

Московский государственный университет имени М. В. Ломоносова

Байсовские методы машинного обучения

Задание 2. EM алгоритм для детектива

Выполнил:

студент 4 курса 417 группы

Бабичев Дмитрий Олегович

Москва, 2017

Содержание

1	Описание модели	2
2	Необходимые формулы	3
3	Краткий анализ	5
3.1	Пункт 1	5
3.2	Пункт 2	7
3.3	Пункт 3	9
3.4	Пункт 4	11
3.5	Пункт 5	12

1 Описание модели

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. Значения h, w в выданных данных указаны в описании задания в anytask.

Введем следующие обозначения:

- $\mathbf{X}_k(i, j)$ – пиксель k -ого изображения;
- $\mathbf{B} \in \mathbb{R}^{H \times W}$ – изображение чистого фона без лица преступника, $\mathbf{B}(i, j)$ – пиксель этого изображения;
- $\mathbf{F} \in \mathbb{R}^{h \times w}$ – изображение лица преступника, $\mathbf{F}(i, j)$ – пиксель этого изображения;
- $\mathbf{d}_k = (d_k^h, d_k^w)$ – координаты верхнего левого угла изображения лица на k -ом изображении (d_k^h – по вертикали, d_k^w – по горизонтали), $\mathbf{d} = (\mathbf{d}_1, \dots, \mathbf{d}_K)$ – набор координат для всех изображений выборки.

Также будем считать 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:

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

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 $\mathbf{A} \in \mathbb{R}^{(H-h+1) \times (W-w+1)}$ следующим образом:

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

где $\mathbf{A}(i, j) = a_{ij}$ – элемент матрицы \mathbf{A} .

В итоге имеем следующую совместную вероятностную модель:

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

2 Необходимые формулы

Попробуем расписать подробно классический ЕМ-алгоритм для этой задачи. В текущей формулировке скрытыми переменными являются координаты левого верхнего угла изображения с лицом. На Е-шаге нужно вычислить их апостериорную оценку $q(\mathbf{d})$:

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

Воспользовавшись формулой условной вероятности $p(A|B) = \frac{p(A,B)}{p(B)}$ и правилом суммирования $p(B) = \int_A p(A,B)dA$, можно получить следующее выражение $p(\mathbf{d}_k \mid \mathbf{X}_k, \boldsymbol{\theta}, \mathbf{A})$:

$$p(\mathbf{d}_k \mid \mathbf{X}_k, \boldsymbol{\theta}, \mathbf{A}) = \frac{p(\mathbf{X}_k, \mathbf{d}_k \mid \boldsymbol{\theta}, \mathbf{A})}{p(\mathbf{X}_k \mid \boldsymbol{\theta}, \mathbf{A})} = \frac{p(\mathbf{X}_k \mid \mathbf{d}_k, \boldsymbol{\theta})p(\mathbf{d}_k \mid \mathbf{A})}{\sum_{i=0}^{H-h} \sum_{j=0}^{W-w} p(\mathbf{X}_k \mid \mathbf{d}_{ij} = (i, j), \boldsymbol{\theta})p(\mathbf{d}_{ij} = (i, j) \mid \mathbf{A})}.$$

При фиксированных параметрах $\boldsymbol{\theta}, \mathbf{A}$ все необходимые для вычислений величины нам известны.

На М-шаге нам нужно вычислить точечные оценки $\boldsymbol{\theta}, \mathbf{A}$:

$$\mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} \mid \boldsymbol{\theta}, \mathbf{A}) \rightarrow \max_{\boldsymbol{\theta}, \mathbf{A}}, \quad \sum_{i=1}^H \sum_{j=1}^W a_{ij} = 1.$$

Распишем подробнее:

$$\begin{aligned} \mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} \mid \boldsymbol{\theta}, \mathbf{A}) &= \mathbb{E}_{q(\mathbf{d})} \sum_{k=1}^K \log p(\mathbf{X}_k, \mathbf{d}_k \mid \boldsymbol{\theta}, \mathbf{A}) = \\ &= \sum_{k=1}^K \mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}_k, \mathbf{d}_k \mid \boldsymbol{\theta}, \mathbf{A}) = \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q^k(\mathbf{d}_{ij}) \log p(\mathbf{X}_k, \mathbf{d}_{ij} \mid \boldsymbol{\theta}, \mathbf{A}) = \\ &= \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k (\log p(\mathbf{X}_k \mid \mathbf{d}_{ij}, \boldsymbol{\theta}) + \log a_{ij}) \end{aligned}$$

Очевидно, что пытаться максимизировать по a_{ij} бесполезно, ведь чем эти коэффициенты больше, тем лучше. Но к счастью, у нас задача с ограничениями в виде равенства. Запишем для нее Лагранжиан:

$$L(\boldsymbol{\theta}, \mathbf{A}, \lambda) = \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k (\log p(\mathbf{X}_k \mid \mathbf{d}_{ij}, \boldsymbol{\theta}) + \log a_{ij}) - \lambda \left(\sum_{i=1}^H \sum_{j=1}^W a_{ij} - 1 \right).$$

