

# **Байесовские постановки задачи обучения распознаванию образов по методу опорных векторов с управляемой селективностью отбора информативных признаков объектов**

**Татарчук Александр Игоревич  
Моттль Вадим Вячеславович  
Вычислительный центр РАН**

## Типовая задача распознавания образов в множествах объектов реального мира

Наблюдателя интересует некоторое множество объектов реального мира:  $\omega \in \Omega$ .

**Пара характеристик:**

$$\begin{cases} \mathbf{x}(\omega) \in \mathbb{X} - \text{наблюдаемая,} \\ y(\omega) \in \mathbb{Y} = \{-1, 1\} - \text{скрытая (индекс класса).} \end{cases}$$

**Обучающая выборка:**

$$\Omega' \subset \Omega: \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}.$$

**Решающая функция:**

$$\hat{y}(\mathbf{x}): \mathbb{X} \rightarrow \mathbb{Y},$$

$$\hat{y}(\mathbf{x}(\omega)) \neq y(\omega) - \text{ошибка.}$$

# Линейный подход к обучению распознавания двух классов

Вектор действительных признаков  $\mathbf{x}(\omega) \in \mathbb{R}^n$  погружает множество объектов в  $\mathbb{R}^n$ .

Евклидова метрика в  $\mathbb{R}^n$ :

$$\rho(\mathbf{x}', \mathbf{x}'') = \|\mathbf{x}' - \mathbf{x}''\| = \left( (\mathbf{x}' - \mathbf{x}'')^T (\mathbf{x}' - \mathbf{x}'') \right)^{1/2}.$$

Разделяющая гиперплоскость:

$$\mathcal{H}(\mathbf{a}, b) = \{ \mathbf{x} \in \mathbb{R}^n : \mathbf{a}^T \mathbf{x} + b = 0 \}, \quad \mathbf{a} \in \mathbb{R}^n.$$

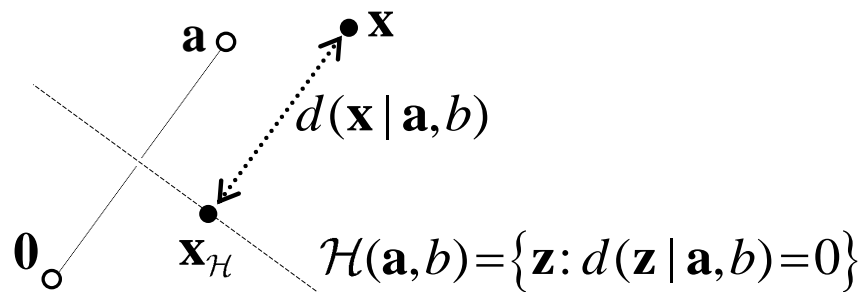
Решающая функция:

$$d(\mathbf{x} | \mathbf{a}, b) = \mathbf{a}^T \mathbf{x} + b \geq 0.$$

Свойство решающей функции:

$$|d(\mathbf{x} | \mathbf{a}, b)| = \rho(\mathbf{x}, \mathbf{x}_{\mathcal{H}}), \quad \mathbf{x}, \mathbf{x}_{\mathcal{H}} \in \mathbb{R}^n, \quad \|\mathbf{a}\| = 1.$$

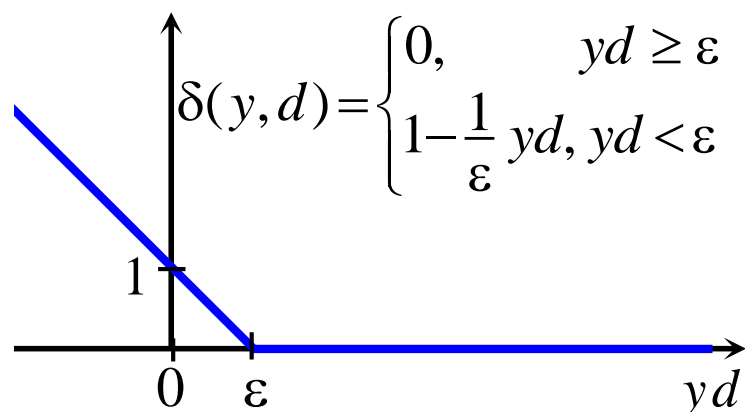
Некоторая конкретизация гипотезы компактности Э.М. Бравермана



