

# **Метод дифференциальной кросс-валидации для выбора уровня сложности обобщенных линейных моделей зависимостей**

**Ангуло Бриан Флориан**

Московский физико-технический институт

**Морозов Алексей Олегович**

Московский физико-технический институт

**Моттль Вадим Вячеславович**

Вычислительный центр РАН

# Обобщенная задача восстановления зависимости в линейном пространстве наблюдений

$\mathbf{x} \in \mathbb{R}^n$  – объекты реального мира, наблюдаемые через вещественные признаки

$y \in \mathbb{Y}$  – скрытое свойство объекта

$y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – неизвестная реально существующая зависимость

$\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$  – множество прецедентов (обучающая совокупность)

$\hat{y} = \hat{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – требуется сконструировать решающее правило, применимое каждому  $\mathbf{x} \in \mathbb{R}^n$ ,  $\hat{y} \approx y$  (приблизительно восстановить скрытую зависимость)

Если  $\mathbb{Y} = \mathbb{R}$  оценивание регрессионной зависимости  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

Если  $\mathbb{Y} = \{-1, 1\}$  двух-классовое распознавание образов  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \{-1, 1\}$

## Обобщенная линейная модель (Generalized Linear Model)

*John Nelder. Generalized Linear Models. Journal of the Royal Statistical Society. Series A, Vol. 135, Issue 3, 1972, pp. 370-384.*

$$z(\mathbf{x}|\mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b : \mathbb{R}^n \rightarrow \mathbb{R}$$

параметры модели:

обобщенная линейная модель зависимости

$\mathbf{a} \in \mathbb{R}^n$  – направляющий вектор,  $b \in \mathbb{R}$  – сдвиг

# Обобщенная задача восстановления зависимости в линейном пространстве наблюдений

$\mathbf{x} \in \mathbb{R}^n$  – объекты реального мира, наблюдаемые через вещественные признаки

$y \in \mathbb{Y}$  – скрытое свойство объекта

$y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – неизвестная реально существующая зависимость

$\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$  – множество прецедентов (обучающая совокупность)

$\hat{y} = \hat{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – требуется сконструировать решающее правило, применимое каждому  $\mathbf{x} \in \mathbb{R}^n$ ,  $\hat{y} \approx y$  (приблизительно восстановить скрытую зависимость)

Если  $\mathbb{Y} = \mathbb{R}$  оценивание регрессионной зависимости  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

Если  $\mathbb{Y} = \{-1, 1\}$  двух-классовое распознавание образов  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \{-1, 1\}$

## Обобщенная линейная модель (Generalized Linear Model)

*John Nelder. Generalized Linear Models. Journal of the Royal Statistical Society. Series A, Vol. 135, Issue 3, 1972, pp. 370-384.*

$$z(\mathbf{x}|\mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b : \mathbb{R}^n \rightarrow \mathbb{R}$$

обобщенная линейная модель зависимости

$$q(y, z) : \mathbb{Y} \times \mathbb{R} \rightarrow \mathbb{R}^+$$

функция связи (штрафа, потерь)

# Обобщенная задача восстановления зависимости в линейном пространстве наблюдений

$\mathbf{x} \in \mathbb{R}^n$  – объекты реального мира, наблюдаемые через вещественные признаки

$y \in \mathbb{Y}$  – скрытое свойство объекта

$y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – неизвестная реально существующая зависимость

$\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$  – множество прецедентов (обучающая совокупность)

$\hat{y} = \hat{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{Y}$  – требуется сконструировать решающее правило, применимое каждому  $\mathbf{x} \in \mathbb{R}^n$ ,  $\hat{y} \approx y$  (приблизительно восстановить скрытую зависимость)

Если  $\mathbb{Y} = \mathbb{R}$  оценивание регрессионной зависимости  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

Если  $\mathbb{Y} = \{-1, 1\}$  двух-классовое распознавание образов  $y = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \{-1, 1\}$

## Обобщенная линейная модель (Generalized Linear Model)

*John Nelder. Generalized Linear Models. Journal of the Royal Statistical Society. Series A, Vol. 135, Issue 3, 1972, pp. 370-384.*

$$z(\mathbf{x}|\mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b : \mathbb{R}^n \rightarrow \mathbb{R}$$

обобщенная линейная модель зависимости

$$q(y, z) : \mathbb{Y} \times \mathbb{R} \rightarrow \mathbb{R}^+$$

функция связи (штрафа, потерь)

