

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

Ангуло Яури Бриан Флориан

Московский физико-технический институт  
Факультет управления и прикладной математики  
Кафедра интеллектуальных систем

Научный руководитель: д.т.н. профессор В.В.Моттль

17 июня 2020 г.

# Цель исследования

## Цель

Разработать подходящий метод выбора оптимальных гиперпараметров для обобщенной линейной модели восстановления зависимостей при селективном отборе признаков.

## Задача

Требуется найти оптимальные гиперпараметры модели в случае когда число признаков больше числа объектов.

## Предложение

Удалить бесконечно малую часть из одного объекта обучающей совокупности, и измерить скорость роста потерь.

# Обобщенный линейный подход

## Обобщенная линейная модель

### Задача восстановления зависимости

$$X = \{(x_j, y_j)_{j=1}^N\}, x_j \in \mathbb{R}^n, y_j \in \mathbb{Y}$$

Требуется сконструировать решающее правило, применимое к каждому  $x \in \mathbb{R}^n$

$$\hat{y} = f(x) : \mathbb{R}^n \rightarrow \mathbb{Y}$$

### Обобщенная линейная модель

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

где  $a$  — направляющий вектор,  $b$  — сдвиг

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

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

## Минимум регуляризованного эмпирического риска

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

где  $\mu$  — гиперпараметр селективности;  
 $\gamma$  — гиперпараметр Ридж-регуляризации;

В качестве регуляризации применяется селективная Ридж-регуляризация<sup>1</sup>.

Результат — подмножество активных (информативных) признаков:

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

# Обобщенный линейный подход

С учетом вышесказанного, что число признаков больше числа объектов, целесообразно использовать LOO.

## Предположение

Из-за селективности признаков нельзя напрямую использовать LOO. Поскольку удаление одного объекта целиком, может привести к изменению подмножества активных признаков.

## Идея

Удаление бесконечно малой части из одного объекта обучающей совокупности, не приведет к изменению подмножества активных (информативных) признаков. Доля объекта, которую удаляем, будет равна  $p$ ,  $p \rightarrow 0$ .

# Обобщенный линейный подход

## Взвешенная задача

Рассмотрим немного более общую задачу.

### Взвешенная задача

Во взвешенной задаче предполагается, что объекты встречаются в обучающей выборке с некоторыми весами  $r_j = 1$ ,  $r_t = 1 - p$ ,  $p \rightarrow 0$

$$J(a|\gamma, \mu) = \gamma \sum_{i=1}^n \left( \begin{array}{l} 2\mu|a_i|, |a_i| \leq \mu \\ \mu^2 + a_i^2, |a_i| > \mu \end{array} \right) + \sum_{j=1}^N q(y_j, r_j x_j^T a) \rightarrow \min(a)$$

Можно показать, что следующая задача эквивалентна предыдущей:

### Взвешенная задача

$$J(a|\gamma, \mu) = \gamma \sum_{i=1}^n \left( \begin{array}{l} 2\mu|a_i|, |a_i| \leq \mu \\ \mu^2 + a_i^2, |a_i| > \mu \end{array} \right) + \sum_{j=1}^N r_j q(y_j, x_j^T a) \rightarrow \min(a)$$

# Обобщенный линейный подход

Решение взвешенной задачи в двойственной форме

Для решения  $J(a|\gamma, \mu)$  удобно представить задачу в разделенной записи, следующим образом:

$$\begin{cases} \gamma \sum_{i=1}^n \begin{pmatrix} 2\mu|a_i|, & |a_i| \leq \mu \\ \mu^2 + a_i^2, & |a_i| > \mu \end{pmatrix} + \sum_{t=1}^N r_j q(y_j, z_j) \rightarrow \min(a \in \mathbb{R}^n, z_j) \\ z_j = a^T x_j, j = 1, \dots, N \end{cases}$$

Решение задачи сводится к поиску седловой точки функции Лагранжа

$$-\frac{1}{2} \lambda^T \left( X_{\mathbb{I}_{\gamma, \mu}}^T X_{\mathbb{I}_{\gamma, \mu}} \right) \lambda + \frac{1}{2\gamma} \sum_{t=1}^N r_j q(y_j, z_j) + z^T \lambda \rightarrow \begin{cases} \min(z) \\ \nabla_{\lambda} = 0 \end{cases}$$

# Численное решение задачи

## Итерационный алгоритм Ньютона

Предыдущая задача является выпуклой гладкой задачей оптимизации, если функция потерь выбрана также выпуклой и гладкой.

$$\hat{\lambda}_{\gamma, \mu} = \arg \min \left\{ \frac{1}{2} \lambda^T \left( X_{\mathbb{I}_{\gamma, \mu}}^T X_{\mathbb{I}_{\gamma, \mu}} \right) \lambda + \sum_{j=1}^N \left[ - \inf_{z \in \mathbb{R}} \left( \frac{1}{2\gamma} q(y_j, z_j) + \lambda_j z_j \right) \right] \right\}$$

Итерационный алгоритм Ньютоновского типа  $\lambda^k$ ,  $k = 0, 1, \dots$  основан на разложении выпуклых функции в ряд Тейлора на каждом шаге в окрестности очередной точки  $\lambda^k$ .

