

**МАШИННОЕ ОБУЧЕНИЕ  
И АНАЛИЗ ДАННЫХ**  
**(Machine Learning and Data Mining)**

Н. Ю. Золотых

<http://www.uic.unn.ru/~zny/ml>



*Лекция 6*

**Борьба с переобучением**

## 6.1. Переобучение

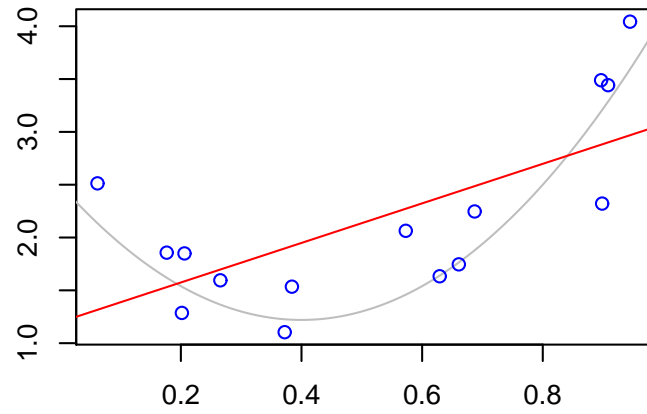
*Переобучение* — это явление, при котором обучающий алгоритм выдает хорошие результаты на обучающей выборке, но имеет очень плохие обобщающие свойства.

Типичное поведение ошибки на обучающей и тестовой выборках с ростом «сложности» модели:

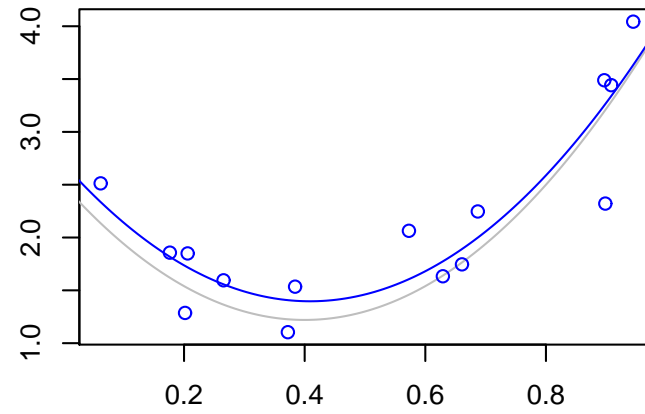
- Обычно на обучающей выборке с ростом сложности модели ошибка уменьшается.
- На тестовой выборке сперва с ростом сложности модели ошибка уменьшается, но с некоторого момента ошибка начинает расти: сказывается синдром переобучения.

Когда модель слишком сложна, она, как правило, хорошо приспособливается к конкретным обучающим данным, улавливая какие-то специфичные для них особенности, но не присущие всей генеральной совокупности, поэтому на тестовых данных ошибка может быть большой.

