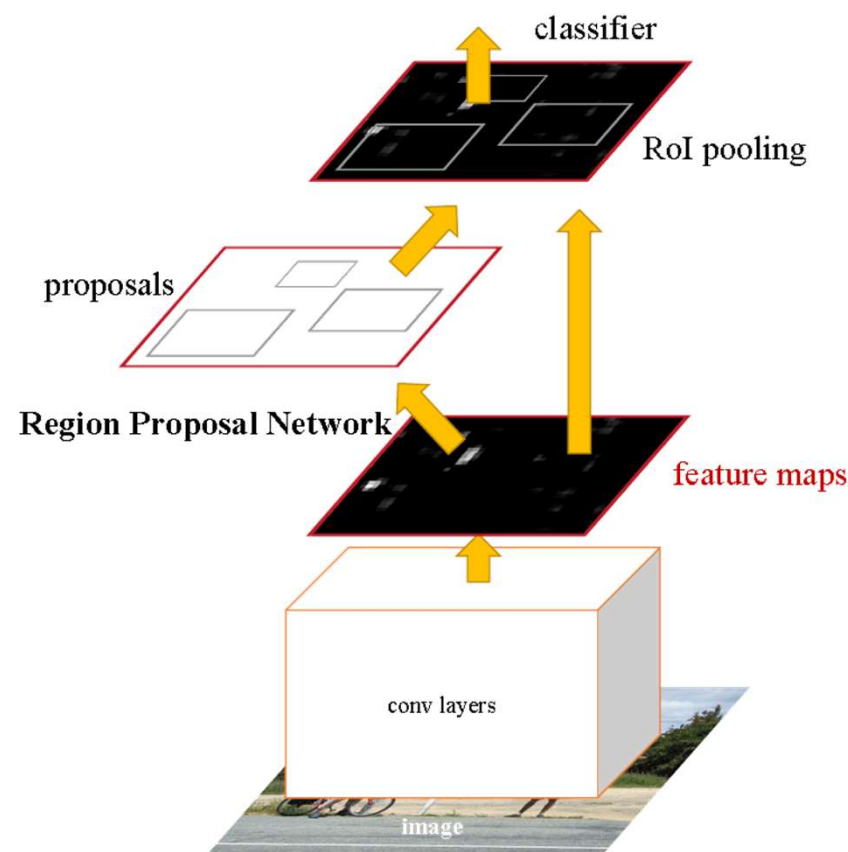
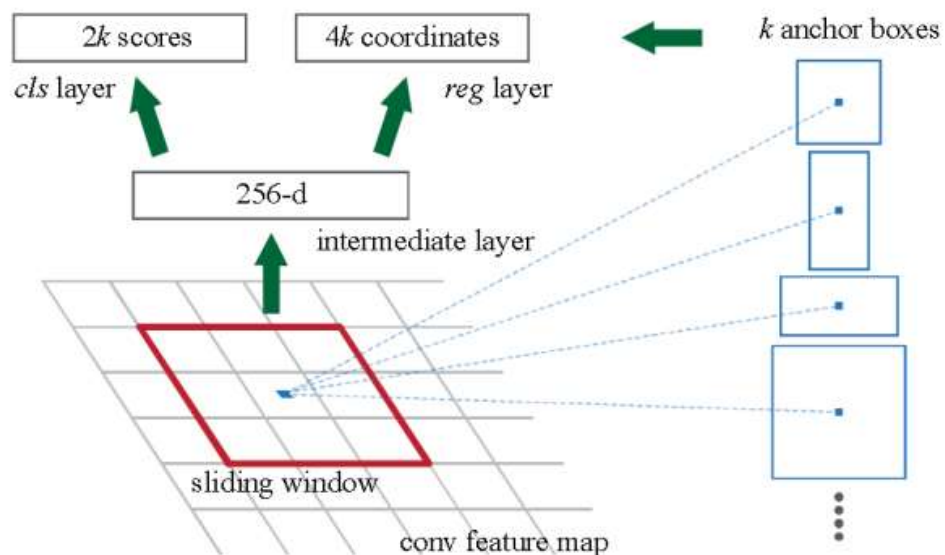
- Игнорировать 😊
 - max, комбинаторный пулинг
- Сглаживать через стохастичность (<http://deepbayes.ru/>)
 - Стохастические активации
- Другие способы сглаживания

Летняя школа
в 2019!

