

# Математика полносвязной сети: геометрия, интерпретируемость, SLAP атака

---

Перминов Андрей

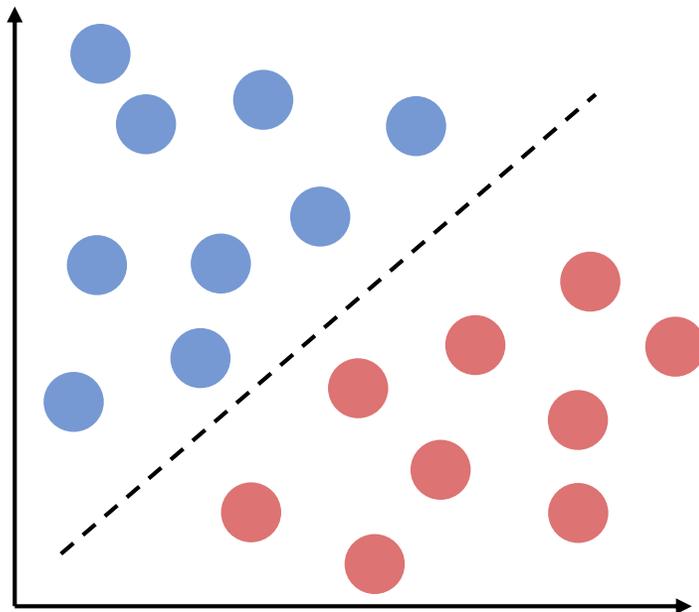
01 октября 2025



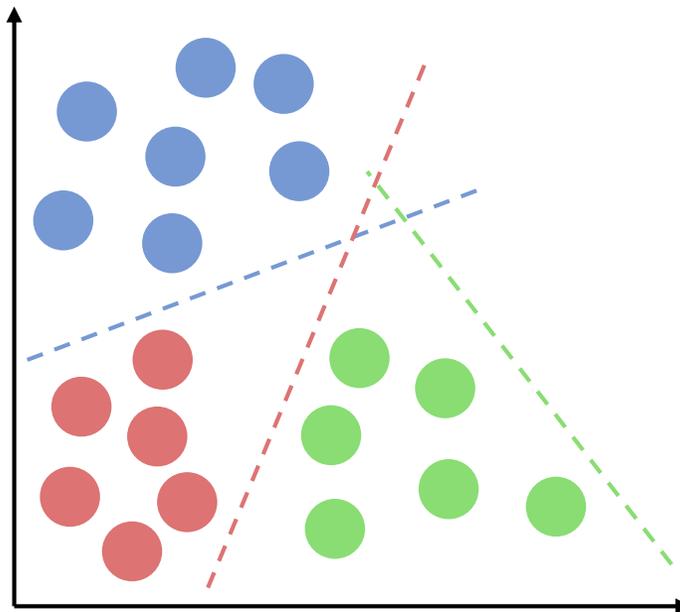
# Задача классификации

**Классификация** – задача, в которой множество **объектов** необходимо **разделить** некоторым образом на **классы**.

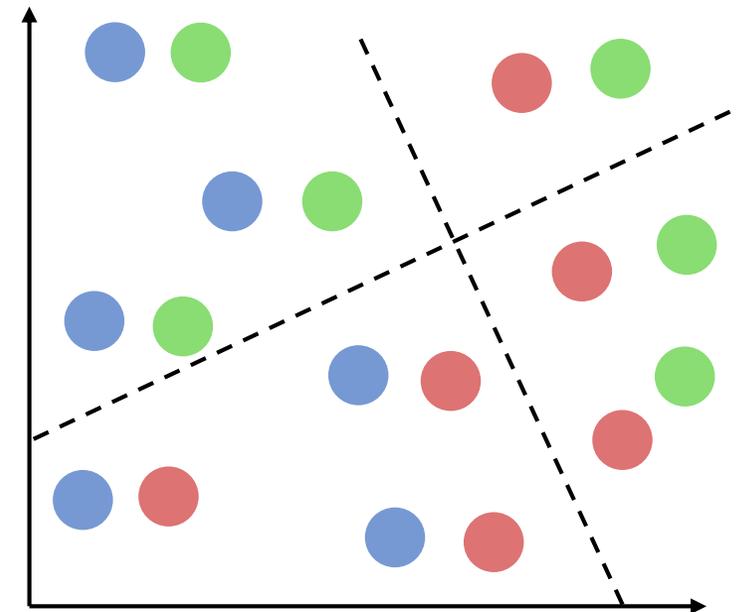
**Выборка** – конечное множество объектов, для которых известно, к каким классам они относятся.



бинарная



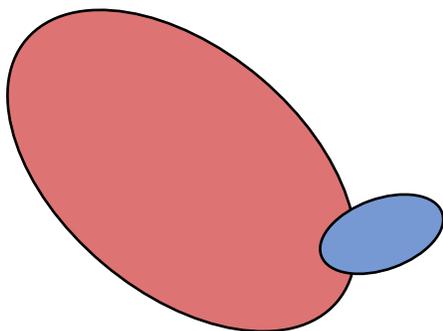
многоклассовая



многометочная

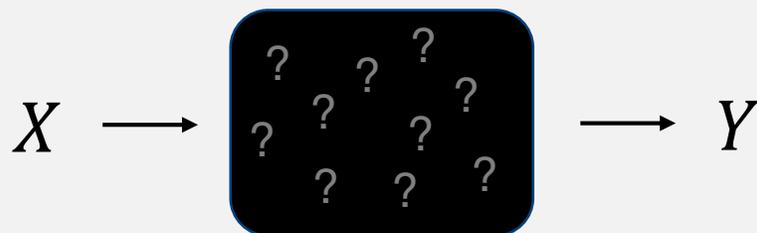
## Дисбаланс классов

- Существенное **различие** априорных вероятностей.
- Деградация точности при классификации **маломощных** классов.
- **Критично** в медицине, финансах, промышленности (редкие события/дефекты).