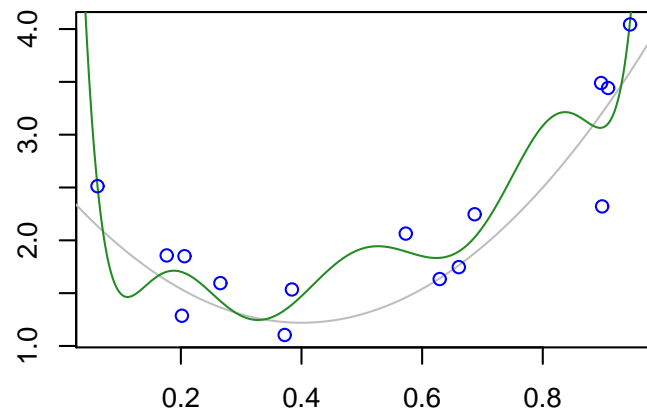
$$y = 8x^2 - 6.4x + 2.5 + N(0, 0.4)$$



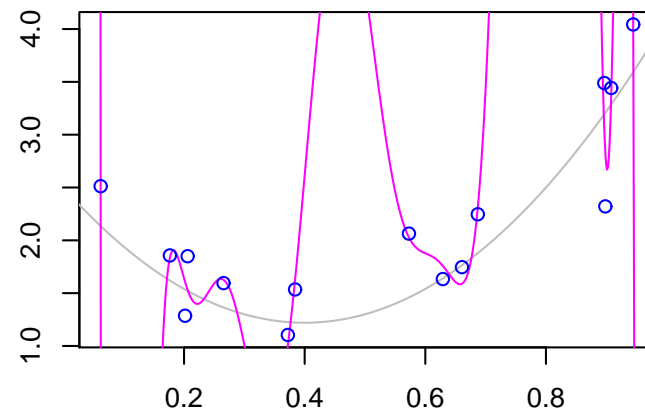
d = 1



d = 2

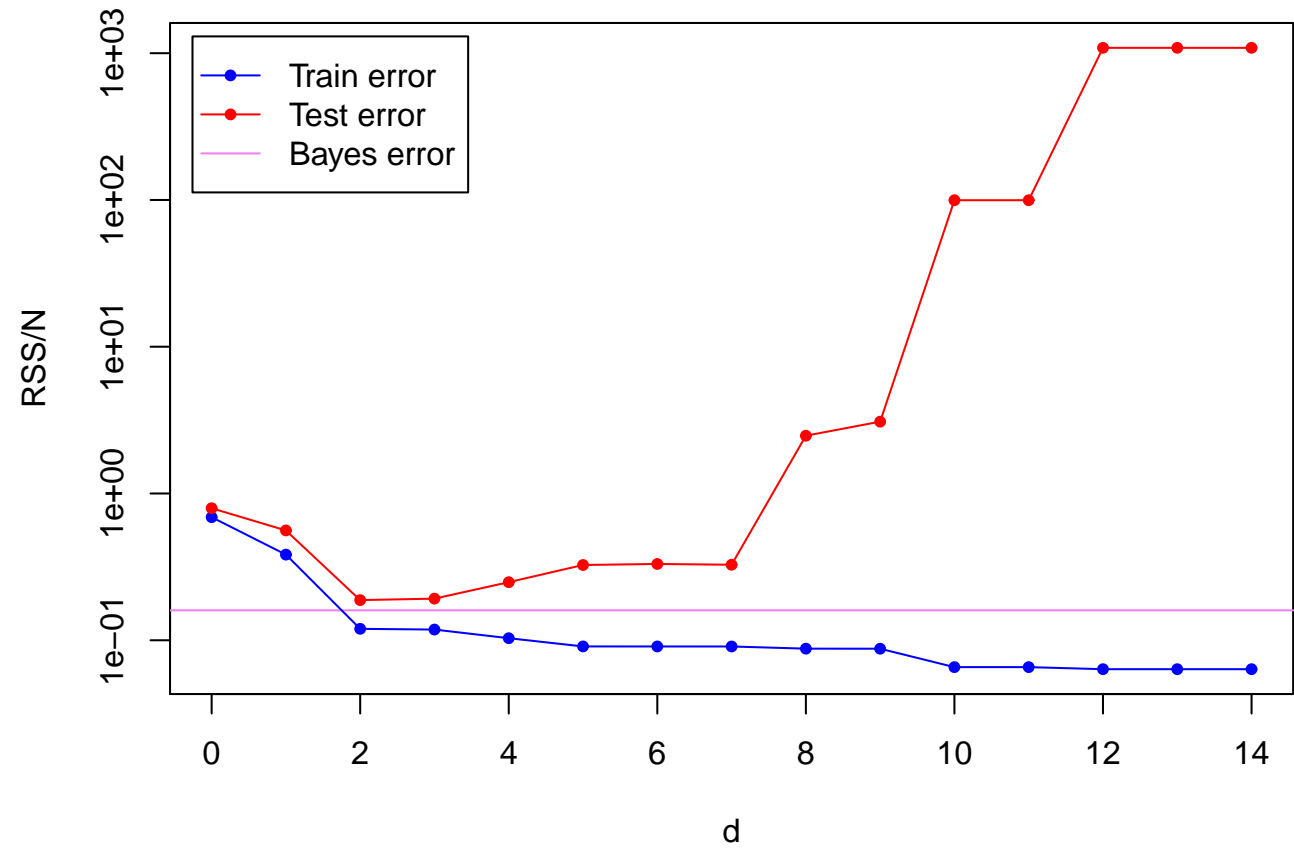


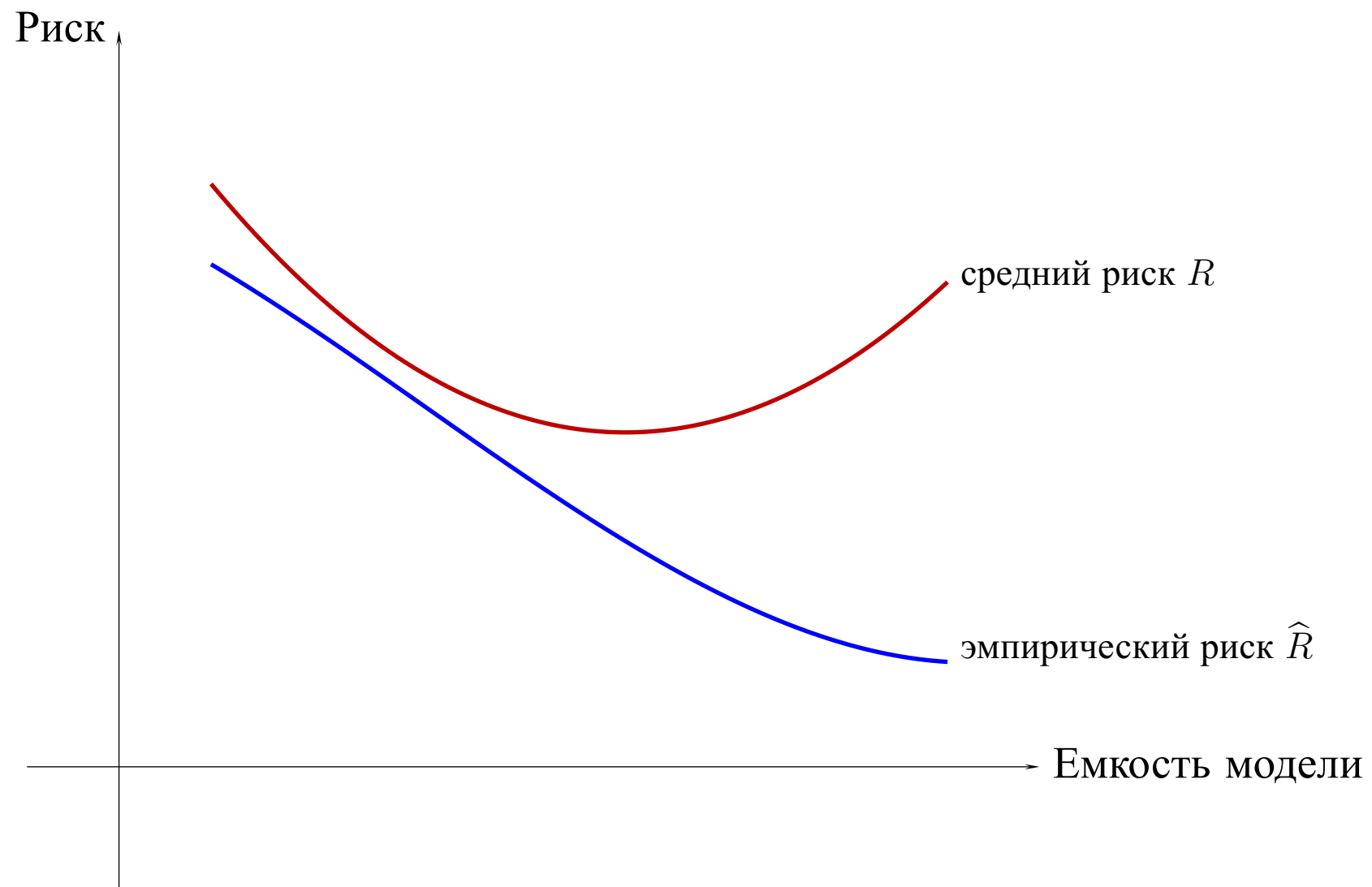
d = 8



d = 13

$RS_{\text{train}}/N_{\text{train}} = 0.1196395$   $RS_{\text{test}(3)}/N_{\text{test}} = 0.1876475$   $\sigma^2 = 0.16$





### 6.1.1. Причины переобучения (в задаче восстановления регрессии)

- *Мультиколлинеарность*:  $\mathbf{X}$  плохо обусловлена (т. е. столбцы матрицы  $\mathbf{X}$  образуют систему, близкую к линейно зависимой).

$$\text{cond}_2 \mathbf{X} = \max_{\beta \neq 0} \frac{\|\mathbf{X}\beta\|_2}{\|\beta\|_2} \bigg/ \min_{\beta \neq 0} \frac{\|\mathbf{X}\beta\|_2}{\|\beta\|_2} = \frac{\max \sqrt{\lambda_j(\mathbf{X}^\top \mathbf{X})}}{\min \sqrt{\lambda_j(\mathbf{X}^\top \mathbf{X})}} = \frac{\max \sigma_j}{\min \sigma_j} = \frac{\sigma_1}{\sigma_d}$$

(отношение максимального сингулярного числа матрицы  $\mathbf{X}$  к минимальному)

- *Наличие неинформативных признаков*
- *Слишком мало данных.*

Даже если все признаки информативны, данных может оказаться слишком мало.



### 6.1.2. Замечание о числе обусловленности

Необходимо смотреть на число обусловленности не для с.л.у., а для задачи L.S.

$$\mathbf{X}\beta \approx \mathbf{y} \quad \Rightarrow \quad (\mathbf{X} + \Delta\mathbf{X})\beta \approx \mathbf{y} + \Delta\mathbf{y}$$

Можно доказать, что при малых возмущениях  $\Delta\mathbf{X}$  и  $\Delta\mathbf{y}$

$$\frac{\|\Delta\beta\|_2}{\|\beta\|_2} \leq \varepsilon \cdot \underbrace{\left( \frac{2 \operatorname{cond}_2 \mathbf{X}}{\cos \theta} + \operatorname{tg} \theta \cdot \operatorname{cond}_2^2 \mathbf{X} \right)}_{\operatorname{cond}_{\text{LS}}(\mathbf{X}, \mathbf{y})} + O(\varepsilon^2),$$

где

$$\varepsilon = \max \left\{ \frac{\|\Delta\mathbf{X}\|_2}{\|\mathbf{X}\|_2}, \frac{\|\Delta\mathbf{y}\|_2}{\|\mathbf{y}\|_2} \right\}, \quad \operatorname{cond}_2 \mathbf{X} = \frac{\sigma_1}{\sigma_d},$$

$\theta$  — угол между  $\mathbf{y}$  и  $L(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d)$ .

Если угол  $\theta$  мал, то  $\operatorname{cond}_{\text{LS}}(\mathbf{X}, \mathbf{y}) \approx \operatorname{cond}_2 \mathbf{X}$

Если  $\theta \rightarrow \pi/2$  мал, то  $\operatorname{cond}_{\text{LS}}(\mathbf{X}, \mathbf{y}) \rightarrow \infty$

## 6.2. Сокращение числа параметров и «усадка» коэффициентов

Рассмотрим некоторые методы по уменьшению числа признаков (уменьшения размерности пространства признаков) и «усадке» коэффициентов.

Зачем?

- *Борьба с переобучением*: чем меньше параметров (и входных переменных) или чем меньше по величине параметры, тем проще модель и снижается возможность переобучения
- *Интерпретация*: чем меньше параметров (и входных переменных), тем проще интерпретировать модель

### 6.2.1. Выбор подмножества признаков

Для каждого  $k \in \{0, 1, \dots, d\}$  найдем подмножество входных параметров мощности  $k$ , для которого  $RSS(\beta)$  минимально.

Большой перебор!

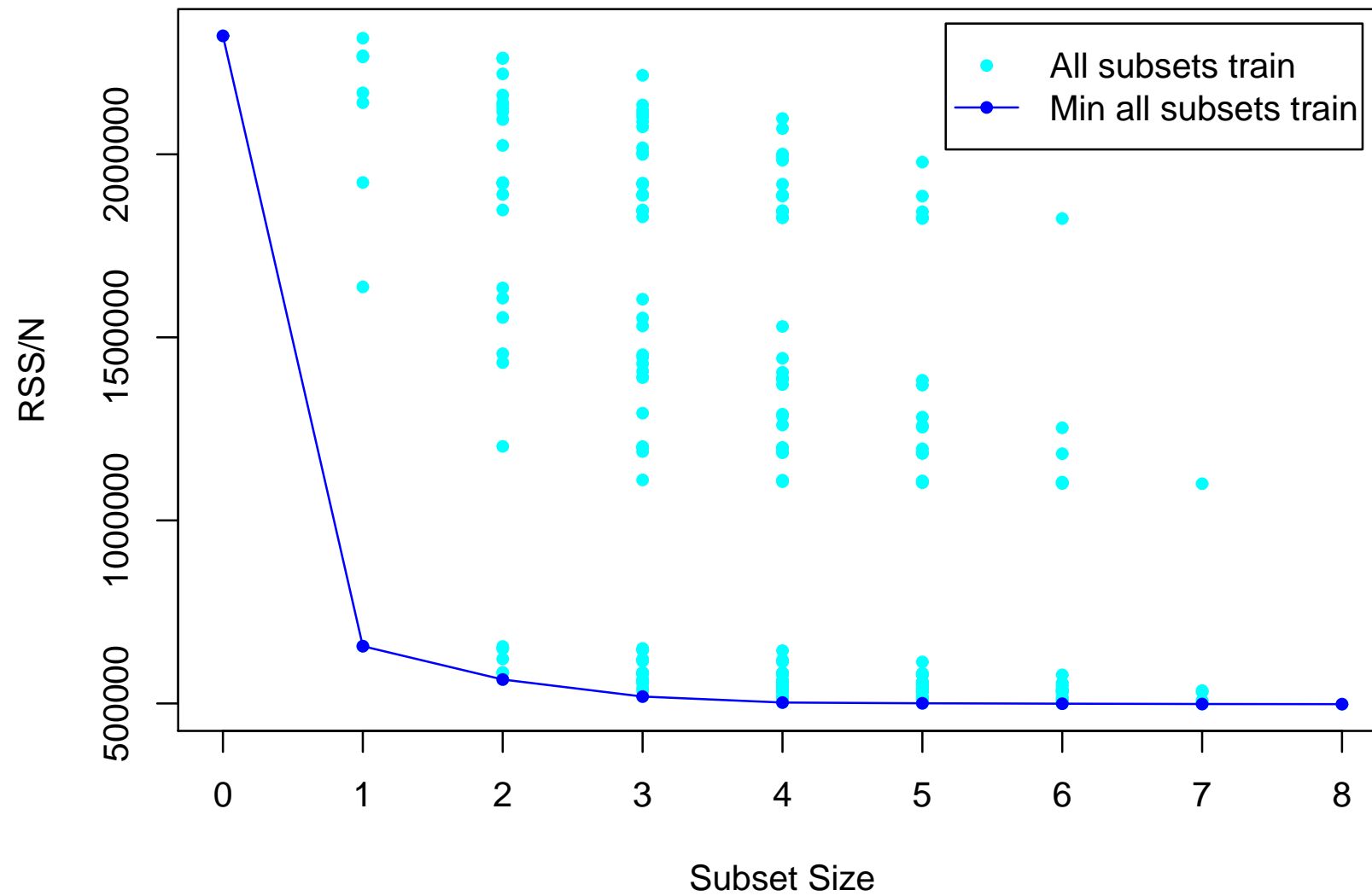
Оптимальное подмножество признаков мощности  $k$  не обязательно содержится в оптимальном подмножестве мощности  $k + 1$