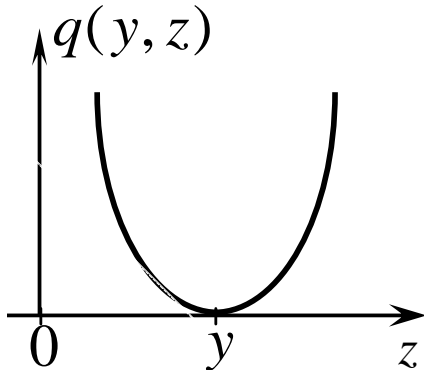
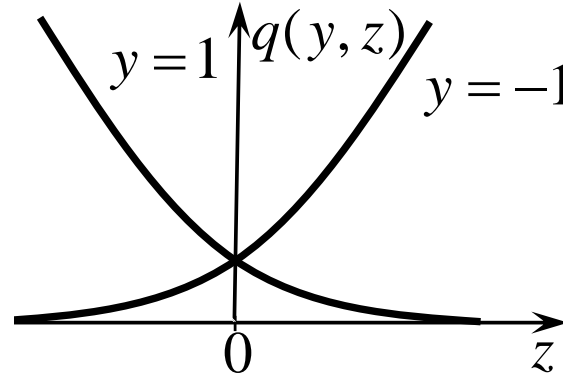
$$\hat{y}(\mathbf{x}|\mathbf{a}, b) = \arg \min_{y \in \mathbb{Y}} q(y, z(\mathbf{x}|\mathbf{a}, b))$$

решающее правило

# Обобщенная линейная модель зависимости

$\begin{cases} z(\mathbf{x} \mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b: \mathbb{R}^n \rightarrow \mathbb{R} \\ q(y, z): \mathbb{Y} \times \mathbb{R} \rightarrow \mathbb{R}^+ \end{cases}$	обобщенный линейный признак объекта, представленного вектором признаков $\mathbf{x} \in \mathbb{R}^n$
	функция связи, выпуклая по $z \in \mathbb{R}$ для любого $y \in \mathbb{Y}$
$\hat{y}(\mathbf{x} \mathbf{a}, b) = \arg \min_{y \in \mathbb{Y}} q(y, z(\mathbf{x} \mathbf{a}, b))$	решающее правило

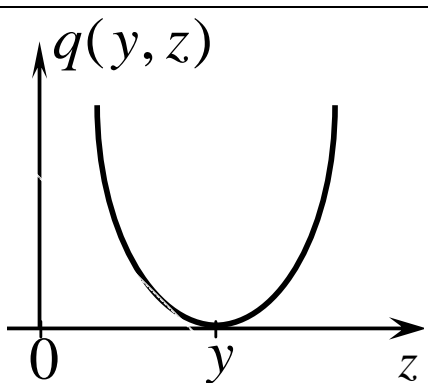
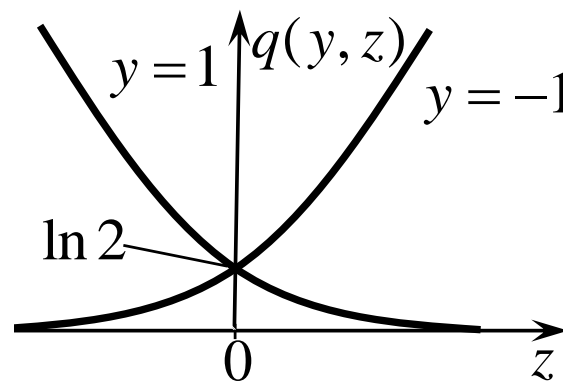
## Частные случаи

Регрессия $\mathbb{Y} = \mathbb{R}$	Двух-классовое распознавание образов $\mathbb{Y} = \{-1, 1\}$
$q(y, z) = (y - z)^2: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$	$\begin{aligned} \lim_{z \rightarrow -\infty} q(y=+1, z) = \infty, & \quad \lim_{z \rightarrow \infty} q(y=+1, z) = 0, \\ \lim_{z \rightarrow -\infty} q(y=-1, z) = 0, & \quad \lim_{z \rightarrow \infty} q(y=-1, z) = \infty. \end{aligned}$
	

