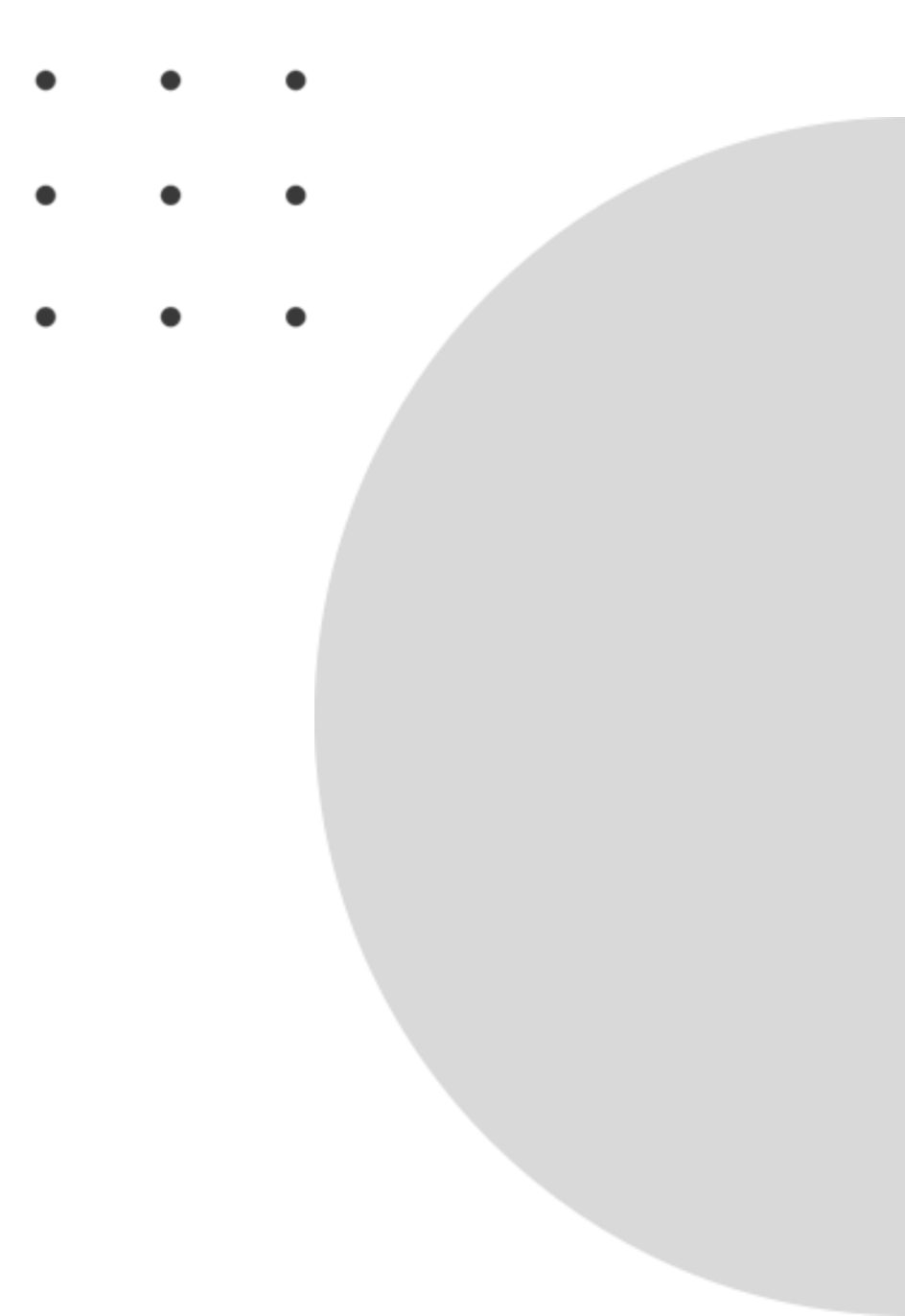


ВВЕДЕНИЕ В МАШИННОЕ ОБУЧЕНИЕ

Лекция №4



План лекции

1. Support Vector Machine
2. Principal component analysis

01

SVM

Линейная классификация

Вспомним, что такое модель линейной классификации

$$X^l = (x_i, y_i)_{i=1}^l, x_i \in R^n, y_i \in \{-1, 1\}$$

$$a(x; w, w_0) = \text{sign}(\langle x, w \rangle - w_0) \quad w \in R^n, w_0 \in R$$

Где функция ошибки:

$$\sum_{i=1}^l [a(x_i; w, w_0) \neq y_i] = \sum_{i=1}^l [M_i(w, w_0) < 0] \rightarrow \min_{w, w_0}.$$

$$M_i(w, w_0) = (\langle x_i, w \rangle - w_0) y_i$$

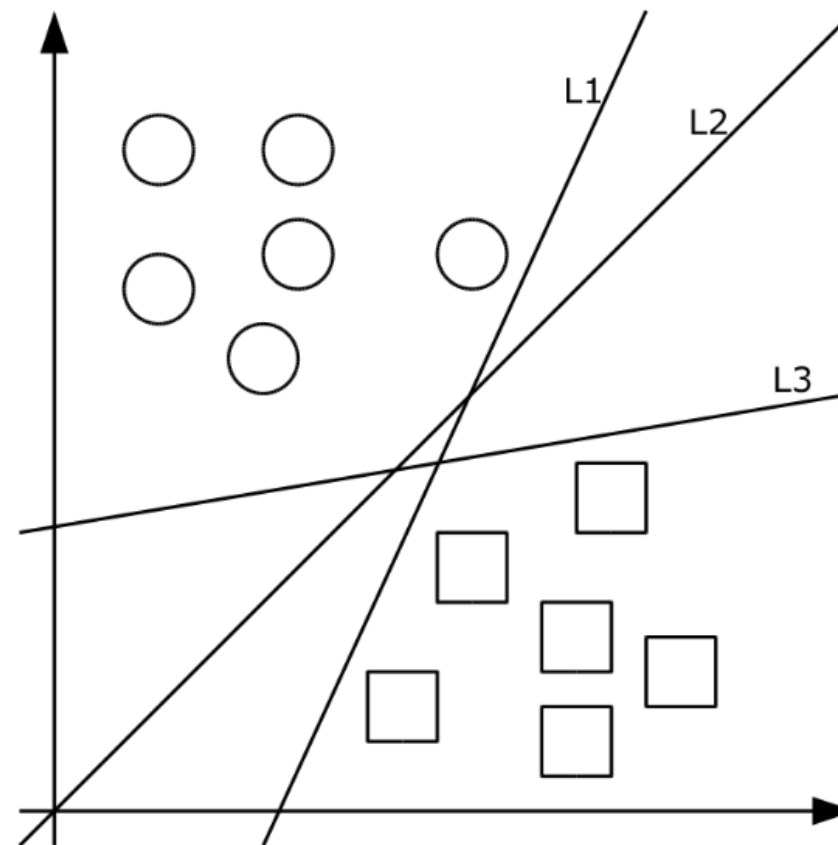
Решение классификации

Вспомним, что такое модель линейной классификации

$$X^l = (x_i, y_i)_{i=1}^l, x_i \in R^n, y_i \in \{-1, 1\}$$

$$\exists w, w_0 : M_i(w, w_0) = y_i(\langle w, x_i \rangle - w_0) > 0, \quad i = 1, \dots, \ell$$

Какое решение l1, l2, l3 лучше всего в данном случае?



Решение классификации

Вспомним, что такое модель линейной классификации

$$X^l = (x_i, y_i)_{i=1}^l, x_i \in R^n, y_i \in \{-1, 1\}$$

$$\exists w, w_0 : M_i(w, w_0) = y_i(\langle w, x_i \rangle - w_0) > 0, \quad i = 1, \dots, \ell$$

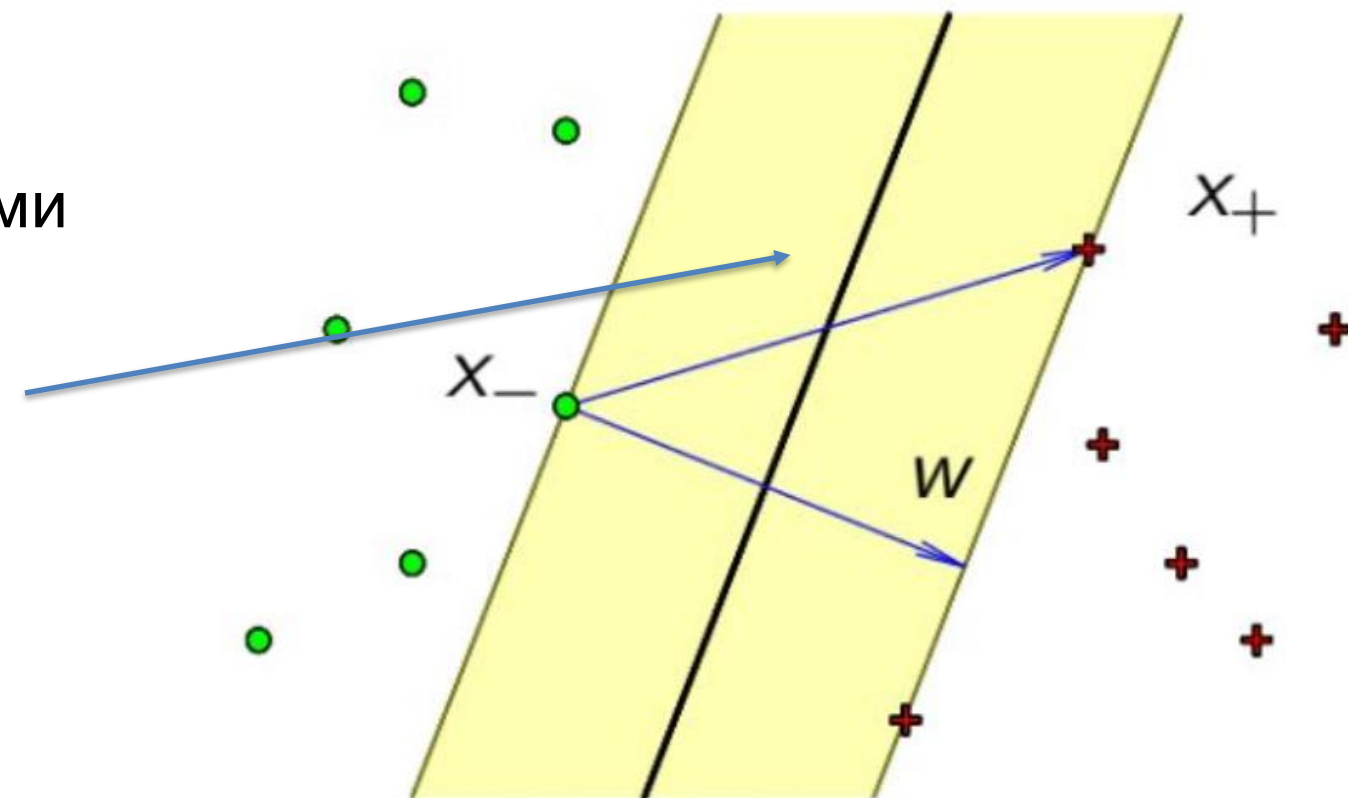
Введем такую нормировку: $\min_{i=1, \dots, \ell} M_i(w, w_0) = 1$