## Задача предсказания цены квартиры

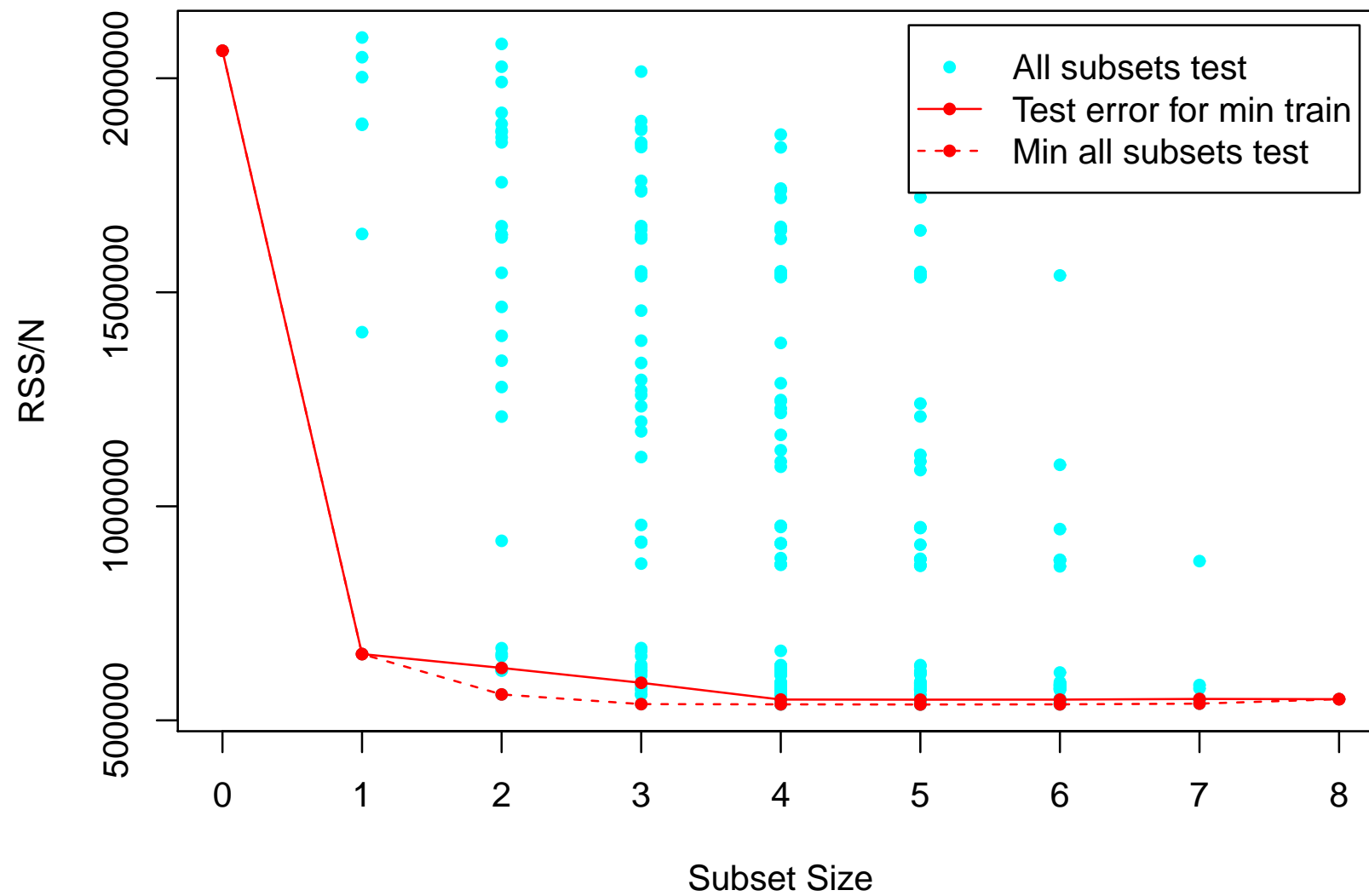
<b>Price</b>	цена квартиры (тыс. руб.)
<b>Date</b>	№ дня, в который квартира выставлена на продажу
<b>Lat</b>	географическая широта объекта недвижимости
<b>Lng</b>	географическая долгота объекта недвижимости
<b>Housing</b>	тип недвижимости (0 — вторичное жилье, 1 — новостройка)
<b>Floors</b>	количество этажей в доме
<b>House</b>	тип дома ( <b>Block</b> — блочный, <b>Brick</b> — кирпичный, <b>Monolithic</b> — монолитный, <b>Panel</b> — панельный, <b>Wooden</b> — деревянный)
<b>Rooms</b>	количество комнат (0 — квартира-студия)
<b>Floor</b>	№ этажа
<b>Area</b>	площадь квартиры (м <sup>2</sup> )

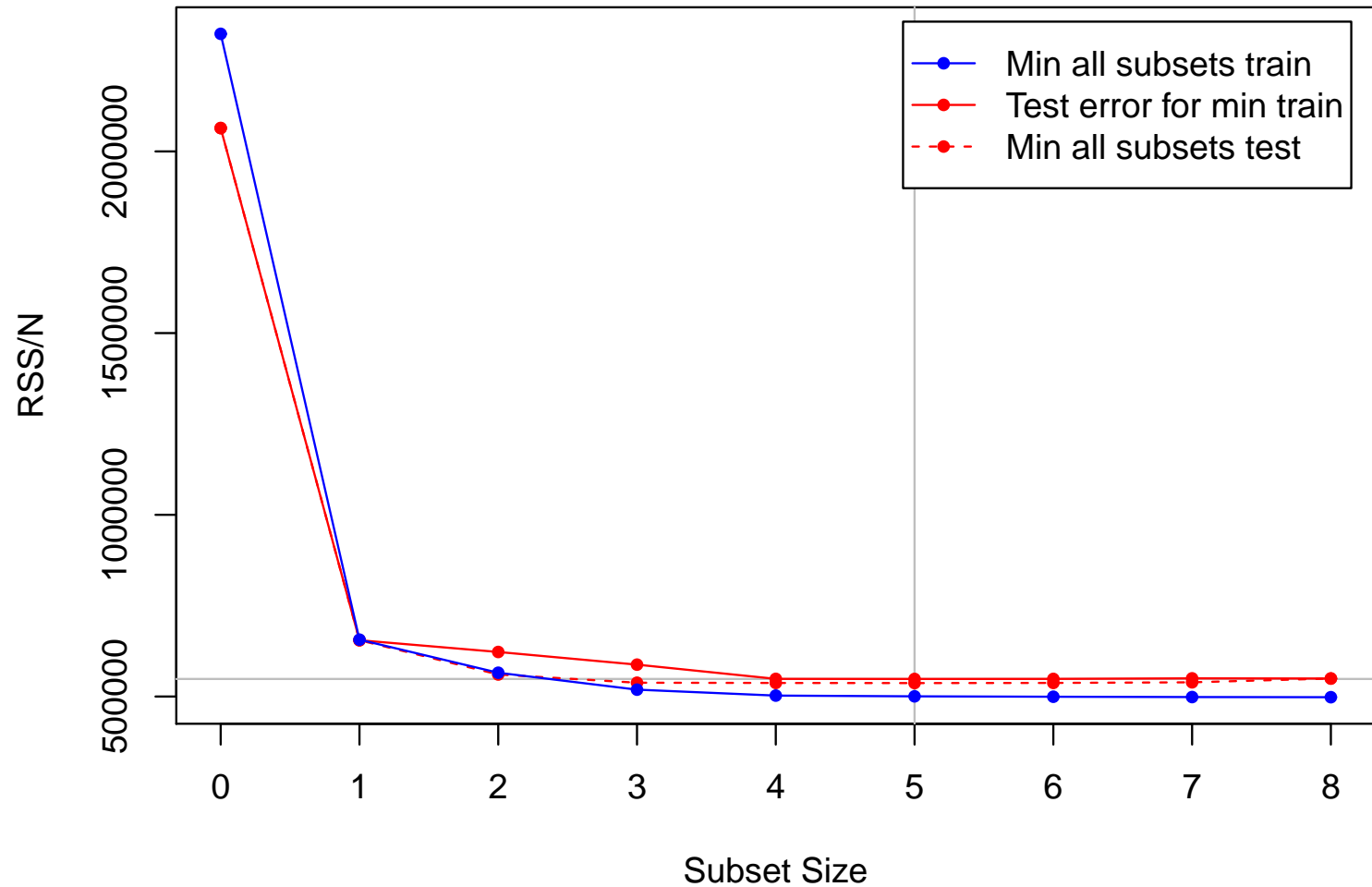
8 входных признаков,  $2^8 = 256$  подмножеств