# Обобщенная линейная модель зависимости

$\begin{cases} z(\mathbf{x} \mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b: \mathbb{R}^n \rightarrow \mathbb{R} \\ q(y, z): \mathbb{Y} \times \mathbb{R} \rightarrow \mathbb{R}^+ \end{cases}$	обобщенный линейный признак объекта, представленного вектором признаков $\mathbf{x} \in \mathbb{R}^n$
	функция связи, выпуклая по $z \in \mathbb{R}$ для любого $y \in \mathbb{Y}$
$\hat{y}(\mathbf{x} \mathbf{a}, b) = \arg \min_{y \in \mathbb{Y}} q(y, z(\mathbf{x} \mathbf{a}, b))$	решающее правило

## Частные случаи

Регрессия $\mathbb{Y} = \mathbb{R}$	Двух-классовое распознавание образов $\mathbb{Y} = \{-1, 1\}$
$q(y, z) = (y - z)^2: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$	$\begin{aligned} \lim_{z \rightarrow -\infty} q(y=+1, z) = \infty, & \quad \lim_{z \rightarrow \infty} q(y=+1, z) = 0, \\ \lim_{z \rightarrow -\infty} q(y=-1, z) = 0, & \quad \lim_{z \rightarrow \infty} q(y=-1, z) = \infty. \end{aligned}$
	

Логистическая регрессия  

$$q(y, z) = \ln[1 + \exp(-yz)]$$

# Общепринятый принцип обучения по прецедентам: Минимум регуляризованного эмпирического риска

Множество прецедентов (обучающая совокупность):  $\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$

Надо выбрать два параметра  $(\mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R})$  линейной модели

Критерий: Минимум потерь в пределах обучающей совокупности

$EmpR(\mathbf{a}, b) =$	$\sum_{j=1}^N q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \rightarrow \min$	эмпирический риск на обуч. совокуп., вместо среднего риска по всем объектам реального мира
-------------------------	---	--

Однако, если  $n > N$ , то существует континуум точек минимума  $(\mathbf{a}, b) \in \mathbb{R}^{n+1}$ .

## Минимум регуляризованного эмпирического риска

$$J(\mathbf{a}, b) = \gamma \mathbf{a}^T \mathbf{a} + \sum_{j=1}^N q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \rightarrow \min(\mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R})$$

простейшая  
ридж регуляризация  
 $\gamma > 0$ ,

$$J(\mathbf{a}, b | \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \rightarrow \min$$

более изощренная  
селективная регуляризация

Параметр селективности  $0 \leq \mu < \infty$ . При увеличении  $\mu$ , штраф  $\mu |a_i|$  «тянет» к нулю коэффициенты при «лишних» признаках, мало способствующих уменьшению эмпирического риска.

Результат – малое подмножество активных признаков:  $\hat{\mathbb{I}}(\gamma, \mu) = \{i: a_{\gamma, \mu, i} \neq 0\} \subseteq \{1, \dots, n\}$

## Задача данного доклада:

# Предложить подходящий способ выбора структурных параметров обобщенной линейной модели зависимости с отбором признаков объектов

Критерий минимума регуляризованного эмпирического риска с отбором признаков содержит два структурных параметра

$$J(\mathbf{a}, b | \gamma, \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \rightarrow \min(\mathbf{a}, b)$$

$\gamma > 0$  – ридж коэффициент  
 $\mu \geq 0$  – коэффициент селективности

Результат работы алгоритма оптимизации (доклады В.В. Моттля и А. Морозова):

$$(\hat{\mathbf{a}}, \hat{b})_{\gamma, \mu} = \arg \min J(\mathbf{a}, b | \gamma, \mu) \quad \text{найденные параметры модели}$$

$$\hat{\mathbb{I}}_{\gamma, \mu} = \{i: \hat{a}_{\gamma, \mu, i} \neq 0\} \subseteq \{1, \dots, n\} \quad \text{и подмножество активных признаков}$$

## Традиционный метод скользящего контроля (Leave-One-Out cross-validation)

Решить задачу обучения  $N$  раз без одного обучающего объекта

$$(\hat{\mathbf{a}}, \hat{b})_{\gamma, \mu}^{(j)} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{t=1}^N q(y_t, \mathbf{a}^T \mathbf{x}_t + b) - q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \right\}$$

Выбрать структурные параметры по критерию  $\sum_{t=1}^N q(y_t, \hat{\mathbf{a}}_{\gamma, \mu}^{(j)T} \mathbf{x}_t + \hat{b}_{\gamma, \mu}^{(j)}) \rightarrow \min_{\gamma, \mu}$  Недостаток – вычисление критерия требует  $N$ -кратного обучения



