

# Московский Государственный Университет

## ЕМ алгоритм для детектива

Выполнил: Курцев Д.В.

Группа: 417

Факультет Вычислительной математики и кибернетики

Кафедра Математических методов прогнозирования

Ноябрь 2022

# Постановка задачи

Дана выборка  $\mathbf{X} = \{\mathbf{X}_k\}_{k=1}^K$  сильно зашумленных черно-белых изображений размера  $H \times W$  пикселей. Каждое из этих изображений содержит один и тот же неподвижный фон и лицо преступника в неизвестных координатах, при этом лицо попадает в любое изображение целиком. Будем считать, что изображение лица имеет прямоугольную форму размера  $h \times w$  пикселей.

Также пусть шум на изображении независимым для каждого пикселя и принадлежащим нормальному распределению  $\mathcal{N}(0, s^2)$ , где  $s$  — стандартное отклонение. Таким образом для одного изображения имеем:

$$p(\mathbf{X}_k | \mathbf{d}_k, \boldsymbol{\theta}) = \prod_{ij} \begin{cases} \mathcal{N}(\mathbf{X}_k(i, j) | \mathbf{F}(i - d_k^h, j - d_k^w), s^2), & \text{если } (i, j) \in \text{faceArea}(\mathbf{d}_k) \\ \mathcal{N}(\mathbf{X}_k(i, j) | \mathbf{B}(i, j), s^2), & \text{иначе} \end{cases},$$

где  $\boldsymbol{\theta} = \{\mathbf{B}, \mathbf{F}, s^2\}$ ,  $\text{faceArea}(\mathbf{d}_k) = \{(i, j) \mid d_k^h \leq i \leq d_k^h + h - 1, d_k^w \leq j \leq d_k^w + w - 1\}$ .

Распределение на неизвестные координаты лица на изображении зададим общим для всех изображений с помощью матрицы параметров  $\mathbf{A} \in \mathbb{R}^{H-h+1, W-w+1}$  следующим образом:

$$p(\mathbf{d}_k | \mathbf{A}) = \mathbf{A}(d_k^h, d_k^w), \quad \sum_{ij} A(i, j) = 1,$$

где  $\mathbf{A}(i, j)$  — элемент матрицы  $\mathbf{A}$ .

В итоге имеем следующую совместную вероятностную модель:

$$p(\mathbf{X}, \mathbf{d} | \boldsymbol{\theta}, \mathbf{A}) = \prod_k p(\mathbf{X}_k | \mathbf{d}_k, \boldsymbol{\theta}) p(\mathbf{d}_k | \mathbf{A}).$$

Требуется решить задачу

$$p(\mathbf{X} | \boldsymbol{\theta}, \mathbf{A}) \rightarrow \max_{\boldsymbol{\theta}, \mathbf{A}}.$$

Для этого предлагается воспользоваться ЕМ-алгоритмом, то есть перейти к следующей задаче оптимизации нижней оценки на логарифм неполного правдоподобия:

$$\mathcal{L}(q, \boldsymbol{\theta}, \mathbf{A}) = \mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} | \boldsymbol{\theta}, \mathbf{A}) - \mathbb{E}_{q(\mathbf{d})} \log q(\mathbf{d}) \rightarrow \max_{q, \boldsymbol{\theta}, \mathbf{A}}$$

$$\mathcal{L}(q, \boldsymbol{\theta}, \mathbf{A}) = \sum_{\mathbf{d}} q(\mathbf{d}) \log p(\mathbf{X}, \mathbf{d} | \boldsymbol{\theta}, \mathbf{A}) - \sum_{\mathbf{d}} q(\mathbf{d}) \log q(\mathbf{d})$$

## Е-шаг

На Е-шаге вычисляется оценка на апостериорное распределение на координаты лица на изображениях:

$$q(\mathbf{d}) = p(\mathbf{d} \mid \mathbf{X}, \boldsymbol{\theta}, \mathbf{A}) = \prod_k p(\mathbf{d}_k \mid \mathbf{X}_k, \boldsymbol{\theta}, \mathbf{A}),$$

По формуле Байеса имеем:

$$p(\mathbf{d}_k \mid \mathbf{X}_k, \boldsymbol{\theta}, \mathbf{A}) = \frac{p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta}, \mathbf{A}) p(\mathbf{d}_k \mid \mathbf{A})}{\sum_{d_k^h} \sum_{d_k^w} p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta}, \mathbf{A}) p(\mathbf{d}_k \mid \mathbf{A})}$$

## М-шаг

На М-шаге вычисляется точечная оценка на параметры  $\boldsymbol{\theta}, \mathbf{A}$ :

$$\mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} \mid \boldsymbol{\theta}, \mathbf{A}) \rightarrow \max_{\boldsymbol{\theta}, \mathbf{A}}.$$

$$\begin{aligned} \mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} \mid \boldsymbol{\theta}, \mathbf{A}) &= \mathbb{E}_{q(\mathbf{d})} \log \prod_k p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta}) p(\mathbf{d}_k \mid \mathbf{A}) = \\ &= \mathbb{E}_{q(\mathbf{d})} \sum_k (\log p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta}) + \log p(\mathbf{d}_k \mid \mathbf{A})) = \\ &= \sum_{\mathbf{d}_k} \sum_k q(\mathbf{d}_k) (\log p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta}) + \log p(\mathbf{d}_k \mid \mathbf{A})) \quad (*) \end{aligned}$$

Выведем точечные оценки на параметры. Для этого продифференцируем полученное выражение и приравняем производные к нулю.