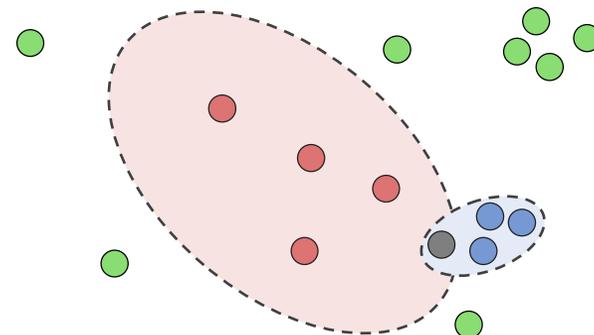
## Интерпретируемость и доверие к модели

- Большинство моделей работают как **чёрный ящик**.
- Снижение **доверия** к модели из-за невозможности анализа принимаемых решений.
- В отличие от деревьев решений, нейросетевые модели **являются плохо интерпретируемыми**.



## Неопределённости в данных

- Несоответствие наблюдений **носителю** распределения и **двусмысленность** в данных.
- Отсутствие **статистической** основы для классификации.
- Необходимость внедрения **отказа** от классификации.



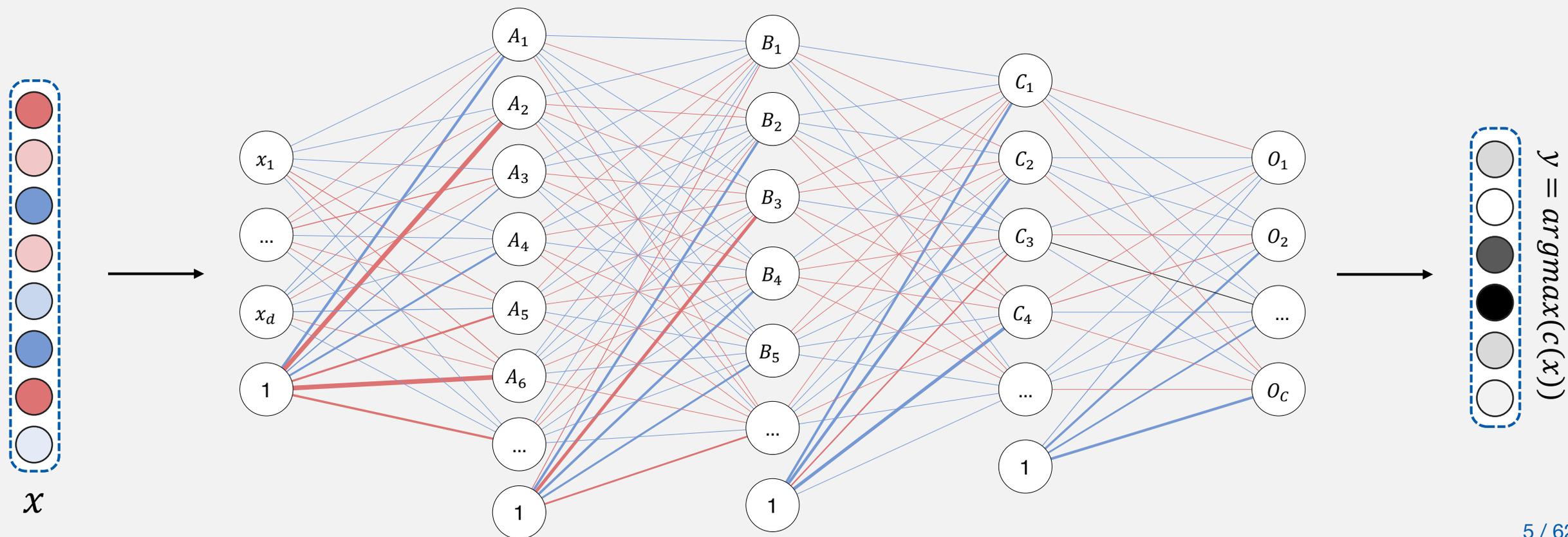


# Многослойный персептрон

Равномерно непрерывная функция  $c(x): R^d \rightarrow R^C$  вида:

$$c(x) = \sigma(W_0 \sigma(\dots \sigma(W_1 x + b_1) \dots) + b_0),$$

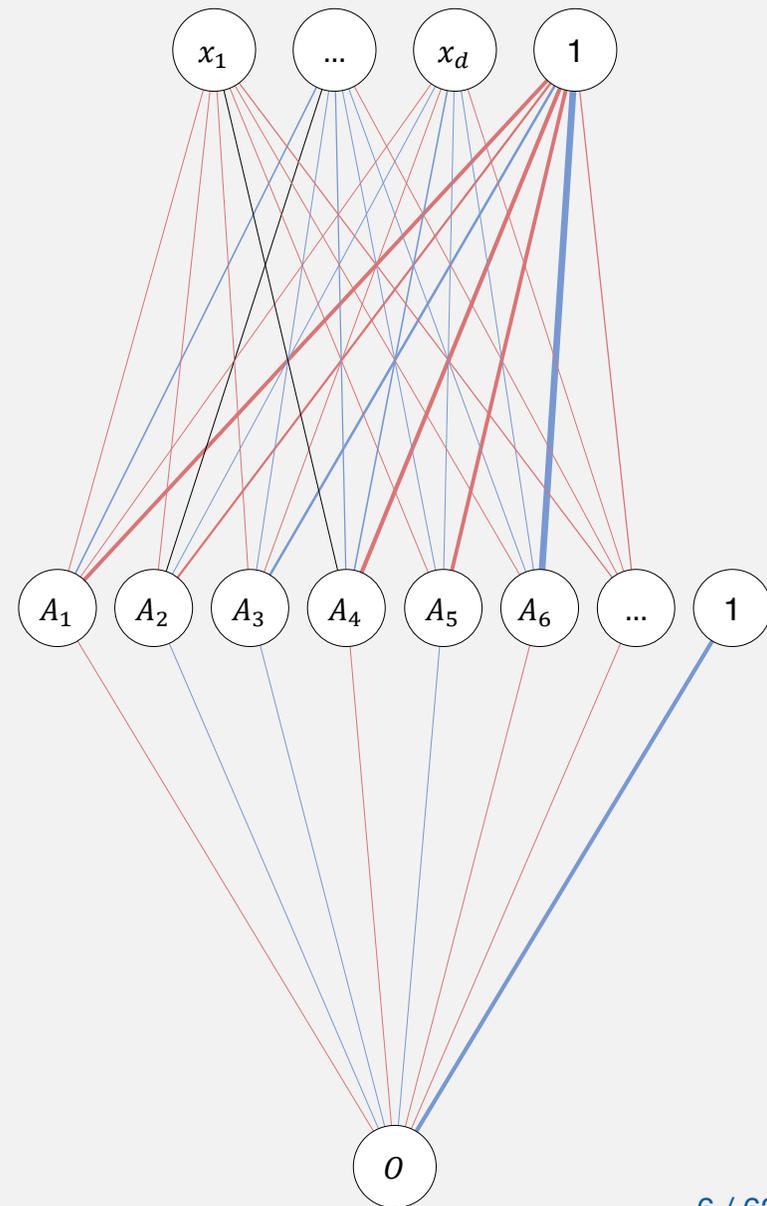
где  $\sigma$  – **нелинейная** функция активации.