Детектор Faster R-CNN

[Ren et al., 2015]

- Одна сеть выдаёт гипотезы объектов (proposals)
- Вторая сеть классифицирует гипотезы
- RoI pooling – недифференцируемый по координатам
- **Игнорируем это!**



Стохастические активации

- Что это такое?
 - $\theta(x)$ – дифференцируемая функция
 - Распределение $p(w | \theta(x))$
 - Лосс $L(w)$ – не дифференцируемый?
- Проход вперёд:
 - Вычисление $\theta(x)$ – параметры распределения
 - Сэмплирование w из $p(w | \theta(x))$
 - Вычисление $L(w)$
- Функция $\mathbb{E}_{w \sim p(w|\theta(x))} L(w)$ – часто дифференцируемая по θ и x
- Градиент получается из log-derivative trick
$$\nabla_{\theta} = \mathbb{E}_{w \sim p(w|\theta)} \nabla_{\theta} [\log p(w|\theta(x))] L(w) \quad \nabla_{\theta} [\log p(w|\theta)] = \frac{\nabla_{\theta} [p(w|\theta)]}{p(w|\theta)}$$
- **Большая дисперсия!**

Обучение: дифференцируемый лосс

- Простой случай:
 - Лосс $L(w)$ – дифференцируемый
 - Распределение $p(w | \theta)$ – «хорошее»
- Репараметризация (если возможна) – самое лучшее решение!
 - Разделение случайности и параметров
 - Представим распределение $p(w | \theta)$ как $g(\theta, \varepsilon)$, $\varepsilon \sim r(\varepsilon)$
 - $z = \mu_\theta(x) + \sigma_\theta(x)\varepsilon$, $\varepsilon \sim r(\varepsilon)$ g – детерминированная функция
 - ε – шум
 - Тогда градиент легко оценить:
$$\nabla_\theta = \nabla_\theta \int p(w|\theta)L(w)dw = \int r(\varepsilon)\nabla_\theta L(g(\theta, \varepsilon))d\varepsilon$$
 - Дисперсия градиента сильно уменьшается

Какие распределения можно репараметризовать?

$p(x y)$	$r(\epsilon)$	$g(\epsilon, y)$
$\mathcal{N}(x \mu, \sigma^2)$	$\mathcal{N}(\epsilon 0, 1)$	$x = \sigma\epsilon + \mu$
$\mathcal{G}(x 1, \beta)$	$\mathcal{G}(\epsilon 1, 1)$	$x = \beta\epsilon$
$\mathcal{E}(x \lambda)$	$\mathcal{U}(\epsilon 0, 1)$	$x = -\frac{\log \epsilon}{\lambda}$
$\mathcal{N}(x \mu, \Sigma)$	$\mathcal{N}(\epsilon 0, I)$	$x = A\epsilon + \mu$, where $AA^T = \Sigma$

Slide credit: Dmitry Vetrov

Нельзя репараметризовать дискретные распределения!

- Категориальное распределение $z \sim \text{Discrete}(\alpha_1, \dots, \alpha_L)$

- **Надо для argmax !**

$$z = (0, 1, 0, \dots, 0)$$

- Релаксация: Gumbel-Softmax [Jang et al., 2017; Maddison et al., 2017]

$$(z_1, \dots, z_L) \sim \text{RelaxedDiscrete}(\alpha_1, \dots, \alpha_L | T)$$

$$z_i = \frac{\exp((\log \alpha_i + G_i)/T)}{\sum_{j=1}^L \exp((\log \alpha_j + G_j)/T)}, \quad G_k \sim \text{Gumbel}$$

$$G_k = -\log(-\log u_k), \quad u_k \sim \text{Uniform}[0, 1]$$

- $\text{RelaxedDiscrete}(\alpha_1, \dots, \alpha_L | T) \xrightarrow{T \rightarrow 0} \text{Discrete}(\alpha_1, \dots, \alpha_L)$

- Трюк: Straight-through estimator

[Hinton et al., 2015 course]

$$\nabla_{\alpha} \approx \nabla_z$$

[Bengio et al. (2013)]

Градиент 1 в выбранную позицию, остальные 0

Недифференцируемый лосс?