Теперь отступ между самыми близкими объектами классов равен 1

$$\{x : -1 \leq \langle w, x \rangle - w_0 \leq 1\}$$

$$\forall x_+ : \langle w, x_+ \rangle - w_0 \geq 1$$

$$\forall x_- : \langle w, x_- \rangle - w_0 \leq -1$$



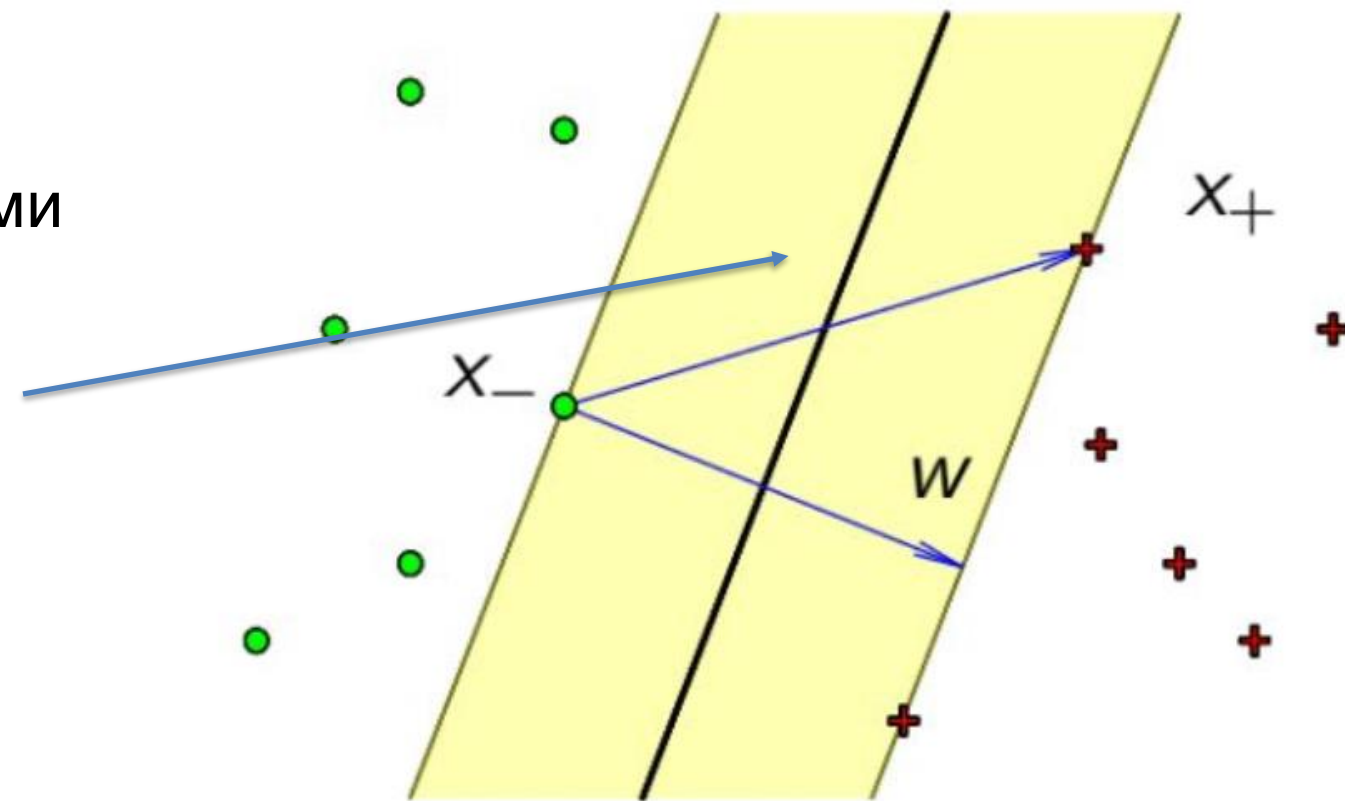
Решение классификации

Теперь отступ между самыми близкими объектами классов равен 1

$$\{x : -1 \leq \langle w, x \rangle - w_0 \leq 1\}$$

$$\forall x_+ : \langle w, x_+ \rangle - w_0 \geq 1$$

$$\forall x_- : \langle w, x_- \rangle - w_0 \leq -1$$



Поставим задачу максимизации такого отступа

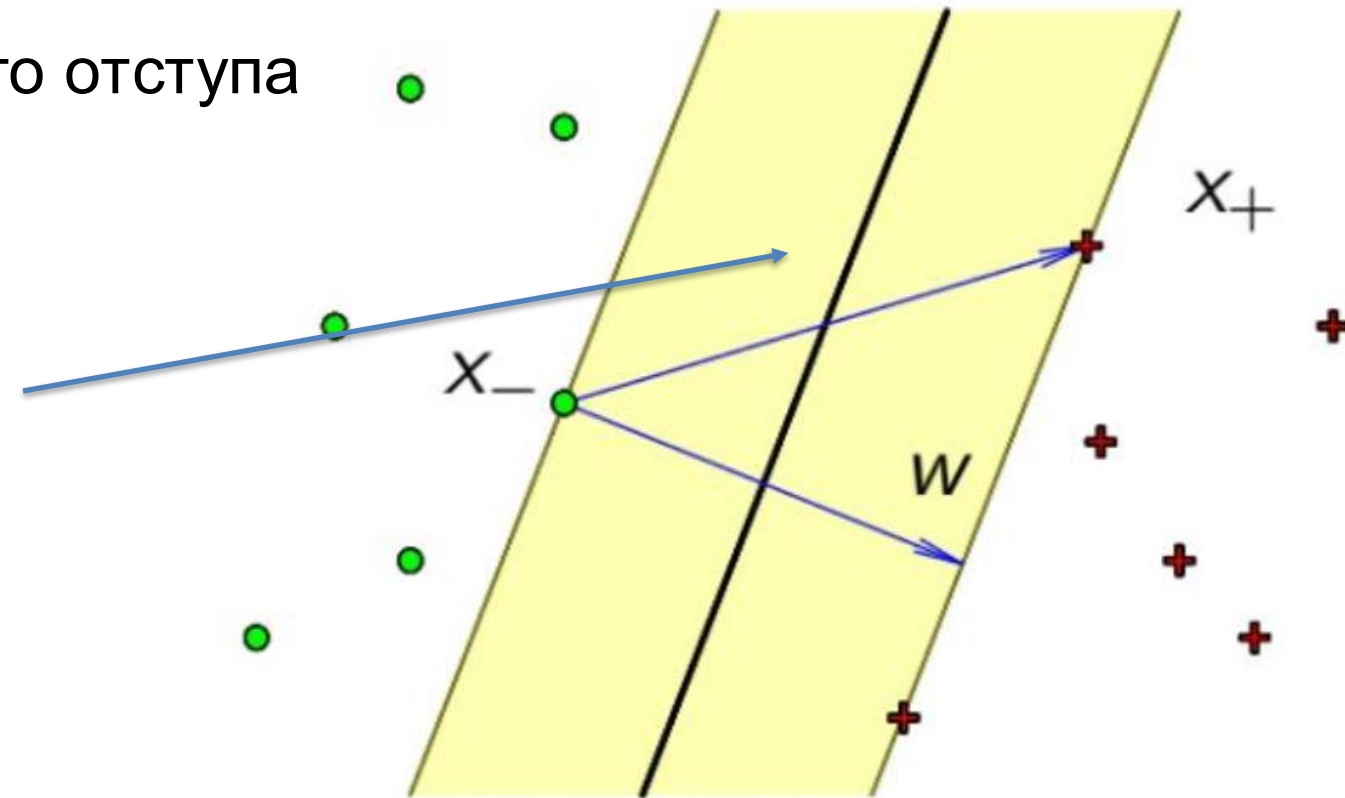
$$\frac{\langle x_+ - x_-, w \rangle}{\|w\|} \geq \boxed{\frac{2}{\|w\|}} \rightarrow \max$$

За счет введенной нормировки M

Идеальный случай

Поставим задачу максимизации такого отступа

$$\begin{cases} \frac{1}{2} \|w\|^2 \rightarrow \min; \\ M_i(w, w_0) \geq 1, \quad i = 1, \dots, \ell. \end{cases}$$



Нелинейная классификация

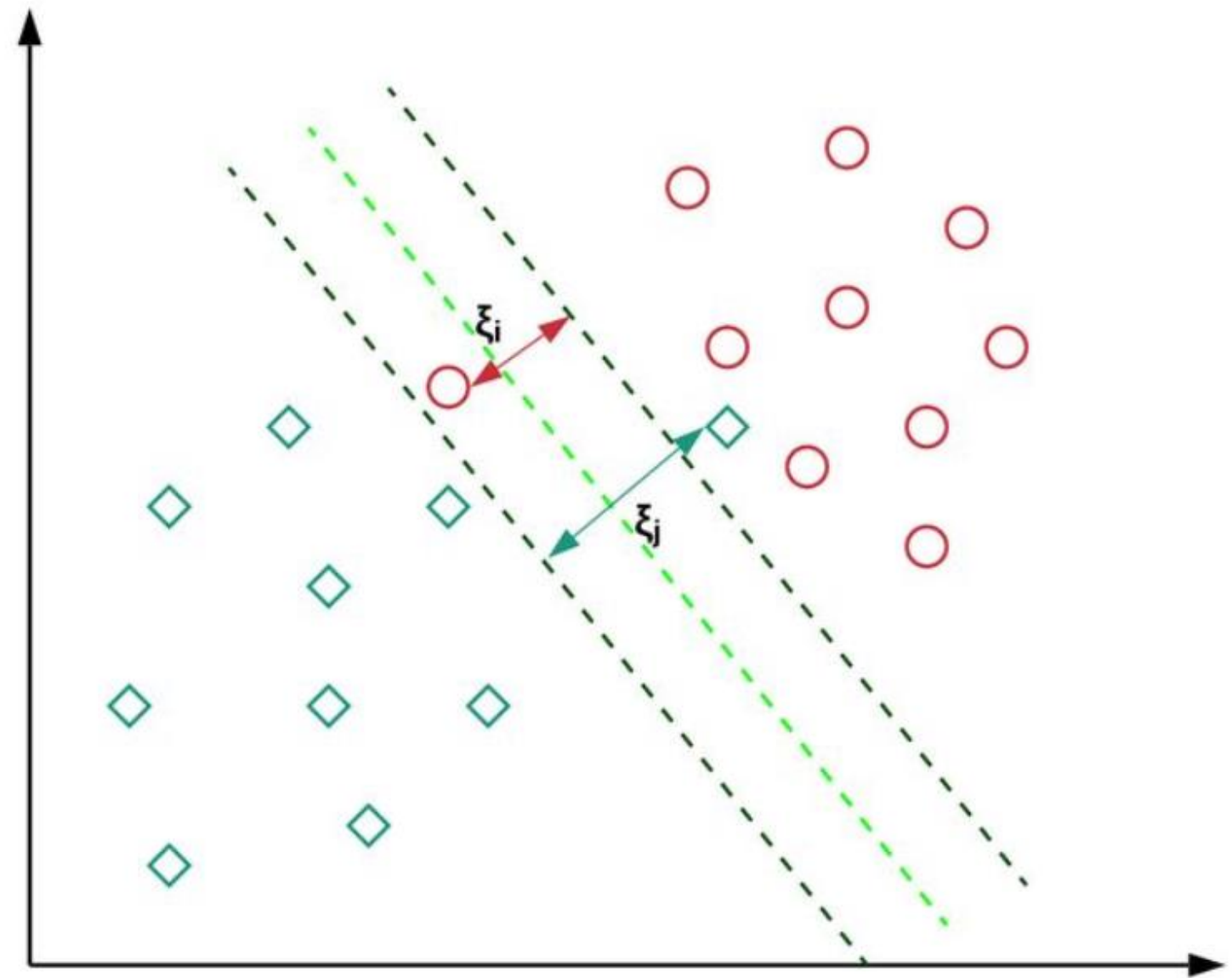
Поставим задачу максимизации такого отступа

$$\begin{cases} \frac{1}{2} \|w\|^2 \rightarrow \min; \\ M_i(w, w_0) \geq 1, \quad i = 1, \dots, \ell. \end{cases}$$

Введем эмпирику

$$M_i(w, w_0) \geq 1 - \xi_i, \quad i = 1, \dots, \ell;$$

Допускаем, что классификатор ошибается



Нелинейная классификация

Оптимизация с ограничениями

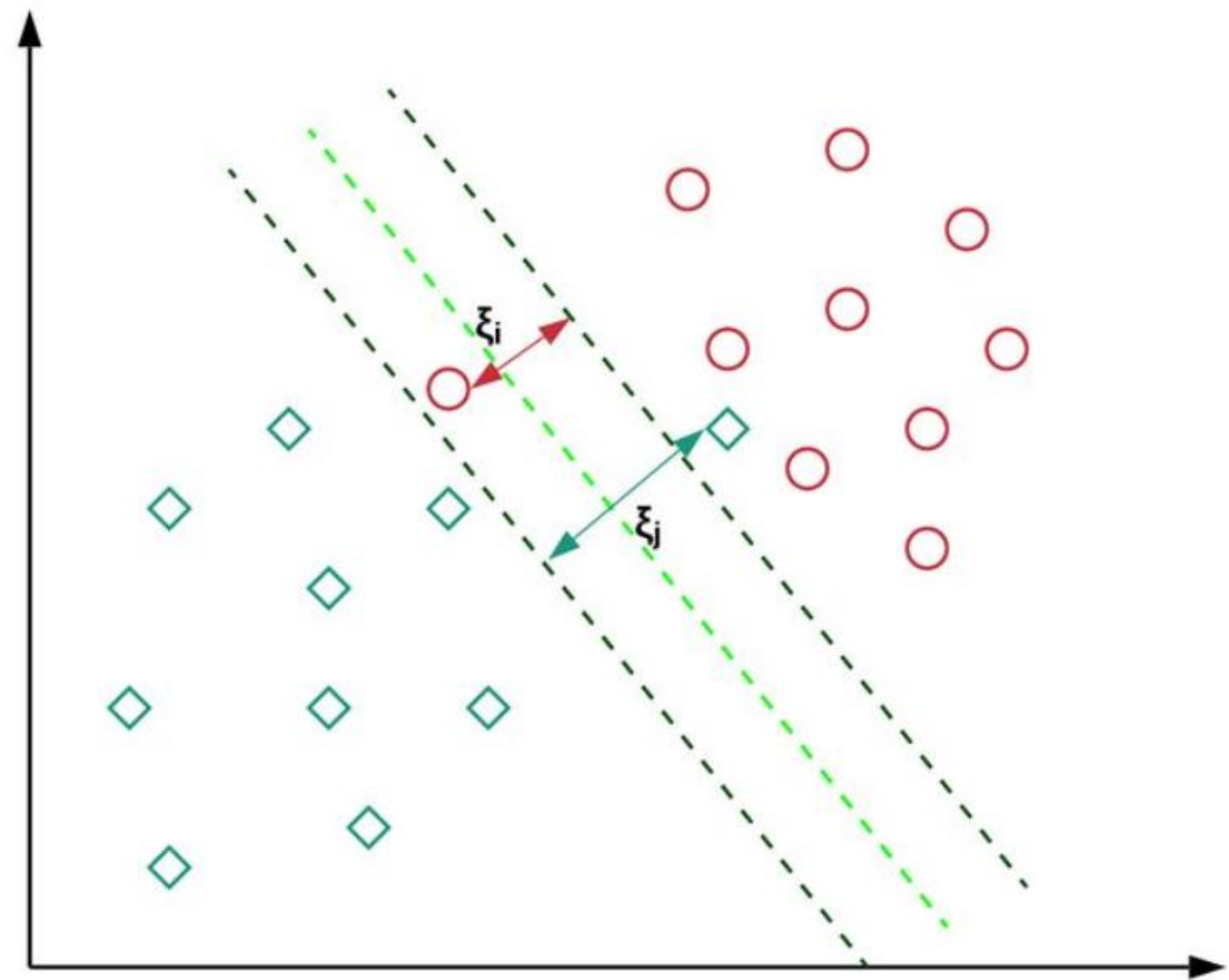
$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min_{w, w_0, \xi}; \\ M_i(w, w_0) \geq 1 - \xi_i, \quad i = 1, \dots, \ell; \\ \xi_i \geq 0, \quad i = 1, \dots, \ell. \end{cases}$$

Эквивалентная задача безусловной оптимизации

$$C \sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2} \|w\|^2 \rightarrow \min_{w, w_0}.$$

Кусочно-линейная

C – скалярная величина, сила регуляризации (гиперпараметр)

