

# **Неявная кросс-валидация для выбора подмножества информативных признаков в задаче обучения распознаванию образов по методу опорных векторов**

**Черноусова Елена Олеговна**

**Левдик Павел Владимирович**

Московский физико-технический институт

**Моттль Вадим Вячеславович**

Вычислительный центр РАН

**Уиндридж, Дэвид**

Университет Суррей, Великобритания

# **Типовая задача восстановления закономерностей в множествах объектов реального мира**

# **Типовая задача восстановления закономерностей в множествах объектов реального мира**

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

## **Типовая задача восстановления закономерностей в множествах объектов реального мира**

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

## Типовая задача восстановления закономерностей в множествах объектов реального мира

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ ,  
то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$

## Типовая задача восстановления закономерностей в множествах объектов реального мира

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ ,  
то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$

Желание наблюдателя:

Иметь инструмент оценивания скрытой характеристики для реальных объектов

$\hat{y}(x): \mathbb{X} \rightarrow \mathbb{Y}; \quad \hat{y}(x) \neq y(x) - \text{ошибка.}$

## Типовая задача восстановления закономерностей в множествах объектов реального мира

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ ,  
то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$

Желание наблюдателя:

Иметь инструмент оценивания скрытой характеристики для реальных объектов  
 $\hat{y}(x): \mathbb{X} \rightarrow \mathbb{Y}; \quad \hat{y}(x) \neq y(x) - \text{ошибка.}$

Обучение по прецедентам:

Подмножество наблюдаемых объектов, для которых измерены значения обеих характеристик  $\Omega^* \subset \Omega: \left\{ (x(\omega_j), y(\omega_j)) = (x_j, y_j), j = 1, \dots, N \right\}$ .

## Типовая задача восстановления закономерностей в множествах объектов реального мира

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ ,  
то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$

Желание наблюдателя:

Иметь инструмент оценивания скрытой характеристики для реальных объектов  
 $\hat{y}(x): \mathbb{X} \rightarrow \mathbb{Y}$ ;  $\hat{y}(x) \neq y(x)$  – ошибка.

Обучение по прецедентам:

Подмножество наблюдаемых объектов, для которых измерены значения обеих характеристик  $\Omega^* \subset \Omega: \left\{ (x(\omega_j), y(\omega_j)) = (x_j, y_j), j = 1, \dots, N \right\}$ .

Задача: Продолжить функцию на все множество  $\Omega$ , так чтобы можно было в дальнейшем оценивать значение рассматриваемой характеристики  $\hat{y}(\omega) = \hat{y}(x(\omega))$  для новых объектов  $\omega \in \Omega \setminus \Omega^*$ .



## Типовая задача восстановления закономерностей в множествах объектов реального мира

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ ,  
то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$

Желание наблюдателя:

Иметь инструмент оценивания скрытой характеристики для реальных объектов  
 $\hat{y}(x): \mathbb{X} \rightarrow \mathbb{Y}$ ;  $\hat{y}(x) \neq y(x)$  – ошибка.

Обучение по прецедентам:

Подмножество наблюдаемых объектов, для которых измерены значения обеих характеристик  $\Omega^* \subset \Omega: \left\{ (x(\omega_j), y(\omega_j)) = (x_j, y_j), j = 1, \dots, N \right\}$ .

Простейшие случаи:

Задача распознавания образов

$\mathbb{Y} = \{y_1, \dots, y_m\}$  – конечное неупорядоченное множество; в частности  $\mathbb{Y} = \{-1, 1\}$ .

Задача восстановления числовой зависимости

$\mathbb{Y} = \mathbb{R}$  – множество действительных чисел.

# **Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности**

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0$ ,  $\int_{\mathbb{Y}} \int_{\mathbb{X}} f^*(x, y) dx dy = 1$ .

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0, \int \int_{\mathbb{Y} \times \mathbb{X}} f^*(x, y) dx dy = 1$ .

Тем не менее, наблюдатель вынужден выбрать оценку  $\hat{y}(x, \mathbf{a}), \mathbf{a} \in \mathbb{R}^n$ , предполагая некоторую функцию потерь  $Loss(y, \hat{y}(x, \mathbf{a})) = q(y, x, \mathbf{a})$ .

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0$ ,  $\int_{\mathbb{Y}} \int_{\mathbb{X}} f^*(x, y) dx dy = 1$ .

Тем не менее, наблюдатель вынужден выбрать оценку  $\hat{y}(x, \mathbf{a})$ ,  $\mathbf{a} \in \mathbb{R}^n$ , предполагая некоторую функцию потерь  $Loss(y, \hat{y}(x, \mathbf{a})) = q(y, x, \mathbf{a})$ .

Средний риск ошибки: 
$$r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx$$

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0$ ,  $\int_{\mathbb{Y}} \int_{\mathbb{X}} f^*(x, y) dx dy = 1$ .

Тем не менее, наблюдатель вынужден выбрать оценку  $\hat{y}(x, \mathbf{a})$ ,  $\mathbf{a} \in \mathbb{R}^n$ , предполагая некоторую функцию потерь  $Loss(y, \hat{y}(x, \mathbf{a})) = q(y, x, \mathbf{a})$ .

Средний риск ошибки:  $r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx$

Естественный способ выбора решающего правила – минимизация среднего риска

$$r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx \rightarrow \min(\mathbf{a})$$

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0$ ,  $\int_{\mathbb{Y}} \int_{\mathbb{X}} f^*(x, y) dx dy = 1$ .

Тем не менее, наблюдатель вынужден выбрать оценку  $\hat{y}(x, \mathbf{a})$ ,  $\mathbf{a} \in \mathbb{R}^n$ , предполагая некоторую функцию потерь  $Loss(y, \hat{y}(x, \mathbf{a})) = q(y, x, \mathbf{a})$ .

Средний риск ошибки:  $r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx$

Естественный способ выбора решающего правила – минимизация среднего риска

$$r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx \rightarrow \min(\mathbf{a})$$

Проклятие неопределенности! Плотность  $f^*(x, y)$  неизвестна наблюдателю!

## Вероятностная интерпретация задачи восстановления зависимости по обучающей совокупности

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

С каждым объектом связана пара характеристик  $\begin{cases} x(\omega) \in \mathbb{X} - \text{наблюдаемая} \\ y(\omega) \in \mathbb{Y} - \text{скрытая} \end{cases}$

Природа случайным образом многократно и независимо выбирает объект  $\omega \in \Omega$ , то есть пару  $(x, y) \in \mathbb{X} \times \mathbb{Y}$ .

Плотность распределения наблюдателю неизвестна  $f^*(x, y) \geq 0$ ,  $\int_{\mathbb{Y}} \int_{\mathbb{X}} f^*(x, y) dx dy = 1$ .

Тем не менее, наблюдатель вынужден выбрать оценку  $\hat{y}(x, \mathbf{a})$ ,  $\mathbf{a} \in \mathbb{R}^n$ , предполагая некоторую функцию потерь  $Loss(y, \hat{y}(x, \mathbf{a})) = q(y, x, \mathbf{a})$ .

Средний риск ошибки:  $r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx$

Естественный способ выбора решающего правила – минимизация среднего риска

$$r_{av}(\mathbf{a}) = r_{av}[\hat{y}(\cdot, \mathbf{a})] = \int_{\mathbb{X}} \int_{\mathbb{Y}} q(x, y, \mathbf{a}) f^*(x, y) dy dx \rightarrow \min(\mathbf{a})$$

Проклятие неопределенности! Плотность  $f^*(x, y)$  неизвестна наблюдателю!

Стандартный выход – минимизация регуляризованного эмпирического риска.



# **Выбор решающего правила по критерию минимума регуляризованного эмпирического риска**

## Выбор решающего правила по критерию минимума регуляризованного эмпирического риска

Два «островка определенности», доступные наблюдателю.

- 1) Обучающая совокупность  $(\mathbf{x}, \mathbf{y}) = \{(x_j, y_j), j = 1, \dots, N\}$

Эмпирический риск выбора решающего правила  $Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) = \left(\frac{1}{N}\right) \sum_{j=1}^N q(y_j, x_j, \mathbf{a}) -$

несмещенная оценка сред среднего риска по обучающей совокупности

Желание наблюдателя – минимизировать хотя бы эмпирический риск

$Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) \rightarrow \min(\mathbf{a})$ .

