

## Task 1

$$P_\theta = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$L(S; \theta) = -\frac{1}{2\sigma^2} \sum_{i=1}^m (x_i - \mu)^2 - m \log(\sigma\sqrt{2\pi})$$

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x_i, \quad \hat{\sigma} = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_i - \hat{\mu})^2}$$

Покажем, что  $\mathbb{E}[\hat{\sigma}^2] = \frac{m-1}{m}\sigma^2$ .

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m (x_i^2 - 2x_i\hat{\mu} + \hat{\mu}^2)\right] = \\ &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x_i^2 - \hat{\mu}^2\right] = \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}[x_i^2] - \mathbb{E}[\hat{\mu}^2] = \\ &= \mathbb{E}[x^2] - \mathbb{E}[\hat{\mu}^2] = \\ &= (\mu^2 + \sigma^2) - (\mu^2 + \mathbb{D}[\hat{\mu}]) = \\ &= \sigma^2 - \frac{1}{m^2} \mathbb{D}\left[\sum_{i=1}^m x_i\right] = \sigma^2 - \frac{1}{m^2} \sum_{i=1}^m \mathbb{D}[x_i] = \sigma^2 - \frac{1}{m}\sigma^2 = \\ &= \frac{m-1}{m}\sigma^2 \end{aligned}$$

Таким образом, оценка максимального правдоподобия для нормально распределенной случайной является смещенной.

## Task 2

Добавим в исходную выборку один элемент класса 1 и один элемент класса 0. Обозначим их  $x_{m+1}, x_{m+2}$  соответственно. Отметим, что речь идет о распределении Бернулли, при этом  $P[X = 1] = \theta$ . Учитывая, что

$$\begin{aligned} L(S; \theta) &= \log\left(\prod_{i=1}^m P_\theta(x_i)\right) = \sum_{i=1}^m \log(P_\theta(x_i)), \quad \hat{\theta} \in \underset{\theta}{\operatorname{argmax}} L(S; \theta); \\ l(\theta; x) &= -\log(P_\theta[x]); \\ \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^m (-\log(P_\theta[x_i])) &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m (\log(P_\theta[x_i])); \\ \hat{\theta} &= \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{m+2} \sum_{i=1}^{m+2} x_i = \frac{1}{m+2} \left(1 + \sum_{i=1}^m x_i\right); \end{aligned}$$

получаем

$$L(S; \theta) = \sum_{i=1}^{m+2} \log(P_\theta(x_i)) = \sum_{i=1}^m \log(P_\theta(x_i)) + \log(\theta) + \log(1 - \theta).$$

Т.е. минимизация регуляризированной целевой функции для исходной выборки эквивалентна минимизации эмперического риска для расширенной выборки.

Получим верхнюю оценку для  $|\hat{\theta} - \theta^*|$ .

$$|\hat{\theta} - \theta^*| \leq |\hat{\theta} - \mathbb{E}[\hat{\theta}]| + |\mathbb{E}[\hat{\theta}] - \theta^*|, \mathbb{E}[\hat{\theta}] = \frac{1+m\theta^*}{m+2}.$$

Тогда

$$|\hat{\theta} - \mathbb{E}[\hat{\theta}]| = \frac{m}{m+2} \left| \frac{1}{m} \sum_{i=1}^m x_i - \theta^* \right|, |\mathbb{E}[\hat{\theta}] - \theta^*| = \left| \frac{1-2\theta^*}{m+2} \right| \leq \frac{1}{m+2}$$

Используя неравенство Хефдинга, получаем

$$\mathbb{P} \left[ |\hat{\theta} - \theta^*| \geq \frac{1}{m+2} + \epsilon \right] = \mathbb{P} \left[ |\hat{\theta} - \theta^*| \geq \epsilon \right] \leq 2e^{-m\epsilon^2}.$$

Тогда с вероятностью  $1 - \delta$

$$|\hat{\theta} - \theta^*| \leq \sqrt{\frac{\log \frac{1}{\delta}}{m}} = O\left(\frac{1}{\sqrt{m}}\right).$$