Перейдем к системе:

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial a_{ij}} = \frac{\sum_{k=1}^K q_{ij}^k}{a_{ij}} - \lambda = 0 \Rightarrow a_{ij} = \frac{\sum_{k=1}^K q_{ij}^k}{\lambda}, \\ \frac{\partial L}{\partial \boldsymbol{\theta}} = 0 \\ \frac{\partial L}{\partial \lambda} = \sum_{i=1}^H \sum_{j=1}^W a_{ij} - 1 = 0 \Rightarrow \sum_{i=1}^H \sum_{j=1}^W a_{ij} = 1, \end{array} \right. \quad \begin{array}{l} \text{распишем позже,} \\ \Rightarrow \lambda = \sum_{k=1}^K \sum_{i=1}^H \sum_{j=1}^W q_{ij}^k \Rightarrow a_{ij} = \frac{\sum_{k=1}^K q_{ij}^k}{K}. \end{array}$$

Вспомним, что $\boldsymbol{\theta} = \{\mathbf{B}, \mathbf{F}, s^2\}$. Трудностей с дифференцированием суммы логарифмов двух нормальных распределений нет, лично у меня возникла проблема с адекватной записью пределов суммирования. Если кратко, то b_{ij} и f_{ij} представляют собой усреднения по всем изображениям и по всевозможным положениям лица на изображениях соответствующих элементов $\mathbf{X}_k(i, j)$ с весами q_{ij}^k . Запишем наш Лагранжиан без независимых от $\boldsymbol{\theta}$ членов (BG - background, FA - face area):

$$\begin{aligned} L(\mathbf{B}, \mathbf{F}, s) &= \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k \log p(\mathbf{X}_k | \mathbf{d}_{ij}, \boldsymbol{\theta}) = \\ &= \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k \left[\left(\sum_{\substack{l,m \in \\ FA(i,j)}} -\log s - \frac{1}{2} 2\pi - \frac{1}{2s^2} (x_{l,m}^k - f_{l-i,m-j})^2 \right) + \right. \\ &\quad \left. + \left(\sum_{\substack{l,m \in \\ BG(i,j)}} -\log s - \frac{1}{2} 2\pi - \frac{1}{2s^2} (x_{l,m}^k - b_{l,m})^2 \right) \right]. \end{aligned}$$

Распишем производные по элементам матриц \mathbf{B}, \mathbf{F} и по s :

$$\begin{cases} \frac{\partial L}{\partial b_{l,m}} = \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} \frac{q_{ij}^k}{s^2} (x_{l,m}^k - b_{l,m}) I[l, m \in BG(i, j)] = 0, \\ \frac{\partial L}{\partial f_{l,m}} = \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} \frac{q_{ij}^k}{s^2} (x_{i+l-1, j+m-1}^k - f_{l,m}) I[l, m \in FA(i, j)] = 0, \\ \frac{\partial L}{\partial s} = \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} \frac{q_{ij}^k}{s} \left(\frac{1}{s^2} \left(\sum_{\substack{l,m \in \\ FA(i,j)}} (x_{l,m}^k - f_{l-i,m-j})^2 + \sum_{\substack{l,m \in \\ BG(i,j)}} (x_{l,m}^k - b_{l,m})^2 \right) - HW \right) = 0 \end{cases} \implies$$

$$\begin{cases} b_{l,m} = \frac{\sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k x_{l,m}^k I[l, m \in BG(i, j)]}{\sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k I[l, m \in BG(i, j)]} \\ f_{l,m} = \frac{\sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k x_{i+l-1, j+m-1}^k I[l, m \in FA(i, j)]}{\sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k I[l, m \in FA(i, j)]} \\ s^2 = \frac{\sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k \left(\sum_{\substack{l,m \in \\ FA(i,j)}} (x_{l,m}^k - f_{l-i,m-j})^2 + \sum_{\substack{l,m \in \\ BG(i,j)}} (x_{l,m}^k - b_{l,m})^2 \right)}{KHW} \end{cases}$$

Осталось вывести лишь формулу нижней оценки $\mathcal{L}(q, \boldsymbol{\theta}, \mathbf{A})$:

$$\begin{aligned} \mathcal{L}(q, \boldsymbol{\theta}, \mathbf{A}) &= \mathbb{E}_{q(\mathbf{d})} \log p(\mathbf{X}, \mathbf{d} | \boldsymbol{\theta}, \mathbf{A}) - \mathbb{E}_{q(\mathbf{d})} \log q(\mathbf{d}) = \\ &= \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k (\log p(\mathbf{X}_k | \mathbf{d}_{ij}, \boldsymbol{\theta}) + \log a_{ij}) - \sum_{k=1}^K \sum_{i=1}^{H-h} \sum_{j=1}^{W-w} q_{ij}^k \log q_{ij}^k, \end{aligned}$$