На обучающей выборке:



На тестовой выборке:





$$d = 5 : \quad \hat{R}_{\text{train}} = \frac{1}{N_{\text{train}}} \text{RSS}_{\text{train}} = 499710.7, \quad \hat{R}_{\text{test}} = \frac{1}{N_{\text{test}}} \text{RSS}_{\text{test}} = 545722.3$$

$$y = -101009 + 2294.1 \text{ Lng} - 152.54 \text{ Housing} + 32.016 \text{ Floors} - 370.43 \text{ Rooms} + 74.944 \text{ Area}$$

## Forward stepwise

находим константную регрессию и ошибку на тестовой выборке  $R_0$

$J \leftarrow \emptyset$

**for**  $k = 1, \dots, d$

$R^* \leftarrow \infty$

**for**  $j \in \{1, \dots, d\} \setminus J$

строим линейную регрессию, используя признаки из  $J \cup \{j\}$

вычисляем  $\widehat{R}_{\text{test}}$  — ошибку этой регрессии на тестовой выборке

**if**  $\widehat{R}_{\text{test}} < R^*$

$j^* \leftarrow j$

$R^* \leftarrow \widehat{R}_{\text{test}}$

$j_k \leftarrow j^*$

$R_k \leftarrow R^*$

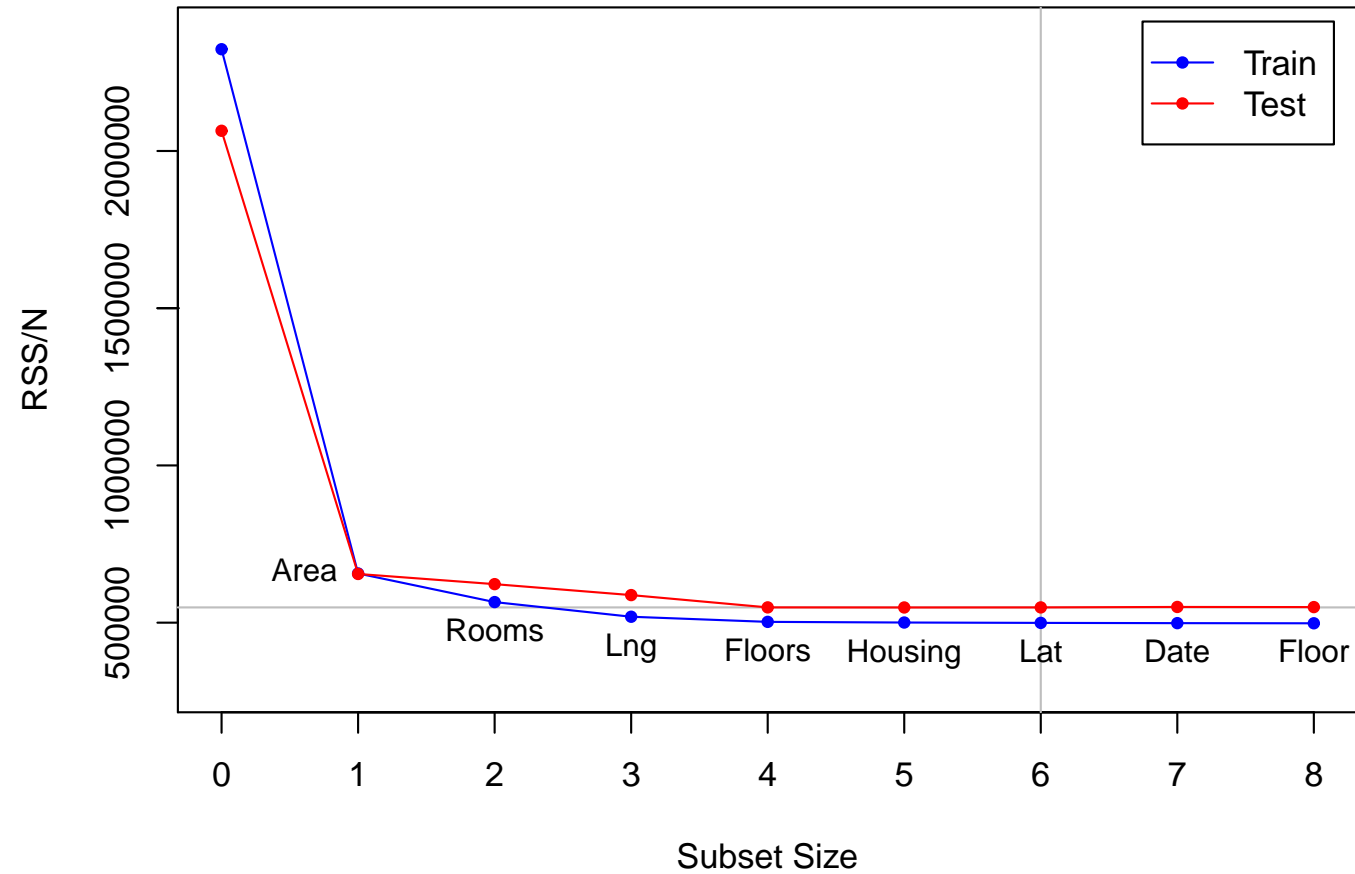
$J \leftarrow J \cup \{j^*\}$

На выходе получаем последовательность номеров признаков  $j_1, j_2, \dots, j_d$  и значения ошибок

$R_0, R_1, \dots, R_d$ , где  $R_k$  — ошибка на тестовой выборке регрессии, построенной по признакам  $j_1, j_2, \dots, j_k$