Что делать?

- Модель:
 - $\theta(x)$ – дифференцируемая функция
 - Распределение $p(w \mid \theta(x))$
 - Лосс $L(w)$ недифференцируемый
- Обучение с подкреплением ☹
 - Политика $p(w \mid \theta(x))$
 - Награда-reward: $-L(w)$
- Log-derivative trick лежит в основе REINFORCE, policy gradients
$$\nabla_{\theta} = \mathbb{E}_{w \sim p(w|\theta)} \nabla_{\theta} [\log p(w|\theta(x))] L(w)$$
 - Борьба с дисперсией! (всеми средствами)

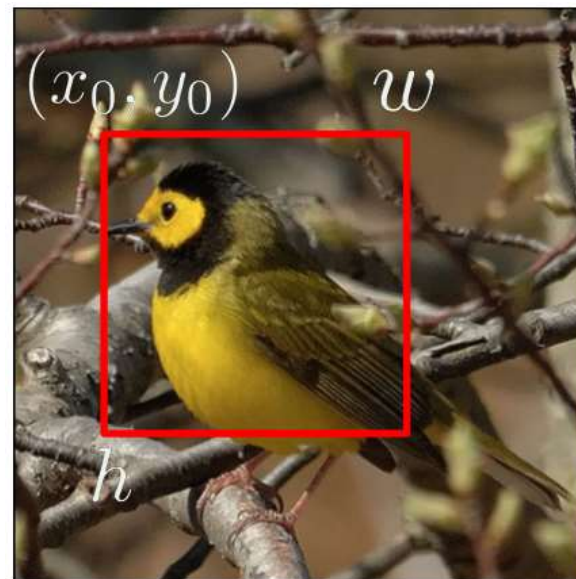
Что делать с дисперсией?

- Модель:
 - $\theta(x)$ – дифференцируемая функция
 - Распределение $p(w \mid \theta(x))$
 - Лосс $L(w)$ недифференцируемый
- Градиент: $\nabla_{\theta} \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} [\log p(w_i \mid \theta(x))] L(w_i)$
- Идея: baseline $\nabla_{\theta} \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} [\log p(w_i \mid \theta(x))] (L(w_i) - b)$
 - Почему?
$$\int p(w \mid \theta) \nabla_{\theta} [\log p(w \mid \theta(x))] b \, dw = \int \nabla_{\theta} p(w \mid \theta) b \, dw = b \nabla_{\theta} \int p(w \mid \theta) \, dw = 0$$
- Идея: дифференцируемый бейзлайн, зависящий от w
 - Компенсировать смещение репараметризацией бейзлайна

Другие способы сглаживания: Spatial Transformer

[Jaderberg et al., 2015]

- Параметрическая модель фрагмента (x_0, y_0, w, h)
- Сеть по картинке выдает параметры
- Нужно вырезать патч
- Можно обучать REINFORCE (но работает не очень) [Mnih et al. , 2014]
- **Spatial Transformer** – лучше
 - **Идея:** билинейная интерполяция дифференцируема
 - «Локальное внимание»
 - Расширение области определения



Slide credit:
Michael Figurnov

Дифференцируемая интерполяция

[Jaderberg et al., 2015]

- Билинейная интерполяция вычисляет цвет в нецелой точке (x, y)
- U – картинка, v – значение в (x, y)
- Основное наблюдение:

$$v = \sum_{i=1}^H \sum_{j=1}^W U_{ij} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$$

- Градиенты

$$\frac{dv}{dU_{ij}} = \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$$

$$\frac{dv}{dx} = \sum_{i=1}^H \sum_{j=1}^W U_{ij} \max(0, 1 - |y - j|) \begin{cases} 0, & \text{if } |x - i| \geq 1 \\ 1, & \text{if } x < i \\ -1, & \text{if } x \geq i \end{cases}$$