## A

Так как у нас имеется ограничение на  $\mathbf{A}$ :  $\sum_{i,j} A(i, j) = 1$ , запишем лагранжиан.  $p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta})$  не зависит от  $\mathbf{A}$ , поэтому не будем его писать, так как данное правдоподобие уйдёт при дифференцировании.

$$\begin{aligned} L &= \sum_{\mathbf{d}} \sum_k q(\mathbf{d}_k) \log \mathbf{A}(\mathbf{d}_k) - \lambda (\sum_{\mathbf{d}} \mathbf{A}(\mathbf{d}) - 1) \\ \frac{\partial L}{\partial \mathbf{A}(\mathbf{d})} &= \frac{\sum_k q(\mathbf{d}_k)}{\mathbf{A}(\mathbf{d})} - \lambda = 0 \quad \Rightarrow \quad \mathbf{A}(\mathbf{d}) = \frac{\sum_k q(\mathbf{d}_k)}{\lambda} \end{aligned}$$

$$\frac{\partial L}{\lambda} = \sum_d \mathbf{A}(\mathbf{d}) - 1 = 0 \Rightarrow \sum_d \mathbf{A}(\mathbf{d}) = 1$$

Тогда получим, что

$$1 = \sum_d \mathbf{A}(\mathbf{d}) = \sum_d \frac{\sum_k q(\mathbf{d}_k)}{\lambda} \Rightarrow \lambda = \frac{1}{\sum_d \sum_k q(\mathbf{d}_k)} = \frac{1}{\sum_k \sum_d q(\mathbf{d}_k)} = \frac{1}{\sum_k 1} = \frac{1}{K}$$

Таким образом оценка на  $\mathbf{A}$  равна:

$$\mathbf{A}(\mathbf{d}) = \frac{1}{K} \sum_k q(\mathbf{d}_k)$$

## F

Аналогично  $\mathbf{A}$ , в (\*) второе слагаемое не зависит от оставшихся параметров, поэтому для сокращения записи не будем его писать. Получим следующее выражение:

$$\begin{aligned} \sum_d \sum_k q(\mathbf{d}_k) \log p(\mathbf{X}_k | \mathbf{d}_k, \boldsymbol{\theta}) &= \sum_d \sum_k q(\mathbf{d}_k) \sum_{i,j} \log \mathcal{N}(\mathbf{X}_k(i, j) | \mu_{ij}, s^2) = \\ &= \sum_d \sum_k q(\mathbf{d}_k) \sum_{i,j} \left[ -\frac{1}{2} \log 2\pi s^2 - \frac{1}{2s^2} (\mathbf{X}_k(i, j) - \mathbf{B}(i, j))^2 \mathbb{1}\{i, j \notin face\} - \right. \\ &\quad \left. - \frac{1}{2s^2} (\mathbf{X}_k(i, j) - \mathbf{F}(i - d_k^h, j - d_k^w))^2 \mathbb{1}\{i, j \in face\} \right] \quad (**) \end{aligned}$$

Для начала продифференцируем (\*\*) по  $\mathbf{F}$ :

$$\begin{aligned} \frac{\partial(**)}{\partial \mathbf{F}(i, j)} &= \sum_d \sum_k q(\mathbf{d}_k) \frac{1}{2s^2} 2 (\mathbf{X}_k(i + d_k^h, j + d_k^w) - \mathbf{F}(i, j)) \mathbf{F}(i, j) = 0 \Rightarrow \\ &\Rightarrow \sum_d \sum_k q(\mathbf{d}_k) \mathbf{F}(i, j) = \sum_d \sum_k q(\mathbf{d}_k) \mathbf{X}_k(i + d_k^h, j + d_k^w) \end{aligned}$$

Таким образом оценка на  $\mathbf{F}$  равна:

$$\mathbf{F}(i, j) = \frac{1}{K} \sum_d \sum_k q(\mathbf{d}_k) \mathbf{X}_k(i + d_k^h, j + d_k^w)$$

## B

Теперь найдём оценку на  $\mathbf{B}$ :

$$\begin{aligned} \frac{\partial(**)}{\partial \mathbf{B}(i, j)} &= \sum_d \sum_k q(\mathbf{d}_k) \frac{1}{2s^2} 2 (\mathbf{X}_k(i, j) - \mathbf{B}(i, j)) \mathbf{B}(i, j) \mathbb{1}\{i, j \notin \text{face}\} = 0 \Rightarrow \\ \Rightarrow \sum_d \sum_k q(\mathbf{d}_k) \mathbf{B}(i, j) \mathbb{1}\{i, j \notin \text{face}\} &= \sum_d \sum_k q(\mathbf{d}_k) \mathbf{X}_k(i, j) \mathbb{1}\{i, j \notin \text{face}\} \end{aligned}$$

Таким образом оценка на  $\mathbf{B}$  равна:

$$\mathbf{B}(i, j) = \frac{\sum_d \sum_k q(\mathbf{d}_k) \mathbf{X}_k(i, j) \mathbb{1}\{i, j \notin \text{face}\}}{\sum_d \sum_k q(\mathbf{d}_k) \mathbb{1}\{i, j \notin \text{face}\}}$$