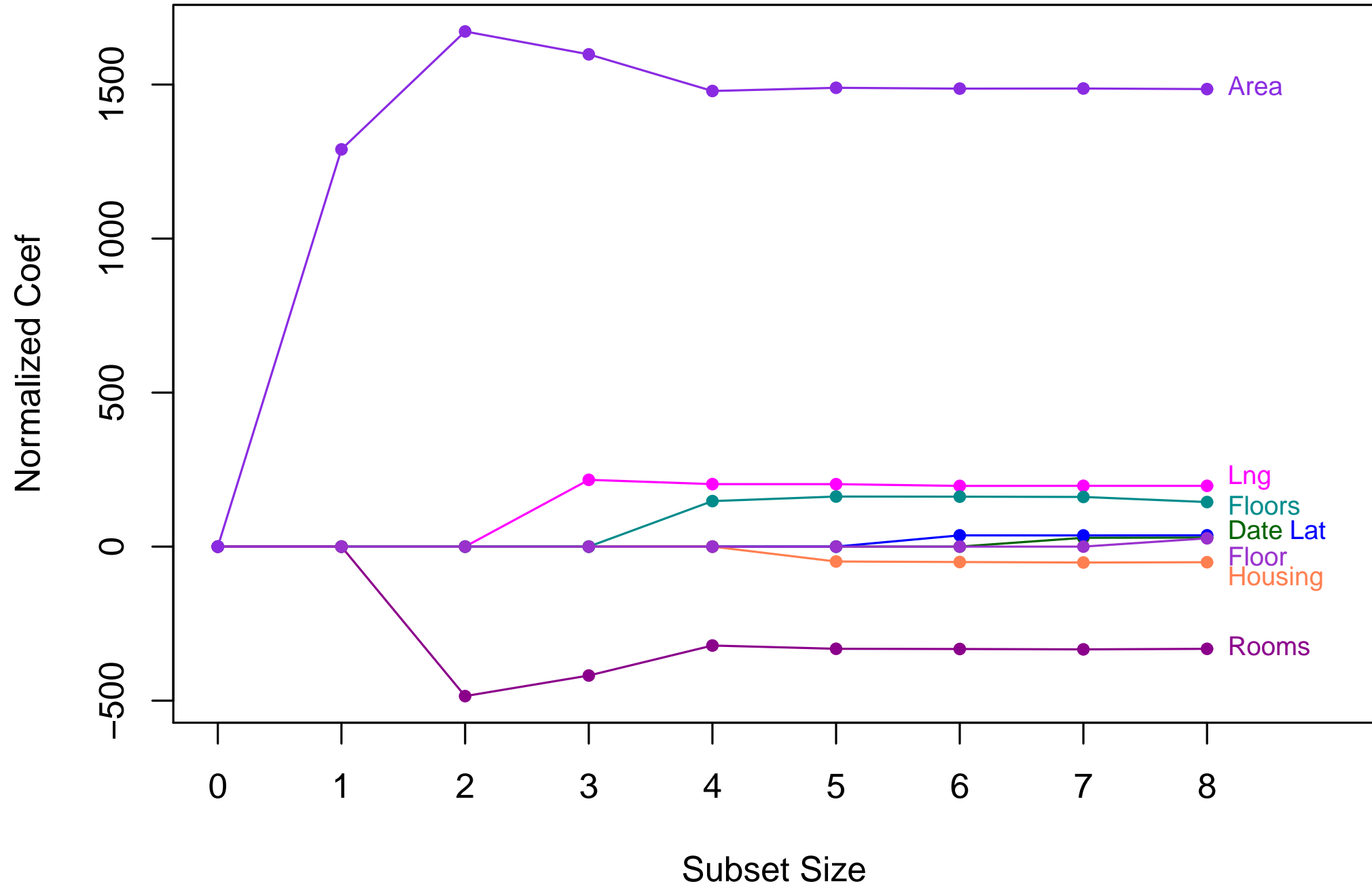


Другая стратегия — начав с рассмотрения всего набора входных переменных, последовательно исключать их из этого набора (backward stepwise).



$$y = -141660 + 772.88 \text{ Lat} + 2229.3 \text{ Lng} - 158.11 \text{ Housing} \\ + 31.941 \text{ Floors} - 371.24 \text{ Rooms} + 74.812 \text{ Area}$$

$$k = 5 : \quad \hat{R}_{\text{train}} = \frac{1}{N_{\text{train}}} \text{RSS}_{\text{train}} = 545813.4, \quad \hat{R}_{\text{test}} = \frac{1}{N_{\text{test}}} \text{RSS}_{\text{test}} = 497700.3$$



## 6.2.2. «Ридж» (гребневая) регрессия (регуляризация)

(*ridge regression*)

Добавим к RSS штрафную функцию, чувствительную к абсолютной величине коэффициентов  $\beta_j$ :

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left( y^{(i)} - \beta_0 - \sum_{j=1}^d x_j^{(i)} \beta_j \right)^2 + \lambda \sum_{j=1}^d \beta_j^2 \right\},$$

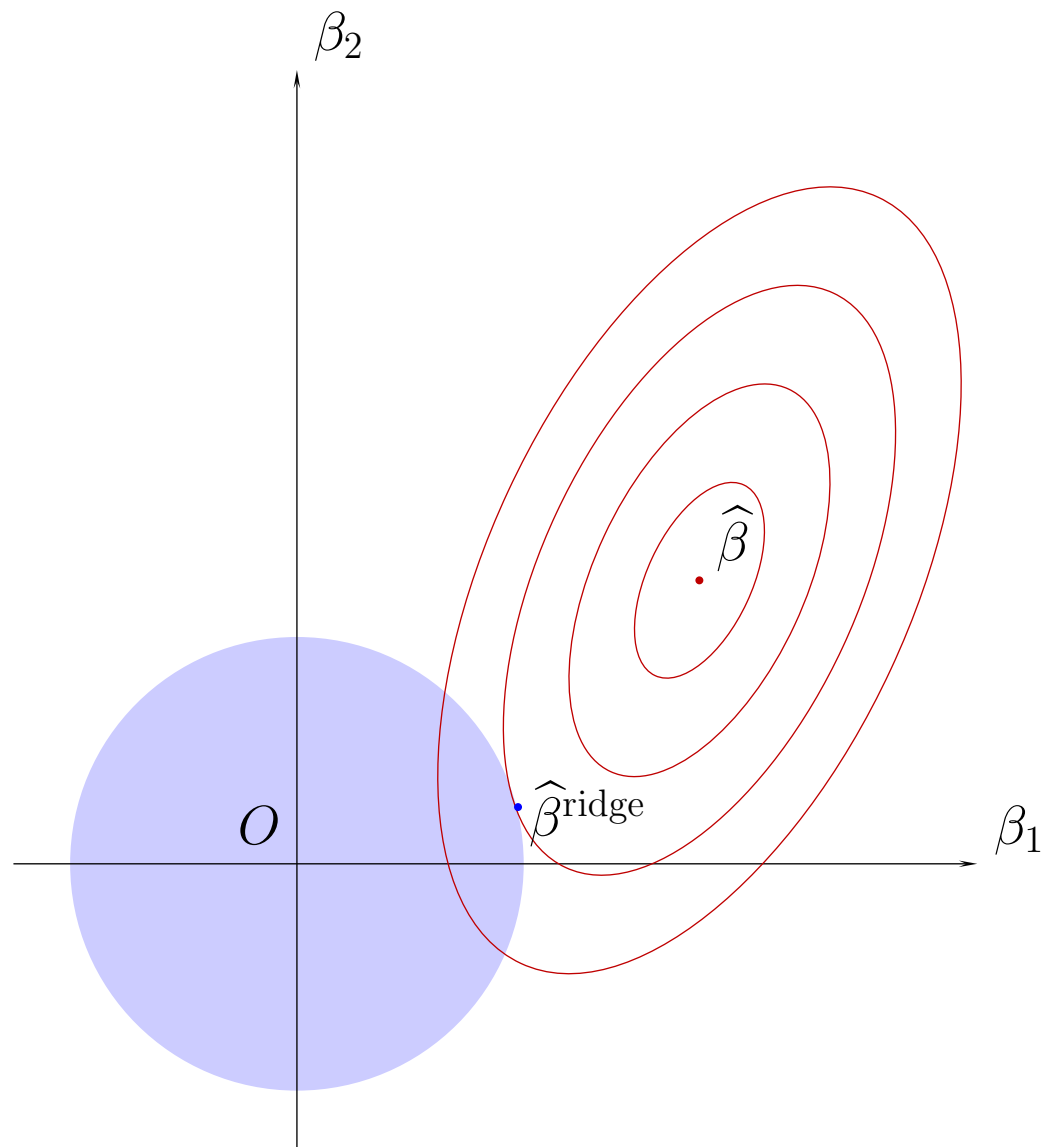
где  $\lambda$  — неотрицательный параметр:

чем больше  $\lambda$ , тем больше чувствительность штрафа к величине коэффициентов

Понятно, что задача эквивалентна следующей задаче условной минимизации:

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left( y^{(i)} - \beta_0 - \sum_{j=1}^d x_j^{(i)} \beta_j \right)^2 \right\}, \text{ при условии } \sum_{j=1}^d \beta_j^2 \leq s.$$

Параметр  $s$  определяется по  $\lambda$ , и наоборот.



Если в линейной регрессионной модели много коррелированных переменных, то параметры модели определяются с трудом и имеют большую дисперсию.

Например,  $X_j$  и  $X_{j'}$  сильно зависимы.

Тогда возможна ситуация:

$\beta_j$  велик и положителен, а  $\beta_{j'}$  велик по абсолютному значению и отрицателен:

так как  $X_j$  и  $X_{j'}$  сильно коррелируют, то вклад  $\beta_j x_j$  компенсируется  $\beta_{j'} x_{j'}$ .

Таким образом, добавляемое ограничение на величину коэффициентов  $\beta_j$  должно предотвратить подобную ситуацию.

1. Коэффициент  $\beta_0$  не участвует в сумме  $\sum_{j=1}^d \beta_j^2$ :

$$\sum_{i=1}^N \left( \beta_0 + \sum_{j=1}^d x_j^{(i)} \beta_j - y^{(i)} \right)^2 \rightarrow \min_{\beta}, \quad \text{при ограничениях} \quad \sum_{j=1}^d \beta_j^2 \leq s$$

2. Как правило,  $y$  центрируют, что позволяет исключить из модели коэффициент  $\beta_0$ :

$$y = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_d x_d$$

Ридж-регрессия:

$$\sum_{i=1}^N \left( \sum_{j=1}^d x_j^{(i)} \beta_j - y^{(i)} \right)^2 \rightarrow \min_{\beta}, \quad \text{при ограничениях} \quad \sum_{j=1}^d \beta_j^2 \leq s$$

3. Как правило,  $X$  центрируют и нормируют на среднеквадратическое отклонение, чтобы вклад в сумму  $\sum_{j=1}^d \beta_j^2$  каждого признака по возможности был бы одинаковым

Коэффициент  $\beta_0$  может быть устранен из штрафной функции, так как он не соответствует ни одной переменной  $x_j$ .

Слагаемое  $\beta_0$  можно выкинуть и из целевой функции.