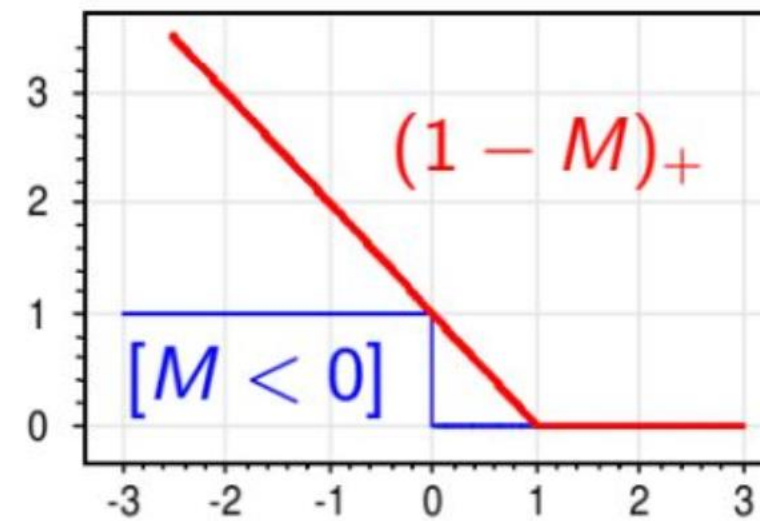


Геометрическая постановка

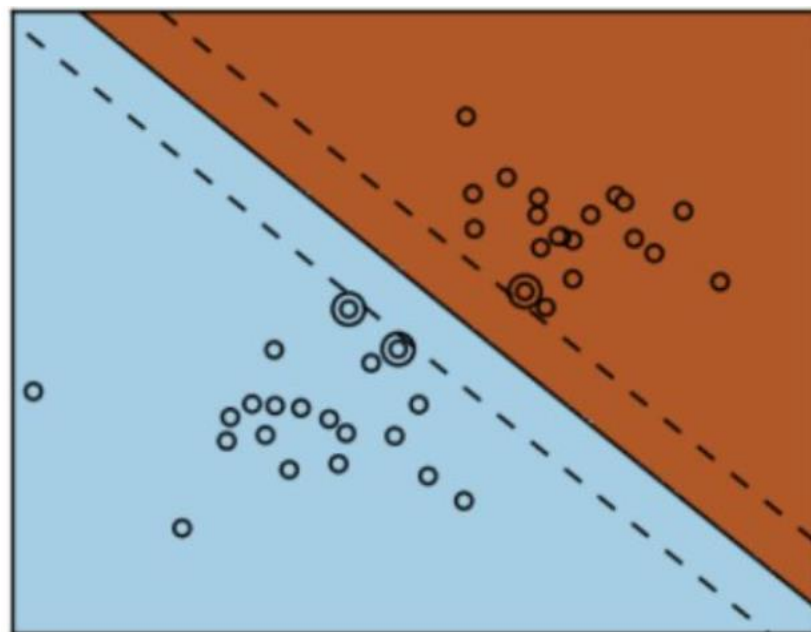
Регуляризация

Оптимизация с ограничениями

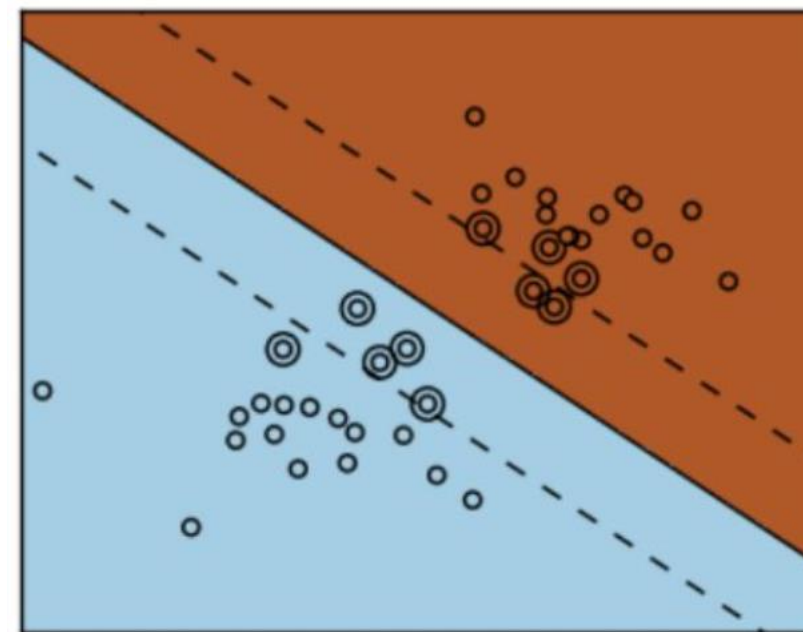
$$Q(w, w_0) = \sum_{i=1}^v [M_i(w, w_0) < 0] \leq$$
$$\leq \underbrace{\sum_{i=1}^l (1 - M_i(w, w_0))_+}_{\text{Approximation}} + \underbrace{\frac{1}{2C} \|w\|^2}_{\text{Regularization}} \rightarrow \min$$



С большая: слабая оптимизация



С маленькая: сильная оптимизация



Нелинейная классификация

Задача квадратичного программирования. Теорема Каруша-Куна-Таккера

$$\begin{cases} f(x) \rightarrow \min_x; \\ g_i(x) \leq 0, & i = 1, \dots, m; \\ h_j(x) = 0, & j = 1, \dots, k. \end{cases}$$

Необходимые условия. Если x — точка локального минимума, то существуют множители $\mu_i, i = 1, \dots, m, \lambda_j, j = 1, \dots, k$:

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial x} = 0, & \mathcal{L}(x; \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i g_i(x) + \sum_{j=1}^k \lambda_j h_j(x); \\ g_i(x) \leq 0; & h_j(x) = 0; \text{ (исходные ограничения)} \\ \mu_i \geq 0; & \text{(двойственные ограничения)} \\ \mu_i g_i(x) = 0; & \text{(условие дополняющей нежёсткости)} \end{cases}$$

Нелинейная классификация

Задача квадратичного программирования. Теорема
Каруша-Куна-Таккера

Функция Лагранжа: $\mathcal{L}(w, w_0, \xi; \lambda, \eta) =$

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^{\ell} \lambda_i (M_i(w, w_0) - 1) - \sum_{i=1}^{\ell} \xi_i (\lambda_i + \eta_i - C),$$

λ_i — переменные, двойственные к ограничениям $M_i \geq 1 - \xi_i$;
 η_i — переменные, двойственные к ограничениям $\xi_i \geq 0$.

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0, & \frac{\partial \mathcal{L}}{\partial w_0} = 0, & \frac{\partial \mathcal{L}}{\partial \xi} = 0; \\ \xi_i \geq 0, & \lambda_i \geq 0, & \eta_i \geq 0, & i = 1, \dots, \ell; \\ \lambda_i = 0 \text{ либо } M_i(w, w_0) = 1 - \xi_i, & i = 1, \dots, \ell; \\ \eta_i = 0 \text{ либо } \xi_i = 0, & i = 1, \dots, \ell; \end{cases}$$

Нелинейная классификация

Задача квадратичного программирования. Теорема
Каруша-Куна-Таккера

Функция Лагранжа: $\mathcal{L}(w, w_0, \xi; \lambda, \eta) =$

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^{\ell} \lambda_i (M_i(w, w_0) - 1) - \sum_{i=1}^{\ell} \xi_i (\lambda_i + \eta_i - C),$$

Необходимые условия седловой точки функции Лагранжа:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^{\ell} \lambda_i y_i x_i = 0 & \implies w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ \frac{\partial \mathcal{L}}{\partial w_0} = - \sum_{i=1}^{\ell} \lambda_i y_i = 0 & \implies \sum_{i=1}^{\ell} \lambda_i y_i = 0; \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = -\lambda_i - \eta_i + C = 0 & \implies \eta_i + \lambda_i = C, \quad i = 1, \dots, \ell. \end{aligned}$$

Понятия опорного вектора

Типизация объектов:

1. $\lambda_i = 0$; $\eta_i = C$; $\xi_i = 0$; $M_i \geq 1$.
— периферийные (неинформативные) объекты.
2. $0 < \lambda_i < C$; $0 < \eta_i < C$; $\xi_i = 0$; $M_i = 1$.
— **опорные** граничные объекты.
3. $\lambda_i = C$; $\eta_i = 0$; $\xi_i > 0$; $M_i < 1$.
— **опорные**-нарушители.

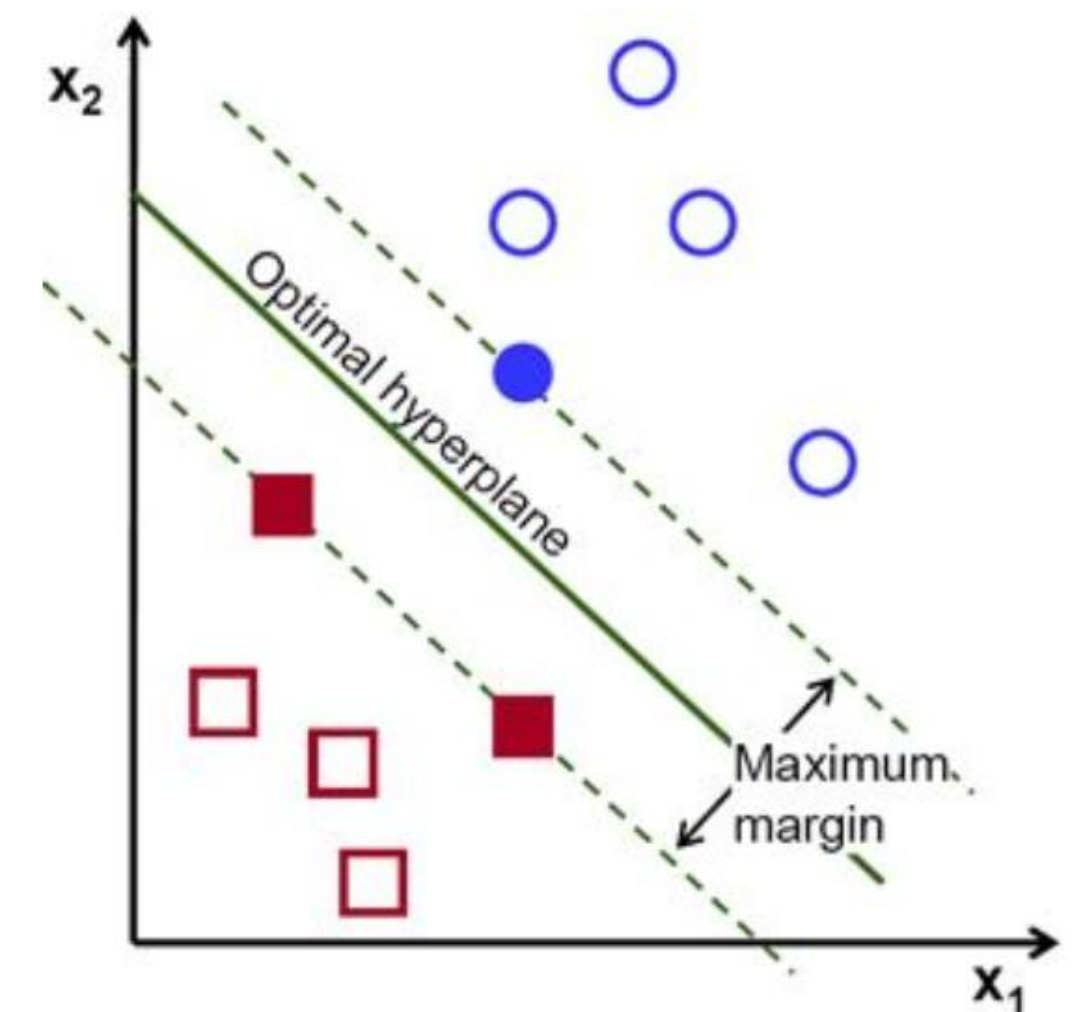
Определение

Объект x_i называется *опорным*, если $\lambda_i \neq 0$.

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i;$$

$$\sum_{i=1}^{\ell} \lambda_i y_i = 0;$$

$$\eta_i + \lambda_i = C, \quad i = 1, \dots, \ell.$$



Понятия опорного вектора

Решение не зависит от элемента выборки i

Типизация объектов:

1. $\lambda_i = 0$; $\eta_i = C$; $\xi_i = 0$; $M_i \geq 1$.
— периферийные (неинформативные) объекты.

2. $0 < \lambda_i < C$; $0 < \eta_i < C$; $\xi_i = 0$; $M_i = 1$.
— **опорные** граничные объекты.

3. $\lambda_i = C$; $\eta_i = 0$; $\xi_i > 0$; $M_i < 1$.
— **опорные**-нарушители.