- 2) Априорные предпочтения относительно параметра решающего правила  
 $V(\mathbf{a}, C) \rightarrow \min(\mathbf{a})$ ,  $C$  – скалярный либо векторный «структурный» параметр.

## Выбор решающего правила по критерию минимума регуляризованного эмпирического риска

Два «островка определенности», доступные наблюдателю.

- 1) Обучающая совокупность  $(\mathbf{x}, \mathbf{y}) = \{(x_j, y_j), j = 1, \dots, N\}$

Эмпирический риск выбора решающего правила  $Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) = \left(\frac{1}{N}\right) \sum_{j=1}^N q(y_j, x_j, \mathbf{a})$  –

несмещенная оценка сред среднего риска по обучающей совокупности

Желание наблюдателя – минимизировать хотя бы эмпирический риск

$Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) \rightarrow \min(\mathbf{a})$ .

- 2) Априорные предпочтения относительно параметра решающего правила

$V(\mathbf{a}, C) \rightarrow \min(\mathbf{a})$ ,  $C$  – скалярный либо векторный «структурный» параметр

Компромисс – минимизация регуляризованного эмпирического риска

Критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}) = \arg \min_{\mathbf{a}} \{V(\mathbf{a}, C) + Q(\mathbf{y}, \mathbf{x}, \mathbf{a})\}$

## Выбор решающего правила по критерию минимума регуляризованного эмпирического риска

Два «островка определенности», доступные наблюдателю.

- 1) Обучающая совокупность  $(\mathbf{x}, \mathbf{y}) = \{(x_j, y_j), j = 1, \dots, N\}$

Эмпирический риск выбора решающего правила  $Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) = \left(\frac{1}{N}\right) \sum_{j=1}^N q(y_j, x_j, \mathbf{a})$  –

несмещенная оценка сред среднего риска по обучающей совокупности

Желание наблюдателя – минимизировать хотя бы эмпирический риск

$Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) \rightarrow \min(\mathbf{a})$ .

- 2) Априорные предпочтения относительно параметра решающего правила

$V(\mathbf{a}, C) \rightarrow \min(\mathbf{a})$ ,  $C$  – скалярный либо векторный «структурный» параметр

Компромисс – минимизация регуляризованного эмпирического риска

Критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}) = \arg \min_{\mathbf{a}} \{V(\mathbf{a}, C) + Q(\mathbf{y}, \mathbf{x}, \mathbf{a})\}$

Как выбрать структурный параметр  $C$ ?

## Выбор решающего правила по критерию минимума регуляризованного эмпирического риска

Два «островка определенности», доступные наблюдателю.

- 1) Обучающая совокупность  $(\mathbf{x}, \mathbf{y}) = \{(x_j, y_j), j = 1, \dots, N\}$

Эмпирический риск выбора решающего правила  $Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) = \left(\frac{1}{N}\right) \sum_{j=1}^N q(y_j, x_j, \mathbf{a})$  –

несмещенная оценка сред среднего риска по обучающей совокупности

Желание наблюдателя – минимизировать хотя бы эмпирический риск

$Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) \rightarrow \min(\mathbf{a})$ .

- 2) Априорные предпочтения относительно параметра решающего правила  
 $V(\mathbf{a}, C) \rightarrow \min(\mathbf{a})$ ,  $C$  – скалярный либо векторный «структурный» параметр

Компромисс – минимизация регуляризованного эмпирического риска

Критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}) = \arg \min_{\mathbf{a}} \{V(\mathbf{a}, C) + Q(\mathbf{y}, \mathbf{x}, \mathbf{a})\}$

Как выбрать структурный параметр  $C$ ?

Выбрать его по критерию минимума среднего риска невозможно

$$\int \int_{\mathbb{Y} \times \mathbb{X}} Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) F^*(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \rightarrow \min(C), \quad F^*(\mathbf{x}, \mathbf{y}) = \prod_{j=1}^N f^*(x_j, y_j)$$

*неизвестное совместное распределение случайной обучающей выборки*

## Выбор решающего правила по критерию минимума регуляризованного эмпирического риска

Два «островка определенности», доступные наблюдателю.

- 1) Обучающая совокупность  $(\mathbf{x}, \mathbf{y}) = \{(x_j, y_j), j = 1, \dots, N\}$

Эмпирический риск выбора решающего правила  $Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) = \left(\frac{1}{N}\right) \sum_{j=1}^N q(y_j, x_j, \mathbf{a})$  –

несмещенная оценка сред среднего риска по обучающей совокупности

Желание наблюдателя – минимизировать хотя бы эмпирический риск

$Q(\mathbf{y}, \mathbf{x}, \mathbf{a}) \rightarrow \min(\mathbf{a})$ .

- 2) Априорные предпочтения относительно параметра решающего правила

$V(\mathbf{a}, C) \rightarrow \min(\mathbf{a})$ ,  $C$  – скалярный либо векторный «структурный» параметр

Компромисс – минимизация регуляризованного эмпирического риска

Критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}) = \arg \min_{\mathbf{a}} \{V(\mathbf{a}, C) + Q(\mathbf{y}, \mathbf{x}, \mathbf{a})\}$

Как выбрать структурный параметр  $C$ ?

Выбрать его по критерию минимума среднего риска невозможно

$$\int \int_{\mathbb{Y} \times \mathbb{X}} Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) F^*(\mathbf{x}, \mathbf{y}) dx dy \rightarrow \min(C), \quad F^*(\mathbf{x}, \mathbf{y}) = \prod_{j=1}^N f^*(x_j, y_j)$$

*неизвестное совместное распределение случайной обучающей выборки*

В этом заключается одна из основных проблем теории обучения.

Именно эта проблема составляет предмет данного доклада.

## Предпосылки исследования: Методы выбора структурного параметра обучения по единственной обучающей совокупности

- **Минимизация структурного риска** (В.Н. Вапник, А.Я. Червоненкис, 1974). *Беспереборный метод*. Основан на понятии VC-dimension и применим только к критериям обучения распознаванию образов.
- **Информационный критерий Акаике** (Hirotugu Akaike, 1974). *Беспереборный метод*. Основан на априорной упорядоченности элементов вектора параметров  $\mathbf{a} = (a_1, \dots, a_\lambda, a_{\lambda+1}=0, \dots, a_n=0)$  и применим только к квадратичным критериям обучения.
- **Скользкий контроль** (М.Н. Вайнгцвайг, 1971). Универсальный *переборный метод*. Последовательное выделение одного объекта обучающей совокупности для контроля ошибки при обучении по остальным объектам. *Высокая вычислительная сложность*.
- **Кросс-валидация – Jackknife** (Bradley Efron, 1982). Более общий *переборный метод*. Основан на многократном делении выборки на обучающую и контрольную части. *Вычислительная сложность может быть еще больше*.

### Задачи данного исследования

- Показать, что идея информационного критерия Акаике для выбора структурного параметра основана на принципе неявной кросс-валидации.
- Разработать беспереборный критерий кросс-валидации для квадратичной задачи оценивания линейной регрессии, в котором классический критерий Акаике являлся бы частным случаем.
- Разработать беспереборный критерий кросс-валидации для задачи обучения распознаванию образов по методу опорных векторов – Support Vector Machine, SVM.