## S

Теперь найдём оценку на  $s$ :

$$\begin{aligned} \frac{\partial(**)}{\partial s} &= \sum_d \sum_k q(\mathbf{d}_k) \sum_{i, j} \left[ -\frac{1}{s} + \frac{1}{s^3} (\mathbf{X}_k(i, j) - \mathbf{B}(i, j))^2 \mathbb{1}\{i, j \notin \text{face}\} - \right. \\ &\quad \left. + \frac{1}{s^3} (\mathbf{X}_k(i, j) - \mathbf{F}(i - d_k^h, j - d_k^w))^2 \mathbb{1}\{i, j \in \text{face}\} \right] = \\ &= \frac{1}{s} (-HWK + \frac{1}{s^2}) \sum_d \sum_k \sum_{i, j} q(\mathbf{d}_k) [(\mathbf{X}_k(i, j) - \mathbf{B}(i, j))^2 \mathbb{1}\{i, j \notin \text{face}\} + \\ &\quad (\mathbf{X}_k(i, j) - \mathbf{F}(i - d_k^h, j - d_k^w))^2 \mathbb{1}\{i, j \in \text{face}\}] = 0 \end{aligned}$$

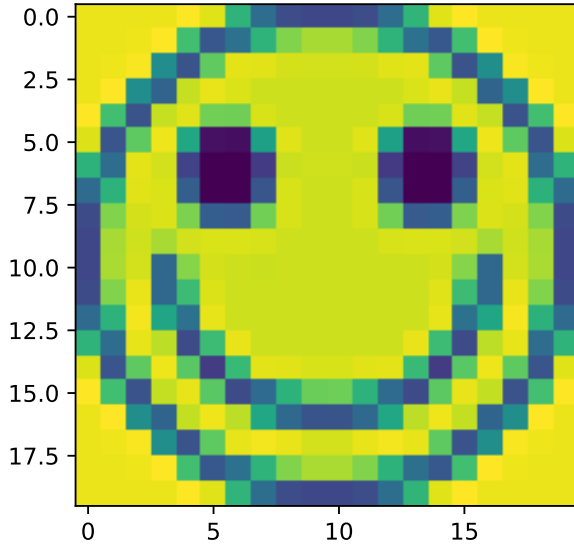
Таким образом оценка на  $s$  равна:

$$\begin{aligned} s^2 &= \frac{1}{HWK} \sum_d \sum_k \sum_{i, j} q(\mathbf{d}_k) [(\mathbf{X}_k(i, j) - \mathbf{B}(i, j))^2 \mathbb{1}\{i, j \notin \text{face}\} + \\ &\quad (\mathbf{X}_k(i, j) - \mathbf{F}(i - d_k^h, j - d_k^w))^2 \mathbb{1}\{i, j \in \text{face}\}] \end{aligned}$$

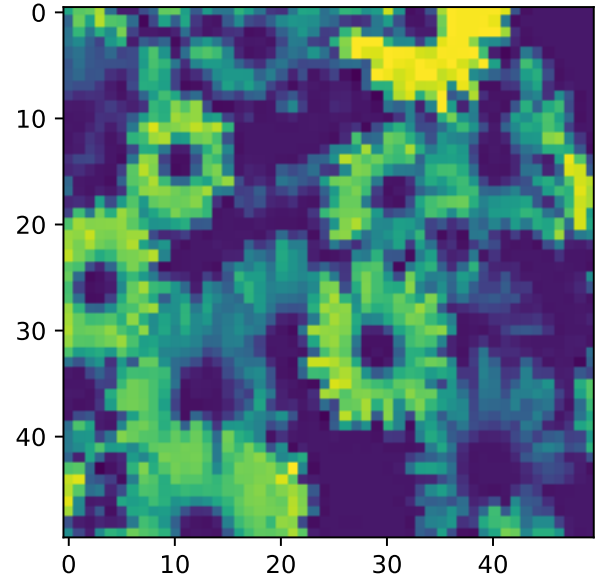
## Эксперименты

Посмотрим, как работает написанный алгоритм. Для этого возьмём произвольный объект - смайлик, который необходимо будет обнаружить на фоне цветочков.

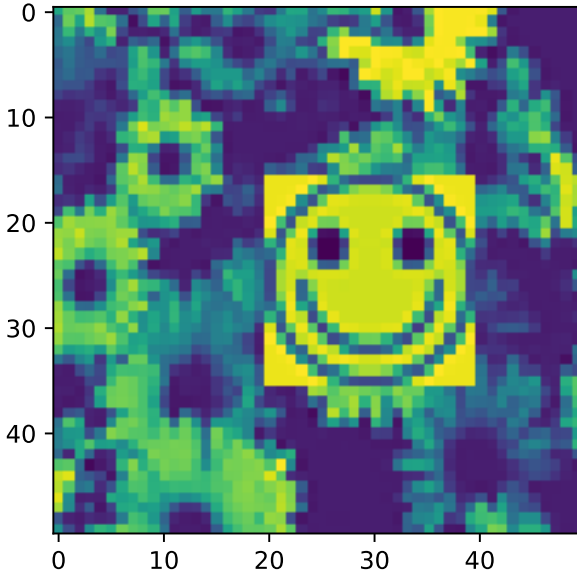
Создадим выборку из размещённых в разных местах на фоне цветочков смайликов. Пример одной картинке с шумом и без представлен ниже.