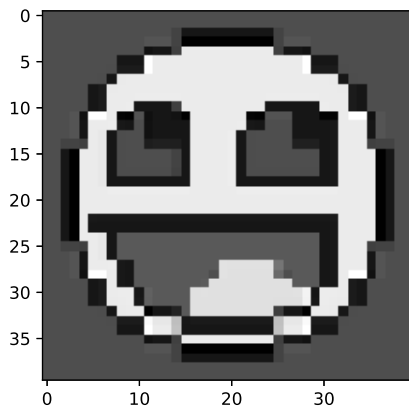
где

$$\log p(\mathbf{X}_k | \mathbf{d}_{ij}, \boldsymbol{\theta}) = \left[\left(\sum_{\substack{l,m \in \\ FA(i,j)}} -\log s - \frac{1}{2} 2\pi - \frac{1}{2s^2} (x_{l,m}^k - f_{l-i,m-j})^2 \right) + \left(\sum_{\substack{l,m \in \\ BG(i,j)}} -\log s - \frac{1}{2} 2\pi - \frac{1}{2s^2} (x_{l,m}^k - b_{l,m})^2 \right) \right].$$

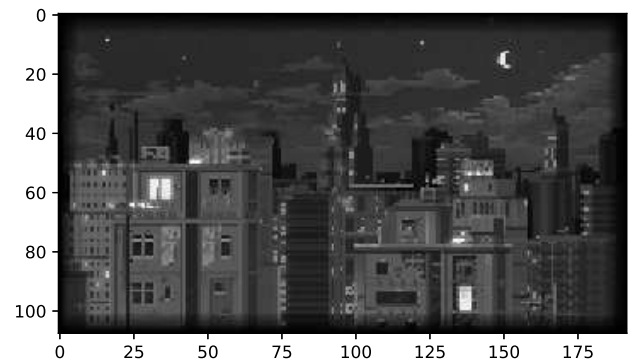
Зная все необходимые распределения вычисление этой оценки не составит труда.

Чтобы получить так называемый hard EM, необходимо преобразовать $q(\mathbf{d})$ так, что для каждого изображения \mathbf{X}_k оценка $q(\mathbf{d}_k)$ принимает значений 1 только в точке максимума апостериорного распределения $p(\mathbf{d}_k | \mathbf{X}_k, \boldsymbol{\theta}, \mathbf{A})$. Можно конечно заного все расписать, введя новое обозначение для этого максимума, но на деле же в общем виде формулы не поменяются, добавится лишь один шаг между E и M - нахождение максимума и приведение каждого $q(\mathbf{d}_k)$ к соответствующему виду.

3 Краткий анализ



face



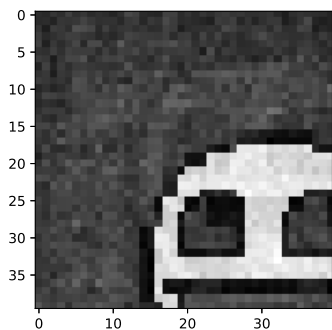
background

Пункт 1

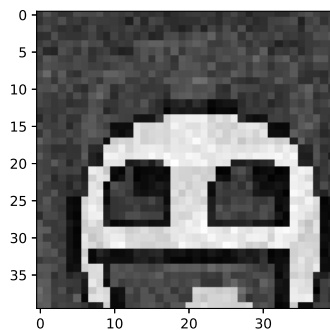
Протестируйте полученный EM алгоритм на сгенерированных данных. Сильно ли влияет начальное приближение на параметры на результаты работы? Стоит ли для данной задачи запускать EM алгоритм из разных начальных приближений?

В текущей реализации для каждого запуска перед первым E-шагом случайным образом генерируются величины \mathbf{F} , \mathbf{B} , \mathbf{A} , s . Посмотрим на результаты работы алгоритма после 12 запусков ($s_{real} = 30$):

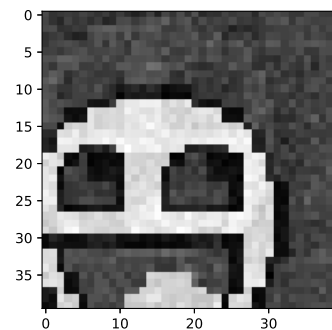
Очень странный результат. После каждого запуска смайлик вполне можно опознать, но вот его расположение меняется. То есть говорить о зависимости результата от начального приближения можно.