Когда решение двойственной задачи найдено  $\hat{\lambda}_{\gamma, \mu} \in \mathbb{R}^N$  компоненты направляющего вектора  $\hat{\mathbb{I}}_{\gamma, \mu}$  находятся независимо друг от друга:

$$\hat{\mathbb{I}}_{\gamma, \mu, i} \begin{cases} 0, & |x_i^T \hat{\lambda}_{\gamma, \mu}| \leq \mu, \\ x_i^T \hat{\lambda}_{\gamma, \mu}, & |x_i^T \hat{\lambda}_{\gamma, \mu}| < \mu. \end{cases} \quad \hat{\mathbb{I}}_{(\lambda|\mu)} = \{i : |x_i^T \hat{\lambda}_{\gamma, \mu}| > \mu\} \subseteq \{1, \dots, n\}$$



# Численное решение задачи

## Итерационный алгоритм Ньютона

В случае линейной регрессии можно получить следующее выражение:

$$W(\boldsymbol{\lambda}|\gamma, \mu) = \boldsymbol{\lambda}^T \left( X_{\mathbb{I}_{\gamma, \mu}}^T X_{\mathbb{I}_{\gamma, \mu}} \right) \boldsymbol{\lambda} + \frac{\gamma}{2} \boldsymbol{\lambda}^T \boldsymbol{\lambda} - Y^T \boldsymbol{\lambda} \rightarrow \min(\boldsymbol{\lambda})$$

Для решения задачи воспользуемся итерационным алгоритмом Ньютона с переменным шагом.

$$\begin{cases} \tilde{\boldsymbol{\lambda}}^{k+1} = \arg \min_{\boldsymbol{\lambda}} \tilde{W}^k(\boldsymbol{\lambda}|\gamma, \mu) \\ \tilde{W}^k(\boldsymbol{\lambda}|\gamma, \mu) = \boldsymbol{\lambda}^T \left( X_{\mathbb{I}^k}^T X_{\mathbb{I}^k} \right) \boldsymbol{\lambda} + \frac{\gamma}{2} \boldsymbol{\lambda}^T \boldsymbol{\lambda} - Y^T \boldsymbol{\lambda} \rightarrow \min(\boldsymbol{\lambda}) \end{cases}$$

Критерий остановки:  $\hat{\mathbb{I}}(\boldsymbol{\lambda}^{k+1}) = \hat{\mathbb{I}}(\boldsymbol{\lambda}^k)$

# Численное решение задачи

## Итерационный алгоритм Ньютона

Поиск седловой точки сводится к системе линейных уравнений.  
Решение системы показывает только наилучшее направление шага

### Линейная регрессия

$$\tilde{\lambda}^{k+1} = \left( \left( X_{\mathbb{I}^k}^T X_{\mathbb{I}^k} \right) + \gamma E \right)^{-1} y$$

### Логистическая регрессия

$$\tilde{\lambda}^{k+1} = \left( \left( X_{\mathbb{I}^k}^T X_{\mathbb{I}^k} \right) + \gamma \left( \hat{G}^k \right)^{-1} \right)^{-1} \tilde{y}^k$$

где  $\tilde{y}^k$  и  $\tilde{G}^k$  — переменные, полученные при квадратичном представлении логистической функции потерь в окрестности точки  $z^k$  ряда Тейлора.

# Численное решение задачи

## Итерационный алгоритм Ньютона

Может случиться так, что длина шага очень большая, тогда следует ее сократить

### Проверка

если  $W(\tilde{\lambda}^{k+1}|\gamma, \mu, r) \leq W(\tilde{\lambda}^k|\gamma, \mu, r)$ ,  $\lambda^{k+1} = \tilde{\lambda}^{k+1}$ ;

если  $W(\tilde{\lambda}^{k+1}|\gamma, \mu, r) > W(\tilde{\lambda}^k|\gamma, \mu, r)$ , —длину следует сократить;