# Основная аппроксимационная теорема

Для любой **непрерывной** функции  $f: [a, b]^d \rightarrow R$  и любого  $\varepsilon > 0$  существует такая полносвязная нейронная сеть  $c(x)$  с **одним скрытым слоем** и нелинейной функцией активации  $\sigma$ , что

$$\sup_{x \in [a, b]^d} |f(x) - c(x)| < \varepsilon.$$

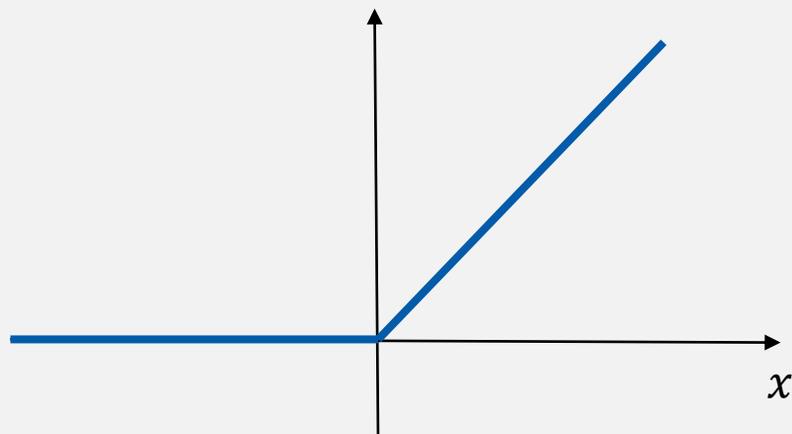


Следует ли из этой теоремы существование персептрона с более чем одним скрытым слоем?

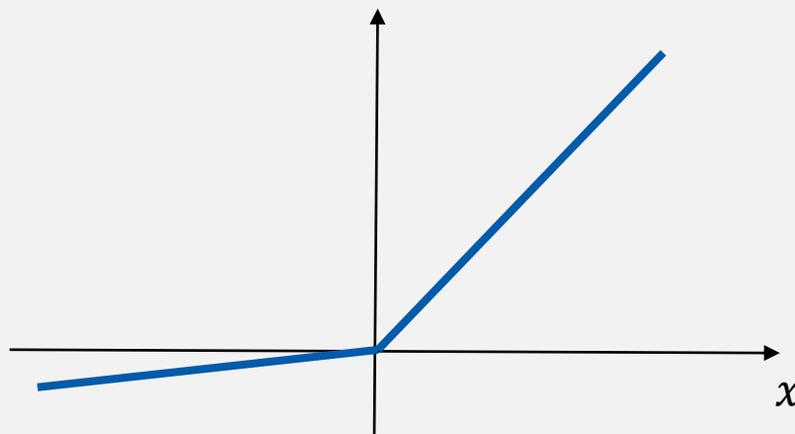
# Кусочно-линейные функции активации

**Кусочно-линейная функция** – функция, линейная на каждом из интервалов, составляющих область определения.

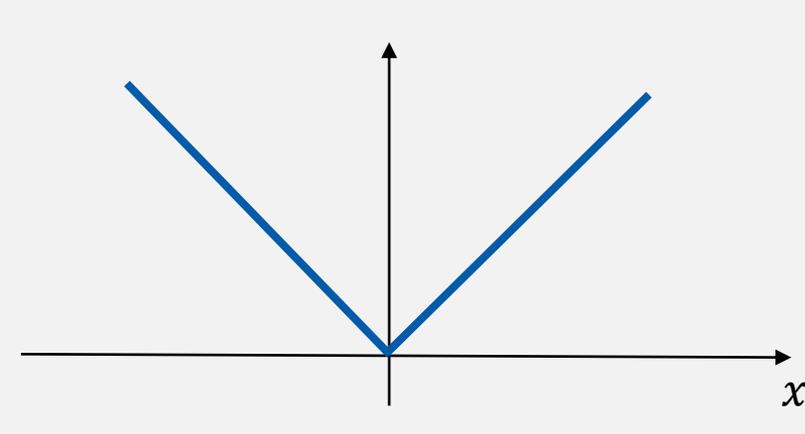
Активно используются в нейросетевых моделях благодаря **быстрому вычислению** значений.