Определение

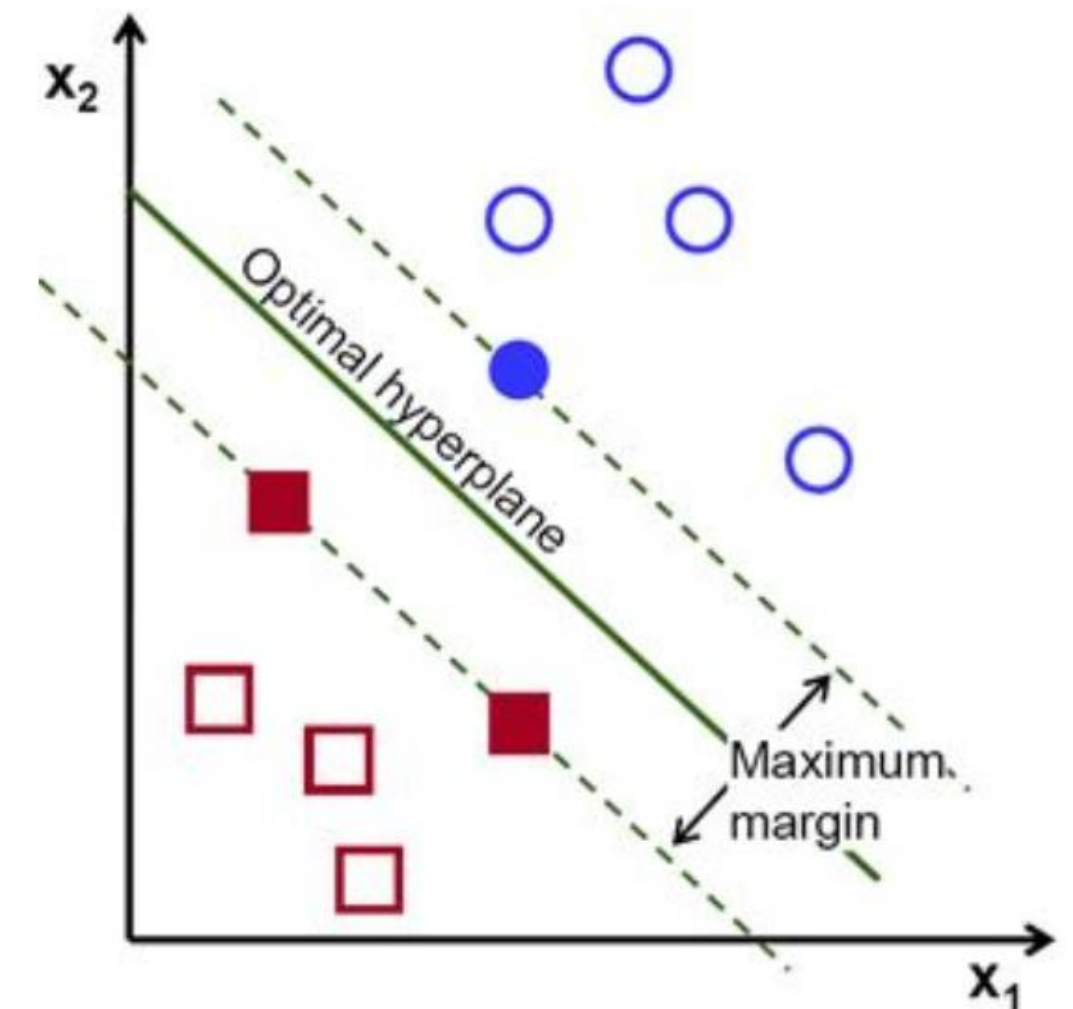
Объект x_i называется *опорным*, если $\lambda_i \neq 0$.

Класс находится за оптимальной гиперплоскостью класса. Нет ошибки
От него не зависит решение

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i;$$

$$\sum_{i=1}^{\ell} \lambda_i y_i = 0;$$

$$\eta_i + \lambda_i = C, \quad i = 1, \dots, \ell.$$



Понятия опорного вектора

Типизация объектов:

1. $\lambda_i = 0$; $\eta_i = C$; $\xi_i = 0$; $M_i \geq 1$.
— периферийные (неинформативные) объекты.
2. $0 < \lambda_i < C$; $0 < \eta_i < C$; $\xi_i = 0$; $M_i = 1$.
— **опорные** граничные объекты. ← Лежат ровно на границе
3. $\lambda_i = C$; $\eta_i = 0$; $\xi_i > 0$; $M_i < 1$.
— **опорные**-нарушители.

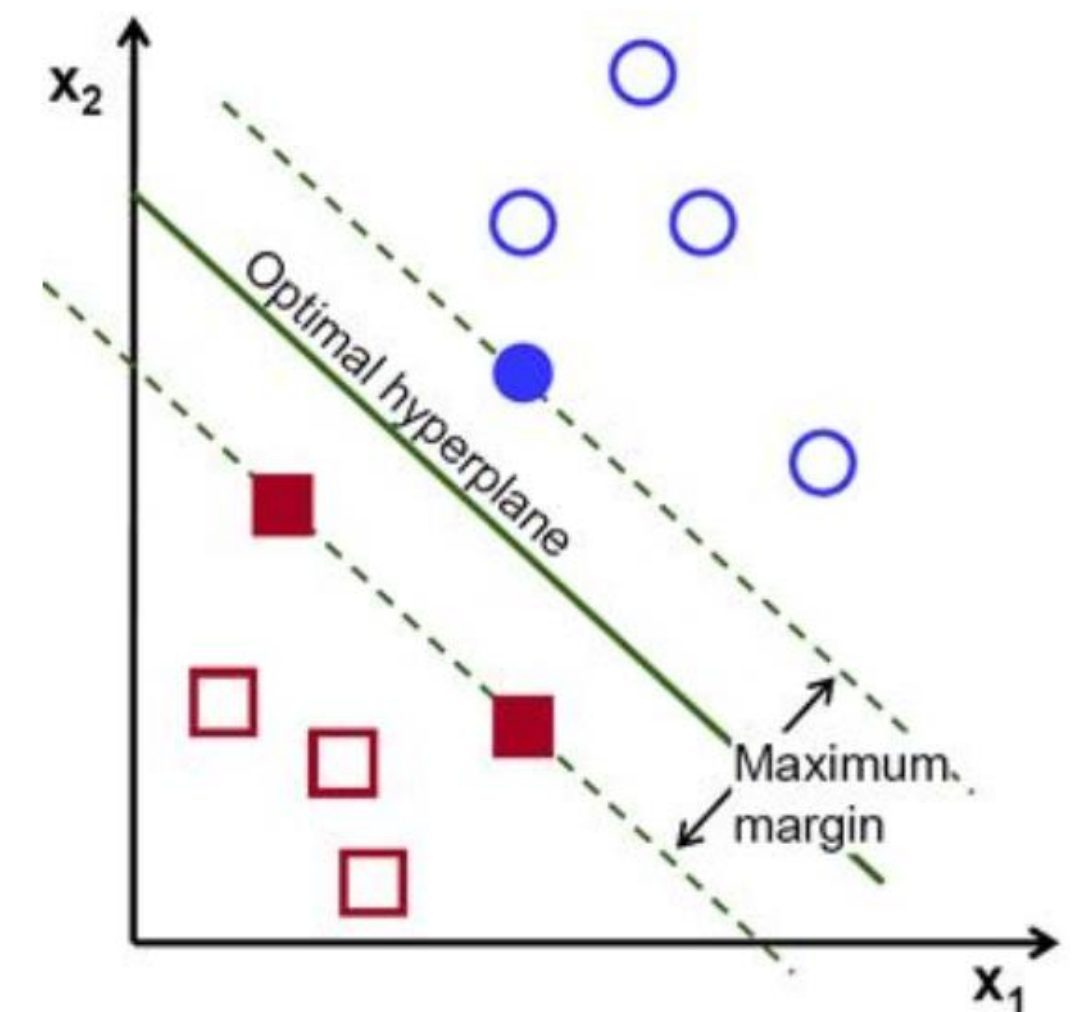
Определение

Объект x_i называется *опорным*, если $\lambda_i \neq 0$.

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i;$$

$$\sum_{i=1}^{\ell} \lambda_i y_i = 0;$$

$$\eta_i + \lambda_i = C, \quad i = 1, \dots, \ell.$$



Понятия опорного вектора

Типизация объектов:

1. $\lambda_i = 0$; $\eta_i = C$; $\xi_i = 0$; $M_i \geq 1$.
— периферийные (неинформативные) объекты.
2. $0 < \lambda_i < C$; $0 < \eta_i < C$; $\xi_i = 0$; $M_i = 1$.
— **опорные** граничные объекты.
3. $\lambda_i = C$; $\eta_i = 0$; $\xi_i > 0$; $M_i < 1$.
— **опорные**-нарушители.

Могут сидеть в «чужом» классе

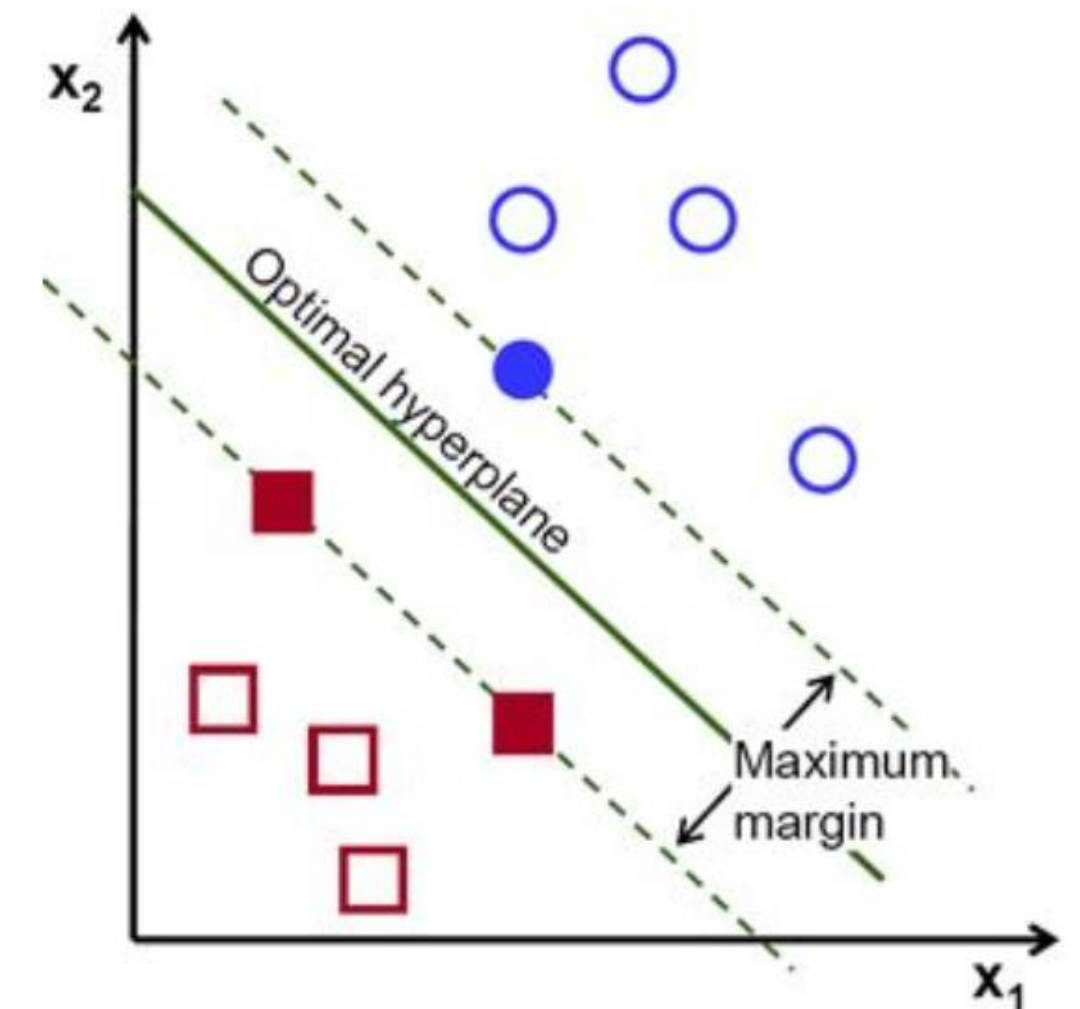
Определение

Объект x_i называется *опорным*, если $\lambda_i \neq 0$.

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i;$$

$$\sum_{i=1}^{\ell} \lambda_i y_i = 0;$$

$$\eta_i + \lambda_i = C, \quad i = 1, \dots, \ell.$$



Двойственная задача

Подставим все ограничения через w и перепишем Лагранжиан через λ_i

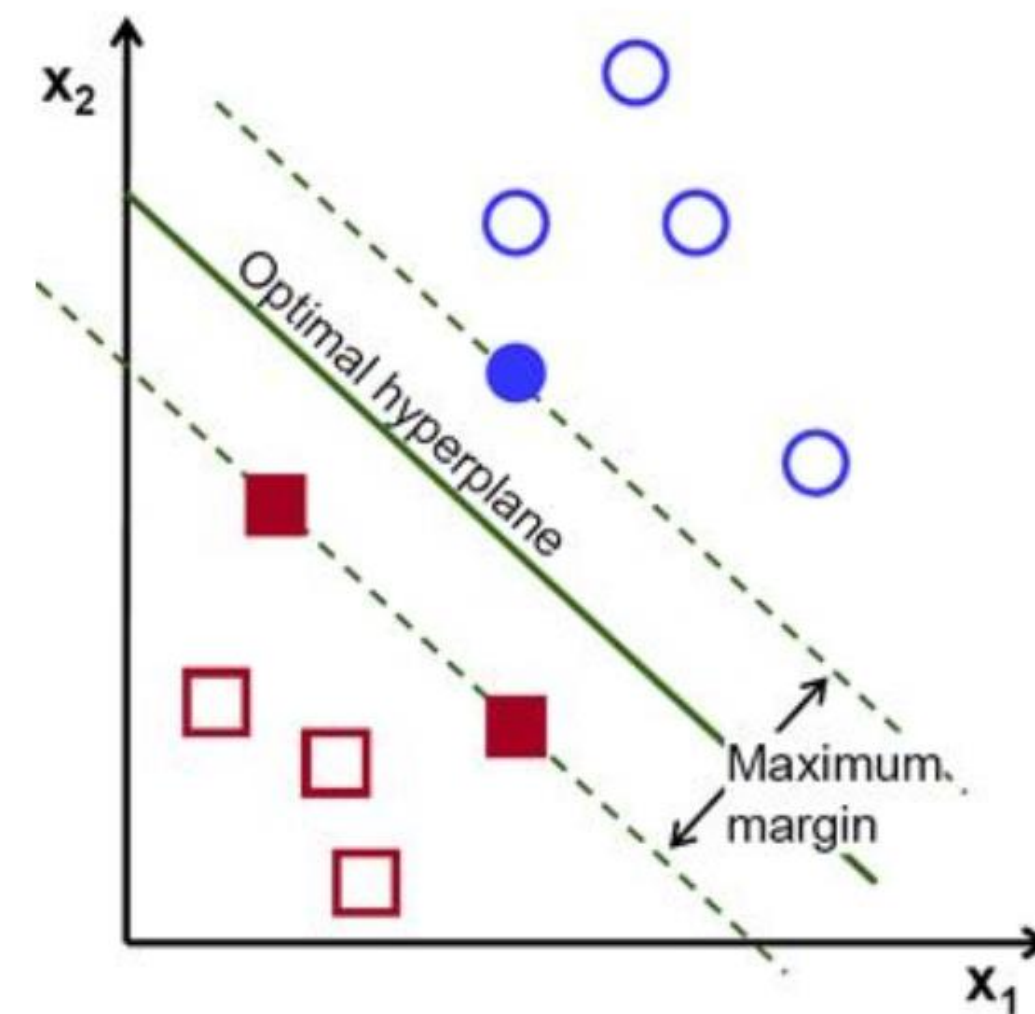
$$\begin{cases} -\mathcal{L}(\lambda) = -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Решение прямой задачи выражается через решение двойственной:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ w_0 = \langle w, x_i \rangle - y_i, \quad \text{для любого } i: \lambda_i > 0, \quad M_i = 1. \end{cases}$$