# Функция потерь: Степень несоответствия значения решающей функции скрытой характеристике объекта

Индекс класса объекта  $y \in \{-1, 1\}$

Метод опорных векторов



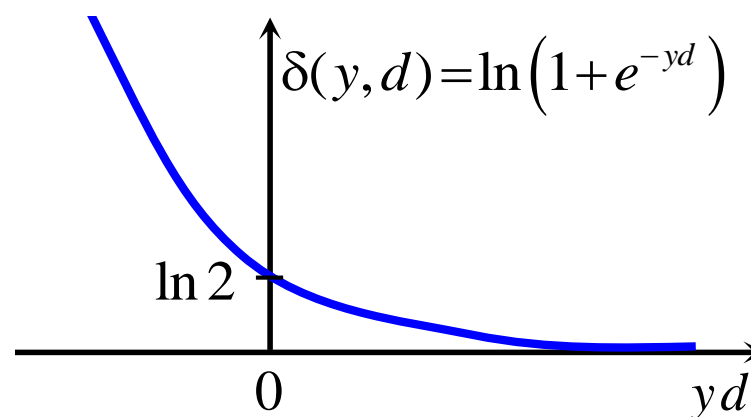
Принцип обучения

Разделение обучающей выборки:

- с наибольшим зазором  $\varepsilon \rightarrow \max$ ,
- с наименьшей суммой потерь

$$\sum_{j=1}^N \delta(y_j, d(\mathbf{x}_j | \mathbf{a}, b)) \rightarrow \min(\mathbf{a}, b).$$

Метод логистической регрессии



} противоречивые требования

## Предложенный невыпуклый критерий обучения

$$\begin{cases} \varepsilon^{-2} + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \boldsymbol{\delta}, \varepsilon), \quad \mathbf{a}^T \mathbf{a} = 1, \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq \varepsilon(1 - \delta_j), \quad \delta_j \geq 0, \quad j = 1, \dots, N, \quad \varepsilon \geq 0. \end{cases}$$

**Утверждение:** Предложенный критерий эквивалентен классическому критерию метода опорных векторов

$$\begin{cases} \mathbf{a}^T \mathbf{a} + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \boldsymbol{\delta}), \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \quad \delta_j \geq 0, \quad j = 1, \dots, N. \end{cases}$$

**Свойства решающей функция для нового объекта:**

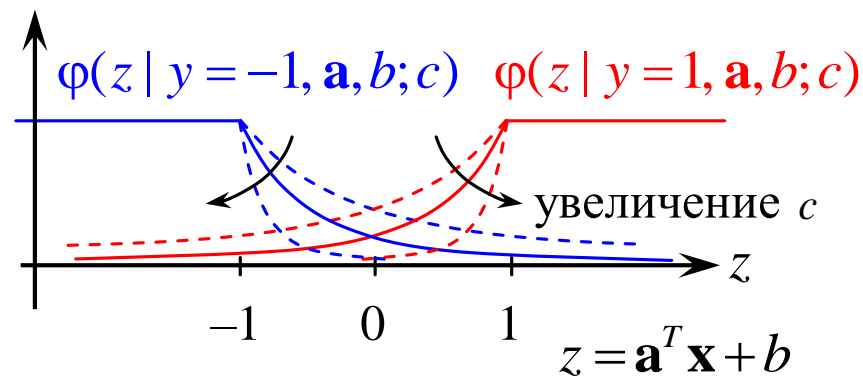
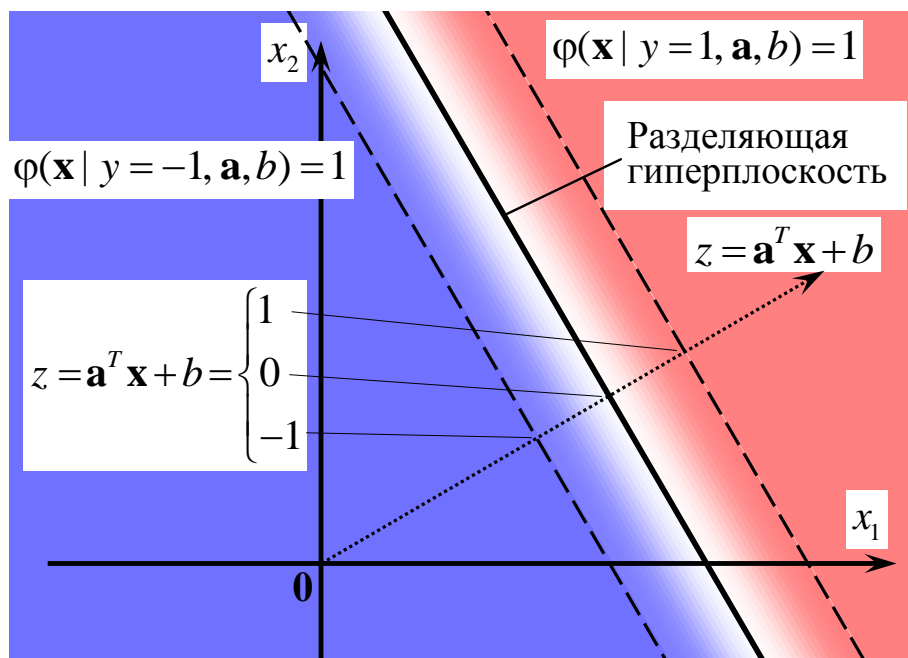
$$d(\mathbf{x} | \hat{\mathbf{a}}, \hat{b}) = \hat{\mathbf{a}}^T \mathbf{x} + \hat{b} = \sum_{j: \hat{\lambda}_j > 0} \hat{\lambda}_j y_j (\mathbf{x}_j^T \mathbf{x}) + b \geq 0$$