## Традиционный метод скользящего контроля (Leave-One-Out Cross-Validation)

Решить задачу обучения  $N$  раз без одного обучающего объекта

$$(\hat{\mathbf{a}}, \hat{b})_{\gamma, \mu}^{(j)} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{t=1}^N q(y_t, \mathbf{a}^T \mathbf{x}_t + b) - q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \right\}$$

Выбрать структурные параметры по критерию  $\sum_{t=1}^N q(y_t, \hat{\mathbf{a}}_{\gamma, \mu}^{(j)T} \mathbf{x}_t + \hat{b}_{\gamma, \mu}^{(j)}) \rightarrow \min_{\gamma, \mu}$  Недостаток – вычисление критерия требует  $N$ -кратного обучения

## Наша идея – принцип дифференциального скользящего контроля (Differential LOO)

1) Удалять из обучающей совокупности не весь очередной объект, а лишь малую его часть, не приводящую к изменению подмножества активных признаков

$$(\hat{\mathbf{a}}, \hat{b})_{\gamma, \mu}(p_j) = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{t=1}^N q(y_t, \mathbf{a}^T \mathbf{x}_t + b) - p_j q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \right\}, p_j \rightarrow 0$$

$$\hat{\mathbb{I}}_{\gamma, \mu}(p_j) = \{i: a_{\gamma, \mu, i}^{(p_j)} \neq 0\} \subseteq \{1, \dots, n\} = \text{const}$$

2) В качестве критерия принять сумму производных потерь на всех объектах

$$DLOO(\gamma, \mu) = \sum_{j=1}^N \frac{\partial}{\partial p_j} q(y_j, \mathbf{x}_j^T \hat{\mathbf{a}}_{\gamma, \mu}(p_j) + \hat{b}_{\gamma, \mu}(p_j)).$$

Такой критерий показывает сумму скоростей роста потерь при «тенденции к удалению» каждого объекта из обучения.

# Наша идея – принцип дифференциального скользящего контроля (**Differential LOO, DLLO**)

1) Удалять из обучающей совокупности не весь очередной объект, а лишь малую его часть, не приводящую к изменению подмножества активных признаков

$$(\hat{\mathbf{a}}, \hat{b})_{\gamma, \mu}(p_j) = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{t=1}^N q(y_t, \mathbf{a}^T \mathbf{x}_t + b) - p_j q(y_j, \mathbf{a}^T \mathbf{x}_j + b) \right\}, p_j \rightarrow 0$$

$$\hat{\mathbb{I}}_{\gamma, \mu}(p_j) = \{i: \hat{a}_{\gamma, \mu, i}(p_j) \neq 0\} \subseteq \{1, \dots, n\} = \text{const}$$

2) В качестве критерия принять сумму производных потерь на всех объектах

$$DLLO(\gamma, \mu) = \sum_{j=1}^N \frac{\partial}{\partial p_j} q(y_j, \mathbf{x}_j^T \hat{\mathbf{a}}_{\gamma, \mu}(p_j) + \hat{b}_{\gamma, \mu}(p_j)).$$

Такой критерий показывает сумму скоростей роста потерь при «тенденции к удалению» каждого объекта из обучения.

**В данном докладе мы применим эту идею к частному случаю селективной модели регрессионной зависимости**

Квадратичная функция связи функция  $q(y_t, \mathbf{a}^T \mathbf{x}_t + b) = (y_t - (\mathbf{a}^T \mathbf{x}_t + b))^2$

Для упрощения выкладок будем полагать  $b = 0$ . Исходная задача:

$$J(\mathbf{a} | \gamma, \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 \rightarrow \min(\mathbf{a}), \mathbf{x}_j, \mathbf{a} \in \mathbb{R}^n$$

## В данном докладе мы применим эту идею к частному случаю селективной модели регрессионной зависимости

Квадратичная функция связи  $q(y_t, \mathbf{x}_t^T \mathbf{a} + b) = (y_t - (\mathbf{x}_t^T \mathbf{a} + b))^2$

Для упрощения выкладок будем полагать  $b = 0$ . Исходная задача:

$$J(\mathbf{a} | \gamma, \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_j - \mathbf{x}_j^T \mathbf{a})^2 \rightarrow \min(\mathbf{a}), \quad \mathbf{x}_t, \mathbf{a} \in \mathbb{R}^n$$

Сначала рассмотрим задачу без селективности

$$\hat{\mathbf{a}}_\gamma = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 \right\}.$$

$$\hat{\mathbf{a}}_\gamma^{(j)} = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{t=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{удалим из обучения один объект } j$$