**Принцип неявной кросс-валидации:  
Первая эвристики наблюдателя**



## Принцип неявной кросс-валидации: Первая эвристики наблюдателя

Еще раз неизвестная совместная плотность  
распределения наблюдаемой и скрытой  
характеристик объекта

$$\underbrace{F^*(\mathbf{x}, \mathbf{y})}_{\text{выборка в целом}} = \prod_{j=1}^N \underbrace{f^*(x_j, y_j)}_{\text{отдельные случайные объекты}}$$

## Принцип неявной кросс-валидации: Первая эвристики наблюдателя

Еще раз неизвестная совместная плотность распределения наблюдаемой и скрытой характеристик объекта

$$\underbrace{F^*(\mathbf{x}, \mathbf{y})}_{\text{выборка в целом}} = \prod_{j=1}^N \underbrace{f^*(x_j, y_j)}_{\text{отдельные случайные объекты}}$$

Представим ее в виде произведения маргинальной плотности распределения наблюдаемой характеристики и условной плотности распределения ненаблюдаемой характеристики:

$$f^*(x, y) = \varphi^*(y | x)g^*(x)$$

$$F^*(\mathbf{x}, \mathbf{y}) = \underbrace{\prod_{j=1}^N \varphi^*(y_j | x_j)}_{\Phi^*(\mathbf{y}|\mathbf{x})} \underbrace{\prod_{j=1}^N g^*(x_j)}_{G^*(\mathbf{x})} = \boxed{\Phi^*(\mathbf{y} | \mathbf{x})} G^*(\mathbf{x})$$

## Принцип неявной кросс-валидации: Первая эвристики наблюдателя

Еще раз неизвестная совместная плотность распределения наблюдаемой и скрытой характеристик объекта

$$\underbrace{F^*(\mathbf{x}, \mathbf{y})}_{\text{выборка в целом}} = \prod_{j=1}^N \underbrace{f^*(x_j, y_j)}_{\text{отдельные случайные объекты}}$$

Представим ее в виде произведения маргинальной плотности распределения наблюдаемой характеристики и условной плотности распределения ненаблюдаемой характеристики:

$$f^*(x, y) = \varphi^*(y | x) g^*(x)$$

$$F^*(\mathbf{x}, \mathbf{y}) = \underbrace{\prod_{j=1}^N \varphi^*(y_j | x_j)}_{\Phi^*(\mathbf{y} | \mathbf{x})} \underbrace{\prod_{j=1}^N g^*(x_j)}_{G^*(\mathbf{x})} = \boxed{\Phi^*(\mathbf{y} | \mathbf{x})} G^*(\mathbf{x})$$

### Первая эвристика наблюдателя:

#### Предположение об условном распределении скрытой характеристики

Неизвестная плотность распределения представляется в виде **неизвестной смеси известных распределений**

$$\Phi^*(\mathbf{y} | \mathbf{x}) = \int_{\mathbb{R}^n} \underbrace{\Phi(\mathbf{y} | \mathbf{x}, \mathbf{a})}_{\substack{\text{известное} \\ \text{параметрическое} \\ \text{семейство}}} \underbrace{\Psi^*(\mathbf{a})}_{\substack{\text{неизвестная} \\ \text{смесь}}} d\mathbf{a}$$

## Принцип неявной кросс-валидации: Первая эвристики наблюдателя

Еще раз неизвестная совместная плотность распределения наблюдаемой и скрытой характеристик объекта

$$\underbrace{F^*(\mathbf{x}, \mathbf{y})}_{\text{выборка в целом}} = \prod_{j=1}^N \underbrace{f^*(x_j, y_j)}_{\text{отдельные случайные объекты}}$$

Представим ее в виде произведения маргинальной плотности распределения наблюдаемой характеристики и условной плотности распределения ненаблюдаемой характеристики:

$$f^*(x, y) = \varphi^*(y | x) g^*(x)$$

$$F^*(\mathbf{x}, \mathbf{y}) = \underbrace{\prod_{j=1}^N \varphi^*(y_j | x_j)}_{\Phi^*(\mathbf{y} | \mathbf{x})} \underbrace{\prod_{j=1}^N g^*(x_j)}_{G^*(\mathbf{x})} = \boxed{\Phi^*(\mathbf{y} | \mathbf{x})} G^*(\mathbf{x})$$

### Первая эвристика наблюдателя:

#### Предположение об условном распределении скрытой характеристики

Неизвестная плотность распределения представляется в виде **неизвестной смеси известных распределений**

$$\Phi^*(\mathbf{y} | \mathbf{x}) = \int_{\mathbb{R}^n} \underbrace{\Phi(\mathbf{y} | \mathbf{x}, \mathbf{a})}_{\substack{\text{известное} \\ \text{параметрическое} \\ \text{семейство}}} \underbrace{\Psi^*(\mathbf{a})}_{\substack{\text{неизвестная} \\ \text{смесь}}} d\mathbf{a}$$

**Принцип незлонамеренности природы.** Наблюдатель правильно «угадал» функцию потерь и класс решающих правил. Природа чаще генерирует пары с низким значением штрафа.

$$\Phi(\mathbf{y} | \mathbf{x}, \mathbf{a}) \propto \exp \left\{ - \underbrace{Q(\mathbf{y}, \mathbf{x}, \mathbf{a})}_{\substack{\text{принятая наблюдателем} \\ \text{функция потерь}}} \right\}$$

↑  
коэффициент не зависит от  $(\mathbf{x}, \mathbf{a})$

**Принцип неявной кросс-валидации:  
Мысленный эксперимент наблюдателя**

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math> <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math> Две мысленные совокупности: <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
--	---

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math> <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>  Две мысленные совокупности:  <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности:  <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math> <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>  Две мысленные совокупности:  <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности:  <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\int \int \int Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\int \int \int Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$



## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math> <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>          Две мысленные совокупности:  <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности:  <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$

Согласно первой эвристике наблюдателя  $Q(\mathbf{x}, \mathbf{y}, \mathbf{a}) = -\ln \Phi(\mathbf{y} | \mathbf{x}, \mathbf{a}) + const$ .

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math> <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>          Две мысленные совокупности:  <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности:  <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$

Согласно первой эвристике наблюдателя  $Q(\mathbf{x}, \mathbf{y}, \mathbf{a}) = -\ln \Phi(\mathbf{y} | \mathbf{x}, \mathbf{a}) + const$ .

Идея скрытой кросс-валидации  $\int [\ln \Phi(\mathbf{y} | \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))] \Phi(\mathbf{y} | \mathbf{x}, \mathbf{a}) d\mathbf{y} \rightarrow \max$

есть максимизация информации Кульбака о распределении  $\Phi(\mathbf{y} | \mathbf{x}, \mathbf{a})$ , содержащейся в его оценке по другой выборке  $\Phi(\mathbf{y} | \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))$ .

Поэтому критерии неявной кросс-валидации уместно называть информационными и рассматривать их как обобщение классической идеи Хиротугу Акаике.

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math>, <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>. Две мысленные совокупности: <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности: <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$

В реальности у наблюдателя имеется единственная обучающая выборка  $(\mathbf{x}, \mathbf{y})$ .

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math>, <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>. Две мысленные совокупности: <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности: <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$

В реальности у наблюдателя имеется единственная обучающая выборка  $(\mathbf{x}, \mathbf{y})$ .

Он может подставить в функцию потерь лишь оценку, вычисленную по той же выборке  $Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}))$ .

## Принцип неявной кросс-валидации: Мысленный эксперимент наблюдателя

<p>Пусть природа разыграла конкретное значение параметра согласно <math>\Psi^*(\mathbf{a})</math>, а также выборку наблюдаемых компонент согласно <math>G^*(\mathbf{x})</math></p>	<p>Затем, дважды применив условное распределение <math>\Phi(\mathbf{y}   \mathbf{x}, \mathbf{a})</math>, получила две реализации <math>\mathbf{y}=(y_1, \dots, y_N) \in \mathbb{R}^N</math>, <math>\tilde{\mathbf{y}}=(\tilde{y}_1, \dots, \tilde{y}_N) \in \mathbb{R}^N</math>. Две мысленные совокупности: <math>(\mathbf{x}, \mathbf{y})</math> – контрольная, <math>(\mathbf{x}, \tilde{\mathbf{y}})</math> – обучающая</p>
<p>Мысленная перекрестная проверка на генеральной совокупности</p>	<p>Если бы наблюдатель знал реализации <math>\mathbf{y}</math> и <math>\tilde{\mathbf{y}}</math>, то мог бы для всякого значения <math>C</math> вычислить оценку <math>\mathbf{a}</math> по обучающей совокупности и подставить в функцию потерь для контрольной совокупности: <math>\hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}) \rightarrow Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x}))</math> – потери на мысленном контроле</p>
<p>Идея скрытой кросс-валидации: минимум математического ожидания потерь</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\mathbf{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \min(C)$

В реальности у наблюдателя имеется единственная обучающая выборка  $(\mathbf{x}, \mathbf{y})$ .