Требуются только **скалярные произведения**  $\mathbf{x}_j^T \mathbf{x}$  нового вектора признаков и **только с опорными векторами.**

# Предложенная вероятностная модель наблюдения объектов в признаковом пространстве

**Предположение 1:** Объекты двух классов  $y = \pm 1$  распределены преимущественно по разные стороны некоторой гиперплоскости в пространстве признаков  $\mathbf{x} \in \mathbb{R}^n$ .

$$\varphi(\mathbf{x}|y, \mathbf{a}, b; c) \propto \begin{cases} 1, & y(\mathbf{a}^T \mathbf{x} + b) \geq 1, \\ \exp[-c(1 - y(\mathbf{a}^T \mathbf{x} + b))], & y(\mathbf{a}^T \mathbf{x} + b) < 1, \end{cases} \quad \mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R}, c > 0.$$



**Предположение 2:** Несобственная совместная априорная плотность параметров гиперплоскости  $(\mathbf{a}, b)$  не зависит от  $b$

$$\Psi(\mathbf{a}, b) \propto \Psi(\mathbf{a}).$$

# Байесовский критерий обучения для пары несобственных плотностей распределения

**Критерий обучения:** Максимизация апостериорной плотности вероятностей параметров гиперплоскости

$$(\hat{\mathbf{a}}, \hat{b} | X, Y; c) = \arg \max_{\mathbf{a}, b} \left[ \ln \Psi(\mathbf{a}) + \sum_{j=1}^N \ln \varphi(\mathbf{x}_j | y_j, \mathbf{a}, b; c) \right].$$

**Теорема 1.** Критерий обучения для заданных плотностей эквивалентен задаче

$$\begin{cases} -\ln \Psi(\mathbf{a}) + c \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \delta), \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{cases}$$

**Свойства критерия:**

Если  $\ln \Psi(\mathbf{a})$  вогнута, то критерий выпуклый и точка минимума  $(\hat{\mathbf{a}}, \hat{b}, \hat{\delta})$  единственна.

Активные ограничения  $y_j (\mathbf{a}^T \mathbf{x}_j + b) = 1 - \delta_j$  определяют опорные векторы  $\mathbf{x}_j$ .

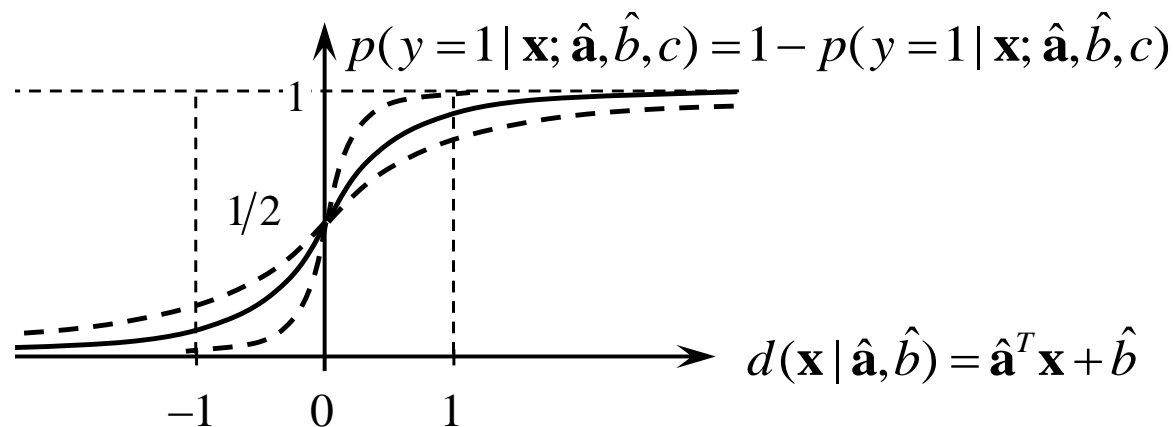
Решающая функция зависит только от опорных векторов.