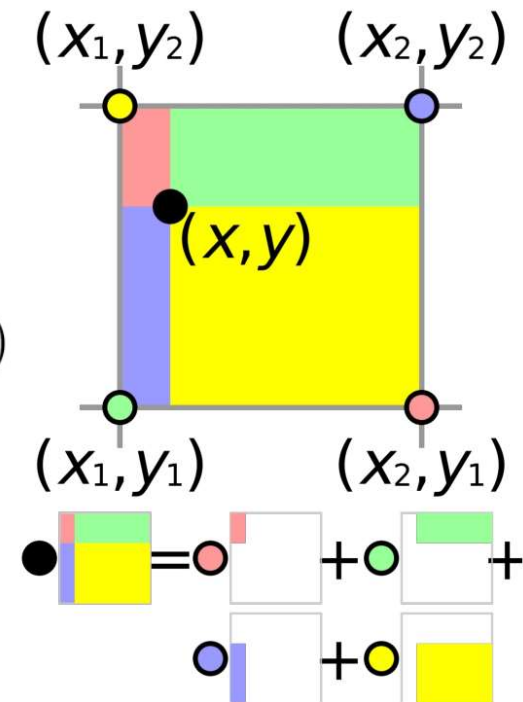
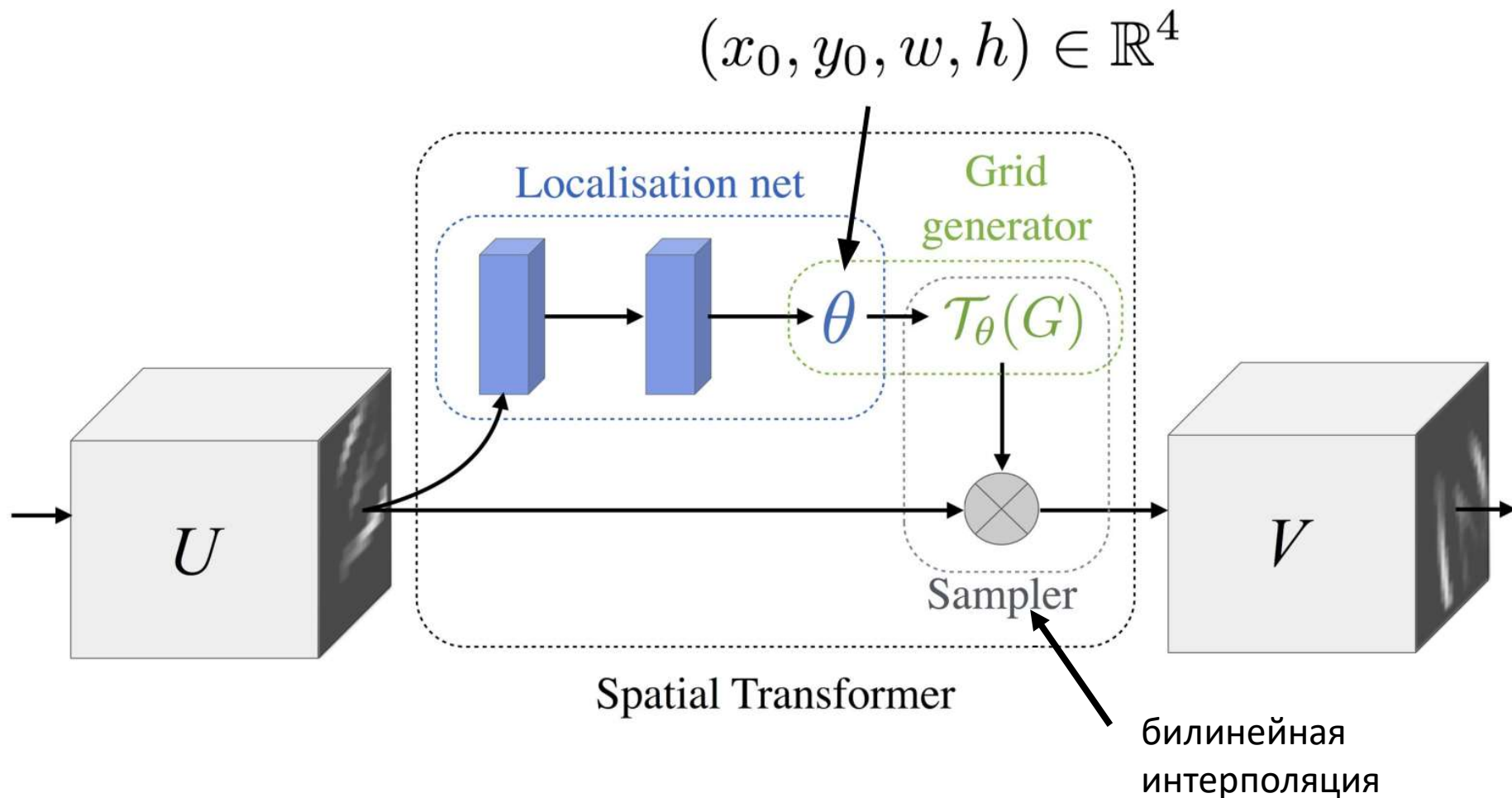