Действительно, решение задачи наименьших квадратов можно разбить на два этапа:

### 1. Центрирование данных:

- каждое  $x_j^{(i)}$  заменяется на  $x_j^{(i)} - \bar{x}_j$  ( $i = 1, 2, \dots, N$ ;  $j = 1, 2, \dots, d$ )
- в качестве  $\beta_0$  выбирается  $\bar{y}$ , где

$$\bar{x}_j = \frac{1}{N} \sum_{i=1}^N x_j^{(i)}, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y^{(i)}.$$

### 2. Гребневая регрессия (без коэффициента $\beta_0$ )

Пусть центрирование произведено, следовательно,  $\mathbf{X}$  имеет  $d$  (а не  $d + 1$ ) столбцов



$$\text{RSS}^{\text{ridge}}(\beta, \lambda) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^\top \beta,$$

$$\text{RSS}^{\text{ridge}}(\beta, \lambda) \rightarrow \min_{\beta}$$

Дифференцируя, находим:

$$\frac{\partial \text{RSS}^{\text{ridge}}}{\partial \beta} = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) + 2\lambda\beta, \quad \frac{\partial^2 \text{RSS}^{\text{ridge}}}{\partial \beta \partial \beta^\top} = 2\mathbf{X}^\top \mathbf{X} + 2\lambda\mathbf{I}.$$

откуда

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \quad (\star)$$

Итак, решение  $\hat{\beta}^{\text{ridge}}$  задачи гребневой линейной регрессии также является функцией, линейно зависящей от  $\mathbf{y}$ .

Если  $\lambda > 0$ , то матрица  $\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I}$  невырождена (положительно определена), даже если столбцы в  $\mathbf{X}$  линейно зависимы

( $\star$ ) представляет собой стандартную *регуляризацию*

Итак, от системы  $\mathbf{X}\beta = \mathbf{y}$  мы перешли к  $(\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})\beta = \mathbf{X}^\top \mathbf{y}$

### 6.2.3. Отступление. Регуляризация (А.Н. Тихонов)

Система  $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})\beta = \mathbf{X}^\top \mathbf{y}$  называется *регуляризованной* к  $\mathbf{X}\beta = \mathbf{y}$ .

$\lambda$  — параметр регуляризации,  $\lambda > 0$ .

Обозначим  $\hat{\beta}(\lambda)$  — (единственное) решение системы  $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})\beta = \mathbf{X}^\top \mathbf{y}$ .

#### Теорема 6.1 (А.Н. Тихонов)

Пусть  $\lambda_n$  — некоторая последовательность,  $\lambda_n > 0$  и  $\lambda_n \rightarrow 0$ .

Тогда  $\hat{\beta}(\lambda_n)$  стремится к нормальному псевдорешению системы  $\mathbf{X}\beta = \mathbf{y}$   
(из всех псевдорешений нормальное псевдорешение имеет минимальную норму).

Если  $\lambda$  близко к нулю, то регуляризованная система близка к вырожденной, поэтому на практике решение  $\hat{\beta}(\lambda)$  будет найдено с большой ошибкой.

Метод регуляризации заключается в решении регуляризованных систем для конечного числа значений  $\lambda_n$  и дальнейшем выборе того решения, для которого норма невязки минимальна.

## 6.2.4. Гребневая регрессия и сингулярное разложение

Выразим  $\hat{\beta}^{\text{ls}}$  и  $\hat{y}^{\text{ls}}$  (без регуляризации), используя *SVD*:

$$\begin{aligned}\hat{\beta}^{\text{ls}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = ((\mathbf{U} \Sigma \mathbf{V}^\top)^\top \mathbf{U} \Sigma \mathbf{V}^\top)^{-1} (\mathbf{U} \Sigma \mathbf{V}^\top)^\top \mathbf{y} = \\ &= (\mathbf{V} \Sigma \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{V}^\top)^{-1} \mathbf{V} \Sigma \mathbf{U}^\top \mathbf{y} = (\mathbf{V}^\top)^{-1} \Sigma^{-2} \mathbf{V}^{-1} \mathbf{V} \Sigma \mathbf{U}^\top \mathbf{y} = \mathbf{V} \Sigma^{-1} \mathbf{U}^\top \mathbf{y}\end{aligned}$$

$$\hat{y}^{\text{ls}} = \mathbf{X} \hat{\beta}^{\text{ls}} = \mathbf{U} \Sigma \mathbf{V}^\top \mathbf{V} \Sigma^{-1} \mathbf{U}^\top \mathbf{y} = \mathbf{U} \mathbf{U}^\top \mathbf{y} = \sum_{j=1}^d \mathbf{u}_j \mathbf{u}_j^\top \mathbf{y}$$

Заметим, что  $\mathbf{u}_j (\mathbf{u}_j^\top \mathbf{y})$  есть проекция вектора  $\mathbf{y}$  на вектор  $\mathbf{u}_j$ , а вся сумма (как мы уже видели) есть проекция вектора  $\mathbf{y}$  на подпространство, натянутое на  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ .

Аналогично, если проводится регуляризация:

$$\hat{y}^{\text{ridge}} = \mathbf{X} \hat{\beta}^{\text{ridge}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{U} \Sigma (\Sigma + \lambda \mathbf{I})^{-1} \Sigma \mathbf{U}^\top \mathbf{y} = \sum_{j=1}^d \mathbf{u}_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{u}_j^\top \mathbf{y}$$

$$\hat{\mathbf{y}}^{\text{ridge}} = \mathbf{U}\Sigma(\Sigma + \lambda\mathbf{I})^{-1}\Sigma\mathbf{U}^{\top}\mathbf{y} = \sum_{j=1}^d \mathbf{u}_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{u}_j^{\top}\mathbf{y}$$

Итак, как и для линейной регрессии, линейная гребневая регрессия вычисляет проекцию  $\mathbf{u}(\mathbf{u}_j^{\top}\mathbf{y})$  вектора  $\mathbf{y}$  на вектор  $\mathbf{u}_j$ , но затем домножает эту проекцию на

$$\frac{\sigma_j^2}{\sigma_j^2 + \lambda} \leq 1.$$

Чем больше  $\sigma_j$ , тем ближе этот множитель к 1.

Чем меньше  $\sigma_j$ , тем ближе этот множитель к 0.

Таким образом, наибольшему уменьшению подвергаются компоненты, соответствующие меньшим сингулярным числам  $\sigma_j$ .