## Апостериорные вероятности классов: для случая равных априорных вероятностей

**Теорема 2:** Для случая равных априорных вероятностей  $q_{-1} = q_1 = 1/2$  апостериорные вероятности классов будут иметь вид сигмоидоподобной функции

$$p(y=1|\mathbf{x}, \hat{\mathbf{a}}, \hat{b}; c) = 1 - p(y=-1|\mathbf{x}, \hat{\mathbf{a}}, \hat{b}; c) = \begin{cases} \left\{1 + \exp(c) \exp[-c(\hat{\mathbf{a}}^T \mathbf{x} + \hat{b})]\right\}^{-1}, & \hat{\mathbf{a}}^T \mathbf{x} + \hat{b} < -1, \\ \left\{1 + \exp[-2c(\hat{\mathbf{a}}^T \mathbf{x} + \hat{b})]\right\}^{-1}, & -1 \leq \hat{\mathbf{a}}^T \mathbf{x} + \hat{b} \leq 1, \\ \left\{1 + \exp(-c) \exp[-c(\hat{\mathbf{a}}^T \mathbf{x} + \hat{b})]\right\}^{-1}, & \hat{\mathbf{a}}^T \mathbf{x} + \hat{b} > 1. \end{cases}$$

**Графическое представление** апостериорных вероятностей классов



Параметр  $c > 0$  управляет крутизной сигмоидоподобного обобщения ступенчатой решающей функции  $\hat{\mathbf{a}}^T \mathbf{x} + \hat{b} \geq 0$ .



# Известные частные случаи критериев обучения

<p>Нормальная плотность</p> <p><b>Метод опорных векторов</b></p>	$\Psi(\mathbf{a}) = \Psi(a_1, \dots, a_n) = \prod_{i=1}^n \mathcal{N}(a_i   0, r) = \frac{1}{r^{n/2} (2\pi)^{n/2}} \exp\left(-\frac{1}{2r} \mathbf{a}^T \mathbf{a}\right)$ $\begin{cases} \mathbf{a}^T \mathbf{a} + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b), & C = 2rc, \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{cases}$
<p>Плотность Лапласа</p> <p><b>Lasso SVM (1-norm SVM)</b></p>	$\Psi(\mathbf{a}   \beta, \mu) = (2r)^{-n/2} \exp\left(- (r/2)^{-1/2} \sum_{i=1}^n  a_i \right).$ $\begin{cases} \sum_{i=1}^n  a_i  + c (r/2)^{1/2} \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \delta), \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{cases}$
<p>Специальная плотность</p> <p>Нормирующая константа (доказано Е. Черноусовой)</p> <p><b>Elastic Net SVM (DrSVM)</b></p>	$\Psi(\mathbf{a}   \beta, \mu) = D^n \exp\left[-\sum_{i=1}^n (\beta a_i^2 + \mu  a_i )\right], \beta \geq 0, \mu \geq 0.$ $D = \left[ 2 \sqrt{\frac{\pi}{\beta}} \exp\left(\frac{\mu^2}{4\beta}\right) \Phi\left(\frac{\mu}{\sqrt{2\beta}}\right) \right]^{-1}, \Phi(u) = \frac{1}{\sqrt{2\pi}} \int_u^{\infty} \exp\left(-\frac{z^2}{2}\right) dz.$ $\begin{cases} \beta \sum_{i=1}^n a_i^2 + \mu \sum_{i=1}^n  a_i  + c \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \delta), \\ y_j (\mathbf{a}^T \mathbf{x}_j + b) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{cases}$

## Первый предлагаемый метод: Метод релевантных признаков

**Принцип:** Селекция признаков выполняется путем совместного байесовского оценивания направляющего вектора и дисперсий его компонент.