Он может подставить в функцию потерь лишь оценку, вычисленную по той же выборке  $Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x}))$ .

Насколько испортится критерий, подлежащий максимизации?

Каким должен быть штраф за использование оценки параметра, вычисленного по той же выборке?

**Штраф за подстановку в функцию потерь оценки параметра,  
вычисленной по той же выборке**

## Штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке

Еще раз идея скрытой кросс-  
валидации:  
минимум двойного  
математического ожидания

$$\begin{aligned} & \int \int \int Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} = \\ & \int \int \int Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}} | \mathbf{x}, \mathbf{a}) \Phi(\tilde{\mathbf{y}} | \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow \\ & \rightarrow \min(C) \end{aligned}$$

## Штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке

<p>Еще раз идея скрытой кросс-валидации: минимум двойного математического ожидания</p>	$\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) F^*(\mathbf{y}, \tilde{\mathbf{y}}, \mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} =$ $\iiint Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) \left\{ \int \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Phi(\tilde{\mathbf{y}}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{x} \rightarrow$ $\rightarrow \min(C)$
<p>Добавление нуля к подынтегральной функции</p>	$Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{\mathbf{y}}, \mathbf{x})) + \left\{ Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) - Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) \right\}$



## Штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке

Еще раз идея скрытой кросс-валидации: минимум двойного математического ожидания	$\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) F^*(y, \tilde{y}, \mathbf{x}) d\tilde{y} dy d\mathbf{x} =$ $\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) \left\{ \int \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{y} dy d\mathbf{x} \rightarrow$ $\rightarrow \min(C)$
Добавление нуля к подинтегральной функции	$Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) + \left\{ Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) - Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \right\}$
Эквивалентная запись критерия минимизации по значению структурного параметра $C$ :	
$C^* = \arg \min_C \left\{ \iint \left[ \int Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \Phi(y   \mathbf{x}, \mathbf{a}) d\mathbf{y} + \right. \right.$ $\left. \iint \left( Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) - Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \right) \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Phi(y   \mathbf{x}, \mathbf{a}) d\tilde{y} dy \right] \Psi^*(\mathbf{a}) G^*(\mathbf{x}) d\mathbf{a} d\mathbf{x} \right\}$	
<hr style="width: 60%; margin-left: 0;"/> функционал от функции потерь $Q(\bullet, \mathbf{x}, \mathbf{a})$ и критерия обучения $\hat{\mathbf{a}}_C(\bullet, \mathbf{x})$	

## Штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке

Еще раз идея скрытой кросс-валидации: минимум двойного математического ожидания	$\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) F^*(y, \tilde{y}, \mathbf{x}) d\tilde{y} dy d\mathbf{x} =$ $\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) \left\{ \int \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{y} dy d\mathbf{x} \rightarrow$ $\rightarrow \min(C)$
Добавление нуля к подинтегральной функции	$Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) + \left\{ Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) - Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \right\}$
Эквивалентная запись критерия минимизации по значению структурного параметра $C$ :	
$C^* = \arg \min_C \left\{ \iint \left[ \int Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \Phi(y   \mathbf{x}, \mathbf{a}) d\mathbf{y} + \right. \right.$ $\left. \iint \left( Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) - Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \right) \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Phi(y   \mathbf{x}, \mathbf{a}) d\tilde{y} dy \right] \Psi^*(\mathbf{a}) G^*(\mathbf{x}) d\mathbf{a} d\mathbf{x} \right\}$	
<hr style="width: 60%; margin-left: 0;"/> <i>функционал от функции потерь <math>Q(\cdot, \mathbf{x}, \mathbf{a})</math> и критерия обучения <math>\hat{\mathbf{a}}_C(\cdot, \mathbf{x})</math></i>	

Для многих типичных функций потерь $Q(y, \mathbf{x}, \mathbf{a})$ и критериев обучения $\hat{\mathbf{a}}_C(y, \mathbf{x})$ , адекватных широкому классу практических задач, функционал во втором слагаемом не зависит от $\mathbf{a}$ .	$\Delta(C, \mathbf{x})$
---	-------------------------

## Штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке

Еще раз идея скрытой кросс-валидации: максимум двойного математического ожидания	$\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) F^*(y, \tilde{y}, \mathbf{x}) d\tilde{y} dy d\mathbf{x} =$ $\iiint Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) \left\{ \int \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Phi(\tilde{y}   \mathbf{x}, \mathbf{a}) \Psi^*(\mathbf{a}) d\mathbf{a} \right\} G^*(\mathbf{x}) d\tilde{y} dy d\mathbf{x} \rightarrow$ $\rightarrow \min(C)$
Добавление нуля к подынтегральной функции	$Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(\tilde{y}, \mathbf{x})) + \left\{ Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) - Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \right\}$
Эквивалентная запись критерия максимизации по значению структурного параметра $C$ :	
$C^* = \arg \min_C \left\{ \iint \left[ \int Q(y, \mathbf{x}, \hat{\mathbf{a}}_C(y, \mathbf{x})) \Phi(y   \mathbf{x}, \mathbf{a}) d\mathbf{y} + \right. \right.$ $\left. \left. + \Delta(C, \mathbf{x}) \right] \Psi^*(\mathbf{a}) G^*(\mathbf{x}) d\mathbf{a} d\mathbf{x} \right\}$	

Для многих типичных функций потерь $Q(y, \mathbf{x}, \mathbf{a})$ и критериев обучения $\hat{\mathbf{a}}_C(y, \mathbf{x})$ , адекватных широкому классу практических задач, функционал во втором слагаемом не зависит от $\mathbf{a}$ .	$\Delta(C, \mathbf{x})$
---	-------------------------

## **Вторая эвристика наблюдателя**

## Вторая эвристика наблюдателя

**Теорема 1.** Критерий скрытой кросс-валидации имеет вид:

$$C^* = \arg \min_C \left\{ \iint [Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) + \Delta(C, \mathbf{x})] F^*(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} \right\}$$

## Вторая эвристика наблюдателя

**Теорема 1.** Критерий скрытой кросс-валидации имеет вид:

$$C^* = \arg \min_C \left\{ \iint [Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) + \Delta(C, \mathbf{x})] F^*(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} \right\}$$

Однако в таком виде критерий по-прежнему неприменим, так как распределение генеральной совокупности  $F^*(\mathbf{x}, \mathbf{y})$  неизвестно.

Идея заключается в замене математического ожидания его несмещенной оценкой по единственной доступной наблюдателю выборке  $\mathbf{y} \in \mathbb{R}^N$ :

$$\hat{C}(\mathbf{y}, \mathbf{x}) = \arg \min_C \left\{ \underbrace{Q(\mathbf{y}, \mathbf{x}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{x})) + \Delta(C, \mathbf{x})}_{\text{эмпирический риск}} \right\}$$

*структурный риск*

## Вторая эвристика наблюдателя

**Теорема 1.** Критерий скрытой кросс-валидации имеет вид:

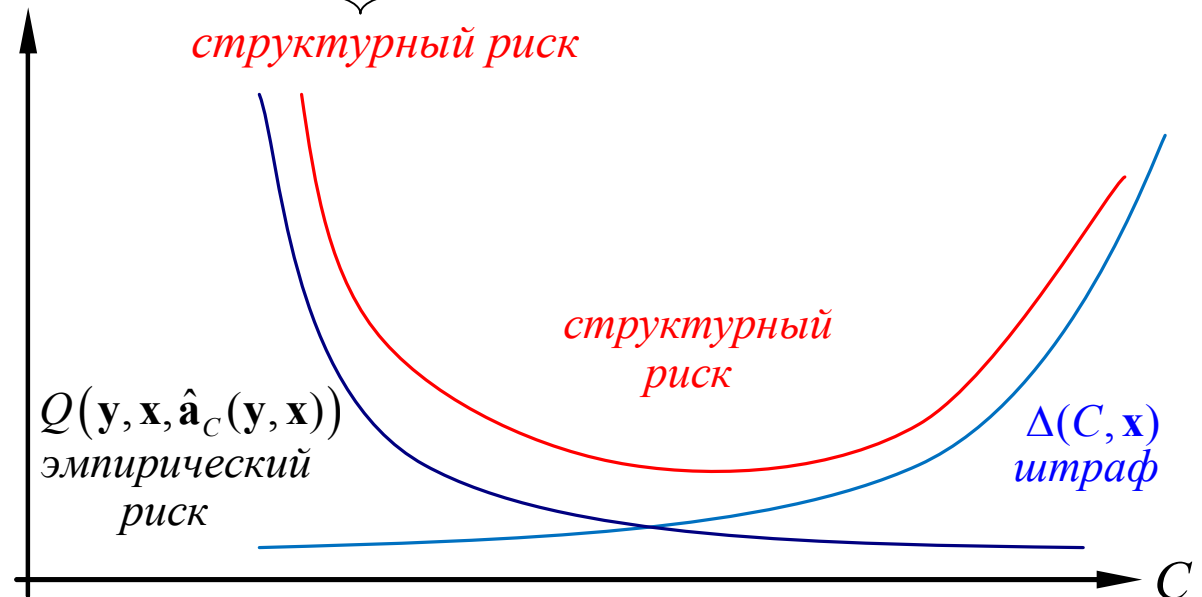
$$C^* = \arg \min_C \left\{ \iint [Q(y, x, \hat{a}_C(y, x)) + \Delta(C, x)] F^*(x, y) dy dx \right\}$$