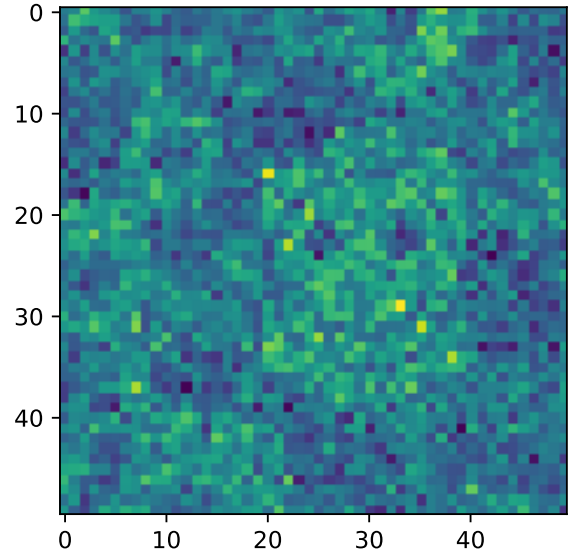
(a) Объект



(b) Фон



(a) Объект на фоне



(b) Зашумлённый объект на фоне

Для начала посмотрим, как на результат работы влияет начальная инициализация параметров алгоритма. Будем использовать 3 стратегии: равномерное на отрезке от нуля до наибольшего значения из выборки -  $\mathcal{U}[0, \max \mathbf{X}_k(i, j)]$ , нормальное, где мат ожиданием является среднее значение из выборки с дисперсией 50 -  $\mathcal{N}(\text{mean} \mathbf{X}_k(i, j), 50)$  и просто константным приближением - среднее. Матрицу  $\mathbf{A}$  во всех случаях отнормируем на сумму.

Для обучения возьмём выборку из 100 изображений зашумлённых нормальных шумом с  $s = 100$ .