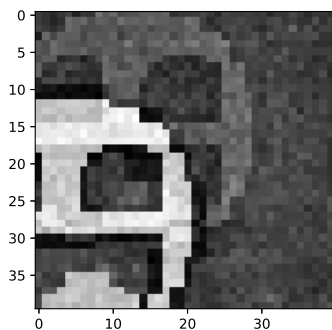
restart 1



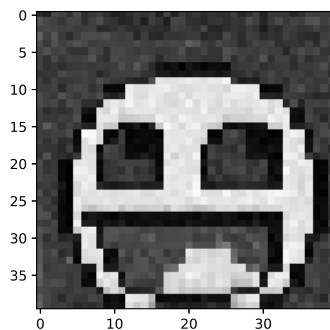
restart 2



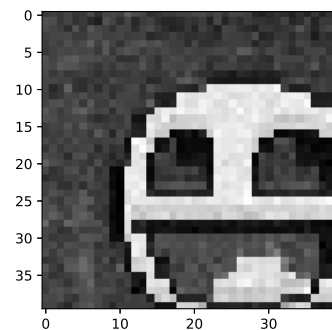
restart 3



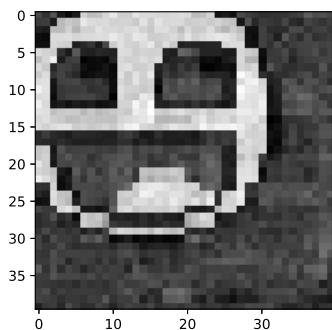
restart 4



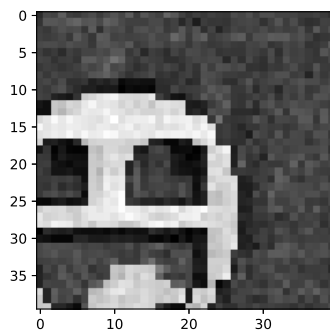
restart 5



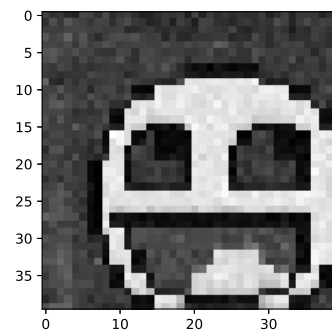
restart 6



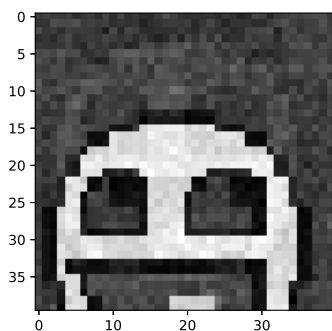
restart 7



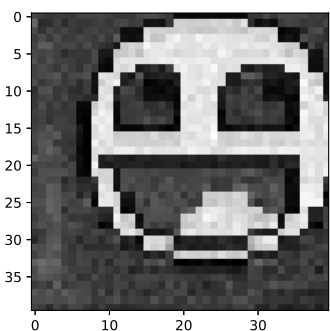
restart 8



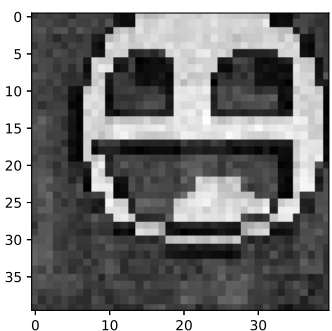
restart 9



restart 10



restart 11

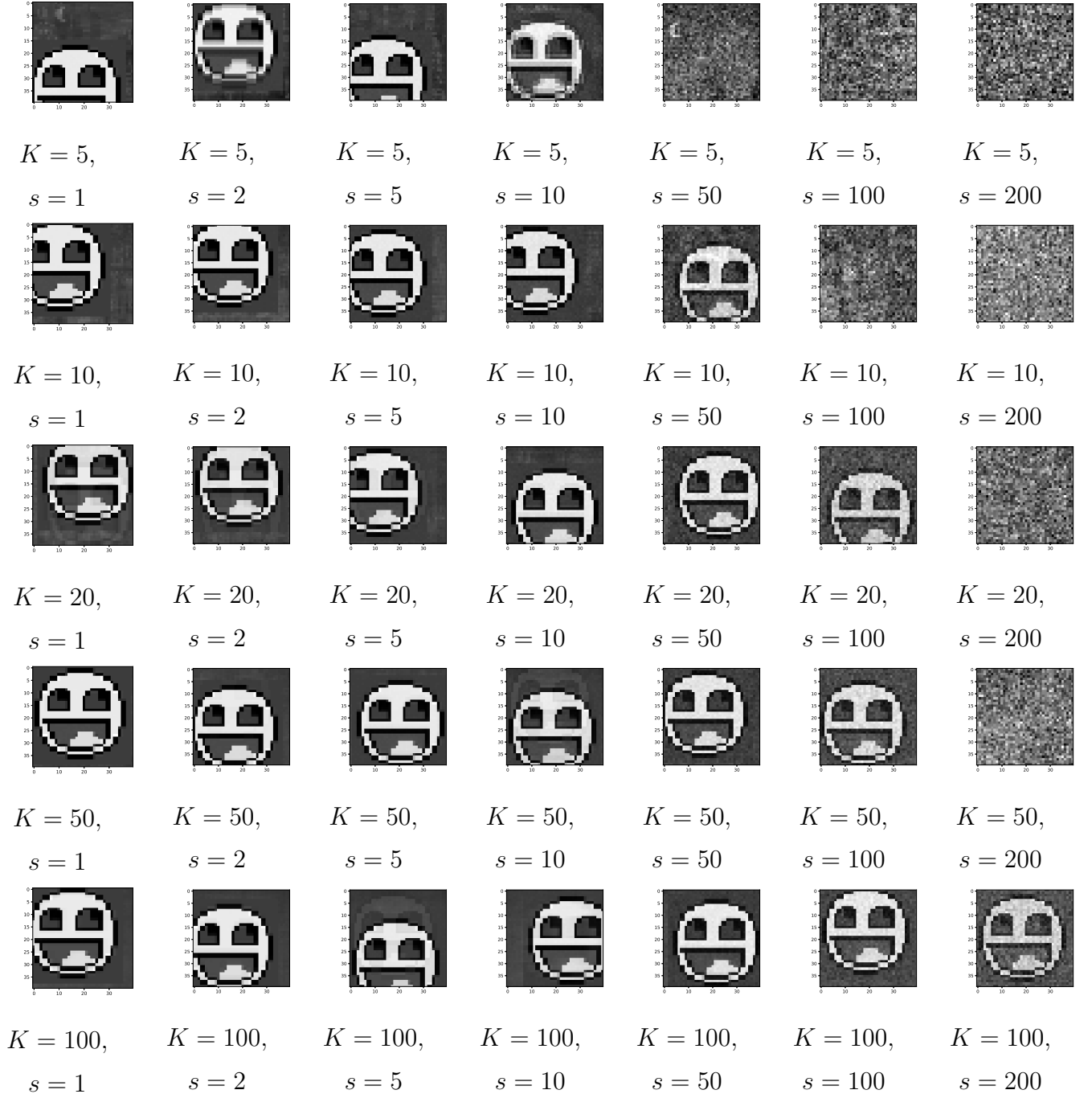


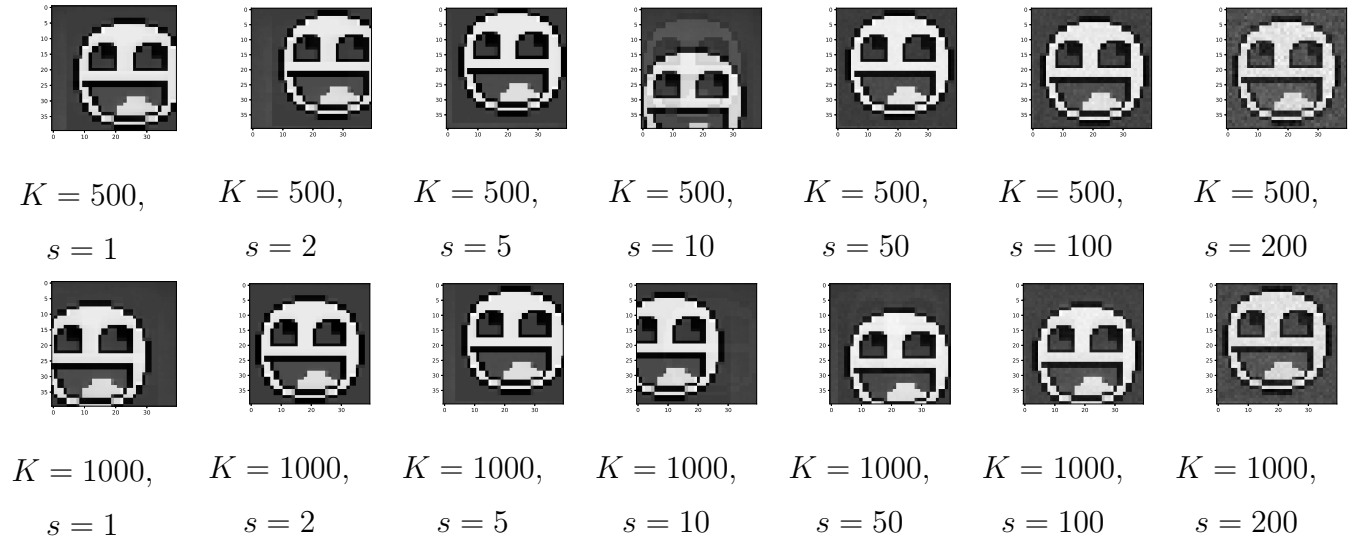