Классический критерий *LOO*

$$Loo(\gamma) = \sum_{j=1}^N \left( y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma^{(j)} \right)^2$$

## В данном докладе мы применим эту идею к частному случаю селективной модели регрессионной зависимости

Квадратичная функция связи  $q(y_t, \mathbf{x}_t^T \mathbf{a} + b) = (y_t - (\mathbf{x}_t^T \mathbf{a} + b))^2$

Для упрощения выкладок будем полагать  $b = 0$ . Исходная задача:

$$J(\mathbf{a} | \gamma, \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_j - \mathbf{x}_j^T \mathbf{a})^2 \rightarrow \min(\mathbf{a}), \quad \mathbf{x}_t, \mathbf{a} \in \mathbb{R}^n$$

### Сначала рассмотрим задачу без селективности

$$\hat{\mathbf{a}}_\gamma = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 \right\}.$$

$$\hat{\mathbf{a}}_\gamma^{(j)} = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{t=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{удалим из обучения один объект } j$$

Классический критерий *LOO* – перебор не требуется

$$Loo(\gamma) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma}{1 - \mathbf{x}_j^T \left( (\mathbf{X}\mathbf{X}^T) + \gamma \mathbf{I}_n \right)^{-1} \mathbf{x}_j} \right)^2$$

M.A. Goldberg, H.A. Cho. Introduction to Regression Analysis. WIT Press, 2004.  
The Sherman-Morrison-Woodbury formula.  
Page 141.

## В данном докладе мы применим эту идею к частному случаю селективной модели регрессионной зависимости

Квадратичная функция связи  $q(y_t, \mathbf{x}_t^T \mathbf{a} + b) = (y_t - (\mathbf{x}_t^T \mathbf{a} + b))^2$

Для упрощения выкладок будем полагать  $b = 0$ . Исходная задача:

$$J(\mathbf{a} | \gamma, \mu) = \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_j - \mathbf{x}_j^T \mathbf{a})^2 \rightarrow \min(\mathbf{a}), \quad \mathbf{x}_t, \mathbf{a} \in \mathbb{R}^n$$

Сначала рассмотрим задачу без селективности

$$\hat{\mathbf{a}}_\gamma = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 \right\}.$$

$$\hat{\mathbf{a}}_\gamma^{(j)} = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{t=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{удалим из обучения один объект } j$$

Классический критерий  $LOO$  – перебор не требуется

$$Loo(\gamma) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma}{1 - \mathbf{x}_j^T \left( (\mathbf{X}\mathbf{X}^T) + \gamma \mathbf{I}_n \right)^{-1} \mathbf{x}_j} \right)^2$$

Но мы не можем так делать из-за присутствия селективной регуляризации, изменяющей размерность векторов признаков.

## Сначала рассмотрим задачу без селективности

$$\hat{\mathbf{a}}_\gamma^{(j)} = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{t=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{удалим из обучения один объект } j$$

Классический критерий *LOO* – перебор не требуется

$$Loo(\gamma) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma}{1 - \mathbf{x}_j^T ((\mathbf{X}\mathbf{X}^T) + \gamma \mathbf{I}_n)^{-1} \mathbf{x}_j} \right)^2 \quad \text{Но мы не можем так делать из-за отсутствия селективной регуляризации, изменяющей размерность векторов признаков.}$$

$$\hat{\mathbf{a}}_\gamma^{(j)} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{Обычный LOO: полное удаление}$$

Вместо полного удаления очередного объекта удалим его малую часть, тогда оценка направляющего вектора изменится мало, и его размерность не изменится.

$$\hat{\mathbf{a}}_\gamma^{(p_j)} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - p_j (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \begin{array}{l} \text{Дифференц. LOO:} \\ p_j \rightarrow 0 \end{array}$$

Остается вычислить скорости роста ошибок  $\partial / \partial p_j (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma^{(p_j)})$ .