Результаты работы алгоритма можно увидеть на Рис. 3 и 4

Как можно заметить, результат работы метода практически не зависит от

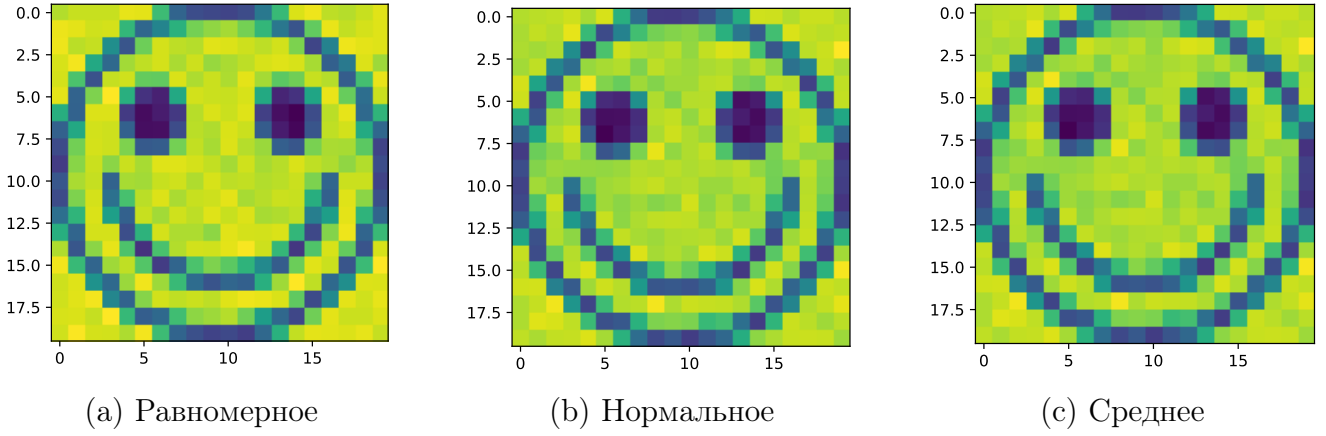


Рис. 3: Результат ЕМ для объекта

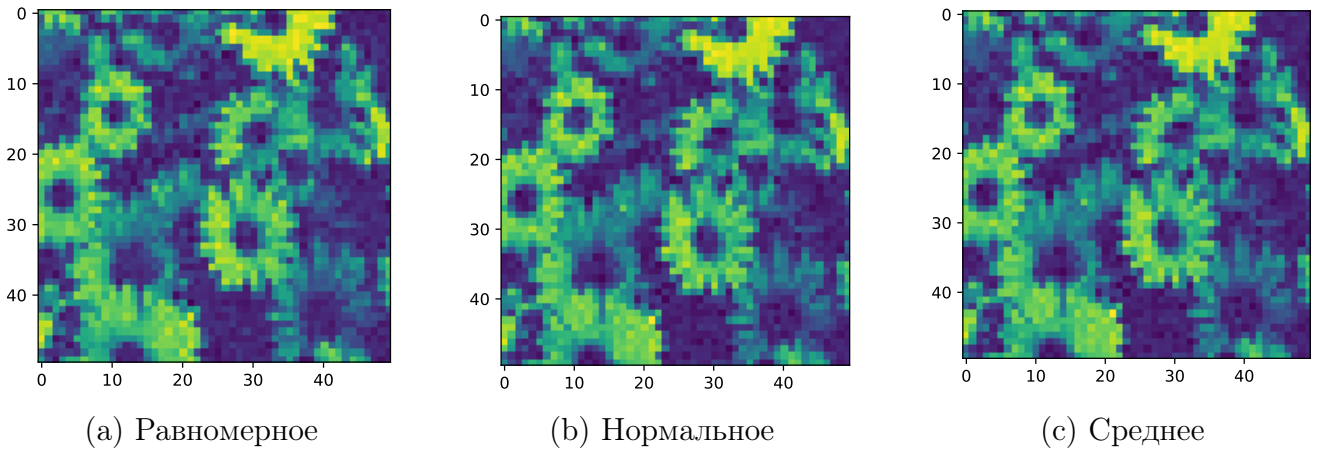


Рис. 4: Результат ЕМ для фона

начальной инициализации. Неполное правдоподобие для равномерной инициализации чуть больше, чем для остальных двух (-1423467 против -1423942). Видно, что для первой стратегии результат получился чуть более точным, чем для оставшихся. Так же заметно, что последние 2 инициализации работают практически одинаково. Вероятно, это связано с тем, что для нормального распределения была выбрана не очень большая дисперсия, поэтому значения параметров не очень сильно отклонились от среднего значения.

Будем всюду далее использовать первую стратегию для задания начальных значений параметров.

Теперь исследуем работу ЕМ алгоритма в зависимости от размера обучающей выборки. Возьмём 30, 50 и 100 зашумлённых изображений.

Результаты работы метода приведены на Рис. 5 и 6. Можно сделать вполне логичный вывод, что чем больше обучающая выборка, тем точнее работает

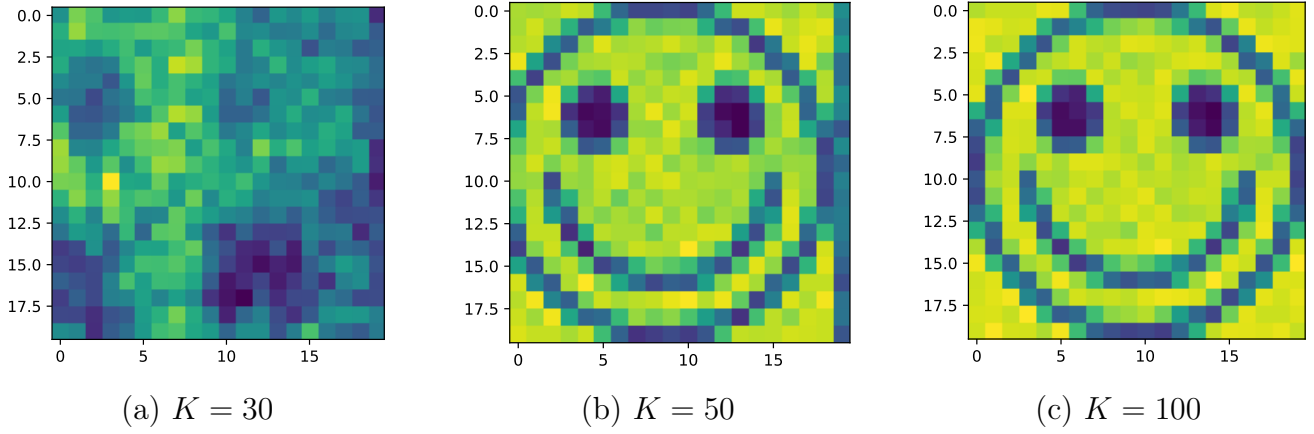


Рис. 5: Результат ЕМ для объекта

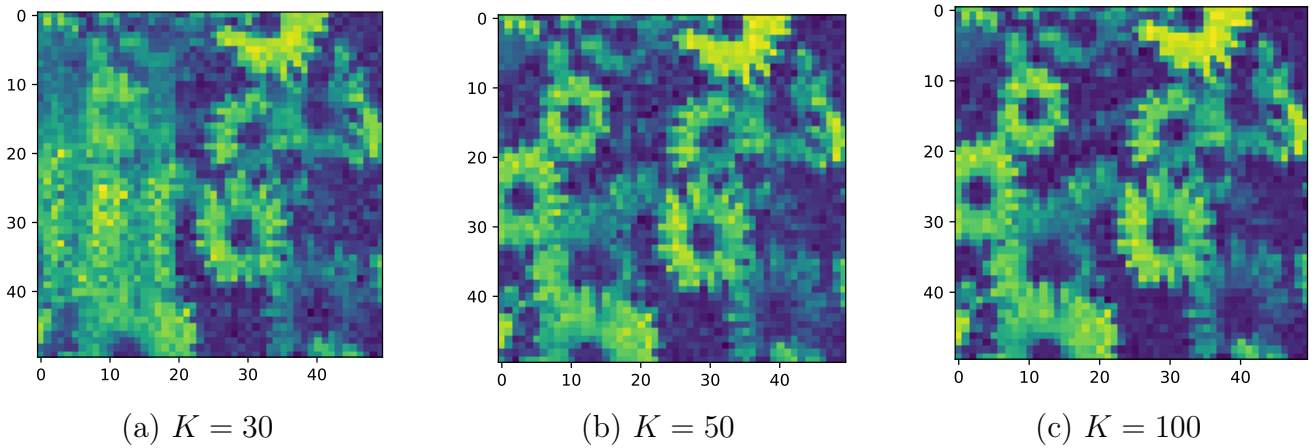


Рис. 6: Результат ЕМ для фона

алгоритм. Видно, что на объёме из 30 изображений, вообще невозможно определить объект и часть фона. Возможно в обучении состояли картинки, на которых смайлик находился на левой стороне изображения. Для выборки объёма 50 ЕМ уже вполне хорошо справляется. Однако полученные картинки чуть более размыты, чем для случая, когда в обучении использовалось 100 изображений.

