# Задача 3, Круглов А.И. Б05-207

## 1) $\hat{\Theta}$ в ridge-регрессии

Вспомним сведения о матр. производной с лекции:

$$f(x) = a^T x \Rightarrow \nabla f = a, \quad f(x) = x^T A x \Rightarrow \nabla f = (A + A^T) x$$

Оптимизируем функционал $F(\Theta) = \|Y - X\Theta\|^2 + \lambda \|\Theta\|^2 =$

$$= (Y - X\Theta)^T (Y - X\Theta) + \lambda \Theta^T \Theta = Y^T Y - 2(Y^T X)\Theta + \Theta^T \cdot$$

$$\cdot (X^T X + \lambda E)\Theta; \quad \text{отметим, что} \quad (X^T X + \lambda E)^T = (X^T X + \lambda E).$$

Тогда $\nabla F(\Theta) = -2 X^T Y + 2(X^T X + \lambda E)\Theta = 0$, откуда

$$X^T Y = (X^T X + \lambda E)\Theta \Rightarrow \hat{\Theta} = (X^T X + \lambda E)^{-1} X^T Y.$$

В МНК было $\hat{\Theta} = (X^T X)^{-1} X^T Y$, а слагаемое $\lambda E$ может

сделать матрицу не вырожденной (обратимой), тривиальный

пример: $X^T X = 0$, $\det(X^T X) = 0$, $\det(X^T X + \lambda E) = \det(\lambda E) \neq 0$

## 2) Шаг град. спуска

GD: $\Theta_{t+1} = \Theta_t - \eta \nabla F(\Theta) \Rightarrow \Theta_{t+1} = \Theta_t + \eta(X^T Y - (X^T X + \lambda E)\Theta_t)$

SGD: тоже, но берем только строки с номерами из $I$ ($I =$

$= \{i_1, ..., i_n\}$, $i_1, ..., i_k \sim U\{1, ..., n\}$ - БАТЧ), т.е. $\Theta_{t+1} = \Theta_t + \eta \frac{n}{N} \cdot$

$$\cdot (X_I^T Y_I - (X_I^T Y_I + \lambda E)\Theta_t)$$

## 3) если признаки не приведены к одинаковому масштабу, то

штраф за величину весов применяется к весам "несправедливо",

неодинаково. Стандартизация позволяет привести столбцы к равному

масштабу (по сети, избавиться от размерности и оставить только информацию о распределении признака), из-за чего ridge-регрессия работает заметно лучше.