Нормальные условные распределения параметров гиперплоскости

$$\Psi(\mathbf{a} | \mathbf{r}) = \prod_{i=1}^n \mathcal{N}(a_i | 0, r_i) = \left( \prod_{i=1}^n r_i \right)^{-1/2} (2\pi)^{-n/2} \exp\left( -\frac{1}{2} \sum_{i=1}^n (1/r_i) a_i^2 \right).$$

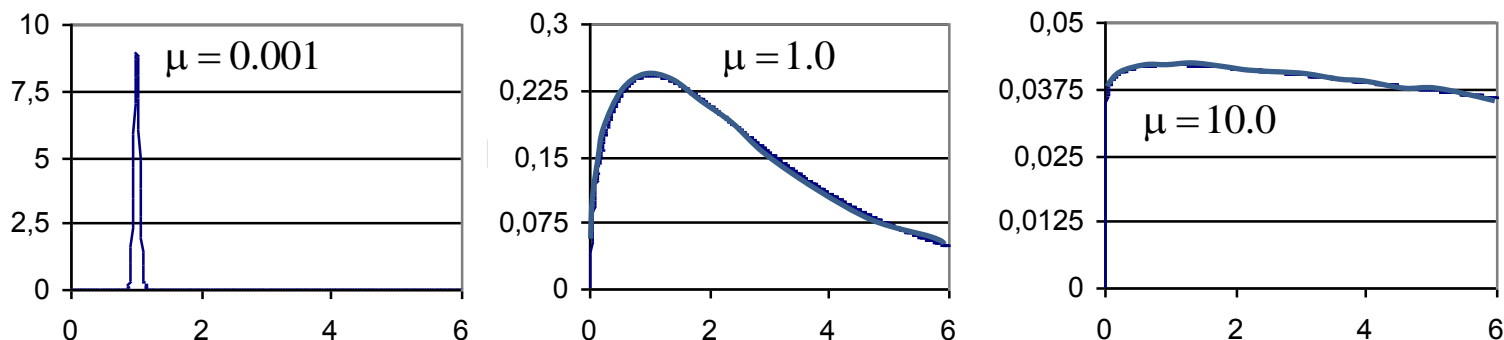
Независимые гамма распределения обратных дисперсий:

$$G(\mathbf{r} | \mu) = \prod_{i=1}^n \gamma(1/r_i | \mu),$$

$$\gamma(1/r_i | \mu) \propto (1/r_i)^{1/2\mu} \exp\left[ -(1/2\mu)(1/r_i) \right],$$

$\mu \geq 0$  – параметр селективности.

Зависимость гамма распределения от параметра селективности



## Критерий обучения

**Теорема 3.** Критерий обучения эквивалентен оптимизационной задаче

$$\begin{cases} \sum_{i=1}^n \left[ (1/r_i) (a_i^2 + (1/\mu)) + ((1/\mu) + 1 + \mu) \ln r_i \right] + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \mathbf{r}, \boldsymbol{\delta}), \\ y_j (\mathbf{a}^T \mathbf{x} + b) \geq 1 - \delta_j, \delta_j \geq 0, j=1, \dots, N, \quad C = 2c. \end{cases}$$

Чем больше дисперсия  $r_i$ , тем больше вес  $i$ -го признака в решающем правиле.

$$d(\mathbf{x} | \hat{\mathbf{a}}, \hat{b}, \hat{\mathbf{r}}) = \sum_{i=1}^n r_i \sum_{j: \lambda_j > 0} y_j \lambda_j x_{ij} x_i + b \stackrel{>}{<} 0.$$

### Итерационный алгоритм обучения

Покоординатная оптимизация двух групп переменных	$\mathbf{r}^k = (r_i^k = 1, i = 1, \dots, n) \Rightarrow (\mathbf{a}^k, b^k), (\delta_j^k, j = 1, \dots, N),$ $\mathbf{a}^k \Rightarrow \mathbf{r}^{k+1}, \mathbf{r}^{k+1} \Rightarrow (\mathbf{a}^{k+1}, b^{k+1}).$ Начальное приближение $\mathbf{r}^0 = \mathbf{1}.$
Шаг 1: $\mathbf{r}^k \Rightarrow (\mathbf{a}^k, b^k)$	$\begin{cases} (\mathbf{a}^k, b^k, \boldsymbol{\delta}^k) = \arg \min_{\mathbf{a}, b, \boldsymbol{\delta}} \sum_{i=1}^n (1/r^k) a_i^2 + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b), \\ y_j (\mathbf{a}^T \mathbf{x} + b) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{cases}$
Шаг 2: $\mathbf{a}^k \Rightarrow \mathbf{r}^{k+1}$	<b>Теорема 4.</b> Новые значения дисперсий $r_i^{k+1} = \arg \min_{r_i} \left[ (1/r_i) (a_i^k)^2 + (1/\mu) + ((1/\mu) + 1 + \mu) \ln r_i \right] = \frac{(a_i^k)^2 + 1/\mu}{\mu + 1 + 1/\mu}$
Критерий останова	$\ \mathbf{a}^k - \mathbf{a}^{k+1}\  < \varepsilon.$ Алгоритм обычно сходится за 10-15 шагов.