Нормированное на объём выборки правдоподобие равняется  $\mathcal{L}_{30} = -14403.3$ ,  $\mathcal{L}_{50} = -14239.2$ ,  $\mathcal{L}_{100} = -14234.6$ . Видно, что чем лучше алгоритм справляется с задачей, тем выше правдоподобие.

Далее посмотрим, как на качество влияет зашумлённость изображений. Ранее рассматривался шум с  $s=100$ . Теперь возьмём обучающую выборку размером из 50 картинок с уровнем зашумлённости  $s \in \{50, 150, 200\}$

Результаты приведены на Рис. 7 и 8. Как можно заметить, чем больше шума в данных, тем сложнее алгоритму найти верное решение. Смайл и фон с



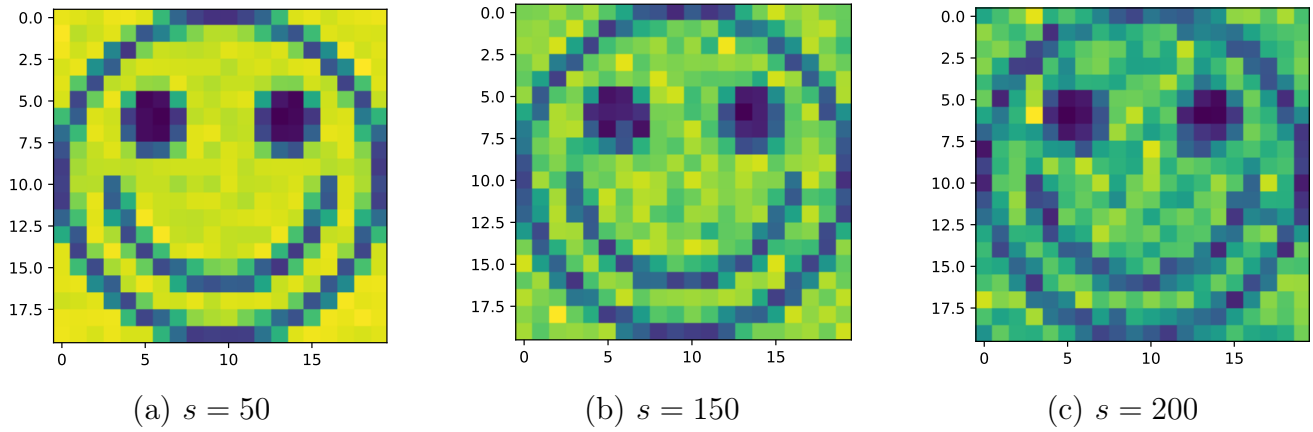


Рис. 7: Результат ЕМ для объекта

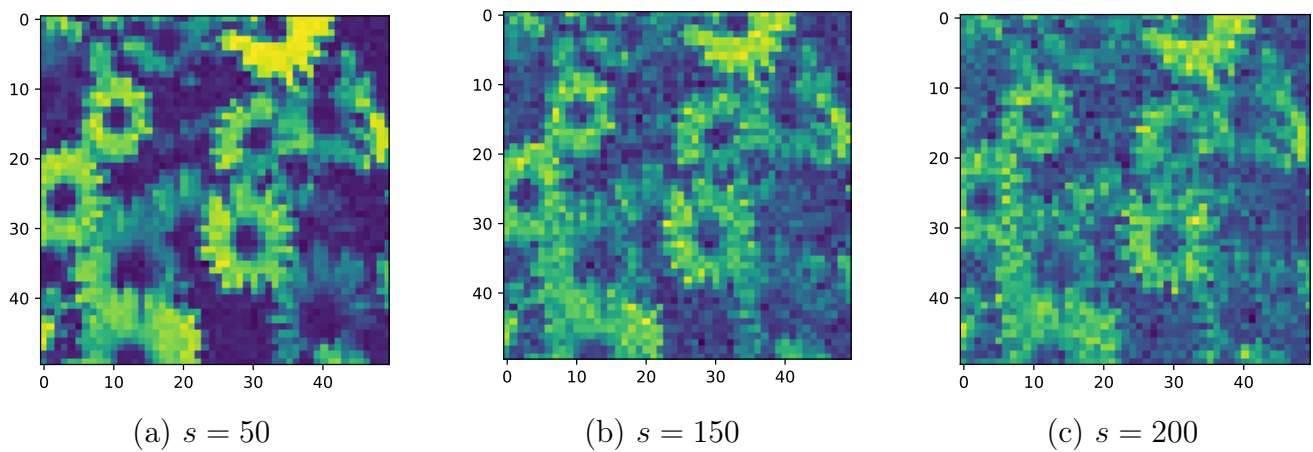


Рис. 8: Результат ЕМ для фона

$s = 200$ , получаются очень размытыми. Зато для  $s = 50$  картинки получаются гораздо чётче, чем был ранее.

Таким образом можно заключить, что чем больше наши данные зашумлены, тем больше должен быть объём обучающей выборки, для корректной работы ЕМ алгоритма

Теперь применим данный алгоритм для зашумлённых снимков преступника. Возьмём данные объёма 100, 300 и 500 изображений.

Результаты представлены на Рис. 9 и 10. Как и ранее здесь можно сделать аналогичный вывод. Чем больше изображений использовалось в обучении, тем точнее получается результат. Маленького объёма обучающей выборки не хватает для того, чтобы найти нужные закономерности в данных и оценить параметры.