### 6.2.5. Лассо (R. Tibshirani)

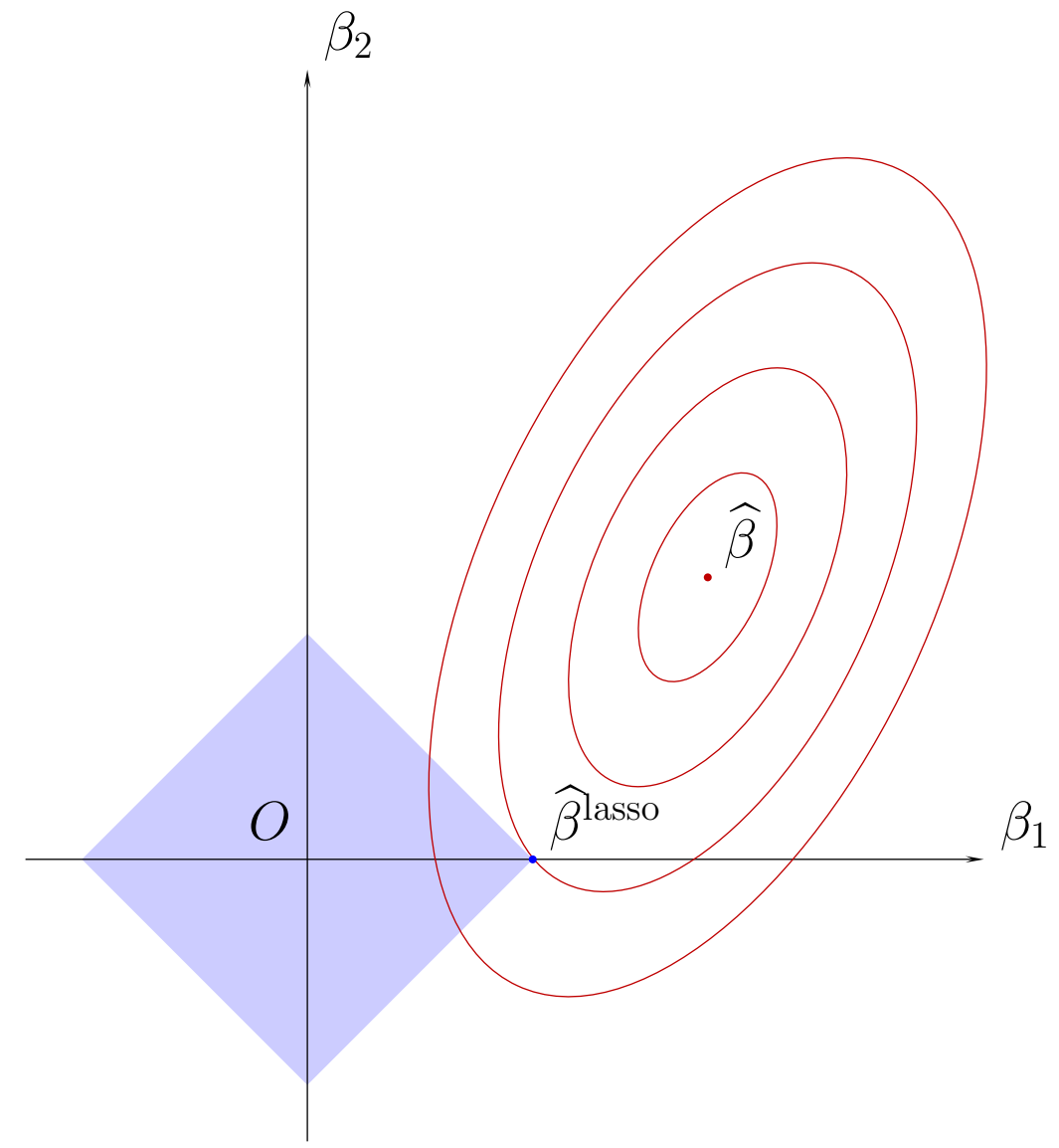
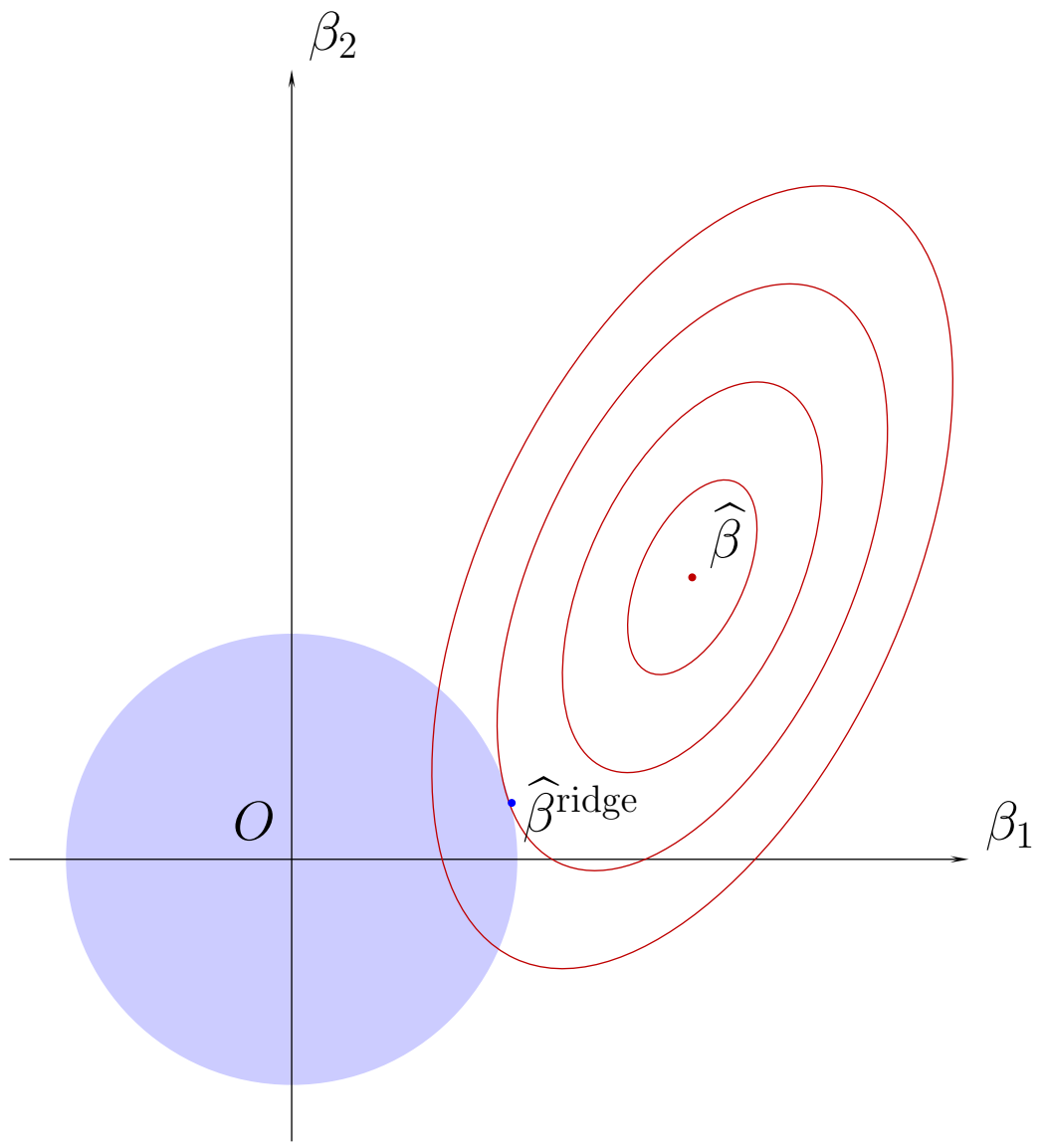
LASSO — «least absolute shrinkage and selection operator»

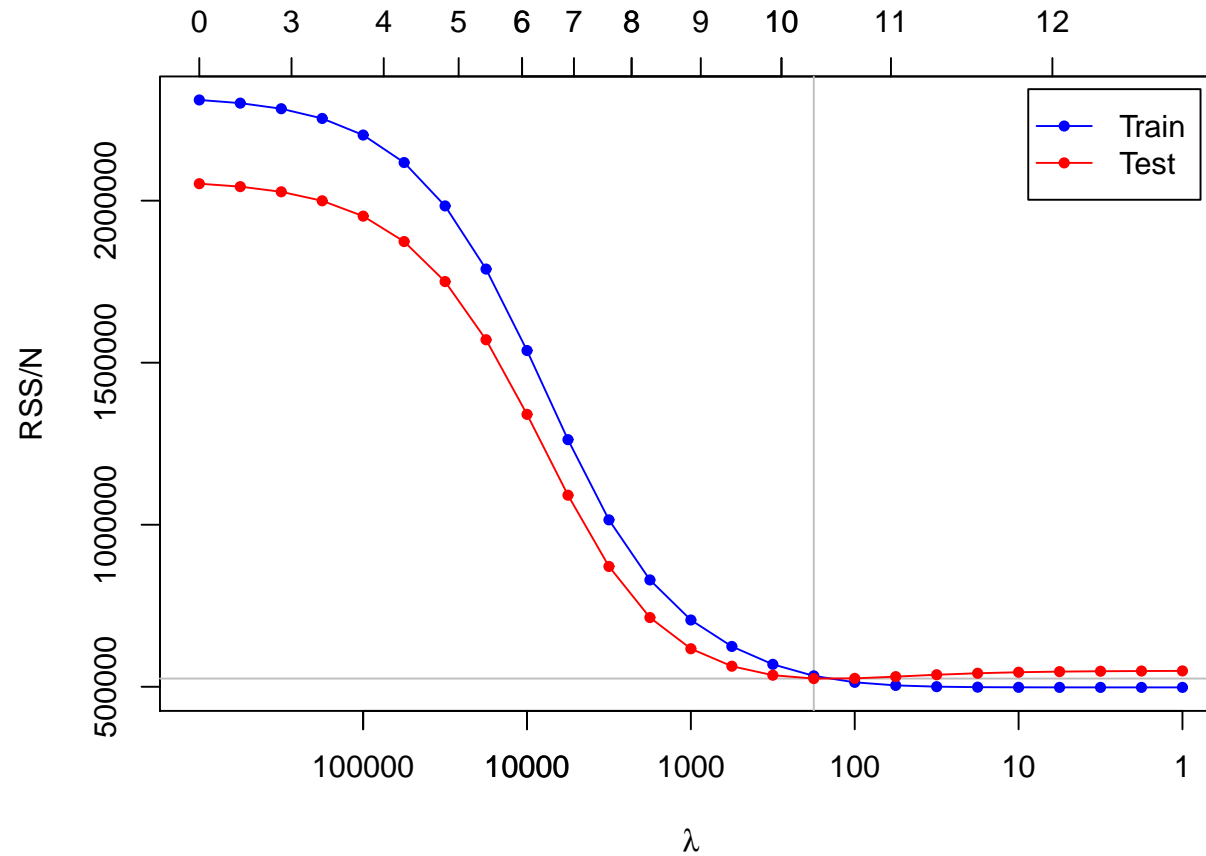
$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left( y^{(i)} - \beta_0 - \sum_{j=1}^d x_j^{(i)} \beta_j \right)^2 \right\}, \text{ при условии } \sum_{j=1}^d |\beta_j| \leq s.$$

Если  $s$  достаточно мало, то в типичной ситуации часть коэффициентов  $\hat{\beta}_j^{\text{lasso}}$  равна в точности 0.

Почему?