Однако в таком виде критерий по-прежнему неприменим, так как распределение генеральной совокупности  $F^*(x, y)$  неизвестно.

Идея заключается в замене математического ожидания его несмещенной оценкой по единственной доступной наблюдателю выборке  $y \in \mathbb{R}^N$ :

$$\hat{C}(y, x) = \arg \min_C \left\{ \underbrace{Q(y, x, \hat{a}_C(y, x))}_{\text{эмпирический риск}} + \Delta(C, x) \right\}$$

Графическое отображение зависимости суммарного структурного риска от параметра регуляризации в предположении, что свобода выбора решающего правила увеличивается при увеличении параметра  $C$ .



# Оценивание числовой зависимости. Линейно-квадратичная модель



# **Оценивание числовой зависимости. Линейно-квадратичная модель**

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$ .

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$ .

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

Искомое линейное решающее правило:  $\hat{y}(\mathbf{x}) = \mathbf{a}^T \mathbf{x} : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ .

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

Искомое линейное решающее правило:  $\hat{y}(\mathbf{x}) = \mathbf{a}^T \mathbf{x}: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ .

Квадратичная функция потерь:  $Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) = \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 = (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a})$

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

Искомое линейное решающее правило:  $\hat{y}(\mathbf{x}) = \mathbf{a}^T \mathbf{x} : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ .

Квадратичная функция потерь:  $Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) = \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 = (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a})$

Напомним общий критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a}} \left\{ \underbrace{V(\mathbf{a}, C)}_{\text{регуляризирующая функция}} + Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) \right\}$

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

Искомое линейное решающее правило:  $\hat{y}(\mathbf{x}) = \mathbf{a}^T \mathbf{x} : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ .

Квадратичная функция потерь:  $Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) = \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 = (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a})$

Напомним общий критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a}} \left\{ \underbrace{V(\mathbf{a}, C)}_{\text{регуляризирующая функция}} + Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) \right\}$

Квадратичная регуляризация:  $V(\mathbf{a}, C) = \mathbf{a}^T \mathbf{B}_C \mathbf{a}$ ,

$\mathbf{B}_C$  – квадратная матрица  $(n \times n)$ , определяемая конкретной прикладной задачей.

## Оценивание числовой зависимости. Линейно-квадратичная модель

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: действительный вектор  $\mathbf{x}(\omega) = (x_1(\omega) \cdots x_n(\omega))^T \in \mathbb{R}^n$ .

Скрытая характеристика: действительное число  $y(\omega) \in \mathbb{R}$

Центрированная и нормированная обучающая совокупность

$(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_N)$  – матрица  $(n \times N)$ ,  $\mathbf{y} \in \mathbb{R}^N$ .

$$\frac{1}{N} \sum_{j=1}^N x_{ij} = 0, \quad \frac{1}{N} \sum_{j=1}^N x_{ij}^2 = 1 \quad \text{для всех } i = 1, \dots, n.$$

Искомое линейное решающее правило:  $\hat{y}(\mathbf{x}) = \mathbf{a}^T \mathbf{x} : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{a} \in \mathbb{R}^n$ .

Квадратичная функция потерь:  $Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) = \sum_{j=1}^N (y_j - \mathbf{a}^T \mathbf{x}_j)^2 = (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a})$

Напомним общий критерий обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a}} \left\{ \underbrace{V(\mathbf{a}, C)}_{\text{регуляризирующая функция}} + Q(\mathbf{y}, \mathbf{X}, \mathbf{a}) \right\}$

Квадратичная регуляризация:  $V(\mathbf{a}, C) = \mathbf{a}^T \mathbf{B}_C \mathbf{a}$ ,

$\mathbf{B}_C$  – квадратная матрица  $(n \times n)$ , определяемая конкретной прикладной задачей.

Итак, квадратичная задача обучения:  $\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a}} \left\{ \mathbf{a}^T \mathbf{B}_C \mathbf{a} + (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a}) \right\}$

Очевидно, что задача обучения сводится к решению системы  $n$  линейных уравнений.



## Оценивание числовой зависимости. Линейно-квадратичная модель

Итак, квадратичная задача обучения для заданного значения параметра регуляризации  $C$ :

$$\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a} \in \mathbb{R}^n} \left\{ \mathbf{a}^T \mathbf{B}_C \mathbf{a} + (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a}) \right\}$$

Напомним общий вид критерия неявной кросс-валидации:

$$\hat{C}(\mathbf{y}, \mathbf{X}) = \arg \min_C \left\{ Q(\mathbf{y}, \mathbf{X}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X})) + \Delta(C, \mathbf{X}) \right\}.$$

## Оценивание числовой зависимости. Линейно-квадратичная модель

Итак, квадратичная задача обучения для заданного значения параметра регуляризации  $C$ :

$$\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a} \in \mathbb{R}^n} \left\{ \mathbf{a}^T \mathbf{B}_C \mathbf{a} + (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a}) \right\}$$

Напомним общий вид критерия неявной кросс-валидации:

$$\hat{C}(\mathbf{y}, \mathbf{X}) = \arg \min_C \left\{ Q(\mathbf{y}, \mathbf{X}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X})) + \Delta(C, \mathbf{X}) \right\}.$$

Напомним также первую эвристику наблюдателя в нашей теории:

### Принцип незлонамеренности природы.

Наблюдатель правильно «угадал» функцию потерь и класс решающих правил. Природа чаще генерирует пары с низким значением штрафа.

$$\Phi(\mathbf{y} | \mathbf{X}, \mathbf{a}) \propto \exp \left\{ - \underbrace{Q(\mathbf{X}, \mathbf{y}, \mathbf{a})}_{\text{принятая наблюдателем функция потерь}} \right\}$$

*принятая наблюдателем  
функция потерь*

## Оценивание числовой зависимости. Линейно-квадратичная модель

Итак, квадратичная задача обучения для заданного значения параметра регуляризации  $C$ :

$$\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a} \in \mathbb{R}^n} \left\{ \mathbf{a}^T \mathbf{B}_C \mathbf{a} + (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a}) \right\}$$

Напомним общий вид критерия неявной кросс-валидации:

$$\hat{C}(\mathbf{y}, \mathbf{X}) = \arg \min_C \left\{ Q(\mathbf{y}, \mathbf{X}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X})) + \Delta(C, \mathbf{X}) \right\}$$

Напомним также первую эвристику наблюдателя в нашей теории:

### Принцип незлонамеренности природы.

Наблюдатель правильно «угадал» функцию потерь и класс решающих правил. Природа чаще генерирует пары с низким значением штрафа.

$$\Phi(\mathbf{y} | \mathbf{X}, \mathbf{a}) \propto \exp \left\{ - \underbrace{Q(\mathbf{X}, \mathbf{y}, \mathbf{a})}_{\text{принятая наблюдателем функция потерь}} \right\}$$

*принятая наблюдателем  
функция потерь*

**Теорема 2.** В методе неявной кросс-валидации для линейно-квадратичной модели штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке, определяется выражением  $\Delta(C, \mathbf{X}) = \text{Tr} \left[ \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \mathbf{B}_C)^{-1} \right]$ .