*ReLU*( $x$ ):  
 $\max(0, x)$



*LeakyReLU*( $x$ ):  
 $\max(0.1 \cdot x, x)$



*Abs*( $x$ ):  
 $|x|$

# Геометрическая интерпретация: первый слой

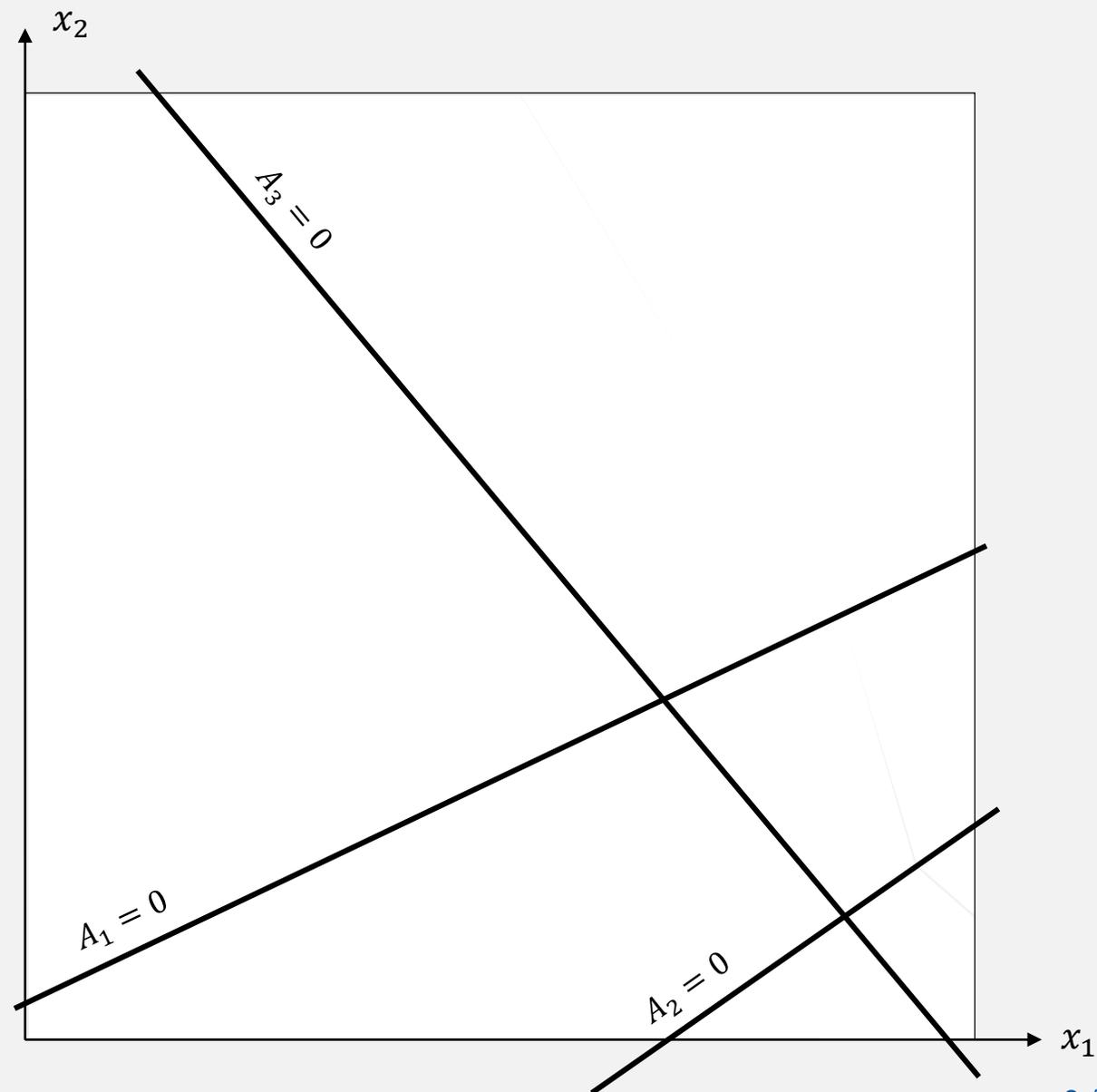
$$A_1 = -0.4x_1 + 0.8x_2$$

$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

**Первый слой** персептрона разбивает входной компакт на непересекающиеся ячейки.



Как определить конкретную ячейку?