## Дифференциальный LOO для квадратичной модели

$$\hat{\mathbf{a}}_\gamma(p_j) = \arg \min \left\{ \gamma \mathbf{a}^T \mathbf{a} + \sum_{j=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 - p_j (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right\} \quad \text{Дифференц. LOO: } p_j \rightarrow 0$$

$\hat{\mathbf{a}}_\gamma = \hat{\mathbf{a}}_\gamma(0)$  – оценка направляющего вектора по всей обучающей совокупности

$\hat{\mathbf{a}}_\gamma^{(j)} = \hat{\mathbf{a}}_\gamma(1)$  – полный LOO,  $p_j=1$ ,  $\hat{\mathbf{a}}_\gamma(p_j)$  – зависимость от удаляемой доли объекта

$$\hat{\mathbf{a}}_\gamma(p_j) = (1 - p_j) \hat{\mathbf{a}}_\gamma + p_j \hat{\mathbf{a}}_\gamma^{(j)} \quad \text{производная } \frac{\partial}{\partial p_j} \hat{\mathbf{a}}_\gamma(p_j) = \hat{\mathbf{a}}_\gamma^{(j)} - \hat{\mathbf{a}}_\gamma$$

Теорема. В этом случае 
$$\frac{\partial}{\partial p_j} (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma(p_j)) = \frac{y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma}{1 - \mathbf{x}_j^T ((\mathbf{X}\mathbf{X}^T) + \gamma \mathbf{I}_n)^{-1} \mathbf{x}_j}$$

Итоговый дифференциальный критерий LOO для квадратичной модели формально совпадает с обычным критерием LOO

$$Dif Loo(\gamma) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_\gamma}{1 - \mathbf{x}_j^T ((\mathbf{X}\mathbf{X}^T) + \gamma \mathbf{I}_n)^{-1} \mathbf{x}_j} \right)^2$$

Вопрос: Как применить этот критерий к селективной модели?

$$J(\mathbf{a} | \gamma, \mu) = \gamma \sum_{i=1}^n \left( \begin{array}{l} 2\mu |a_i|, |a_i| \leq \mu \\ \mu^2 + a_i^2, |a_i| > \mu \end{array} \right) + \sum_{j=1}^N (y_t - \mathbf{x}_t^T \mathbf{a})^2 \rightarrow \min(\mathbf{a})$$

# Дифференциальный LOO для селективной квадратичной модели

Пусть задача решена при конкретных значениях структурных параметров  $(\gamma, \mu)$

$$\hat{\mathbf{a}}_{\gamma, \mu} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N \left( y_j - \sum_{i=1}^n a_i x_{j,i} \right)^2 \right\} \quad \hat{\mathbf{a}}_{\gamma, \mu} = (\hat{a}_{\gamma, \mu, 1} \cdots \hat{a}_{\gamma, \mu, n})$$

доклады В.В. Моттля  
и А.О. Морозова

$$\hat{\mathbb{I}}_{\gamma, \mu} = \{i : \hat{a}_{\gamma, \mu, i} \neq 0\} \subseteq \{1, \dots, n\} - \text{подмножество } \hat{n}_{\gamma, \mu} = |\hat{\mathbb{I}}_{\gamma, \mu}| \leq n \text{ активных признаков}$$

$$\mathbf{x}_{\gamma, \mu, j} = (x_{j,i}, i \in \hat{\mathbb{I}}_{\gamma, \mu}) \in \mathbb{R}^{\hat{n}_{\gamma, \mu}} - \text{активные признаки} \quad \mathbf{X}_{\gamma, \mu} = (\mathbf{x}_{\gamma, \mu, 1} \cdots \mathbf{x}_{\gamma, \mu, N}) - \text{матрица } (\hat{n}_{\gamma, \mu} \times N)$$

Свойство решения:

$$\hat{\mathbf{a}}_{\gamma, \mu} = \arg \min \left\{ \gamma \sum_{i \in \hat{\mathbb{I}}_{\gamma, \mu}} a_i^2 + \sum_{j=1}^N \left( y_j - \sum_{i \in \hat{\mathbb{I}}_{\gamma, \mu}} a_i x_{j,i} \right)^2 \right\}$$

Если запомнить подмножество активных признаков  $\hat{\mathbb{I}}_{\gamma, \mu} = \{i : \hat{a}_{\gamma, \mu, i} \neq 0\}$ , то вектор  $\hat{\mathbf{a}}_{\gamma, \mu}$  легко восстановить

Это эквивалентная модель зависимости, в которой роль структурных параметров вместо пары  $(\gamma, \mu)$  играет пара  $(\mathbb{I}_{\gamma, \mu}, \gamma)$ . Идея дифференциального критерия LOO для исходной селективной модели выражается условием

$$Dif Loo(\gamma, \mu) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}_{\gamma, \mu}, \gamma}}{1 - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \left( (\mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}} \mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}}^T) + \gamma \mathbf{I}_{\hat{n}_{\gamma, \mu}} \right)^{-1} \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}} \right)^2 \rightarrow \min(\gamma, \mu)$$