## Оценивание числовой зависимости. Линейно-квадратичная модель

Итак, квадратичная задача обучения для заданного значения параметра регуляризации  $C$ :

$$\hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{a} \in \mathbb{R}^n} \left\{ \mathbf{a}^T \mathbf{B}_C \mathbf{a} + (\mathbf{y} - \mathbf{X}^T \mathbf{a})^T (\mathbf{y} - \mathbf{X}^T \mathbf{a}) \right\}$$

Напомним общий вид критерия неявной кросс-валидации:

$$\hat{C}(\mathbf{y}, \mathbf{X}) = \arg \min_C \left\{ Q(\mathbf{y}, \mathbf{X}, \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X})) + \Delta(C, \mathbf{X}) \right\}$$

Напомним также первую эвристику наблюдателя в нашей теории:

<p><b>Принцип незлонамеренности природы.</b> Наблюдатель правильно «угадал» функцию потерь и класс решающих правил. Природа чаще генерирует пары с низким значением штрафа.</p>	$\Phi(\mathbf{y}   \mathbf{X}, \mathbf{a}) \propto \exp \left\{ - \underbrace{Q(\mathbf{X}, \mathbf{y}, \mathbf{a})}_{\substack{\text{принятая наблюдателем} \\ \text{функция потерь}}} \right\}$
---	---

**Теорема 2.** В методе неявной кросс-валидации для линейно-квадратичной модели штраф за подстановку в функцию потерь оценки параметра, вычисленной по той же выборке, определяется выражением  $\Delta(C, \mathbf{X}) = \text{Tr} \left[ \mathbf{X}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \mathbf{B}_C)^{-1} \right]$ .

**Критерий неявной кросс-валидации для линейно-квадратичной модели:**

$\hat{C}(\mathbf{y}, \mathbf{X}) = \arg \min_C \left\{ (\mathbf{y} - \mathbf{X}^T \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X}))^T (\mathbf{y} - \mathbf{X}^T \hat{\mathbf{a}}_C(\mathbf{y}, \mathbf{X})) + \text{Tr} \left[ \mathbf{X}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \mathbf{B}_C)^{-1} \right] \right\}$	$0 \leq \text{Tr} \left[ \mathbf{X}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \mathbf{B}_{\hat{C}(\mathbf{y}, \mathbf{X})})^{-1} \right] \leq n$ <p>эффективная дробная размерность оцененного вектора коэффициентов</p>
---	--

## **Обучение распознаванию образов. Метод опорных векторов**

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

## Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

## Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: индекс класса  $y(\omega) \in \{-1, 1\}$ .

## Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: Индекс класса  $y(\omega) \in \{-1, 1\}$

Параметрическое семейство  
разделяющих гиперплоскостей

$$\mathbf{a}^T \mathbf{x} + b \begin{cases} > 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = 1 \\ < 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = -1 \end{cases}$$

$$\mathbf{x} = (x_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, \mathbf{a} = (a_i, i \in \hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}) \in \mathbb{R}^{\hat{n}}, b \in \mathbb{R}, |\hat{\mathbb{I}}| = \hat{n}$$



## Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: индекс класса  $y(\omega) \in \{-1, 1\}$

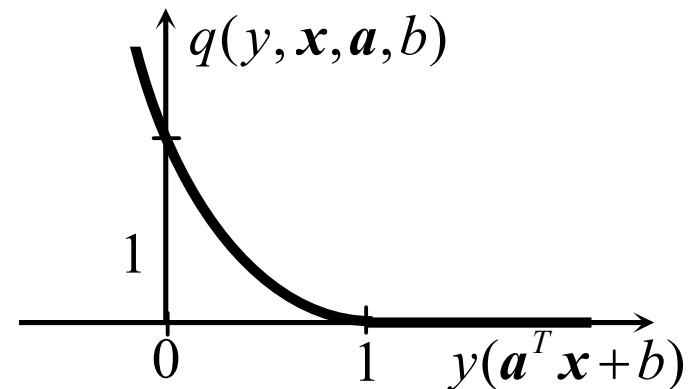
Параметрическое семейство  
разделяющих гиперплоскостей

$$\mathbf{a}^T \mathbf{x} + b \begin{cases} > 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = 1 \\ < 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = -1 \end{cases}$$

$$\mathbf{x} = (x_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, \mathbf{a} = (a_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, b \in \mathbb{R}, |\hat{\mathbb{I}}| = \hat{n}$$

Функция потерь, применимая  
к любому объекту  $(\mathbf{x}, y)$

$$q(y, \mathbf{x}, \mathbf{a}, b) = \begin{cases} 0, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) > 1, \\ (1 - y(\mathbf{a}^T \mathbf{x} + b))^2, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) < 1, \end{cases}$$



## Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: Индекс класса  $y(\omega) \in \{-1, 1\}$

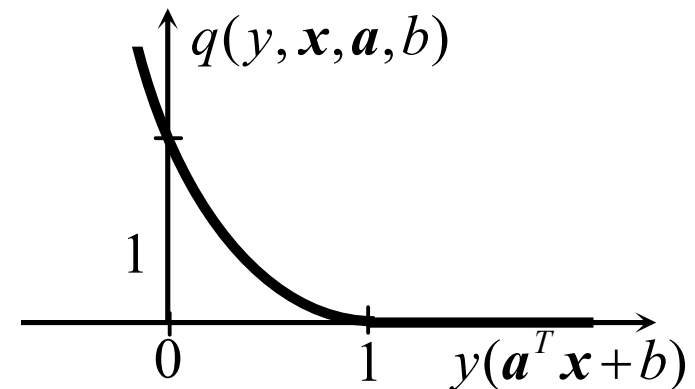
Параметрическое семейство  
разделяющих гиперплоскостей

$$\mathbf{a}^T \mathbf{x} + b \begin{cases} > 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = 1 \\ < 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = -1 \end{cases}$$

$$\mathbf{x} = (x_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, \mathbf{a} = (a_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, b \in \mathbb{R}, |\hat{\mathbb{I}}| = \hat{n}$$

Функция потерь, применимая  
к любому объекту  $(\mathbf{x}, y)$

$$q(y, \mathbf{x}, \mathbf{a}, b) = \begin{cases} 0, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) > 1, \\ (1 - y(\mathbf{a}^T \mathbf{x} + b))^2, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) < 1, \end{cases}$$



Критерий обучения (квадратичная версия SVM):

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{\mathbf{a}, b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

# Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: Индекс класса  $y(\omega) \in \{-1, 1\}$

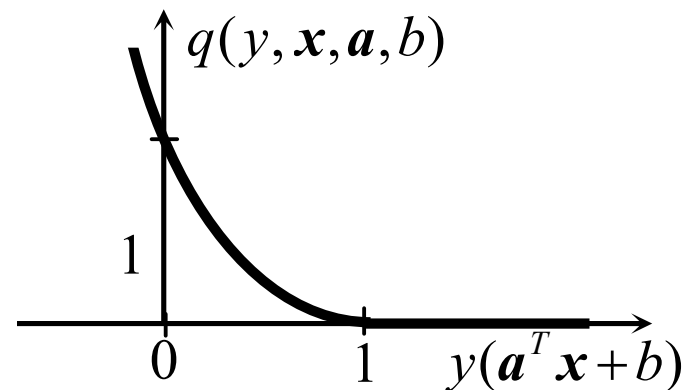
Параметрическое семейство  
разделяющих гиперплоскостей

$$\mathbf{a}^T \mathbf{x} + b \begin{cases} > 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = 1 \\ < 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = -1 \end{cases}$$

$$\mathbf{x} = (x_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, \mathbf{a} = (a_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, b \in \mathbb{R}, |\hat{\mathbb{I}}| = \hat{n}$$

Функция потерь, применимая  
к любому объекту  $(\mathbf{x}, y)$

$$q(y, \mathbf{x}, \mathbf{a}, b) = \begin{cases} 0, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) > 1, \\ (1 - y(\mathbf{a}^T \mathbf{x} + b))^2, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) < 1, \end{cases}$$



Критерий обучения (квадратичная версия SVM):

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{\mathbf{a}, b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

# Обучение распознаванию образов. Метод опорных векторов

Некоторое множество реально существующих объектов  $\omega \in \Omega$ .