restart 12

Пункт 2

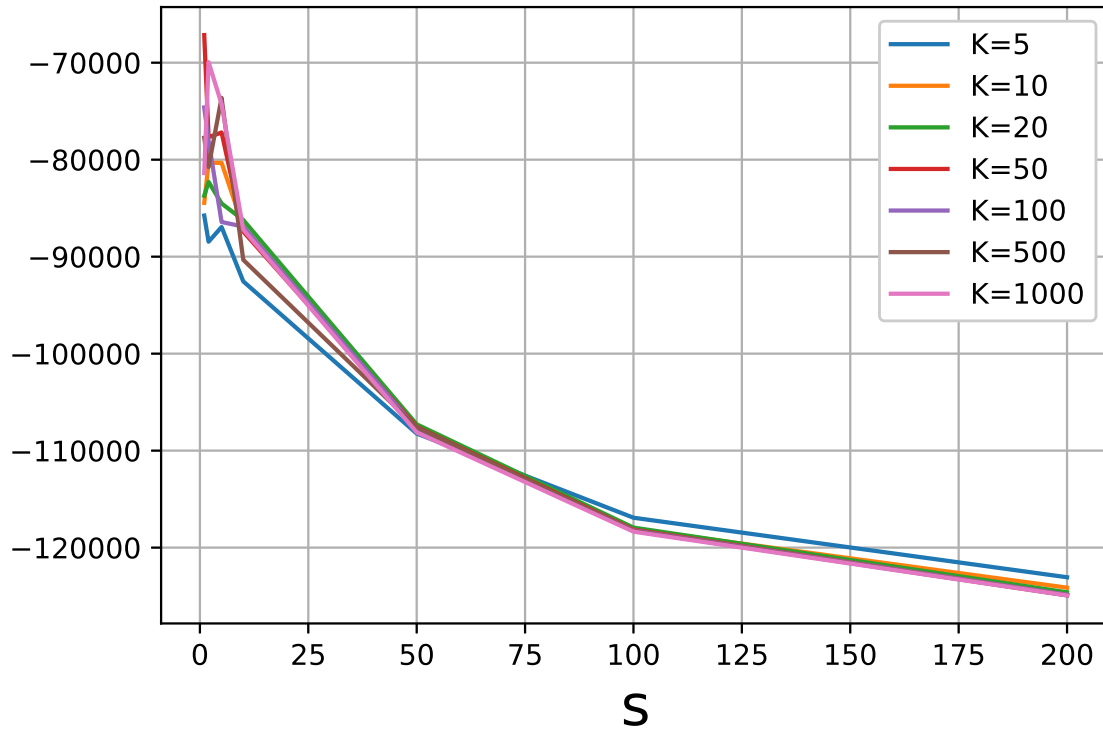
Запустите ЕМ алгоритм на сгенерированных выборках разных размеров и с разным уровнем зашумления. Как изменения в обучающей выборке влияют на результаты работы (получаемые F , B и $L(q, \theta, \mathbf{A})$)? При каком уровне шума ЕМ-алгоритм перестает выдавать вменяемые результаты? В данном пункте учтите, что для сравнения значения $L(q, \theta, \mathbf{A})$ для выборок разного размера стоит нормировать его на объем выборки.

Результат работы алгоритма для разных размеров выборки и разных дисперсий шума:





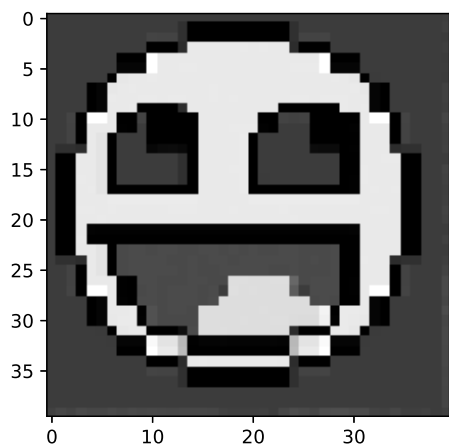
Если не учитывать такие ограничения адекватности дисперсии как максимальное возможное значение яркости пикселя (256), то наблюдаемые результаты вполне логичны: любой аддитивный независимый одинаково распределенный шум можно подавить имея достаточно большую для этого выборку. Убедимся в падении значения $L(q, \theta, \mathbf{A})$ от уровня шума:



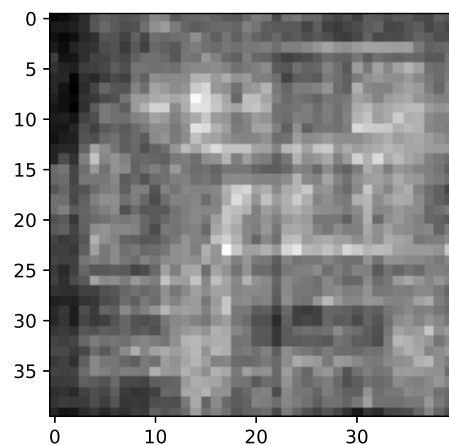
Пункт 3

Сравните качество и время работы ЕМ и hard ЕМ на сгенерированных данных. Как Вы думаете, почему разница в результатах работы так заметна?