# Геометрическая интерпретация: первый слой

$$A_1 = -0.4x_1 + 0.8x_2$$

$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

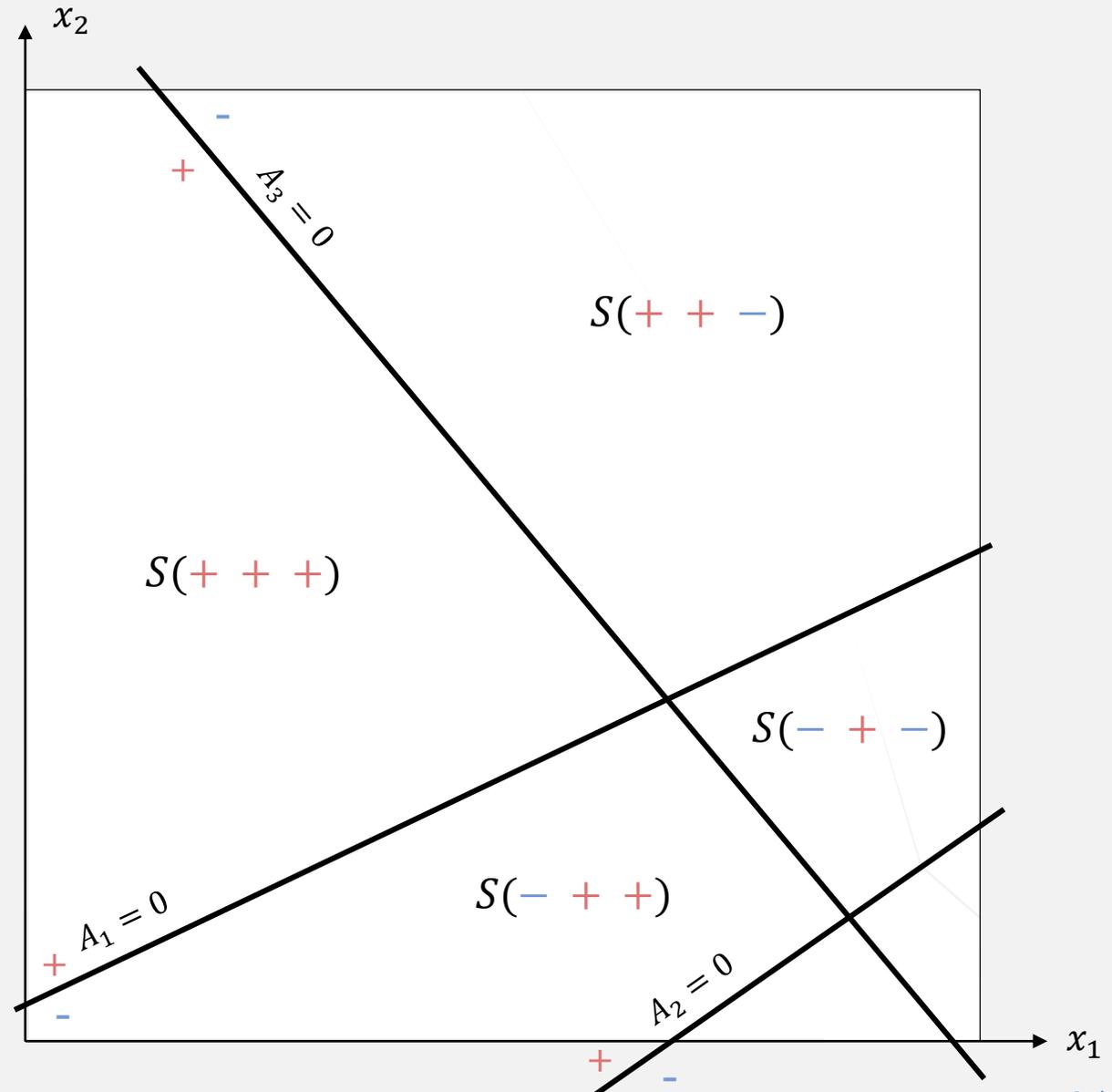
$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

**Первый слой** персептрона разбивает входной компакт на непересекающиеся ячейки.

Каждая ячейка определяется **набором знаков** соответствующих неравенств нейронов  $A_1, \dots, A_k$ .

Что будет на выходном слое?



# Геометрическая интерпретация: выходной слой

$$A_1 = -0.4x_1 + 0.8x_2$$

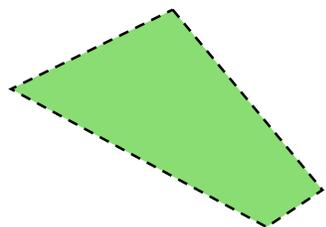
$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

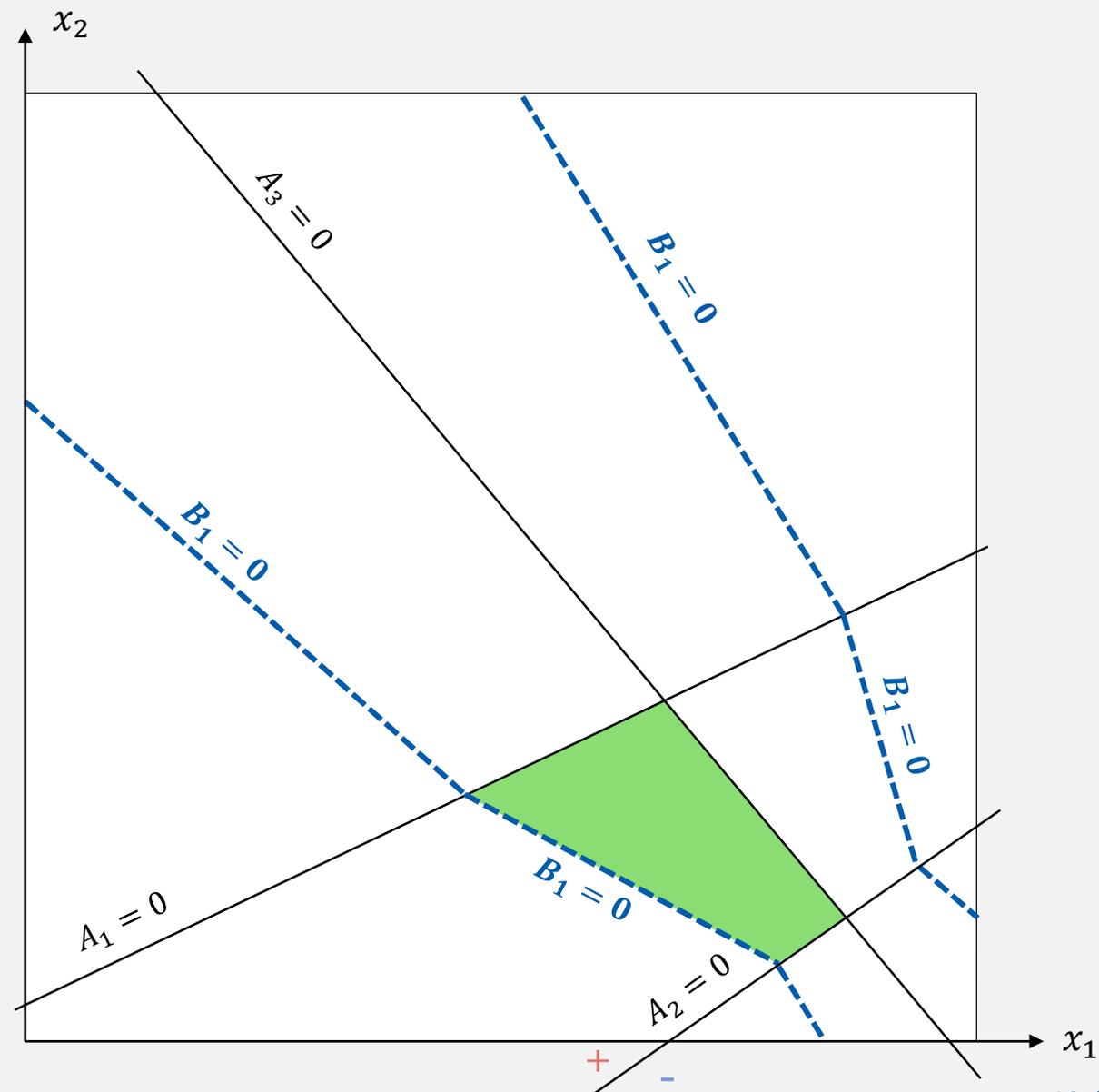
**Выходной слой** персептрона разбивает полученные предыдущим слоем ячейки.