Наблюдаемая характеристика: Исходный вектор признаков  $\mathbf{x}(\omega) = (x_i(\omega), i = 1, \dots, n)$ .

Подмножество «полезных» признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ .

Скрытая характеристика: Индекс класса  $y(\omega) \in \{-1, 1\}$

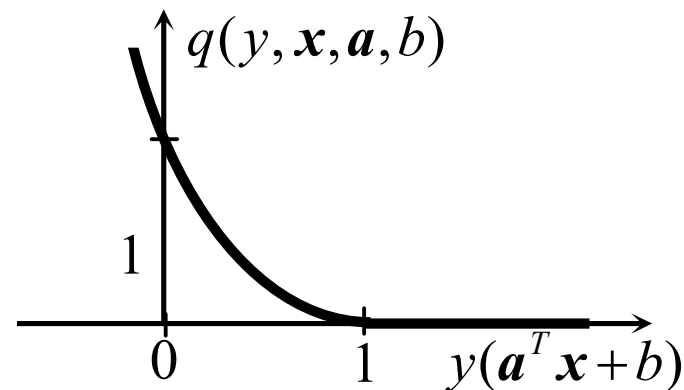
Параметрическое семейство  
разделяющих гиперплоскостей

$$\mathbf{a}^T \mathbf{x} + b \begin{cases} > 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = 1 \\ < 0 \Rightarrow \hat{y}(\mathbf{x} | \mathbf{a}, b) = -1 \end{cases}$$

$$\mathbf{x} = (x_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, \mathbf{a} = (a_i, i \in \hat{\mathbb{I}}) \in \mathbb{R}^{\hat{n}}, b \in \mathbb{R}, |\hat{\mathbb{I}}| = \hat{n}$$

Функция потерь, применимая  
к любому объекту  $(\mathbf{x}, y)$

$$q(y, \mathbf{x}, \mathbf{a}, b) = \begin{cases} 0, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) > 1, \\ (1 - y(\mathbf{a}^T \mathbf{x} + b))^2, & \text{if } y(\mathbf{a}^T \mathbf{x} + b) < 1, \end{cases}$$



Критерий обучения (квадратичная версия SVM):

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{\mathbf{a}, b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

Как выбрать структурные параметры по единственной обучающей совокупности?

## Обучение распознаванию образов. Метод опорных векторов

Еще раз критерий обучения (квадратичная версия SVM):

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{a, b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \quad \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$ .

# Обучение распознаванию образов. Метод опорных векторов

Еще раз критерий обучения (квадратичная версия SVM): Структурные параметры:

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{a,b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

1) подмножество признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ ,  $\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$   
 2) коэффициент баланса  $C$

Объединим параметры дискриминантной гиперплоскости:  $\mathbf{a} = (\mathbf{a}, b) \in \mathbb{R}^{\hat{n}+1}$ ,  $\mathbf{x} = (\mathbf{x}, 1) \in \mathbb{R}^{\hat{n}+1}$

Критерий обучения:

$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}},C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}},C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}.$$

# Обучение распознаванию образов. Метод опорных векторов

Еще раз критерий обучения (квадратичная версия SVM): Структурные параметры:

$$\begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = \arg \min_{\mathbf{a}, b} \left[ \underbrace{\mathbf{a}^T \mathbf{a} + C \sum_{j: y_j(\mathbf{a}^T \mathbf{x}_j + b) < 1} (y_j - (\mathbf{a}^T \mathbf{x}_j + b))^2}_{\text{эмпирический риск}} \right]$$

1) подмножество признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$ ,  $\hat{n} = |\hat{\mathbb{I}}| \leq n$ ,  $\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$

2) коэффициент баланса  $C$

Объединим параметры дискриминантной гиперплоскости:  $\mathbf{a} = (\mathbf{a}, b) \in \mathbb{R}^{\hat{n}+1}$ ,  $\mathbf{x} = (\mathbf{x}, 1) \in \mathbb{R}^{\hat{n}+1}$

Критерий обучения:  $\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}},C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases}$   $\mathbf{B}_{\hat{\mathbb{I}},C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}$   $((\hat{n}+1) \times (\hat{n}+1))$

Решение:

**активные ограничения**

1) подмножество опорных объектов  $\begin{cases} \hat{\mathbb{J}}_{\hat{\mathbb{I}},C} = \{j \in \mathbb{J} : y_j \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) = 1 - \hat{\xi}_j, \hat{\xi}_j > 0\}, \\ \hat{\xi}_j^2 = (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}},C})^2. \end{cases}$

2) параметры гиперплоскости  $\hat{\mathbf{a}}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) = \begin{pmatrix} \hat{\mathbf{a}}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \\ \hat{b}_{\hat{\mathbb{I}},C}(\mathbf{X}, \mathbf{y}) \end{pmatrix} = (\hat{\mathbf{X}}_{\hat{\mathbb{I}},C} \hat{\mathbf{X}}_{\hat{\mathbb{I}},C}^T + \mathbf{B}_{\hat{\mathbb{I}},C})^{-1} \hat{\mathbf{X}}_{\hat{\mathbb{I}},C} \hat{\mathbf{y}}_{\hat{\mathbb{I}},C}$

Здесь:  $\hat{\mathbf{X}}_{\hat{\mathbb{I}},C} = \begin{pmatrix} \mathbf{x}_{j_1} \cdots \mathbf{x}_{j_{\hat{N}_{\hat{\mathbb{I}},C}}} \\ 1 \cdots 1 \end{pmatrix} ((\hat{n}+1) \times \hat{N}_{\hat{\mathbb{I}},C})$ ,  $\hat{\mathbf{y}}_{\hat{\mathbb{I}},C} = (y_{j_1} \cdots y_{j_{\hat{N}_{\hat{\mathbb{I}},C}}})^T$ ,  $\hat{N}_{\hat{\mathbb{I}},C} = |\hat{\mathbb{J}}_{\hat{\mathbb{I}},C}|$

# Обучение распознаванию образов. Метод опорных векторов

Критерий обучения (квадратичная версия SVM):

$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

Подмножество опорных объектов

$$\begin{cases} \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C} = \{j \in \mathbb{J} : y_j \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = 1 - \hat{\xi}_j, \hat{\xi}_j > 0\}, \\ \hat{\xi}_j^2 = (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C})^2. \end{cases}$$

Дополнительные обозначения:

$$\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} \mathbf{x}_{j_1} \cdots \mathbf{x}_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}} \\ 1 \cdots 1 \end{pmatrix}_{(\hat{n}+1) \times \hat{N}_{\hat{\mathbb{I}}, C}}, \quad \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} = (y_{j_1} \cdots y_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}})^T, \quad \hat{N}_{\hat{\mathbb{I}}, C} = |\hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}|.$$

Оптимальная гиперплоскость полностью определяется подмножеством опорных объектов

$$\hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = \operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right) =$$

$$\operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \left\| \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} - \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T \mathbf{a} \right\|^2 \right) = (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C}$$



# Обучение распознаванию образов. Метод опорных векторов

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \quad \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

Критерий обучения (квадратичная версия SVM):

$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \quad \xi_j \geq 0, \quad j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Подмножество опорных объектов

$$\begin{cases} \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C} = \left\{ j \in \mathbb{J} : y_j \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = 1 - \hat{\xi}_j, \quad \hat{\xi}_j > 0 \right\}, \\ \hat{\xi}_j^2 = (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C})^2. \end{cases}$$

Дополнительные обозначения:

$$\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} \mathbf{x}_{j_1} & \dots & \mathbf{x}_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}} \\ 1 & \dots & 1 \end{pmatrix}_{(\hat{n}+1) \times \hat{N}_{\hat{\mathbb{I}}, C}}, \quad \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} = (y_{j_1} \dots y_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}})^T, \quad \hat{N}_{\hat{\mathbb{I}}, C} = |\hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}|$$

Оптимальная гиперплоскость полностью определяется подмножеством опорных объектов

$$\hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = \operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right) =$$

$$\operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \left\| \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} - \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T \mathbf{a} \right\|^2 \right) = (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C}$$