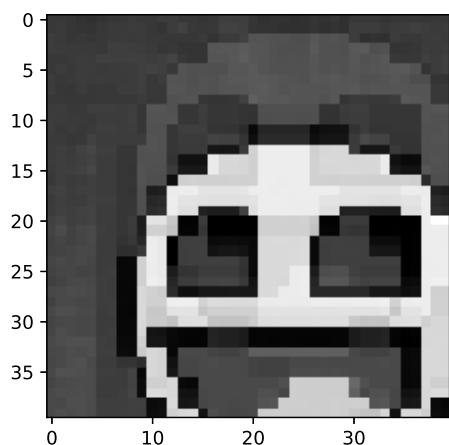
Ну чтож, посмотрим (фиксируем размер выборки и будем менять уровень шума):



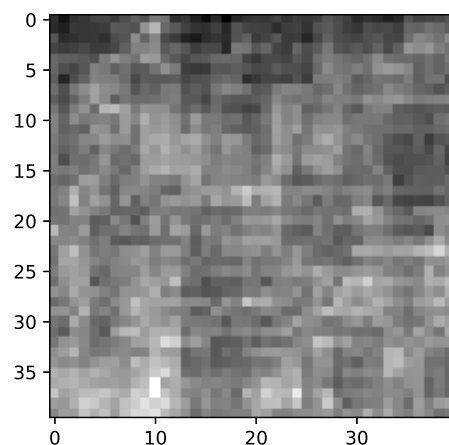
$K = 50, s = 1, \text{full}$



$K = 50, s = 1, \text{hard}$

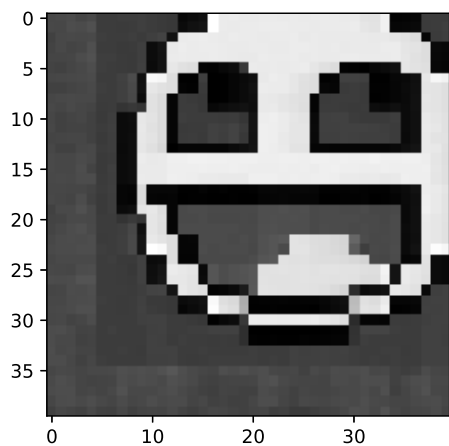


$K = 50, s = 2, \text{full}$

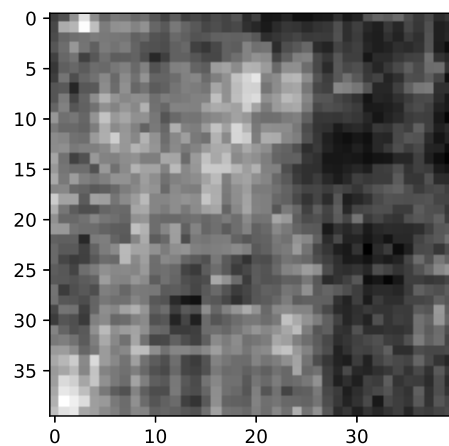


$K = 50, s = 2, \text{hard}$

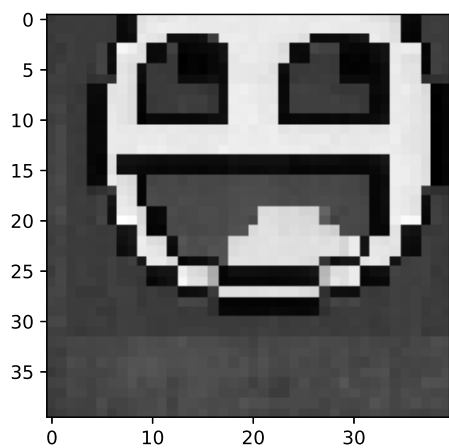
Либо в код закралась ошибка (но я вроде проверял), либо недостаточно большая выборка. Есть ощущение, что основная проблема hard-ЕМ – уменьшение данных, по которым ведется усреднение для расчета s в $h \cdot w$ раз. Этим и вызвано неспособность восстановить лицо при заданных параметрах. Было бы время, то я бы попробовал провести эксперимент с $K = 2000$ и $s = 10$.