Для нахождения подходящей длины шага Ньютона используем метод золотого сечения

### Одномерная оптимизация

$$\tau^{k+1} = \arg \min W[(1 - \tau)\hat{\lambda}^k + \tau\tilde{\lambda}^{k+1}|\gamma, \mu], 0 < \tau < 1,$$

$$\hat{\lambda}^{k+1} = \tau^{k+1}\hat{\lambda}^k + (1 - \tau^{k+1})\tilde{\lambda}^{k+1}$$

# Дифференциальный LOO

## Формула Вудбери

Пусть исходная двойственная задача решена:

### Вектор-решение

$$\hat{\lambda}_{\gamma, \mu}^1 = \underbrace{\left( X_{\mathbb{I}_{\gamma, \mu}}^T X_{\mathbb{I}_{\gamma, \mu}} + \gamma E \right)}_D^{-1} y$$

Решение возмущенной двойственной задачи когда вес  $j$ -го объекта  $r_j = 1 - \rho$

### Вектор-решение с $R^{1j}$

$$\hat{\lambda}_{\gamma, \mu}^{1j} = \left[ \left( X_{\mathbb{I}_{\gamma, \mu}}^T X_{\mathbb{I}_{\gamma, \mu}} + \gamma E \right) - \rho 1^j v_{\mathbb{I}_{\gamma, \mu}, j}^T \right]^{-1} \left[ y - \rho y_j 1^j \right]$$

где  $1^j$  - вектор-столбец:  $1_t = 0, t \in \{1, \dots, N\} \setminus j, 1_j = 1$

# Дифференциальный LOO

## Формула Вудбери

Решение задачи с заглушенным объектом, вес у которого равен  $1 - \rho$ , выражается следующим образом.

Решение задачи с заглушенным объектом

$$\hat{\lambda}_{\gamma,\mu}^{1j} = \hat{\lambda}_{\gamma,\mu} - D_{\mathbb{I}_{\gamma,\mu,j}} \frac{\rho(y_j - v_{\mathbb{I}_{\gamma,\mu,j}}^T \hat{\lambda}_{\gamma,\mu})}{1 - \rho v_{\mathbb{I}_{\gamma,\mu,j}}^T D_{\mathbb{I}_{\gamma,\mu,j}}}$$

В терминах  $z$  решение следующее

Решение в терминах обобщенных признаков

$$\hat{z}_{\gamma,\mu,j}^{1j} = \hat{z}_{\gamma,\mu,j} - \frac{\rho v_{\mathbb{I}_{\gamma,\mu,j}}^T D_{\mathbb{I}_{\gamma,\mu,j}}}{1 - \rho v_{\mathbb{I}_{\gamma,\mu,j}}^T D_{\mathbb{I}_{\gamma,\mu,j}}} (y_{\gamma,\mu,j} - \hat{z}_{\gamma,\mu,j})$$

# Дифференциальный LOO

## Вывод формулы Вудбери

С учетом формулы Вудбери получаем

Отклонение  $\hat{z}_j^{1P}$  от целевой переменной  $y_j$

$$y_j - \hat{z}_{\gamma, \mu}^{1P} = \frac{y_j - \hat{z}_{\gamma, \mu, j}}{1 - \rho \mathbf{v}_{\mathbb{I}_{\gamma, \mu, j}}^T \mathbf{D}_{\mathbb{I}_{\gamma, \mu, j}}}$$

Линейную регрессию можно представить как

Линейная регрессия

$$q(y_j, \hat{z}_{\gamma, \mu}^{1P}) = \left( \frac{y_j - \hat{z}_{\gamma, \mu, j}}{1 - \rho \mathbf{v}_{\mathbb{I}_{\gamma, \mu, j}}^T \mathbf{D}_{\mathbb{I}_{\gamma, \mu, j}}} \right)^2$$

# Дифференциальный LOO

## Критерий

Наш критерий *DiffLOO* построен как среднее скоростей роста потерь для каждого отдельного объекта

## Общий критерий дифференциального LOO

$$\text{DiffLOO}(\gamma, \mu) = \frac{1}{N} \sum_{j=1}^N \frac{\partial}{\partial p} q(y_j, \hat{z}_{\gamma, \mu, j}^{1p})$$

Для линейной регрессии получаем

## Критерий дифференциального LOO

$$\text{DifLOO}(\gamma, \mu) = \frac{1}{N} \sum_{j=1}^N \left( y_j - \hat{z}_{\gamma, \mu, j} \right)^2 \left( \mathbf{v}_{\mathbb{I}_{\gamma, \mu, j}}^T \mathbf{D}_{\mathbb{I}_{\gamma, \mu, j}} \right)$$