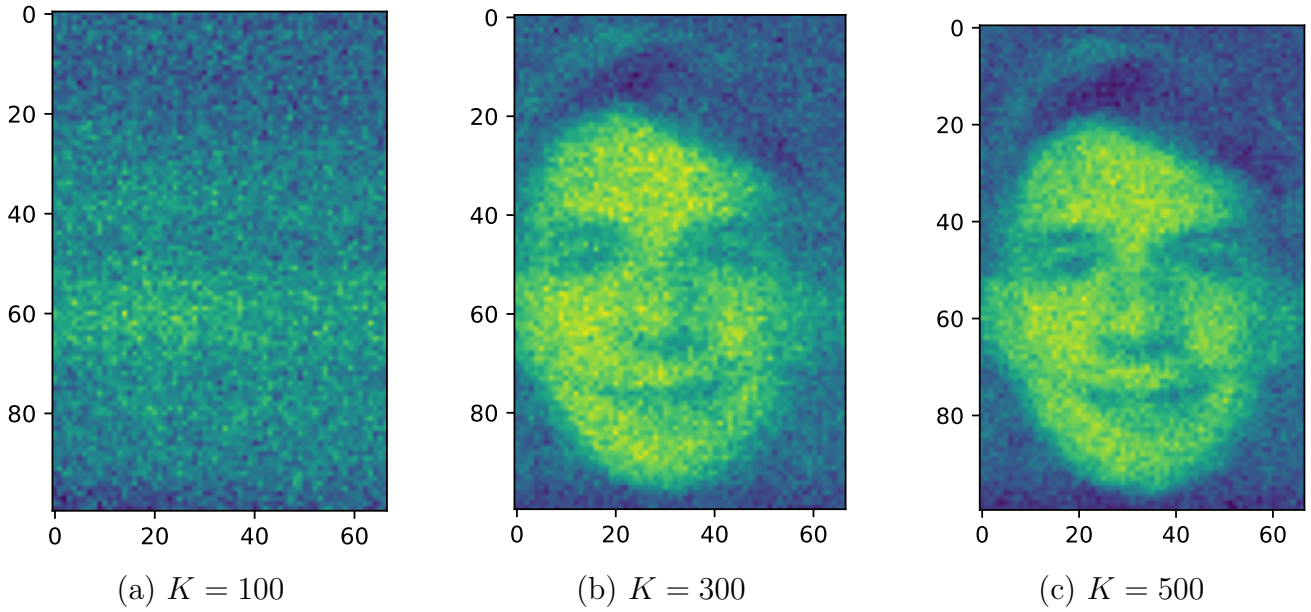


Рис. 9: Результат ЕМ для фото

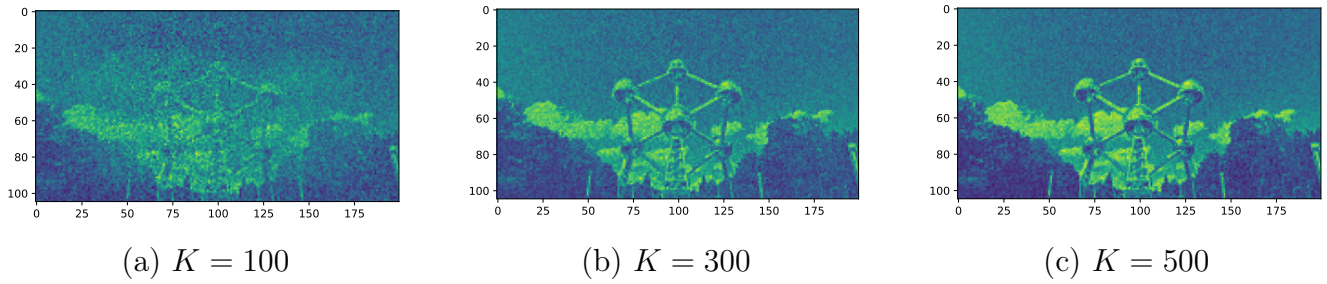


Рис. 10: Результат ЕМ для фона

Нормированное на объём выборки правдоподобие равняется  $\mathcal{L}_{100} = -128937.8$ ,  $\mathcal{L}_{300} = -128874.7$ ,  $\mathcal{L}_{500} = -128887.6$ . Опять же видим, что чем больше объём выборки, тем больше правдоподобие и соответственно качество изображений.

Посмотрим на качество и время работы ЕМ и hard ЕМ этих данных. (На сгенерированных результаты получаются практически одинаковыми).