Линейный классификатор с признаками $f_i(x) = \langle x, x_i \rangle$:

$$a(x) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right).$$



Двойственная задача

Подставим все ограничения через w и перепишем Лагранжиан через λ_i

$$\begin{cases} -\mathcal{L}(\lambda) = -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

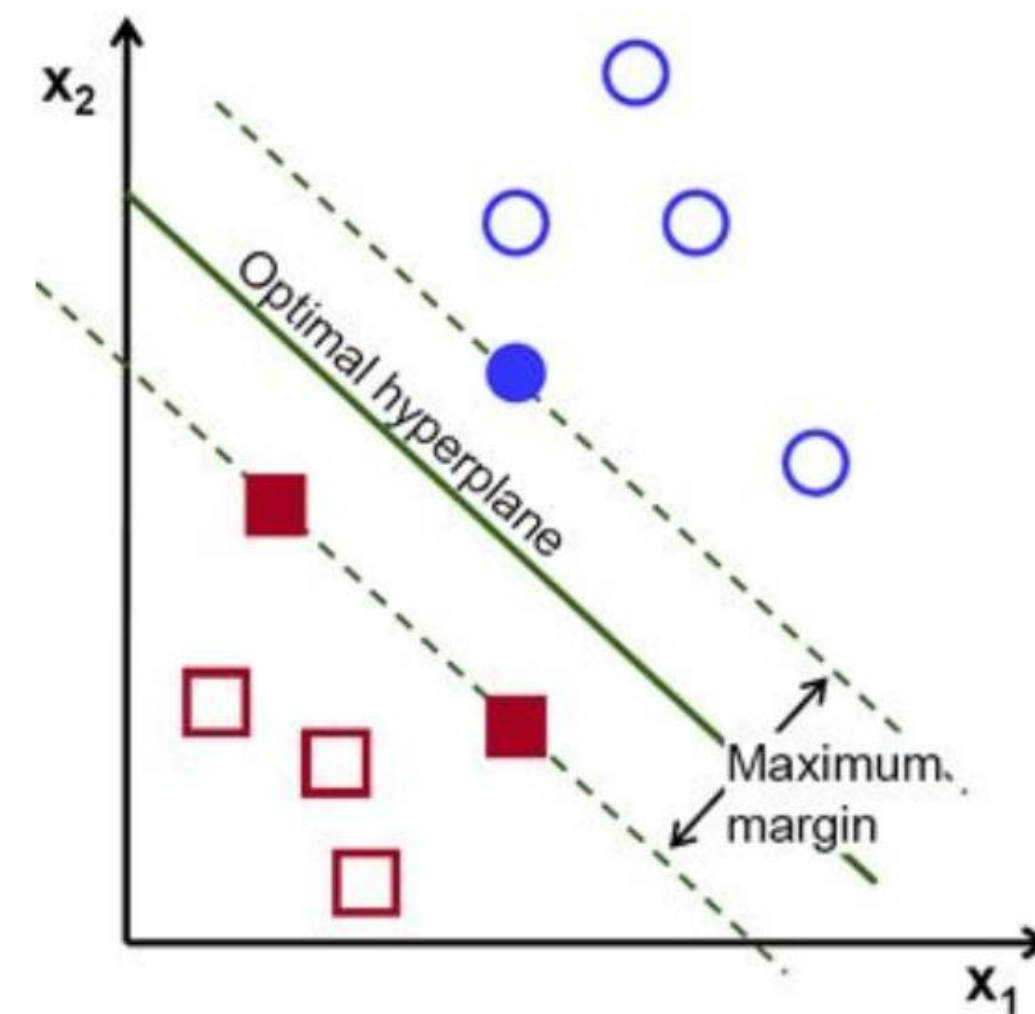
Матрица Грамма (скалярное произведение)

Решение прямой задачи выражается через решение двойственной:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ w_0 = \langle w, x_i \rangle - y_i, \quad \text{для любого } i: \lambda_i > 0, \quad M_i = 1. \end{cases}$$

Линейный классификатор с признаками $f_i(x) = \langle x, x_i \rangle$:

$$a(x) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right).$$



Двойственная задача

Подставим все ограничения через w и перепишем Лагранжиан через λ_i

$$\begin{cases} -\mathcal{L}(\lambda) = -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Матрица Грамма (скалярное произведение)

Куб с ребром C

Выпуклая задача

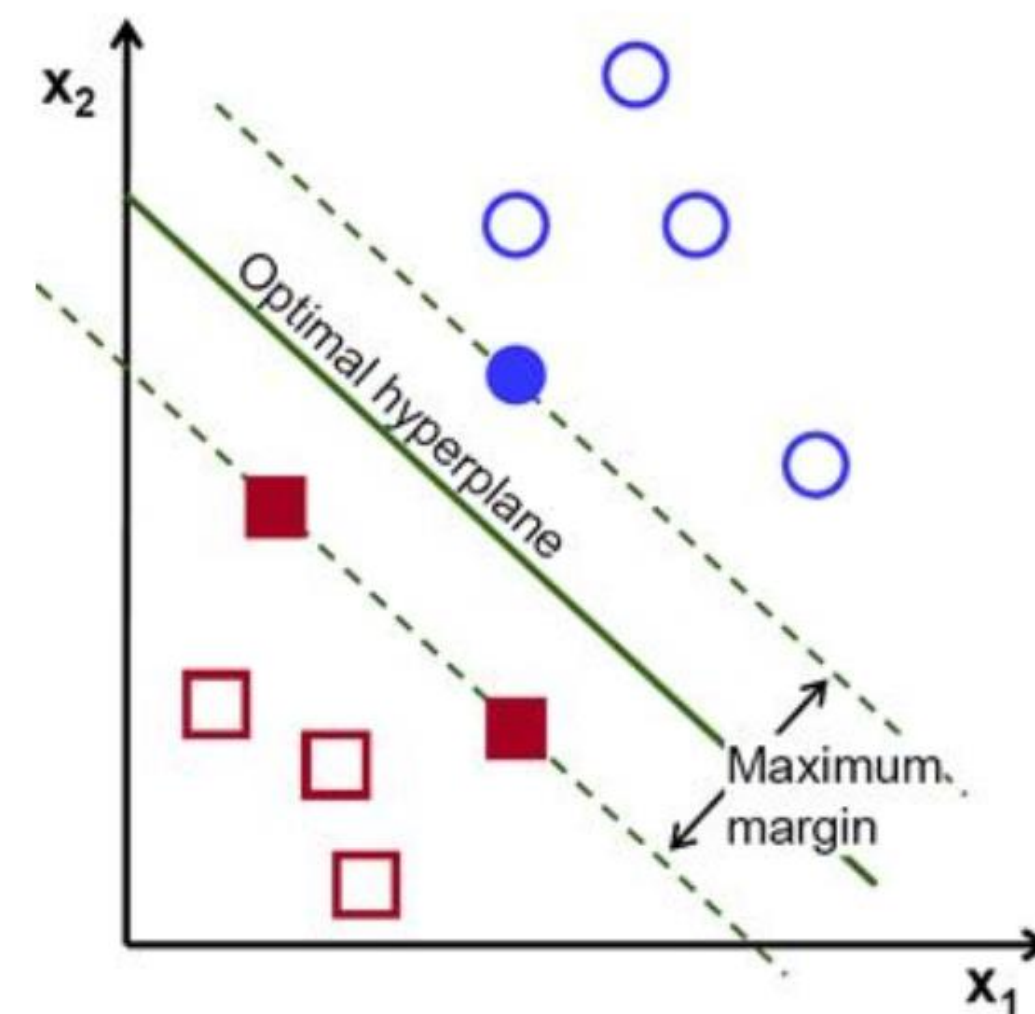
Гиперплоскость

Решение прямой задачи выражается через решение двойственной:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ w_0 = \langle w, x_i \rangle - y_i, \quad \text{для любого } i: \lambda_i > 0, \quad M_i = 1. \end{cases}$$

Линейный классификатор с признаками $f_i(x) = \langle x, x_i \rangle$:

$$a(x) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right).$$



Двойственная задача

Подставим все ограничения через w и перепишем Лагранжиан через λ_i

$$\begin{cases} -\mathcal{L}(\lambda) = -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Матрица Грамма (скалярное произведение)

Куб с ребром C

Выпуклая задача

Гиперплоскость

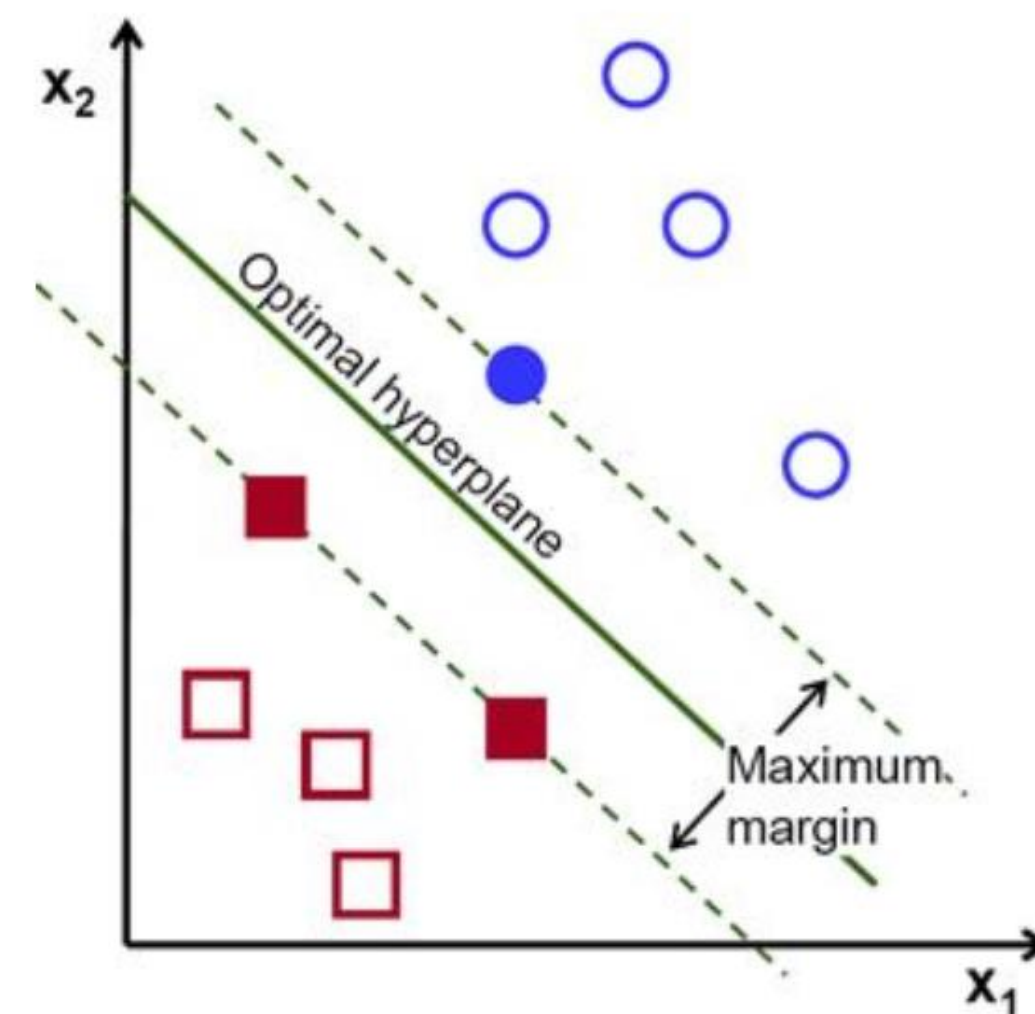
Решение прямой задачи выражается через решение двойственной:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ w_0 = \langle w, x_i \rangle - y_i, \quad \text{для любого } i: \lambda_i > 0, \quad M_i = 1. \end{cases}$$

На практике опорных векторов не так много

Линейный классификатор с признаками $f_i(x) = \langle x, x_i \rangle$:

$$a(x) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right).$$



Двойственная задача

Подставим все ограничения через w и перепишем Лагранжиан через λ_i

$$\begin{cases} -\mathcal{L}(\lambda) = -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Матрица Грамма (скалярное произведение)

Куб с ребром C

Выпуклая задача

Гиперплоскость

Решение прямой задачи выражается через решение двойственной:

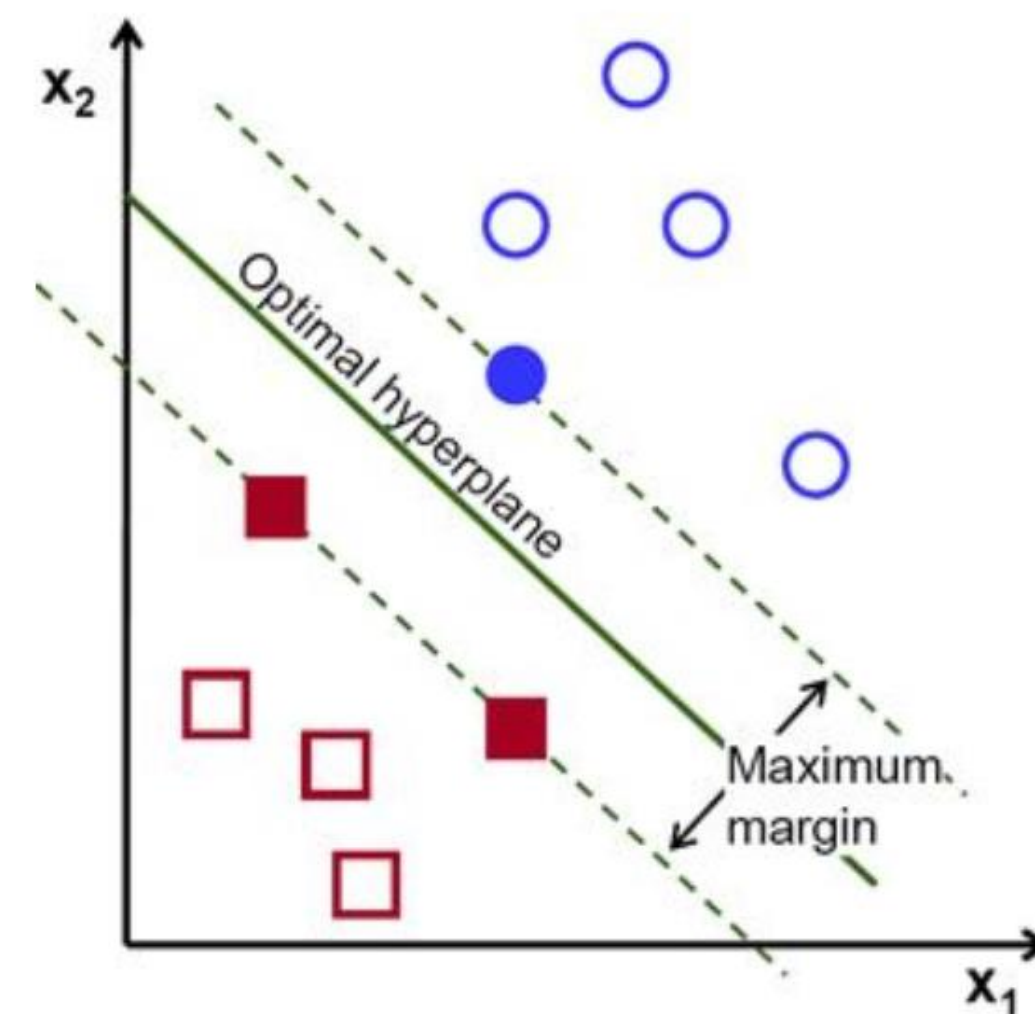
$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i; \\ w_0 = \langle w, x_i \rangle - y_i, \quad \text{для любого } i: \lambda_i > 0, \quad M_i = 1. \end{cases}$$

На практике опорных векторов не так много

Линейный классификатор с признаками $f_i(x) = \langle x, x_i \rangle$:

Задача линейна относительно λ_i \longrightarrow $a(x) = \text{sign} \left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0 \right).$

Близость объекта до «опорного»



Трюк ядра

Идея: заменить $\langle x, x' \rangle$ нелинейной функцией $K(x, x')$.

Переход к спрямляющему пространству,
как правило, более высокой размерности: $\psi: X \rightarrow H$.

Определение

Функция $K: X \times X \rightarrow \mathbb{R}$ — ядро, если $K(x, x') = \langle \psi(x), \psi(x') \rangle$ при некотором $\psi: X \rightarrow H$, где H — гильбертово пространство.

Теорема

Функция $K(x, x')$ является ядром тогда и только тогда, когда она симметрична: $K(x, x') = K(x', x)$; и неотрицательно определена:

$$\int_X \int_X K(x, x') g(x) g(x') dx dx' \geq 0 \text{ для любой } g: X \rightarrow \mathbb{R}.$$

Трюк ядра

$$a(x) = \text{sign}\left(\sum_{i=1}^{\ell} \lambda_i y_i \langle x, x_i \rangle - w_0\right).$$

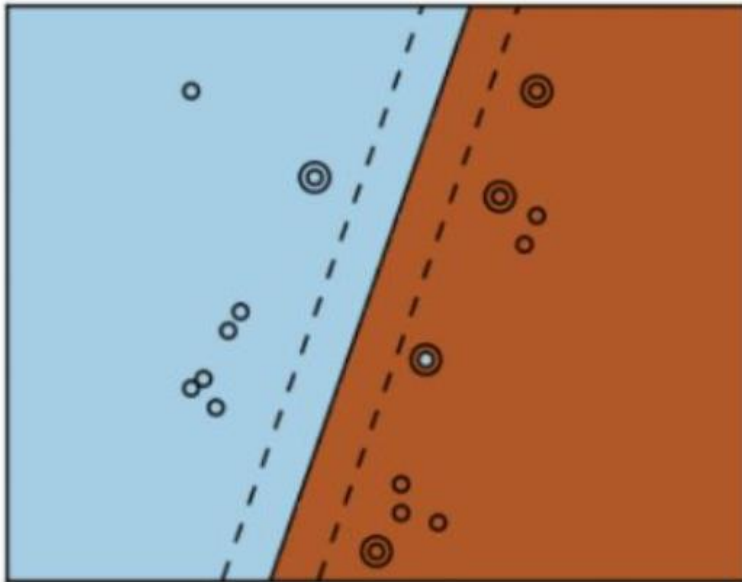
$K(x, x') = \langle x, x' \rangle^2$ - quadratic

$K(x, x') = \langle x, x' \rangle^d$ - polynomial with degree d

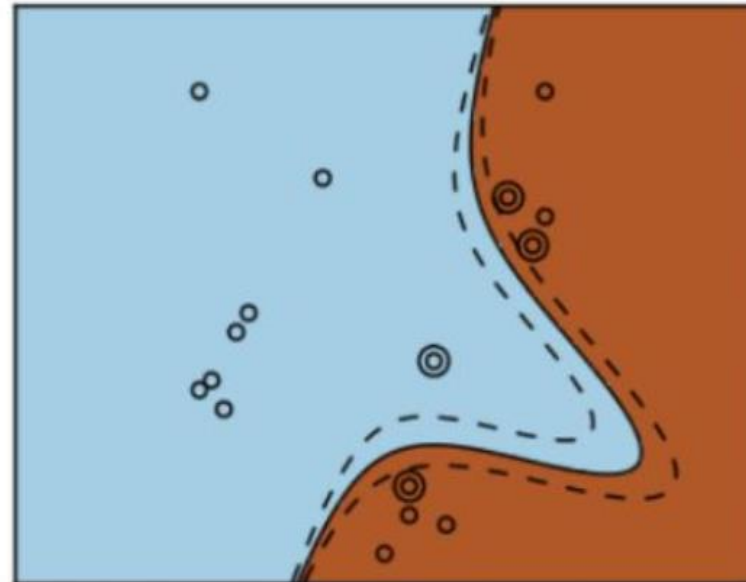
$K(x, x') = (\langle x, x' \rangle + 1)^d$ - polynomial with degree $\leq d$

$K(x, x') = \exp(-\gamma \|x - x'\|^2)$ - Radial Basis Functions (RBF) kernel

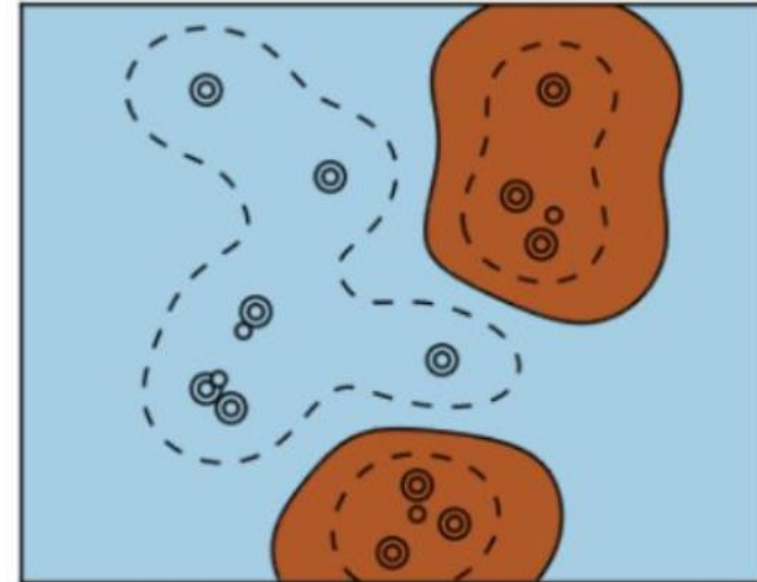
$\langle x, x' \rangle$



$(\langle x, x' \rangle + 1)^d, d=3$



$\exp(-\gamma \|x - x'\|^2)$



02

PCA

Задача понижения размерности

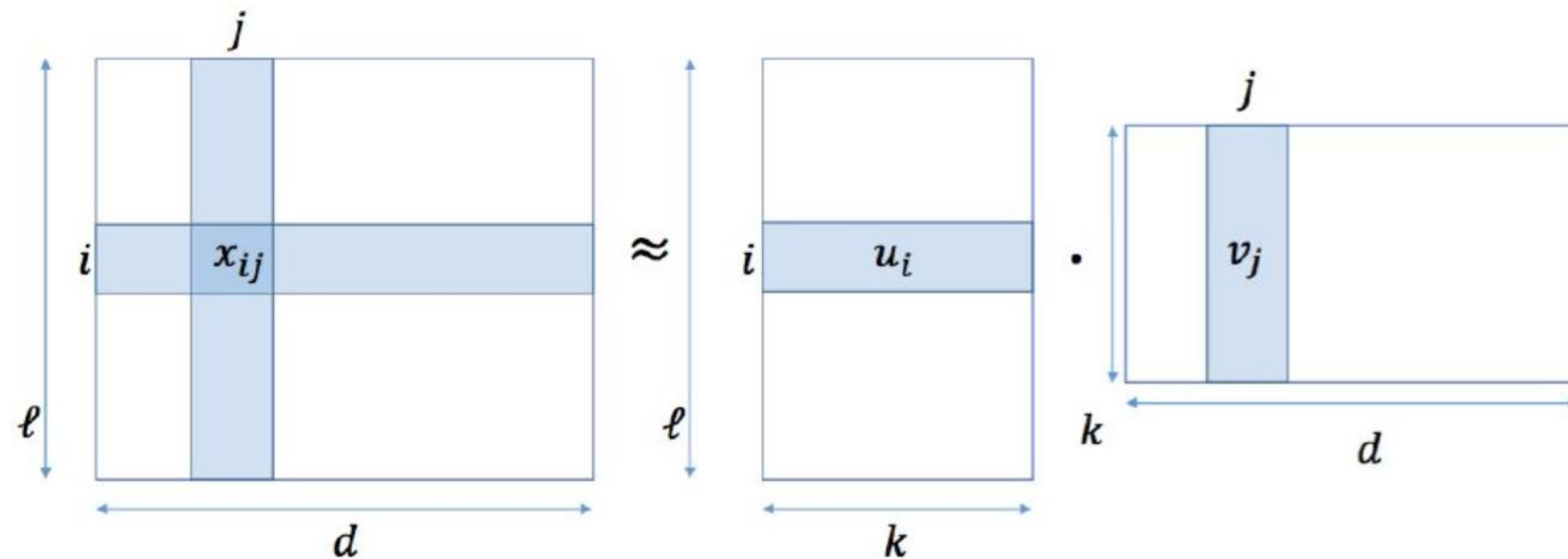
- В ML мы часто работаем с многомерными данными.
- Сотни или тысячи функций
- Их трудно визуализировать
- Обучение занимает много времени
- Некоторые модели хуже справляются с многомерными разреженными входными данными

Композиция матриц

Разложение на множители матриц меньшего ранга

$$X_{l,d} \approx U_{l,k} \cdot V_{k,d}^T$$

$$\|X - UV^T\| \rightarrow \min$$

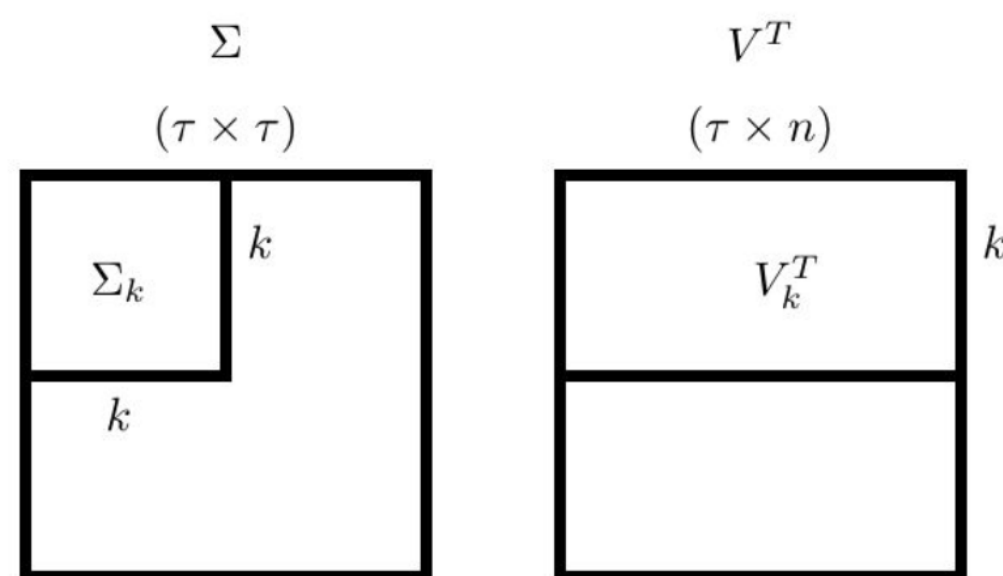
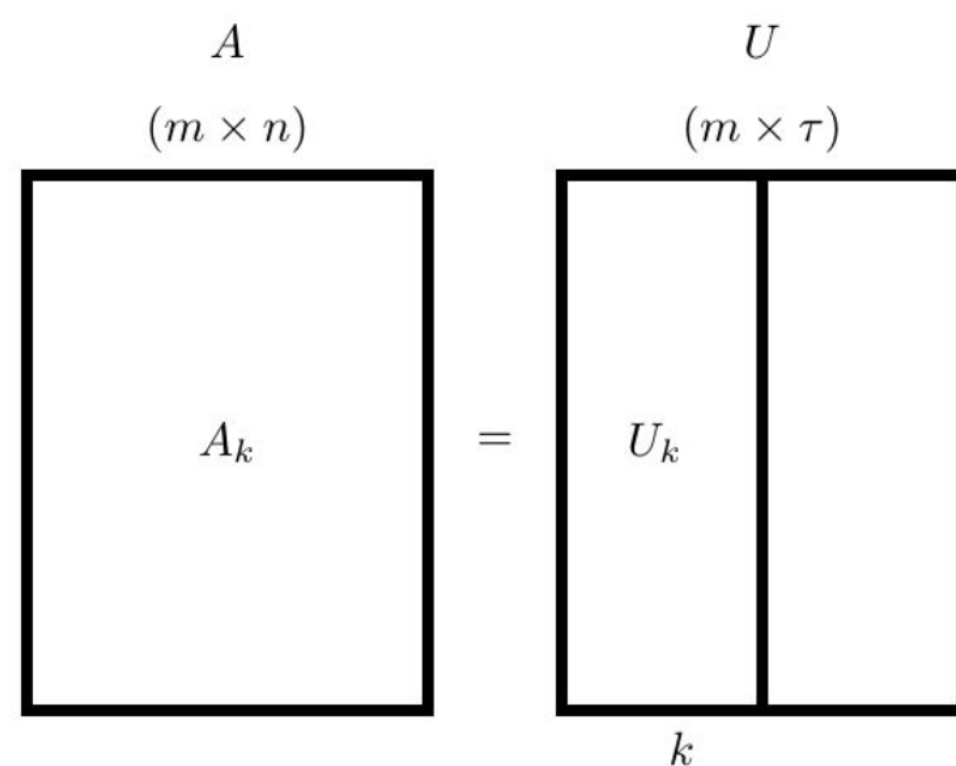