## Эксперимент

Регрессоры  $\mathbf{x}_t = (x_{t,i}, i=1, \dots, n), t=1, \dots, N=251$ :

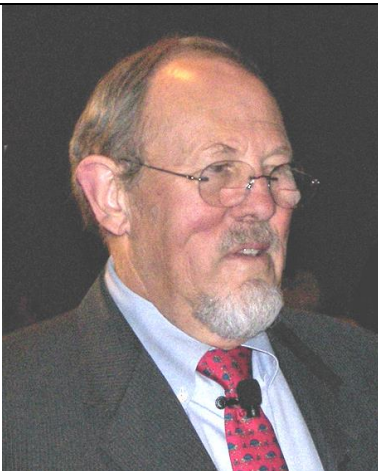
Множество временных рядов месячных доходностей  $n=650$  биржевых ценных бумаг на Нью-Йоркской фондовой бирже чуть больше, чем за два года.

Наблюдаемый сигнал  $y_t, t=1, \dots, N=251$ :

Временной ряд доходностей инвестиционного портфеля, построенного как вложение капитала в равных долях в  $n^*=13$  неизвестных ценных бумаг в множестве  $\{1, \dots, n\}$ .

Требуется найти подмножество активных ценных бумаг.

Регрессионная модель:  $y_t \cong \sum_{i=1}^n a_i x_{t,i}$  – модель Шарпа,  $a_i$  – доли вложения капитала.



**Уильям Шарп**, почетный Professor of Finance  
[Graduate School of Business, Stanford University](#)

Лауреат Нобелевской премии по экономике 1990 г.  
 за ряд достижений, в частности, за

**Returns Based Style Analysis –RBSA**

оценивание стиля и эффективности инвестиционных фондов

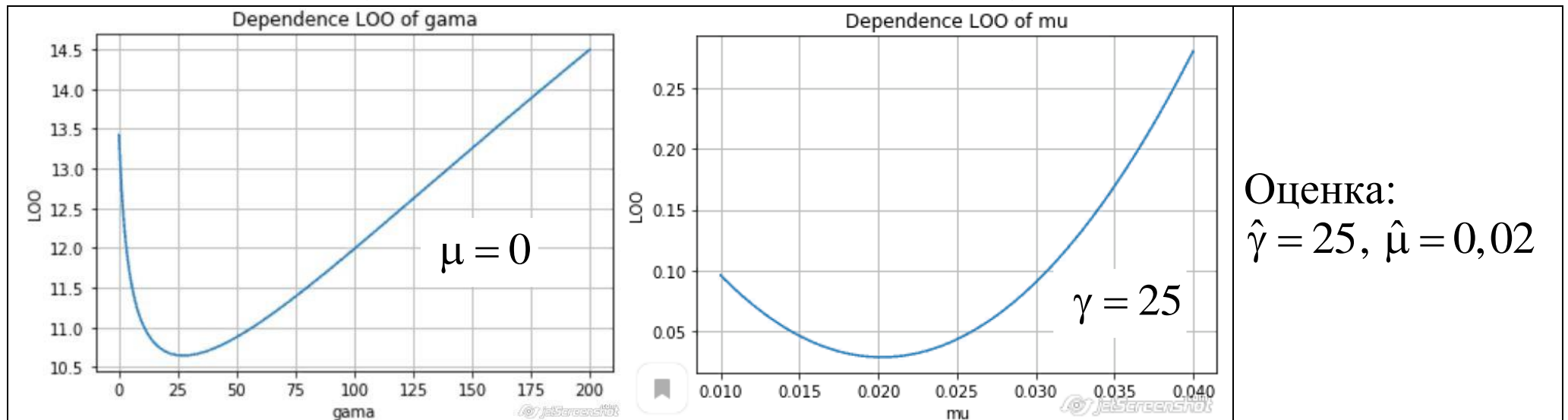
Мы рассматриваем дополнительную задачу Factor Search – оценивание фактического состава инвестиционного портфеля в большом множестве биржевых активов. У Шарпа задачи Factor Search не было.