(При гребневой регрессии такого, как правило, не происходит)



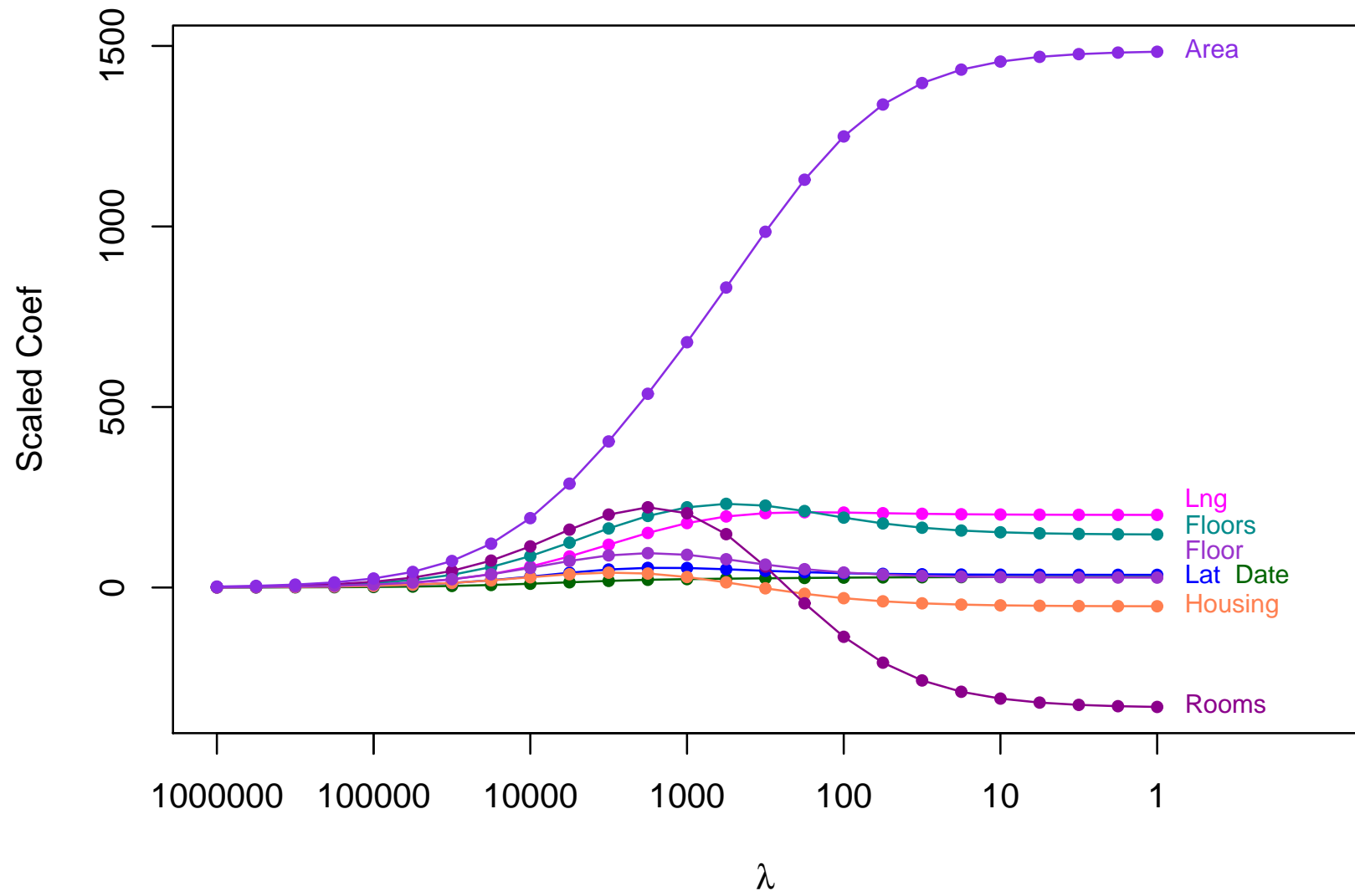


$$\lambda = 177.8279 = 10^{2.25}$$

$$y = 3226.1 + 0.30465 \text{ Date} + 940.46 \text{ Lat} + 2314.9 \text{ Lng} - 53.574 \text{ Housing} \\ + 41.265 \text{ Floors} - 48.654 \text{ Rooms} + 13.105 \text{ Floor} + 56.764 \text{ Area}$$

$$\hat{R}_{\text{train}} = \frac{1}{N_{\text{train}}} \text{RSS}_{\text{train}} = 533652.8, \quad \hat{R}_{\text{test}} = \frac{1}{N_{\text{test}}} \text{RSS}_{\text{test}} = 525341.3$$



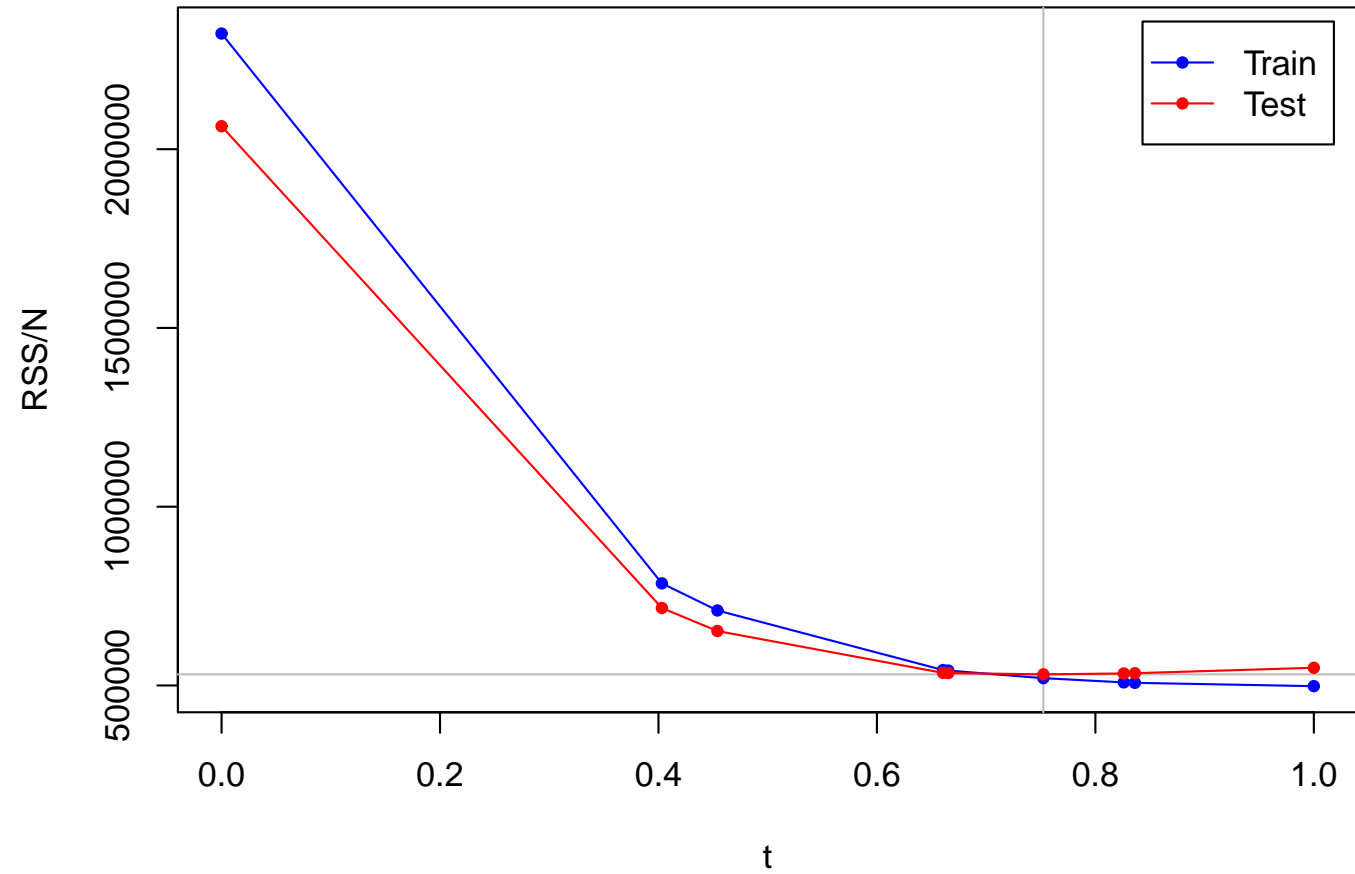


Лассо. Введем параметр

$$t = \frac{s}{\sum_{j=1}^d |\hat{\beta}_j|}$$

При  $t = 0$  все  $\hat{\beta}_j^{\text{lasso}}$  равны нулю.

При  $t = 1$  имеем  $\hat{\beta}_j^{\text{lasso}} = \hat{\beta}_j$  ( $j = 1, 2, \dots, d$ ).



$t = 0.7523712$

$$y = 3226.1 + 2083.4 \text{ Lng} + 34.318 \text{ Floors} - 118.41 \text{ Rooms} + 2.3186 \text{ Floor} + 63.356 \text{ Area}$$

$$\hat{R}_{\text{train}} = \frac{1}{N_{\text{train}}} \text{RSS}_{\text{train}} = 520436.5, \quad \hat{R}_{\text{test}} = \frac{1}{N_{\text{test}}} \text{RSS}_{\text{test}} = 531102.9$$