## Второй предлагаемый метод: Метод опорных признаков

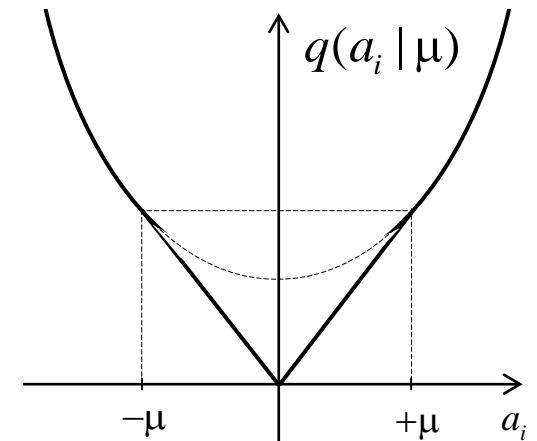
**Принцип:** Априорная плотность  $\Psi(\mathbf{a})$  есть комбинация распределения Лапласа при небольших значениях нормы  $|a| \leq \mu$  и нормального распределения при больших  $|a_i| > \mu$ .

Априорное распределение независимых компонент направляющего вектора

$$\psi(a_i | \mu) \propto \exp(-q(a_i | \mu)), \quad i = 1, \dots, n,$$

$$q(a_i | \mu) = \begin{cases} 2\mu |a_i|, & |a_i| \leq \mu, \\ \mu^2 + a_i^2, & |a_i| > \mu. \end{cases}$$

$\mu \geq 0$  – параметр селективности.



**Теорема 5.** Критерий обучения эквивалентен оптимизационной задаче

$$\begin{cases} J_{SFM}(\mathbf{a}, b, \delta | C, \mu) = 2\mu \sum_{|a_i| \leq \mu} |a_i| + \sum_{|a_i| > \mu} (\mu^2 + a_i^2) + C \sum_{j=1}^N \delta_j \rightarrow \min(\mathbf{a}, b, \delta), \\ y_j (\mathbf{a}^T \mathbf{x} + b) \leq 1 - \delta_j, \quad \delta_j \geq 0, \quad j = 1, \dots, N. \end{cases}$$

Задача выпуклой оптимизации.

# Двойственная задача оптимизационная задача обучения

**Теорема 5.** Критерий обучения имеет двойственную оптимизационную задачу

$$\begin{cases} L(\lambda, \xi | C, \mu) = \sum_{j=1}^N \lambda_j - \sum_{i \in I} \sum_{i=2}^n (1/2) \xi_i \rightarrow \max(\lambda, \xi), \\ \xi_i \geq 0, \quad \xi_i \geq \left( \sum_{j=1}^N y_j x_{ij} \hat{\lambda}_j \right)^2 - \mu^2, \quad i \in I = \{1, \dots, n\}, \\ \sum_{j=1}^N y_j \lambda_j = 0, \quad 0 \leq \lambda_j \leq C/2, \quad j = 1, \dots, N. \end{cases}$$

Решение  $(\hat{\lambda}_1, \dots, \hat{\lambda}_N)$  определяет разбиение множества компонент направляющего:

$$\begin{cases} \hat{a}_i = \sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{ij}, & i \in I^+ = \left\{ i \in I : \left( \sum_{j=1}^N y_j x_{ij} \hat{\lambda}_j \right)^2 - \mu^2 > 0 \right\}, \\ \hat{a}_i = \hat{\eta}_i \sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{ij}, & i \in I^0 = \left\{ i \in I : \left( \sum_{j=1}^N y_j x_{ij} \hat{\lambda}_j \right)^2 - \mu^2 = 0 \right\}, \\ \hat{a}_i = 0, & i \in I^- = \left\{ i \in I : \left( \sum_{j=1}^N y_j x_{ij} \hat{\lambda}_j \right)^2 - \mu^2 < 0 \right\}. \end{cases} \begin{array}{l} \text{опорные признаки} \\ \text{удаленные признаки} \end{array}$$