Результат идентификации селективной модели

$$\hat{\mathbf{a}}_{\gamma, \mu} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N \left( y_j - \sum_{i=1}^n a_i x_{j,i} \right)^2 \right\}, \quad n = 650, \quad N = 251,$$

по критерию дифференциального скользящего контроля

$$Dif Loo(\gamma, \mu) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \hat{\mathbf{a}}_{\hat{\gamma}, \hat{\mu}, \gamma}}{1 - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \left( (\mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}} \mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}}^T) + \gamma \mathbf{I}_{\hat{n}_{\gamma, \mu}} \right)^{-1} \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}} \right)^2 \rightarrow \min(\gamma, \mu)$$



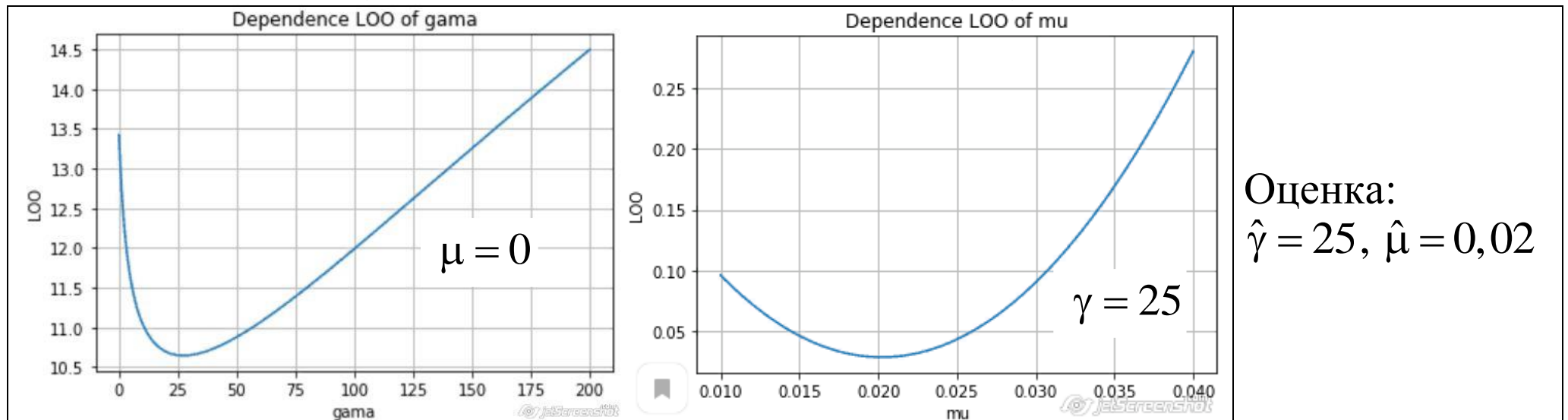
Найденное подмножество  $\hat{\mathbb{I}}_{\hat{\gamma}, \hat{\mu}} = \{i : \hat{a}_{\hat{\gamma}, \hat{\mu}, i} \neq 0\}$  содержит  $\hat{n}_{\hat{\gamma}, \hat{\mu}} = 16$  биржевых активов, содержащее  $n^* = 13$  истинных.

Результат идентификации селективной модели

$$\hat{\mathbf{a}}_{\gamma, \mu} = \arg \min \left\{ \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu |a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{j=1}^N \left( y_j - \sum_{i=1}^n a_i x_{j,i} \right)^2 \right\}, \quad n = 650, \quad N = 251,$$

по критерию дифференциального скользящего контроля

$$Dif Loo(\gamma, \mu) = \sum_{j=1}^N \left( \frac{y_j - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \hat{\mathbf{a}}_{\hat{\gamma}, \hat{\mu}, \gamma}}{1 - \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}^T \left( (\mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}} \mathbf{X}_{\hat{\mathbb{I}}_{\gamma, \mu}}^T) + \gamma \mathbf{I}_{\hat{n}_{\gamma, \mu}} \right)^{-1} \mathbf{x}_{\hat{\mathbb{I}}_{\gamma, \mu}, j}} \right)^2 \rightarrow \min(\gamma, \mu)$$



Найденное подмножество  $\hat{\mathbb{I}}_{\hat{\gamma}, \hat{\mu}} = \{i : \hat{a}_{\hat{\gamma}, \hat{\mu}, i} \neq 0\}$  содержит  $\hat{n}_{\hat{\gamma}, \hat{\mu}} = 16$  биржевых активов, содержащее  $n^* = 13$  истинных.

*Можно сказать, что мы нашли иголку в стоге сена!*

# Благодарности

Работа поддержана грантами РФФИ:

19-37-90159-Аспиранты (А.О. Морозов),

17-07-00436-а (В.В. Моттль).

**Спасибо за внимание!**

**Вопросы?**