Для логистической регрессии получаем

Полученный критерий дифференциального LOO

$$\text{DifLOO}(\gamma, \mu) = \frac{1}{N} \sum_{j=1}^N \left( \tilde{y}_j - \mathbf{v}_{\mathbb{I}_{\gamma, \mu, j}}^T \hat{\boldsymbol{\lambda}}_{\gamma, \mu} \right)^2 \left( \mathbf{v}_{\mathbb{I}_{\gamma, \mu, j}}^T \tilde{\mathbf{D}}_{\mathbb{I}_{\gamma, \mu, j}} \right)$$



В качестве набора регрессоров применяются  $n = 650$  временных рядов месячных доходностей биржевых ценных бумаг на Нью-Йорской фондовой бирже, каждый длиной чуть больше 20 лет

$$X_t = (x_{t,i}, i = 1, \dots, N) \in \mathbb{R}^n, n = 650$$

Наблюдаемый сигнал, состоящий из  $N = 251$ .

$$y_t, t = 1, \dots, N = 251$$

Временной ряд доходностей инвестиционного портфеля, построенного как вложение капитала в равных долях в  $n^* = 13$  неизвестных ценных бумаг. В качестве модели используется регрессионная модель Шарпа

$$y_t \cong \sum_{i=1}^n a_i x_{t,i}$$

где  $a_i$  - это доли вложения капитала.

# Эксперименты

Результаты экспериментов для линейной регрессии

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

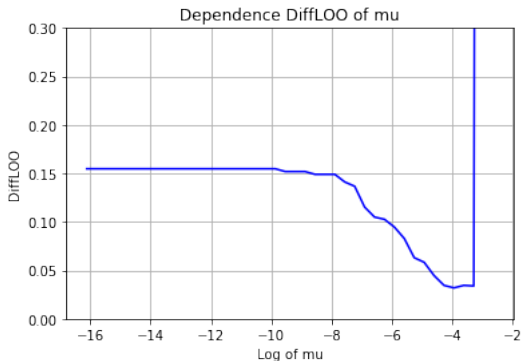


Figure: Зависимость дифференциального LOO от параметра  $\mu$

В качестве выборки взята искусственно сгенерированная выборка, состоящая из 400 признаков и 200 объектов, в которой только 2 признака образуют оптимальное подмножество признаков, а остальные признаки — шумовые.

$$X_t = (x_{t,i=1,\dots,N}) \in \mathbb{R}^n, \quad n = 400$$

Целевая переменная, состоящая из  $N = 200$ .

$$y_t \in \{1, -1\}, \quad t = 1, \dots, N = 200$$

Оптимальные признаки заданы из нормального распределения с  $\mathbb{E} = 0$  и  $\mathbb{D} = 1$ .

# Эксперименты

## Результаты экспериментов для логистической регрессии

Найденное подмножество  $\hat{\mathbb{I}}_{\gamma, \mu}$  в точности совпадает с оптимальным подмножеством.

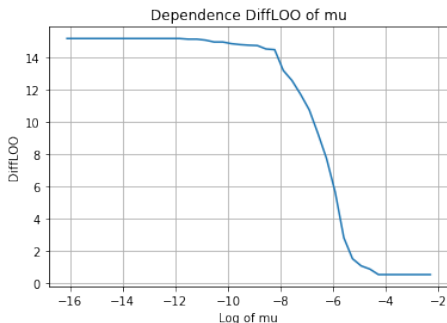


Figure: Зависимость дифференциального LOO от параметра  $\mu$

## Публикация результатов

- 1 Ангуло Б., Морозов А.О., Моттль В.В. Метод дифференциальной кросс-валидации для выбора уровня сложности обобщенных линейных моделей зависимостей. Математические методы распознавания образов: Тезисы докладов 19-ой Всероссийской конференции с международным участием, г. Москва, 2019 г.
- 2 B. Angulo, A. Morozov, V. Mottl, A. Tatarchuk, O. Krasotkina. Differential Leave-One-Out Cross-Validation for Hyperparameter Optimization in Generalized Linear Dependence Models. 25-th International Conference on Pattern Recognition ICPR- 2021. Milan, Italy, January 10-15, 2021. (To appear)

- Сформулирован и подобран алгоритм численного решения для выбора активных признаков;
- Разработан метод дифференциальной кросс-валидации для выбора оптимальных гиперпараметров модели;
- Получен критерий качества решения;
- Продемонстрирована правильность и корректность работы метода на примере приклавных задач.

## планы

- Изменить метод на общий случай когда каждый объект будет учитываться с разным весом
- Изменить метод на общий случай когда вместо скалярных произведений можно было бы взять ядро
- Усовершенствовать имеющее численное решение на эти случаи