**Внутри** каждой ячейки персептрон  $s(x)$  является **линейной функцией**:



$$S(- + + -)$$

$$B_1 = -2.31x_1 - 4.11x_2 + 2.25$$



# Геометрическая интерпретация: выходной слой

$$A_1 = -0.4x_1 + 0.8x_2$$

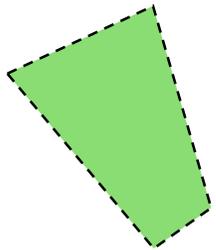
$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

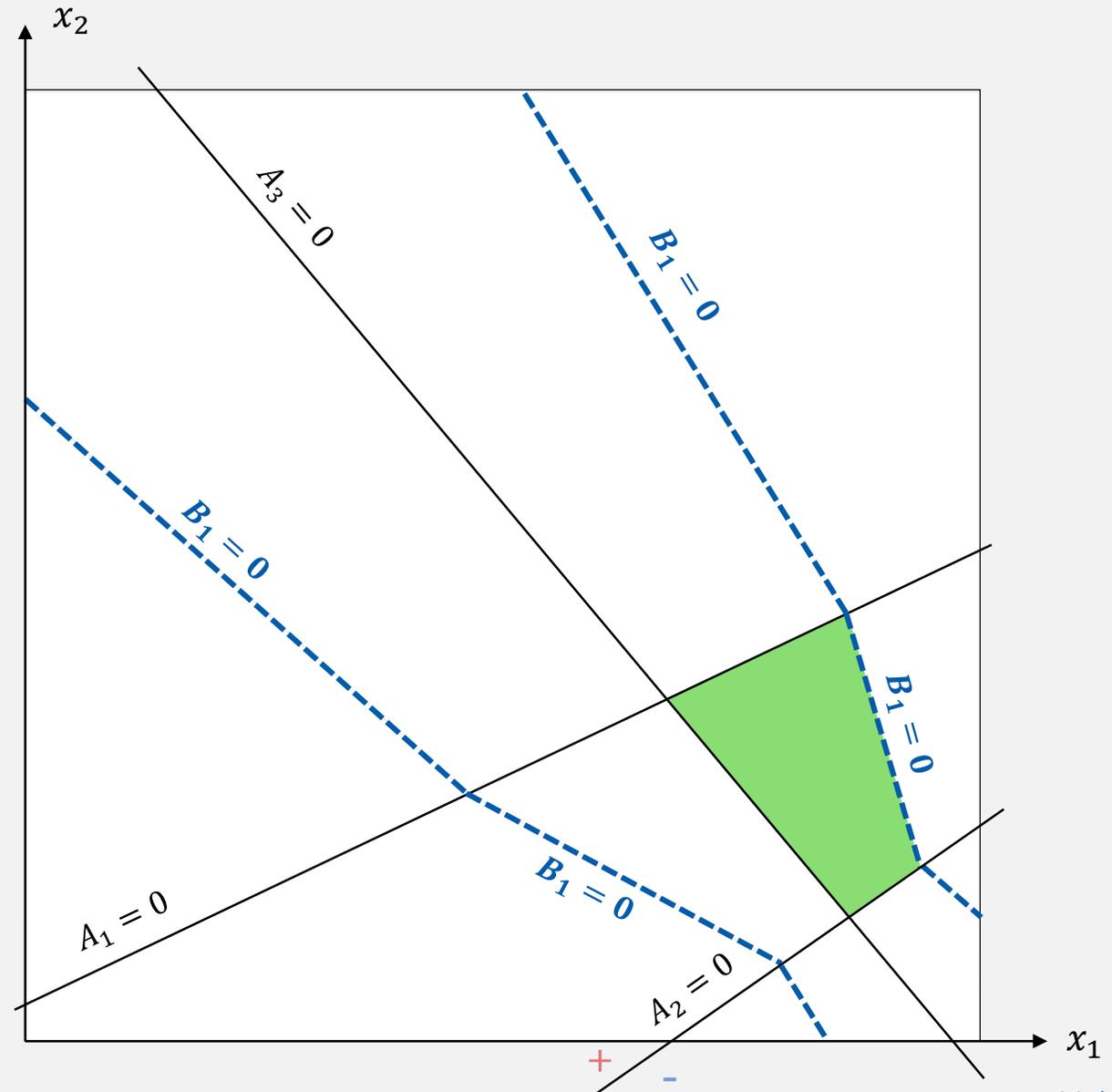
**Выходной слой** персептрона разбивает полученные предыдущим слоем ячейки.