### Task 3

Шаг  $M$  метода *soft k-means* заключается в максимизации ожидаемого логарифма правдоподобия.

$$\{c \in \mathbb{R}^k : \sum_{i=1}^k c_i = 1, c_i \geq 0\}.$$

$$\max_{\theta} \sum_{i=1}^m \sum_{y=1}^k P_{\theta(t)} [Y=y|X=x_i] \left( \log(c_y) - \frac{1}{2} \|x_i - \mu_y\|^2 \right).$$

$$\mu_y = \frac{\sum_{i=1}^m P_{\theta(t)} [Y=y|X=x_i] x_i}{\sum_{i=1}^m P_{\theta(t)} [Y=y|X=x_i]}.$$

Обозначим  $P_{\theta(t)}$  как  $P$ , а также  $\nu_y = \sum_{i=1}^m P[Y=y|X=x_i]$ , тогда получаем, что исходная максимизация эквивалентна следующей задаче оптимизации:

$$\max_{c \in \mathbb{R}^k} \sum_{y=1}^k \nu_y \log(c_y), c_y \geq 0, \sum_y c_y = 1.$$

Пусть  $c^* = \frac{\nu_y}{\sum_y \nu_y}$ . Учитывая, что  $\nu_y \geq 0$  и  $\sum_y c_y^* = 1$ , тогда  $c^*$  - вектор вероятности.

$$\begin{aligned} -D_{RE}(c^* \parallel c) &= \sum_y c_y^* \log\left(\frac{c_y}{c_y^*}\right) = \sum_y c_y^* (\log(c_y) - \log(c_y^*)) = \\ &= Z_1 \left( \sum_y \nu_y \log(c_y) + Z_2 \right), \end{aligned}$$

т.е. минимизация  $D_{RE}(c^* \parallel c)$  при условии  $c_y \geq 0, \sum_y c_y = 1$  эквивалентна исходной задаче.

Так как  $D_{RE}[P \parallel P_{\theta}] = \sum P[x] \log \left( \frac{P[x]}{P_{\theta}[x]} \right)$  и true risk минимален, а  $L(x, \theta)$  достигает максимума при  $\hat{P}_{\theta}^x = P$ .

Найдем данный оптимальный параметр  $\theta^*$ .

$$\begin{aligned}
L(x, \theta) &= \sum_{i=1}^m \sum_{y=1}^k P_{\theta^{(t)}} [Y = y | X = x_i] \left( \log(c_y) - \frac{1}{2} \|x_i - \mu_y\|^2 \right) \\
\frac{\partial L}{\partial \mu_y} &= \sum_{i=1}^m P_{\theta^{(t)}} [Y = y | X = x_i] \|x_i - \mu_y\| = 0 \\
\mu_y &= \frac{\sum_{i=1}^m P_{\theta^{(t)}} [Y = y | X = x_i] x_i}{\sum_{i=1}^m P_{\theta^{(t)}} [Y = y | X = x_i]}.
\end{aligned}$$

Воспользуемся правилом множителей Лагранжа

$$\begin{aligned}
&L(x, \theta) + \lambda \left( \sum_{y=1}^k c_y - 1 \right) \\
\sum_{y=1}^k \left( \frac{\partial L}{\partial c_y} + \lambda \right) &= 0 \Rightarrow \lambda = - \sum_{y'=1}^k \sum_{i=1}^m P_{\theta^{(t)}} [Y = y' | X = x_i] \\
c_y &= \frac{\sum_{i=1}^m P_{\theta^{(t)}} [Y = y | X = x_i]}{\sum_{y'=1}^k \sum_{i=1}^m P_{\theta^{(t)}} [Y = y' | X = x_i]} \Rightarrow c^* = \frac{\nu_y}{\sum_y \nu_y}
\end{aligned}$$

Т.е. функция правдоподобия достигает максимума при  $\theta^* = \{\mu^*, c^*\}$ .