Image credit:
wikipedia

Slide credit:
Michael Figurnov

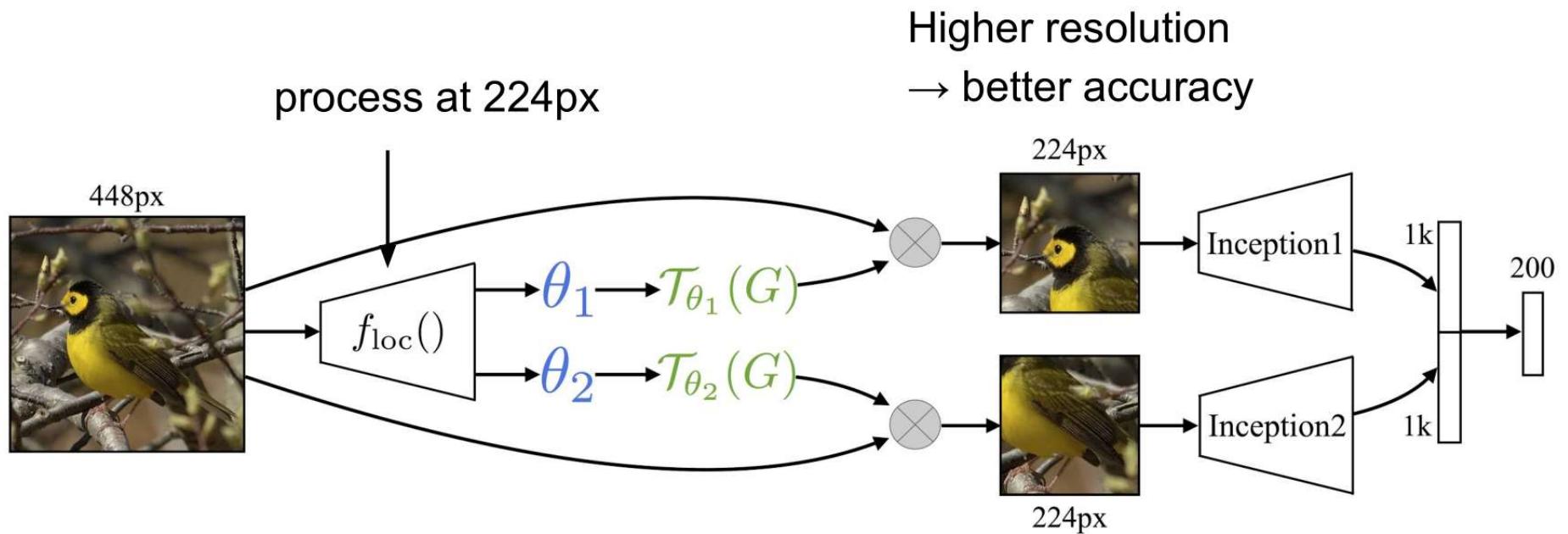
Другие способы сглаживания: Spatial Transformer

[Jaderberg et al., 2015]