Декомпозиция матриц SVD

Сингулярное разложение (Singular value decomposition)

$$A = U \Sigma V^T$$

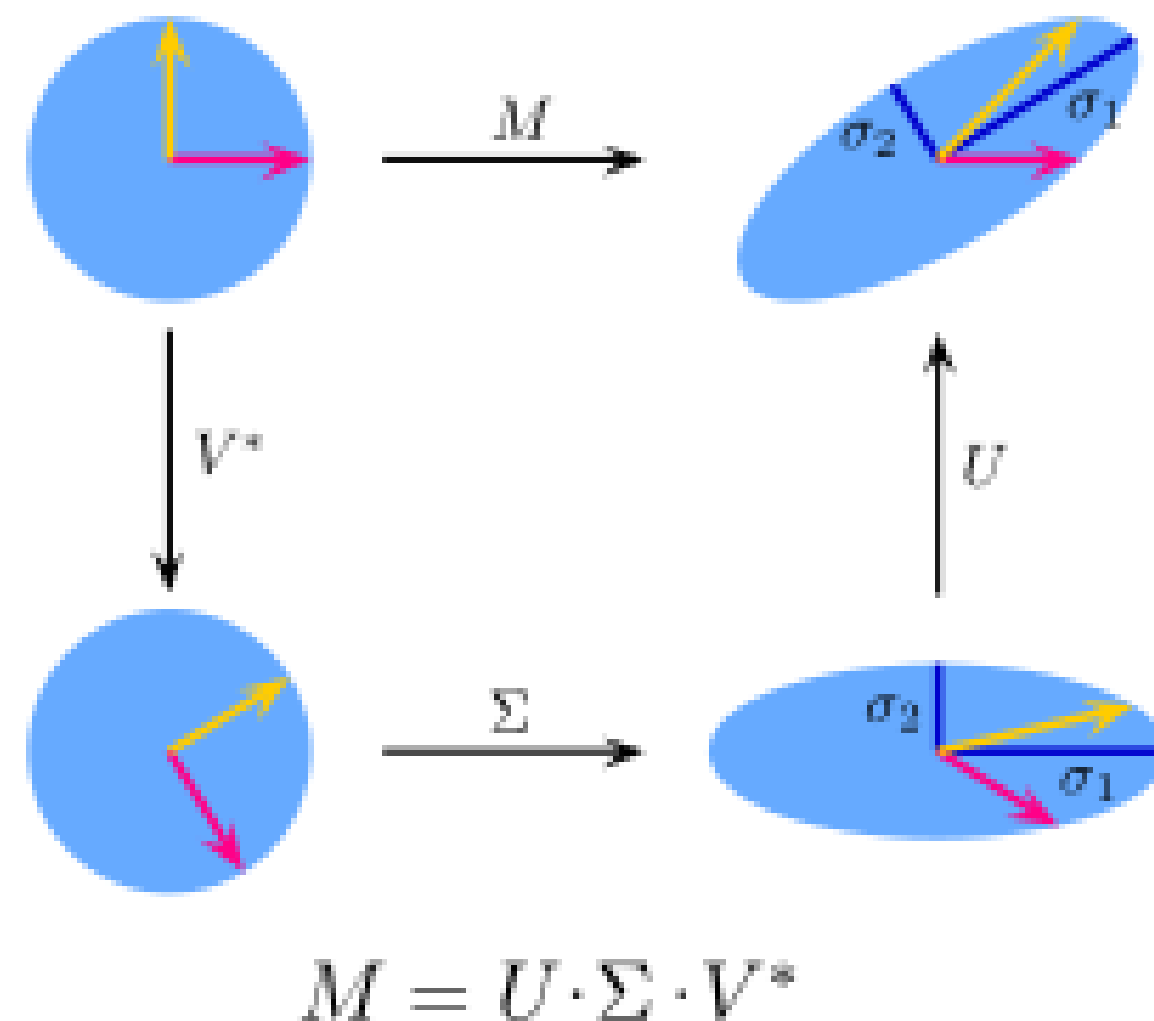
$$A_k = U_k \Sigma_k V_k^T = (U_k \Sigma_k) V_k^T = U_k (\Sigma_k V_k^T)$$



$$U \in \mathbb{R}^{m \times r}; U U^T = I$$

$$V \in \mathbb{R}^{n \times r}; V V^T = I$$

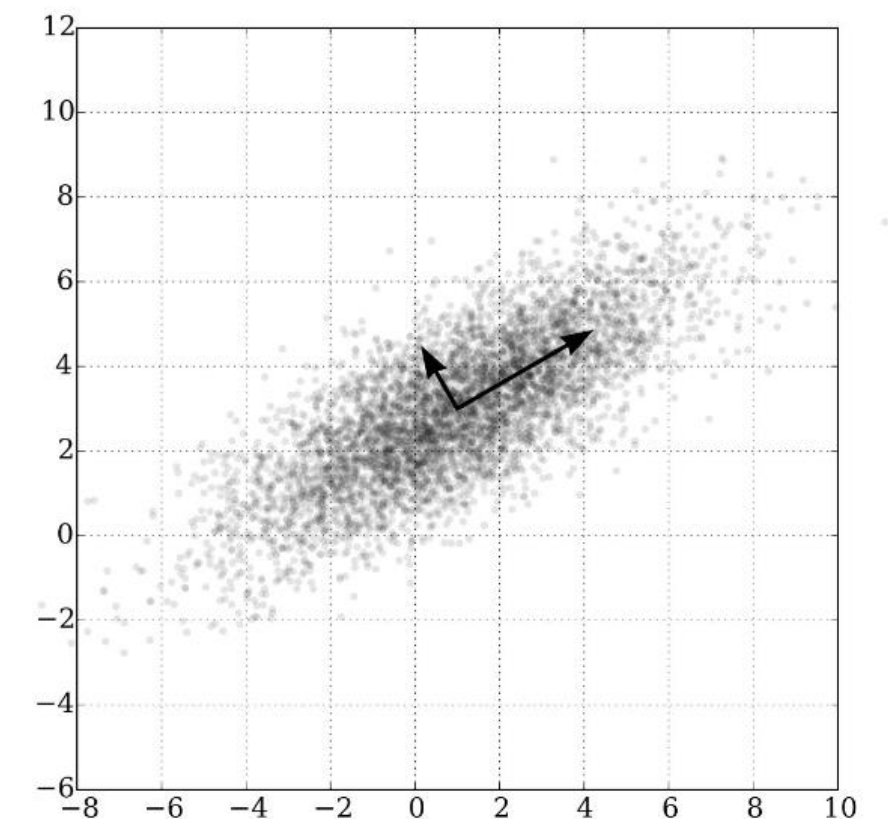
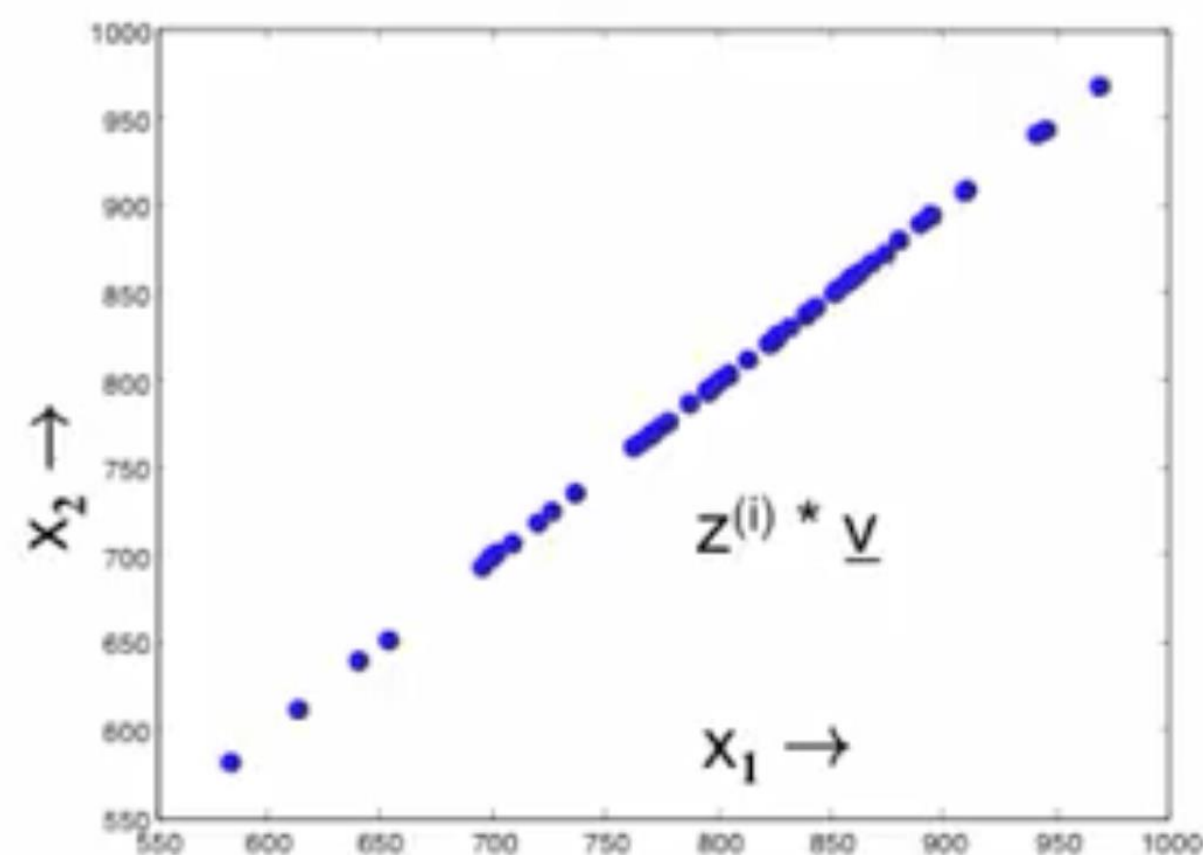
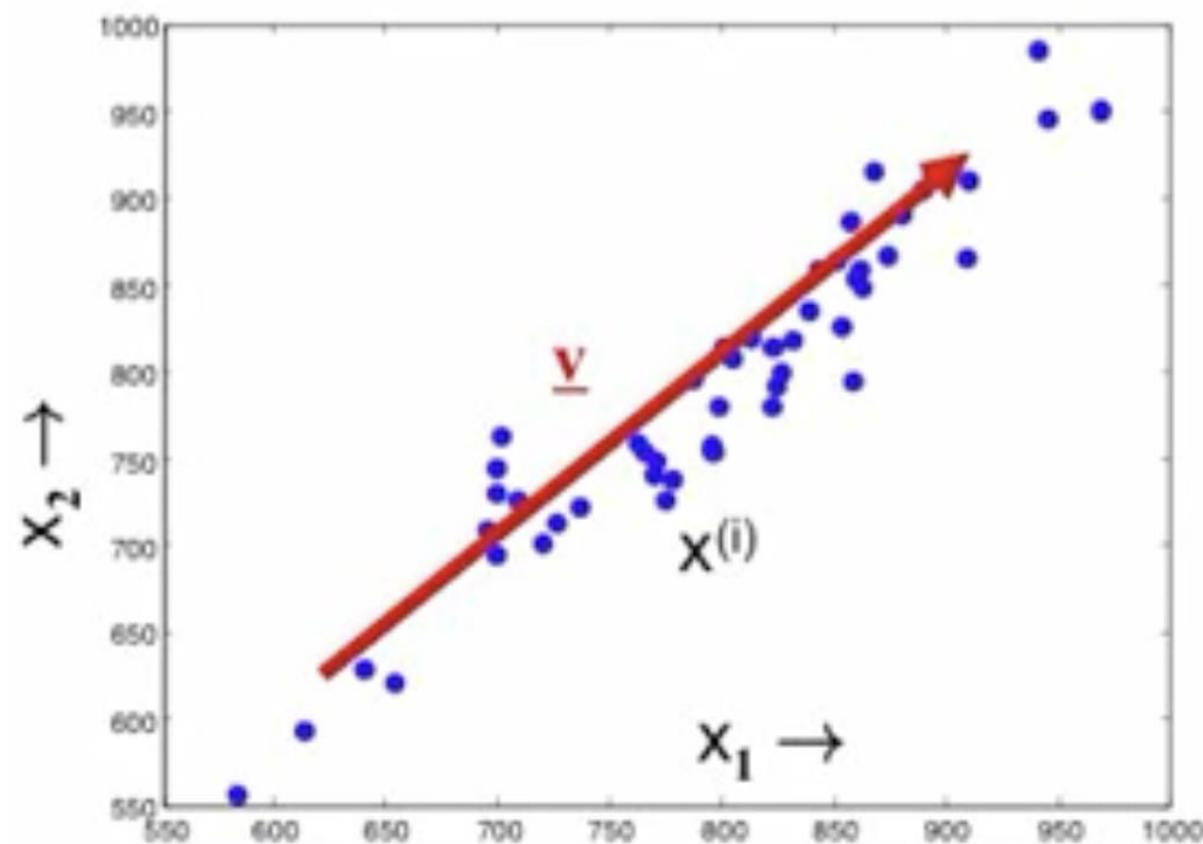
$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r); r = \text{rank}(M)$$