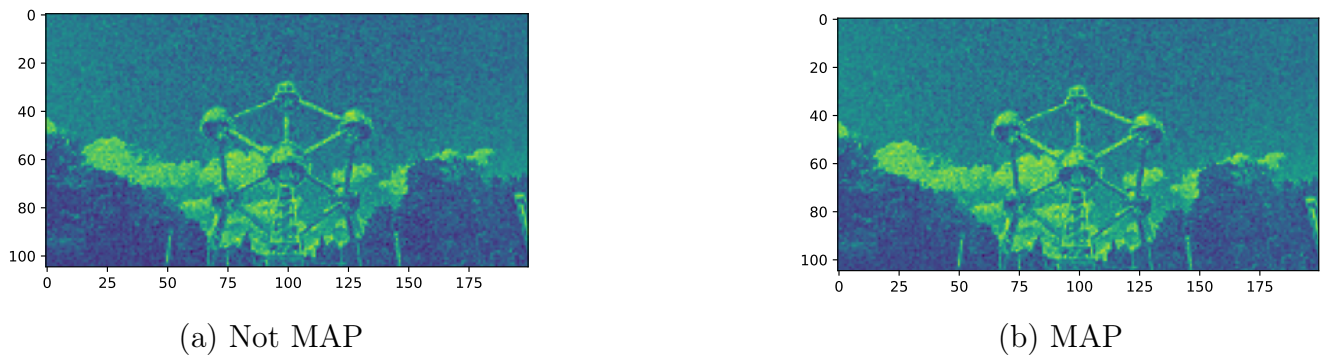


Рис. 11: Результат ЕМ для фона

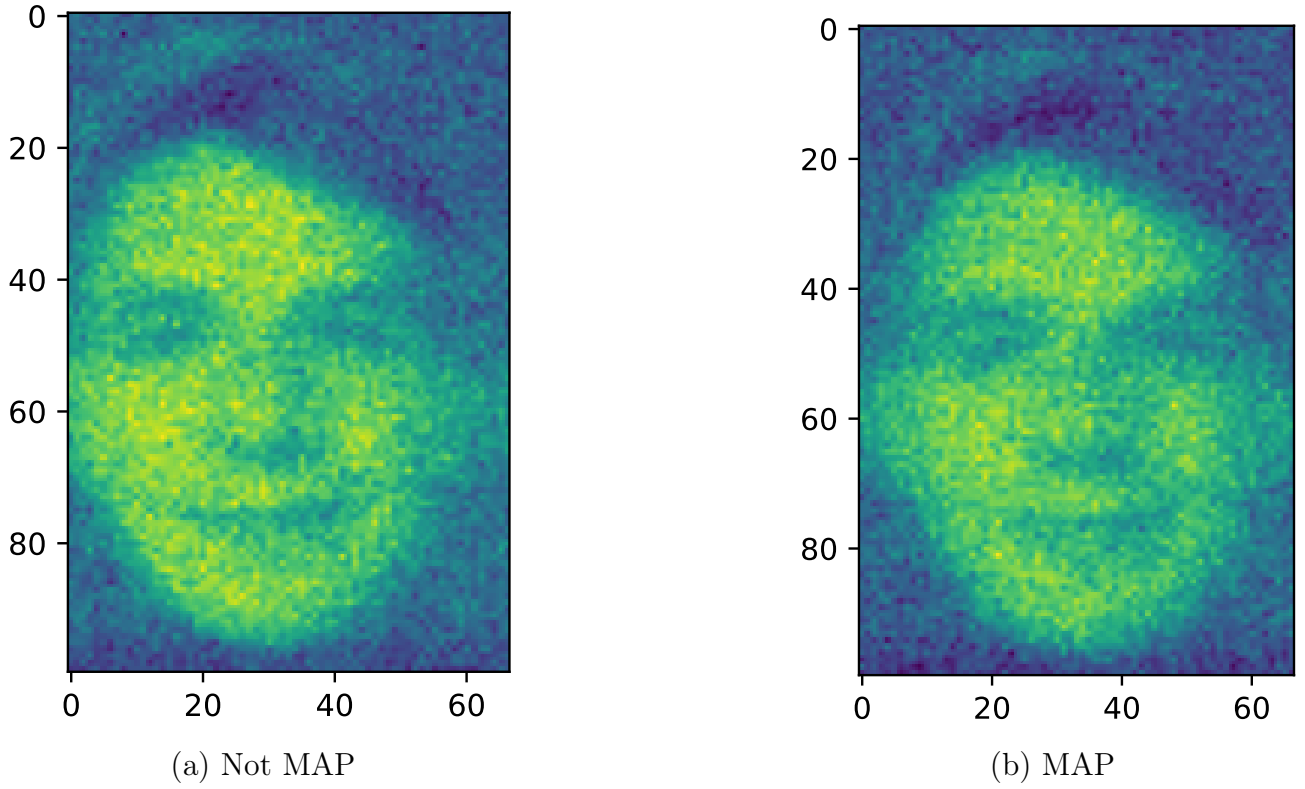


Рис. 12: Результат ЕМ для фото

Как можно заметить из Рис. 11 и 12 классическая реализация ЕМ алгоритма работает чуть лучше. Его модернизация (hard ЕМ) даёт более размытое изображение. Это можно объяснить тем, что мы делаем грубую аппроксимацию на латентные переменные. В связи с чем падает качество (правдоподобие меньше):  $\mathcal{L}_{EM} = -128874.7$ ,  $\mathcal{L}_{hard} = -128877.2$ . Однако hard ЕМ при сопоставимом качестве работы сходится гораздо быстрее:  $t_{EM} = 1min\ 48sec$ ,  $t_{hard} = 22sec$ .

Из модификаций можно предложить выход не по  $\mathcal{L}$ , а по норме  $F$ . То есть в качестве критерия останова использовать следующее условие:  $\|F^{(t+1)} - F^{(t)}\| \leq tol$ .

Так же усовершенствуя идею hard ЕМ, можно выбирать не 1 точку, а  $k$  наиболее вероятных точек и преобразовать  $q(\mathbf{d}_k)$  так, чтобы вся вероятностная масса была сосредоточена в них.

Р.с. Ник очень невнимательный. Максимально невнимательный...