**Внутри** каждой ячейки персептрон  $c(x)$  является **линейной функцией**:



$$S(- + - -)$$

$$B_1 = 3.99x_1 + 1.35x_2 - 4.05$$



# Геометрическая интерпретация: выходной слой

$$A_1 = -0.4x_1 + 0.8x_2$$

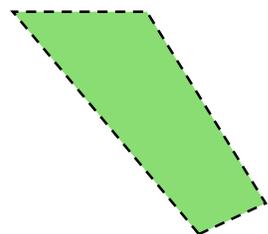
$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

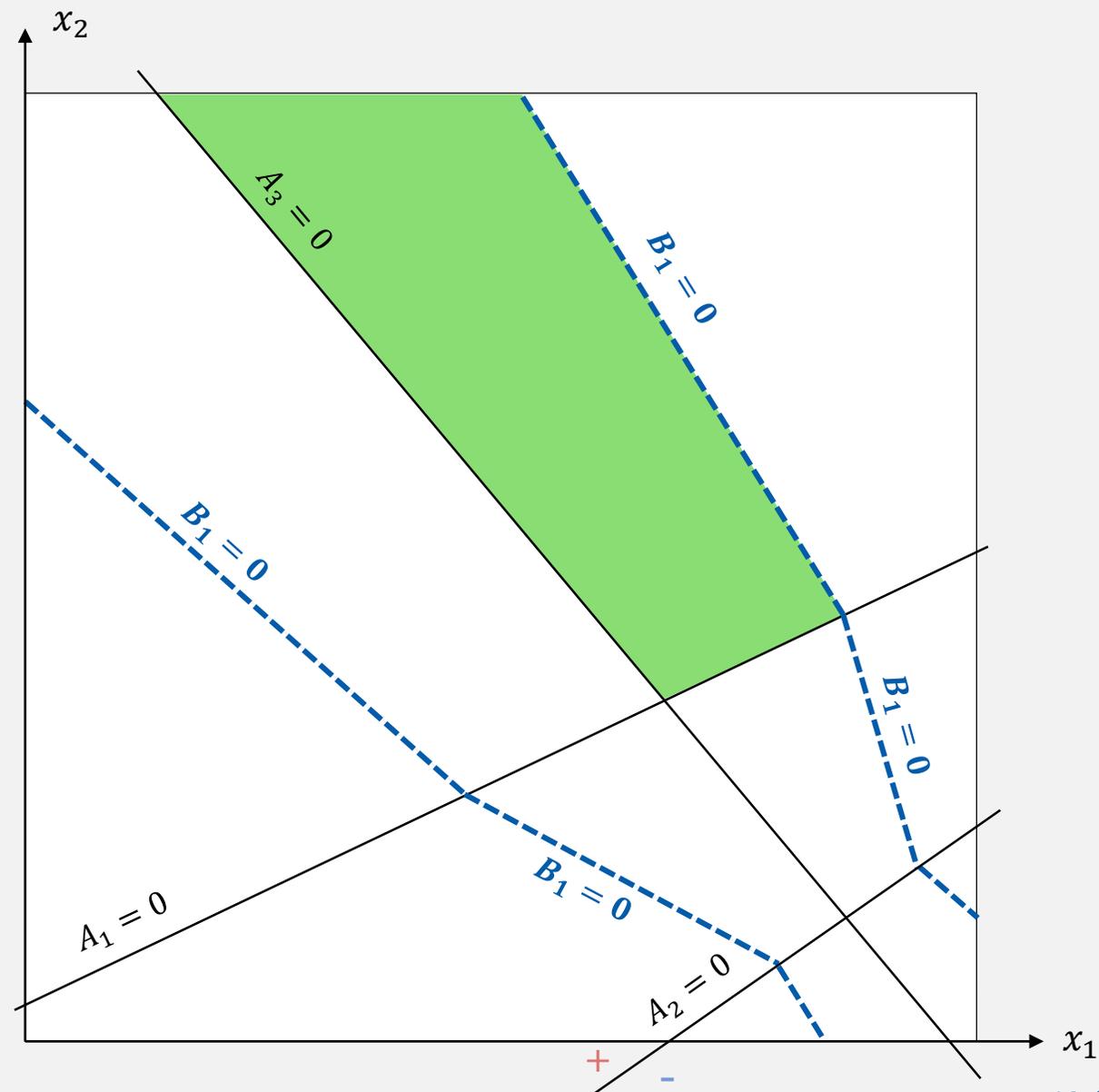
**Выходной слой** персептрона разбивает полученные предыдущим слоем ячейки.

**Внутри** каждой ячейки персептрон  $s(x)$  является **линейной функцией**:



$$S(+ + - -)$$

$$B_1 = 3.51x_1 + 2.31x_2 - 4.05$$



# Геометрическая интерпретация: выходной слой

$$A_1 = -0.4x_1 + 0.8x_2$$

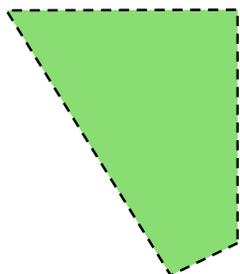
$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