**Теорема 6.** Итоговое решение о классе нового объекта имеет вид

$$d(\mathbf{x} | \hat{\lambda}, \eta_i, i \in I^0, b) = \underbrace{\sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j \left[ \underbrace{\sum_{i \in I^+} x_{ij} x_i + \sum_{i \in I^0} \hat{\eta}_i x_{ij} x_i}_{\text{опорные признаки}} \right]}_{\text{опорные векторы}} + \hat{b} \begin{matrix} > \\ < \end{matrix} 0,$$

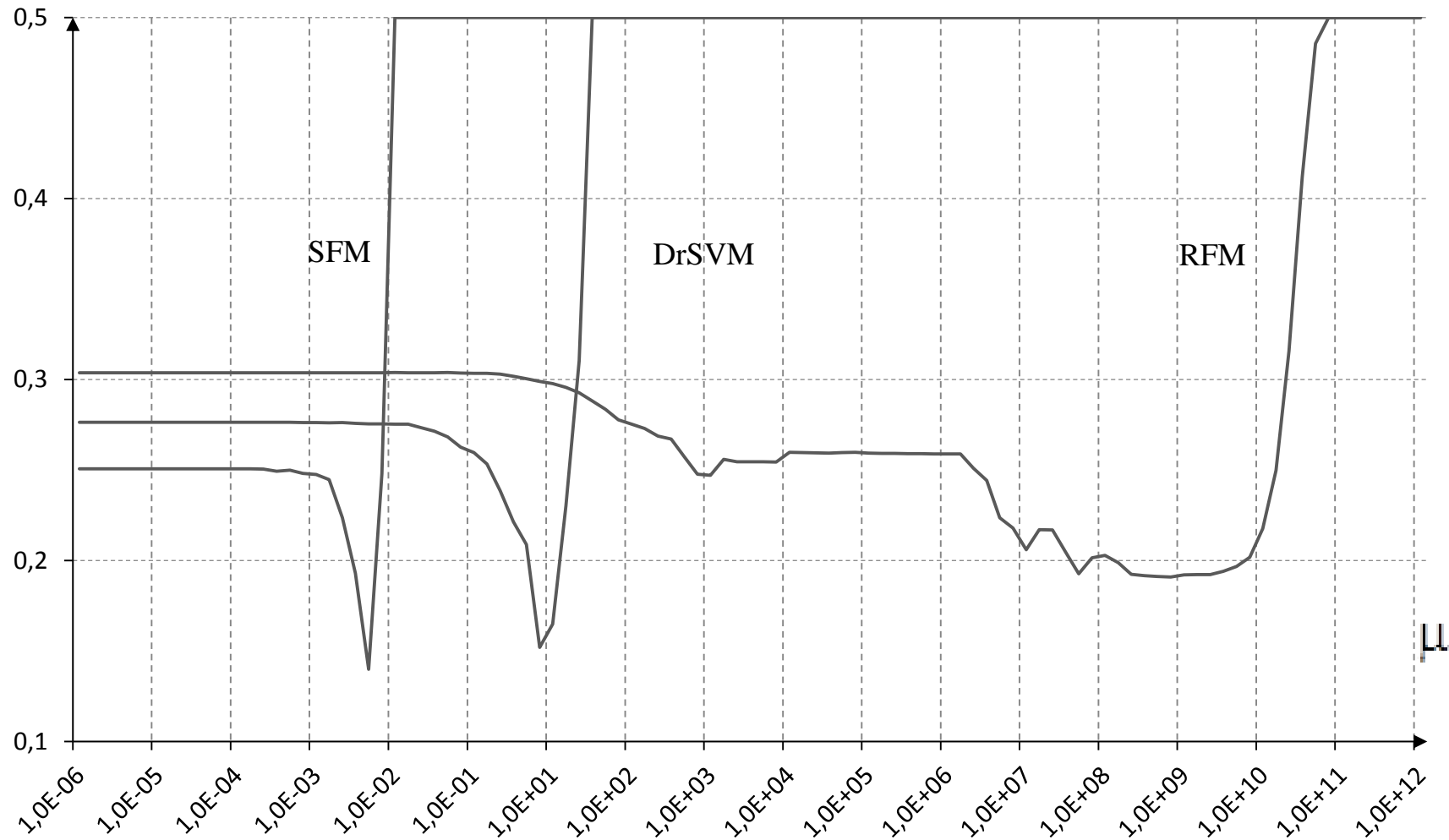
где коэффициенты  $\hat{\eta}_i$  и  $\hat{b}$  определяются задачей линейного программирования.

