Круглов А.И. Б05-202 ML2

**1** $L(\theta) = \sum_{i=1}^{n} w_i (y_i - x_i^T \theta)^2 \rightarrow \min_{\theta}$ , НАЙТИ РЕШ. В МАТРИЧНОМ ВИДЕ

Пусть $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$ - МАТРИЦА ПРИЗНАКОВ , $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ - ВЕКТОР ТАРГЕТОВ , $W = \begin{pmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & & w_n \end{pmatrix}$ - ДИАГ. МАТРИЦА ВЕСОВ

ТОГДА В МАТРИЧНОМ ВИДЕ: $L(\theta) = (y - X\theta)^T W (y - X\theta) \rightarrow \min_{\theta}$

$\frac{\partial L}{\partial \theta} = \left(\frac{\partial}{\partial \theta}(y-X\theta)^T\right) W(y-X\theta) + (y-X\theta)^T W\left(\frac{\partial}{\partial \theta}(y-X\theta)\right) = -X^T W(y-X\theta) +$

$+ (y-X\theta)^T W(-X) = -2X^T W(y-X\theta)$

min : $\frac{\partial L}{\partial \theta} = 0$ , т.е. $X^T W y = X^T W X \hat{\theta}$ . ВЫРАЗИМ $\hat{\theta} = (X^T W X)^{-1} X^T W y$

По полученной формуле также ВИДНО, что МНК - это частный случай при

$w_1 = \ldots = w_n = 1$ , т.е. $\hat{\theta}_{МНК} = (X^T X)^{-1} X^T y$ , что было на лекции.

Ответ: $\hat{\theta} = (X^T W X)^{-1} X^T W y$

**2** $F(\theta) = -\log L_y(\theta) + \lambda \|\theta\|_2^2 = -\ell_y(\theta) + \lambda \theta^T \theta \quad \Rightarrow \quad \nabla F = -\nabla \ell_y + \nabla(\lambda \theta^T \theta) = -\nabla \ell_y + 2\lambda \theta,$

$\nabla \ell_y$ ВЫВЕДЕН НА ЛЕКЦИИ : $\nabla \ell_y = X^T(y - S(\theta))$ , где $S(\theta) = (\sigma(x_1^T \theta), \ldots, \sigma(x_n^T \theta))^T$

GD: ИНИЦИАЛИЗИРУЕМ ВЕСА $\theta_0$ и НА КАЖДОМ ШАГЕ БУДЕМ ОБНОВЛЯТЬ ПО ПРАВИЛУ:

$\theta_{t+1} = \theta_t - \eta \nabla F(\theta_t) = \theta_t - \eta(2\lambda\theta_t - X^T(y - S(\theta_t)))$ , где $\eta$ - РАЗМЕР ШАГА

(learning rate) , $\ominus \eta \nabla F$ - потому что $F \rightarrow \min$

SGD: ПОХОЖЕ НА GD, НО КАЖДЫЙ ШАГ ЗАДЕЙСТВУЕТ БАТЧ $I = \{i_1, \ldots, i_m\} \in U[1, \ldots, n]$ :

$\theta_{t+1} = \theta_t - \eta \cdot \frac{n}{m}(-X_I^T(y_I - S(\theta_t))) - \eta \cdot 2\lambda\theta_t$

IRLS: ЗАДАЁМ НАЧАЛЬНЫЕ ВЕСА $\theta_0$ и БУДЕМ ИХ ОБНОВЛЯТЬ. НА КАЖДОМ ШАГЕ ВЫЧИСЛЯЕМ

$W = \text{diag}(\sigma(x^T\theta)(1 - \sigma(x^T\theta)))$ и $\theta_{t+1} = \theta_t - (X^T W(\theta_t) X)^{-1}(-X^T(y - S(\theta_t)))$ ,

что подробнее выводилось на лекции

ГРАДИЕНТНЫЙ СПУСК ХОРОШО РАБОТАЕТ ДЛЯ МАЛЫХ ОБЪЁМОВ ДАННЫХ, т.к. требует много вычислений. Поэтому НА ПРАКТИКЕ ЧАЩЕ ОБРАЩАЮТСЯ К SGD или IRLS