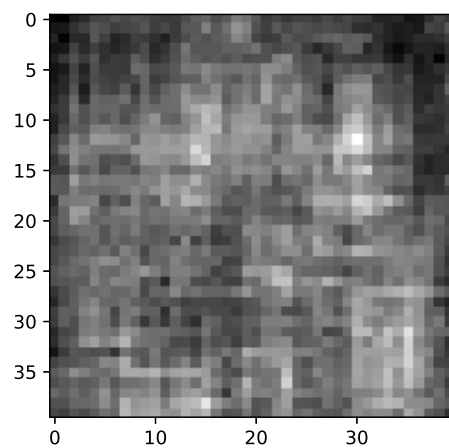
$K = 50, s = 5, \text{ full}$



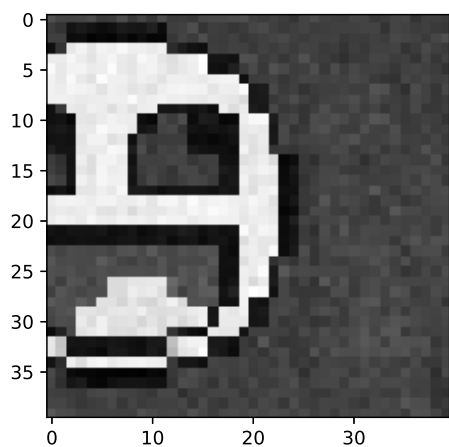
$K = 50, s = 5, \text{ hard}$



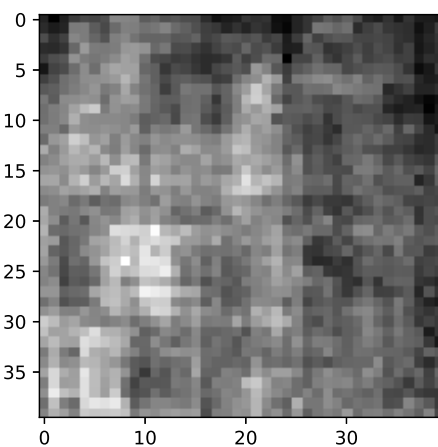
$K = 50, s = 10, \text{ full}$



$K = 50, s = 10, \text{ hard}$



$K = 50, s = 50, \text{ full}$

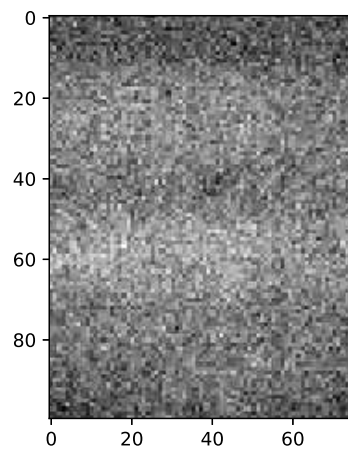


$K = 50, s = 50, \text{ hard}$

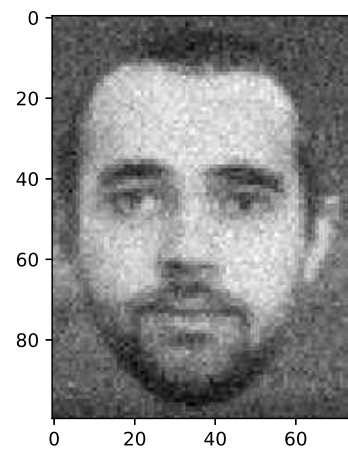
Пункт 4

Примените ЕМ-алгоритм к данным с зашумленными снимками преступника. Приведите результаты работы алгоритма на выборках разного размера.

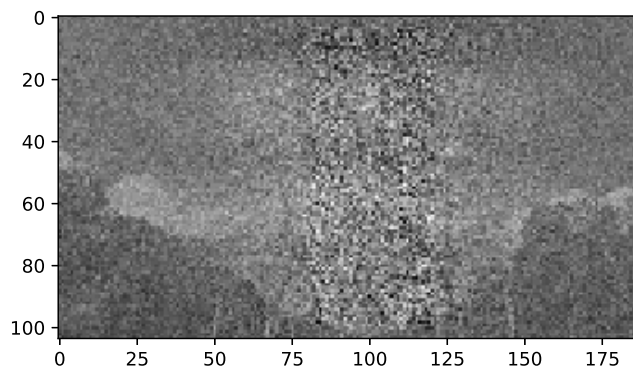
К сожалению я приступил к заданию в тот момент, когда были выложены данные в виду 100 картинок. Так что у меня есть всего лишь 2 результата:



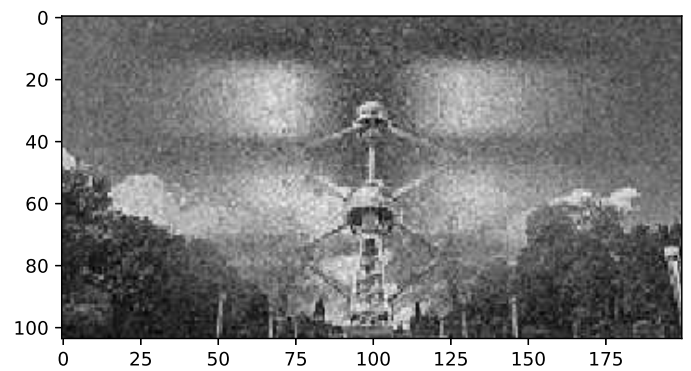
100 картинок



1000 картинок



100 картинок



1000 картинок

Пункт 5

Предложите какую-нибудь модификацию полученного ЕМ алгоритма, которая бы работала на данной задаче качественнее и/или быстрее.

Первое, что пришло в голову: попытаться более грамотно подбирать начальные приближения. В частности, можно найти среднее изображение $\hat{\mathbf{X}}$ как простое арифметическое среднее. Далее для каждого \mathbf{X}_k рассчитать квадрат разности со средним \mathbf{D}_k . В нашей задаче гарантируется, что дисперсия у фона и лица одинакова, а значит минимальное значение $\hat{d} = \min_{k,i,j} \mathbf{D}_k$ будет хорошей оценкой для s^2 . И это значение можно пересчитывать не на каждой итерации, а, например, на каждой 10-ой и в конце. Далее для каждого \mathbf{D}_k мы можем найти координаты окон размера $h \times w$ с максимальной суммой и усреднить $\{\mathbf{X}_k\}$ отдельно по окнам и отдельно по пространству вокруг них. Хотя это уж очень похоже на М-шаг.