# Результаты экспериментального исследования: Структура модельного эксперимента

Модель данных	<p>100-мерное признаковое пространство <math>\mathbf{x} = (x_1, \dots, x_{100})</math>. Наблюдения сгенерированы независимо по многомерному нормальному закону распределения с математическими ожиданиями</p> <p style="text-align: center;"> <i>первый класс</i> <span style="margin-left: 200px;"><i>второй класс</i></span> </p> $\mathbf{m}_+ = \left( \underbrace{0.5, \dots, 0.5}_5, \underbrace{0, \dots, 0}_{95} \right)^T \in \mathbb{R}^{100} \text{ и } \mathbf{m}_- = \left( \underbrace{-0.5, \dots, -0.5}_5, \underbrace{0, \dots, 0}_{95} \right)^T \in \mathbb{R}^{100},$ <p style="text-align: center;"> <i>информативные признаки</i> <span style="margin-left: 100px;"><i>шумовые признаки</i></span> <span style="margin-left: 200px;"><i>информативные признаки</i></span> <span style="margin-left: 100px;"><i>шумовые признаки</i></span> </p> <p>единичной ковариационной матрицей <math>\Sigma = \mathbf{I}_{n \times n} = \begin{pmatrix} 1 &amp; 0 &amp; \dots &amp; 0 \\ 0 &amp; 1 &amp; \dots &amp; 0 \\ \dots &amp; \dots &amp; \dots &amp; \dots \\ 0 &amp; 0 &amp; \dots &amp; 1 \end{pmatrix}</math>.</p> <p><b>Все признаки независимы. Информативные признаки равнозначны.</b></p>
Обучающая совокупность	<p><math>N = N_+ = N_- = 50 + 50 = 100</math> векторов <math>\mathbf{x}_j</math>, <math>j = 1, \dots, N</math> двух классов <math>y_j = \pm 1</math>, выбранных случайно и независимо.</p>
Обучение классификатора	<p>Оценивание параметров гиперплоскости <math>(a_1, \dots, a_{100}, b)</math> независимо по 100 обучающим совокупностям для каждого набора пар значений параметра селективности <math>10^{-6} \leq \mu \leq 10^{12}</math> и разделимости <math>10^{-6} \leq C \leq 10^{12}</math>.</p>
Оценивание качества	<p>Тестовая совокупность <math>N_{test} = 50\,000 + 50\,000 = 100\,000</math> случайных векторов.</p>

# Модельный эксперимент (продолжение)

Доля ошибочно классифицированных объектов тестовой выборки



Пример эффекта регуляризации методов на одной из модельных обучающих совокупностей для возрастающих значений параметра селективности  $\mu$  при фиксированных оптимальных значениях параметра делимости  $\hat{C}$ .

## Итоги серии модельных экспериментов

Метод обучения		Минимальная средняя ошибка на тестовой совокупности по 100 обучающим совокупностям			
		Информативные признаки линейно независимы		Только совместно информативные признаки обеспечивают приемлемую точность	
		Шумовые признаки независимы	Шумовые признаки линейно зависимы	Шумовые признаки независимы	Шумовые признаки линейно зависимы
Предложенные	SFM	<b>0.1495</b>	<b>0.1429</b>	0.3140	0.2364
	RFM	0.1797	0.1611	<b>0.2325</b>	<b>0.2212</b>
Существующие	DrSVM	0.1523	0.1489	0.3140	0.3406
	SVM	0.2353	0.1620	0.4359	0.2364
	SVM на 5 заведомо информативных признаках	<b>0.1430</b>	<b>0.1397</b>	<b>0.2238</b>	<b>0.2177</b>
	«Оракульная» ошибка	<b>0.1320</b>	<b>0.132</b>	<b>0.2150</b>	<b>0.215</b>



# Основные результаты работы

1. Общая математическая постановка задачи обучения двухклассовому распознаванию образов в линейном пространстве признаков на основе количественного измерения расстояния между вектором признаков объекта и разделяющей гиперплоскостью.
2. Вероятностная модель наблюдения объектов в пространстве векторов признаков относительно фиксированной разделяющей гиперплоскости.
3. Два семейства априорных вероятностных моделей направляющего вектора разделяющей гиперплоскости, отражающих стратегии отбора признаков на основе взвешивания всех исходно заданных признаков (feature weighting) и на основе жесткого выбора их подмножества (feature subset selection).
4. Комплекс байесовских методов и алгоритмов оценивания разделяющей гиперплоскости, реализующих принцип обучения с заданной селективностью отбора признаков объектов.