Для сравнения:

Критерий обучения в линейно-квадратичной модели

$$\hat{\mathbf{a}}_C(\mathbf{X}, \mathbf{y}) = \operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_C \mathbf{a} + \left\| \mathbf{y} - \mathbf{X}^T \mathbf{a} \right\|^2 \right) = (\mathbf{X} \mathbf{X}^T + \mathbf{B}_C)^{-1} \mathbf{X} \mathbf{y}$$

# Обучение распознаванию образов. Метод опорных векторов

Критерий обучения (квадратичная версия SVM):

$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{0} \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

Подмножество опорных объектов

$$\begin{cases} \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C} = \left\{ j \in \mathbb{J} : y_j \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = 1 - \hat{\xi}_j, \hat{\xi}_j > 0 \right\}, \\ \hat{\xi}_j^2 = (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C})^2. \end{cases}$$

Дополнительные обозначения:

$$\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} \mathbf{x}_{j_1} & \dots & \mathbf{x}_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}} \\ 1 & \dots & 1 \end{pmatrix}_{(\hat{n}+1) \times \hat{N}_{\hat{\mathbb{I}}, C}}, \quad \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} = (y_{j_1} \dots y_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}})^T, \quad \hat{N}_{\hat{\mathbb{I}}, C} = |\hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}|$$

Оптимальная гиперплоскость полностью определяется подмножеством опорных объектов

$$\begin{aligned} \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) &= \operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} (y_j - \mathbf{x}_j^T \mathbf{a})^2 \right) = \\ &= \operatorname{argmin} \left( \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \left\| \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} - \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T \mathbf{a} \right\|^2 \right) = (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C}. \end{aligned}$$

Центральная идея, сближающая модель SVM и линейную квадратичную модель:

Рассматривать подмножество опорных объектов  $\hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}$  как вторичный структурный параметр, полностью определяемый исходными структурными параметрами  $(\hat{\mathbb{I}}, C)$

# Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Критерий обучения:

$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases}$$

$$\mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Структурные параметры:

1) подмножество признаков

$$\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}, \hat{n} = |\hat{\mathbb{I}}| \leq n,$$

$$\mathbf{a}, \mathbf{x}_j \in \mathbb{R}^{\hat{n}}$$

2) коэффициент баланса  $C$

Подмножество опорных объектов, найденное в результате обучения

$$\begin{cases} \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C} = \{j \in \mathbb{J} : y_j \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C}(\mathbf{X}, \mathbf{y}) = 1 - \hat{\xi}_j, \hat{\xi}_j > 0\}, \\ \hat{\xi}_j^2 = (y_j - \mathbf{x}_j^T \hat{\mathbf{a}}_{\hat{\mathbb{I}}, C})^2. \end{cases}$$

Дополнительные обозначения:

$$\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} \mathbf{x}_{j_1} \cdots \mathbf{x}_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}} \\ 1 \cdots 1 \end{pmatrix}_{(\hat{n}+1) \times \hat{N}_{\hat{\mathbb{I}}, C}}, \quad \hat{\mathbf{y}}_{\hat{\mathbb{I}}, C} = (y_{j_1} \cdots y_{j_{\hat{N}_{\hat{\mathbb{I}}, C}}})^T, \quad \hat{N}_{\hat{\mathbb{I}}, C} = |\hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}|.$$

Критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра, использующий аналогию с линейно-квадратичной моделью:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}, C} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}}, C} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}}, C} \left\{ \sum_{j \in \mathbb{J}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

Напомним критерий обучения: 
$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}, C} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

Напомним критерий обучения: 
$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \hat{\mathbb{J}}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \hat{\mathbb{J}}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \hat{\mathbb{J}} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Алгоритм решения задачи обучения расскажет Павел Левдик в следующем докладе.

## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}, C} \left\{ \sum_{j \in \hat{\mathbb{I}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

Напомним критерий обучения: 
$$\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases} \quad \mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$$

Алгоритм решения задачи обучения расскажет Павел Левдик в следующем докладе.

Опыт показывает, что значение параметра  $C$  следует принимать как можно бóльшим.

Пусть  $C = \text{const} \rightarrow \infty$ .

**Теорема.** 
$$\lim_{C \rightarrow \infty} \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \lim_{C \rightarrow \infty} \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right] = \min \left\{ \hat{n}, \hat{N}_{\hat{\mathbb{I}}, C} \right\} + 1$$

## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\begin{pmatrix} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{pmatrix} = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}, C} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

Напомним критерий обучения:  $\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ y_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases}$   $\mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$

Алгоритм решения задачи обучения расскажет Павел Левдик в следующем докладе.

Опыт показывает, что значение параметра  $C$  следует принимать как можно бóльшим.

Пусть  $C = \text{const} \rightarrow \infty$ .

**Теорема.**  $\lim_{C \rightarrow \infty} \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \lim_{C \rightarrow \infty} \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right] = \min \{ \hat{n}, \hat{N}_{\hat{\mathbb{I}}, C} \} + 1$ .

Критерий неявной кросс-валидации для фиксированного большого значения  $C$ :

$$\hat{\mathbb{I}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \left( \min \{ \hat{n}, \hat{N}_{\hat{\mathbb{I}}, C} \} + 1 \right) \right\}.$$



## Критерий неявной кросс-валидации выбора структурных параметров в квадратичной модели SVM

Итак, критерий неявной кросс-валидации для подмножества опорных объектов как вторичного структурного параметра:

$$\left( \begin{array}{c} \hat{\mathbb{I}}(\mathbf{y}, \mathbf{X}) \\ \hat{C}(\mathbf{y}, \mathbf{X}) \end{array} \right) = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}, C} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) \right\}, \quad \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right]$$

Напомним критерий обучения:  $\begin{cases} \mathbf{a}^T \mathbf{B}_{\hat{\mathbb{I}}, C} \mathbf{a} + \sum_{j \in \mathbb{J}} \xi_j^2 \rightarrow \min(\mathbf{a}, \xi_j, j \in \mathbb{J}), \\ \mathbf{y}_j \mathbf{x}_j^T \mathbf{a} \geq 1 - \xi_j, \xi_j \geq 0, j \in \mathbb{J} = \{1, \dots, N\}, \end{cases}$   $\mathbf{B}_{\hat{\mathbb{I}}, C} = \begin{pmatrix} (1/C) \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{pmatrix}_{((\hat{n}+1) \times (\hat{n}+1))}$

Алгоритм решения задачи обучения расскажет Павел Левдик в следующем докладе.

Опыт показывает, что значение параметра  $C$  следует принимать как можно бóльшим.

Пусть  $C = \text{const} \rightarrow \infty$ .

**Теорема.**  $\lim_{C \rightarrow \infty} \Delta(\hat{\mathbb{I}}, C, \mathbf{X}) = \lim_{C \rightarrow \infty} \text{Tr} \left[ \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T (\hat{\mathbf{X}}_{\hat{\mathbb{I}}, C} \hat{\mathbf{X}}_{\hat{\mathbb{I}}, C}^T + \mathbf{B}_{\hat{\mathbb{I}}, C})^{-1} \right] = \min \{ \hat{n}, \hat{N}_{\hat{\mathbb{I}}, C} \} + 1$

Критерий неявной кросс-валидации для фиксированного большого значения  $C$ :

$$\hat{\mathbb{I}}_C(\mathbf{y}, \mathbf{X}) = \arg \min_{\hat{\mathbb{I}} \subseteq \mathbb{I}} \left\{ \sum_{j \in \hat{\mathbb{J}}_{\hat{\mathbb{I}}, C}} \hat{\xi}_{\hat{\mathbb{I}}, C, j}^2 + \left( \min \{ \hat{n}, \hat{N}_{\hat{\mathbb{I}}, C} \} + 1 \right) \right\}.$$

Конечно, огромное число  $2^n$  всех подмножеств в множестве признаков  $\hat{\mathbb{I}} \subseteq \mathbb{I} = \{1, \dots, n\}$  делает невозможной «наивную» реализацию этого критерия.

В следующем докладе Павел Левдик расскажет, как обойтись вообще без перебора.

Спасибо за внимание!