Principal component analysis PCA

Построим новый базис $f(z)$: $[x_1, x_2, \dots, x_n] = z * v = z * [v_1, v_2, \dots, v_n]$
где v будет направлением максимальной дисперсии по x_1, x_2

Выберем новые оси дисперсия по осям была максимальной. Это эквивалентно поиску ортогонального базиса в пространстве дисперсий, то есть решению $XX^T e = \lambda e$
где e – ортогональный базис и λ – собственные значения матрицы ковариации XX^T



Алгоритм PCA

PCA

Large Table

X1	X2	X3	X4	X5
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*

Covariance matrix

*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*
*	*	*	*	*

Eigenstuff

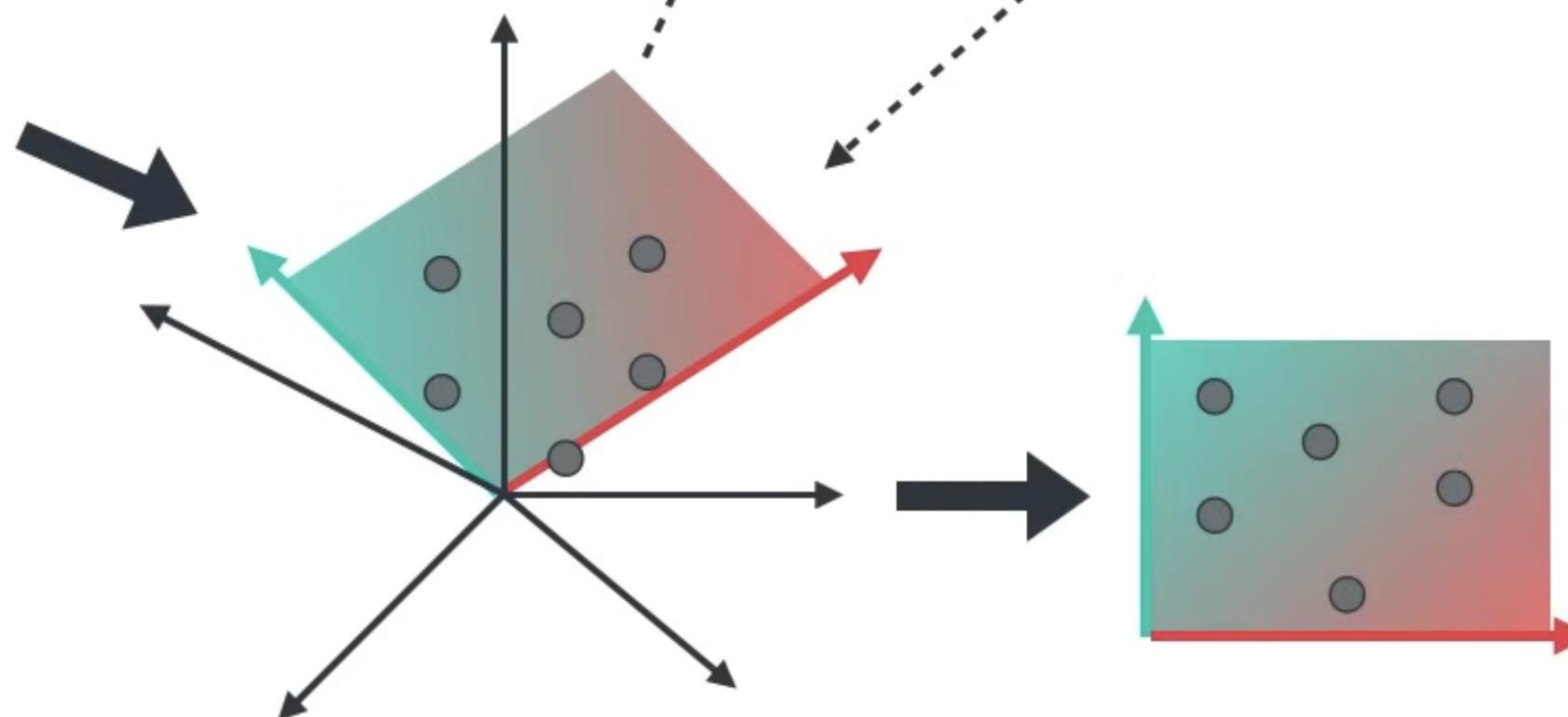
V_1	λ_1
V_2	λ_2
V_3	λ_3
V_4	λ_4
V_5	λ_5

Big

Small

Отсортированы
по значению
собственного
числа

Top@2 – самых
значимых
компонент



5D Plot

2D Plot

SVD и PCA

Factorization of Data Matrix F :

$$F_{N \times M} = U_{N \times N} \cdot \Sigma_{N \times M} \cdot V^T_{M \times M}$$

Covariance Matrix: $R = FF^T$

$$R = U \cdot \Sigma \cdot \underbrace{V^T \cdot V}_I \cdot \Sigma^T \cdot U^T$$

$$R = U \cdot \Lambda \cdot U^T$$

Factorization of Data Matrix F :

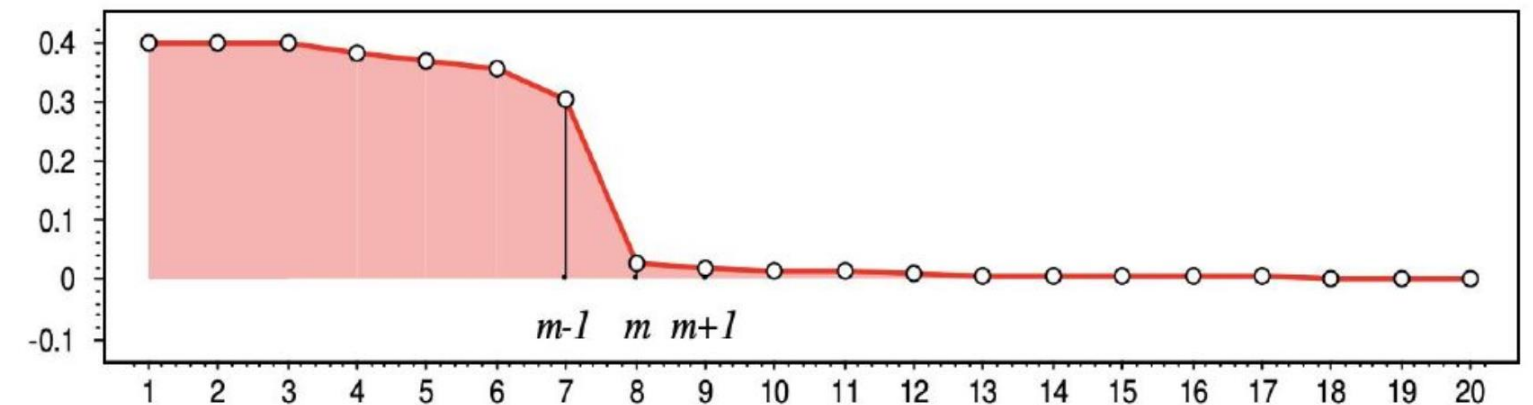
$$R_{N \times N} = U_{N \times N} \cdot \Lambda_{N \times N} \cdot U^T_{N \times N}$$

Eigenvalues:

$$\Lambda_{N \times N} = \begin{bmatrix} \lambda_1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & \lambda_3 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \lambda_4 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & \lambda_K & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & \lambda_{K+1} & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \lambda_N \end{bmatrix} \quad \lambda_i = \sigma_i^2$$

Понижение размерности PCA

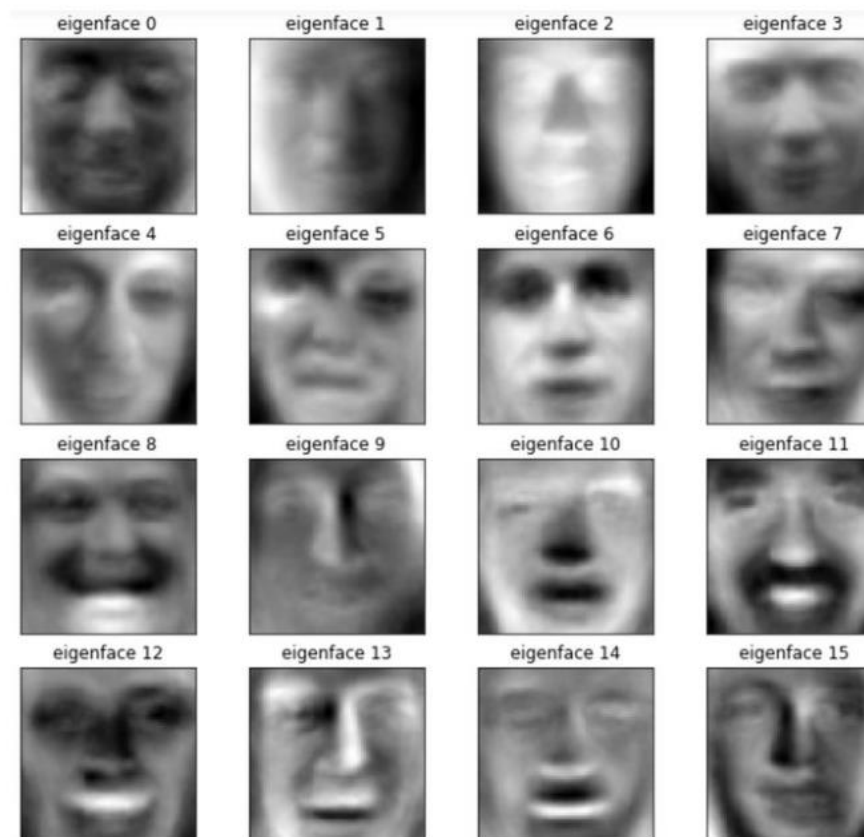
- Часто данные содержат помехи и неинформативны
- Удалим компоненты с низкой дисперсией в PCA



$$E_m = \frac{\|GU^T - F\|^2}{\|F\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\lambda_1 + \dots + \lambda_n} \leq \varepsilon.$$

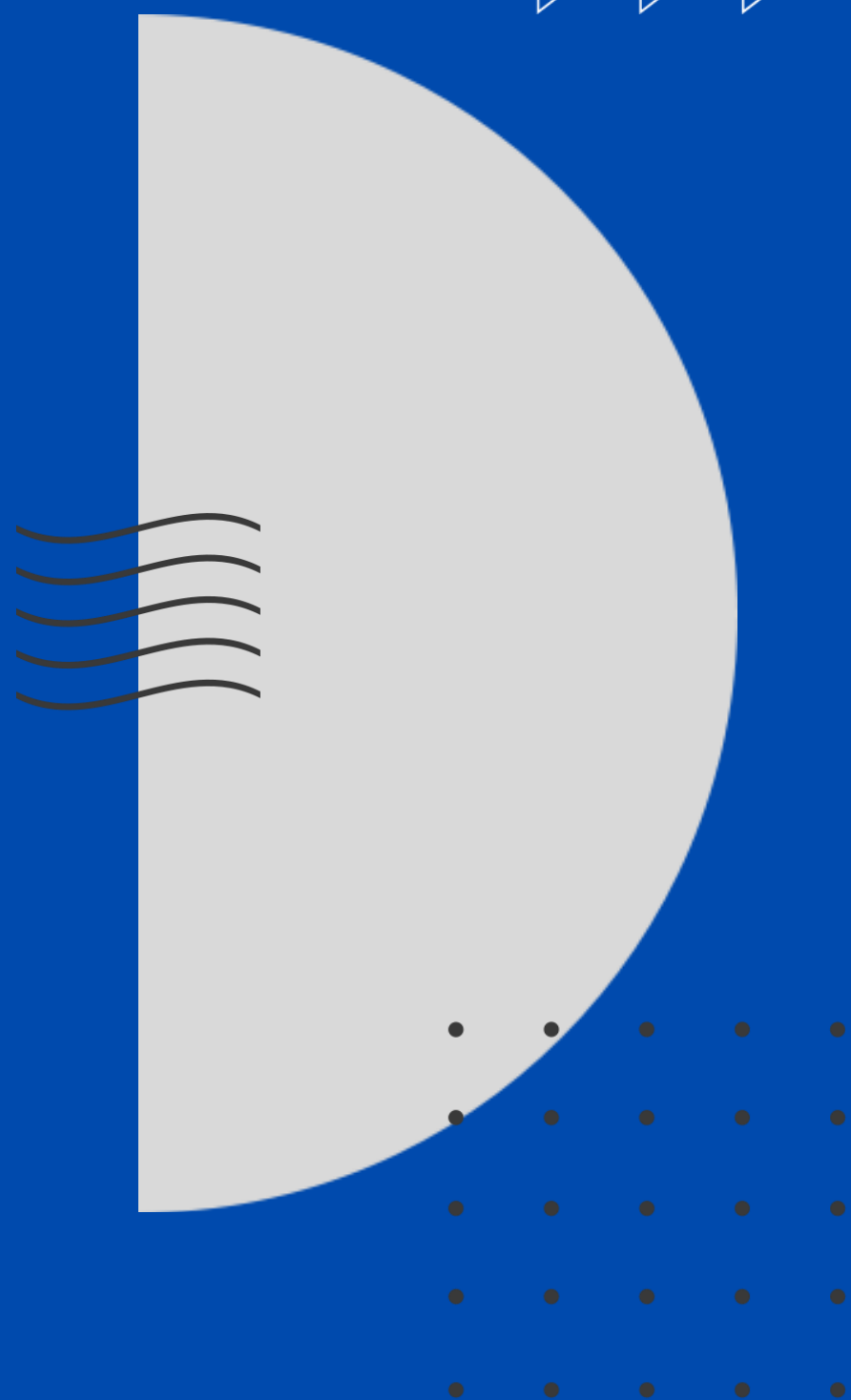


Top16



Top50





Место для ваших
вопросов