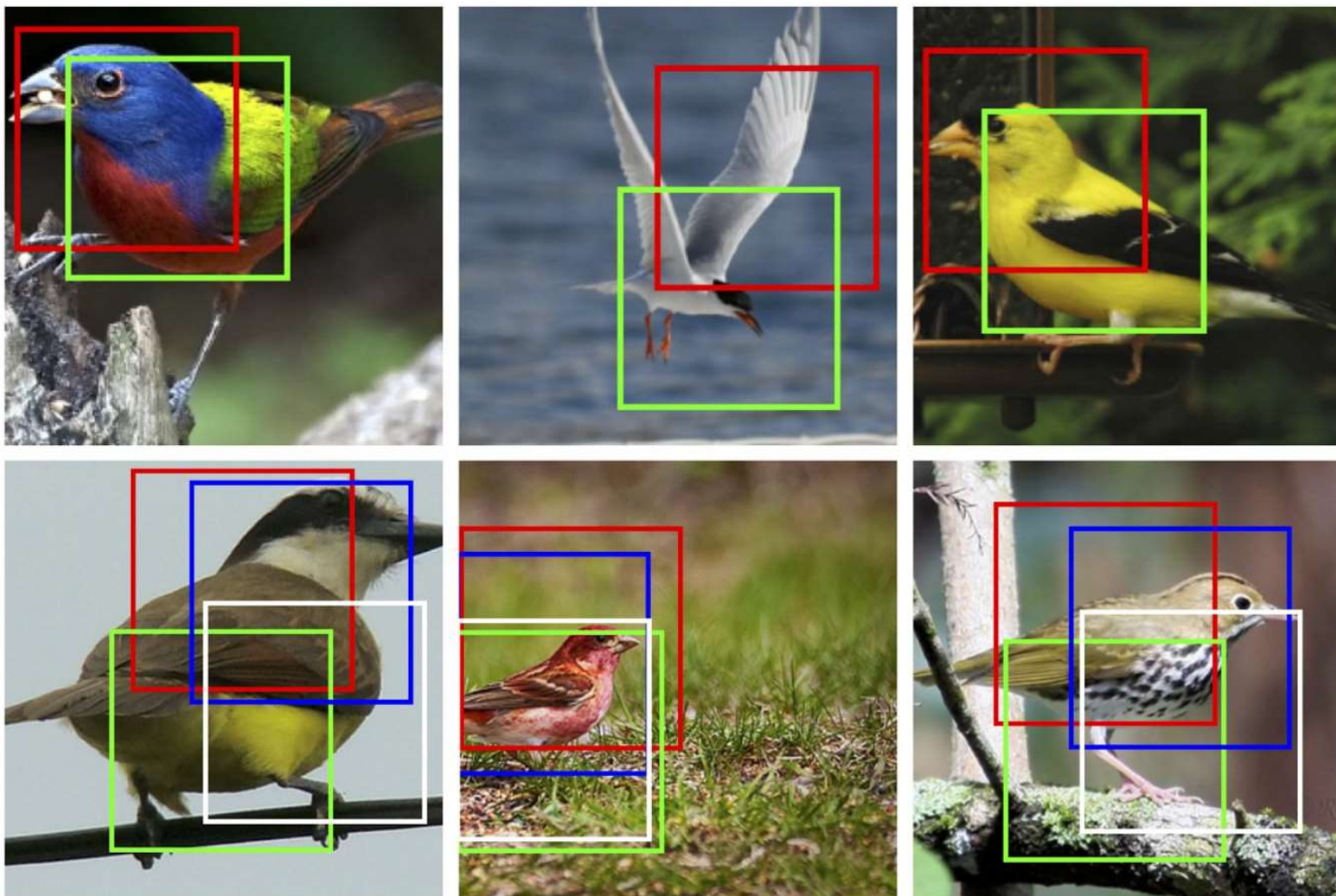
Spatial Transformer для классификации птиц

[Jaderberg et al., 2015]



Spatial Transformer для классификации птиц

[Jaderberg et al.,
2015]



Заключение

- Разные способы встраивать алгоритмы в нейросети
 - Структурный пулинг
 - Итерации алгоритма => слои нейросети
 - Прямое (аналитическое) дифференцирование по входу
- С недифференцируемыми функциями можно бороться
 - Основной способ – сведение к дифференцируемым
 - Игнорировать недифференцируемую операцию
 - Стохастические активации (репараметризация – хорошо работает)
 - Совсем недифференцируемо – RL
 - Все сложно, нестабильно, долго, но иногда возможно