**Выходной слой** персептрона разбивает полученные предыдущим слоем ячейки.

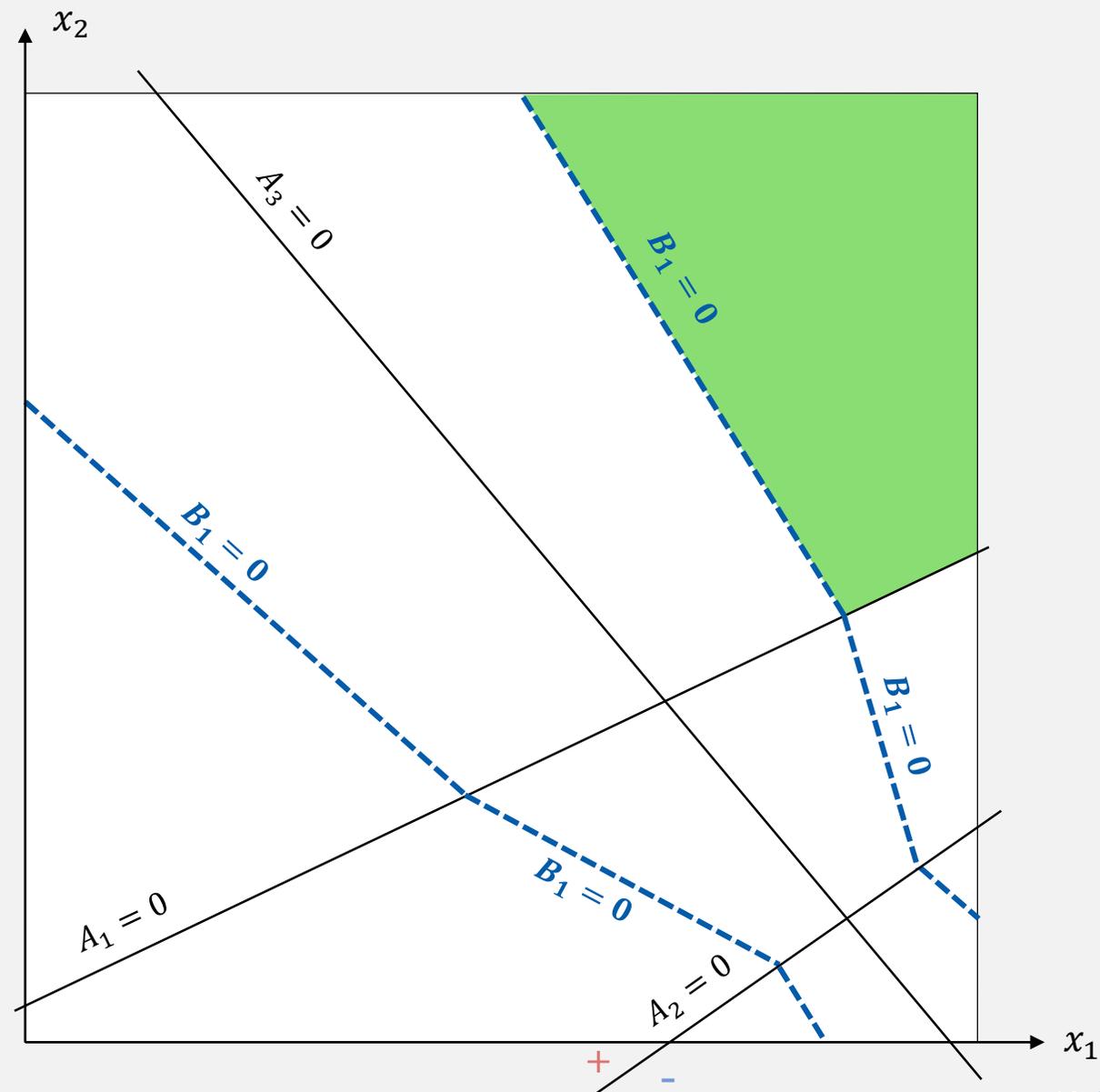
**Внутри** каждой ячейки персептрон  $s(x)$  является **линейной функцией**:



$$S(+ + - +)$$

$$B_1 = 3.51x_1 + 2.31x_2 - 4.05$$

Что делать, если слоёв больше одного?



# Бинарный байесовский классификатор

$(X, Y)$  – вектор из **распределения**  $P$ ;

$X \in [0, 1]^d$  – признаки из **распределения**  $P_x$  с **носителем**  $S$ ;

$Y \in \{\pm 1\}$  – **бинарные** метки;

$f: [0, 1]^d \rightarrow \{\pm 1\}$  – **дискриминантная** функция.

**Решение задачи регрессии:**

$$g(x) = \mathbb{E}(Y | X = x), \quad x \in [0, 1]^d$$

**Классический байесовский классификатор:**

$$s(x) = \begin{cases} +1, & g(x) > 0 \text{ и } x \in S \\ \pm 1, & g(x) = 0 \text{ или } x \notin S \\ -1, & g(x) < 0 \text{ и } x \in S \end{cases}$$

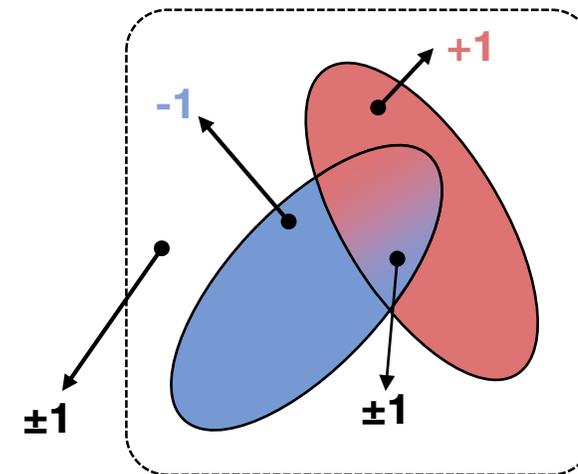
## Общая задача

**Классификация:**

$$\mathbb{P}(Y \neq f(X)) \rightarrow \min_f$$

**Регрессия:**

$$\mathbb{E}(Y - f(X))^2 \rightarrow \min_f$$



# Модификация байесовского классификатора

Добавим данные с **нулевыми** метками и **равномерно распределёнными** на  $[0, 1]^d$  признаками («фон»):

**Смесь двух распределений:**

$$P_\alpha = (1 - \alpha)P + \alpha\mathcal{P}, \quad \alpha \in (0, 1),$$

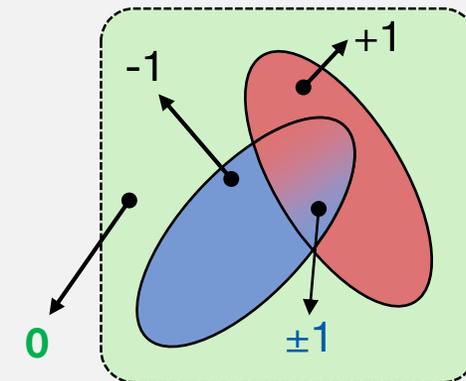
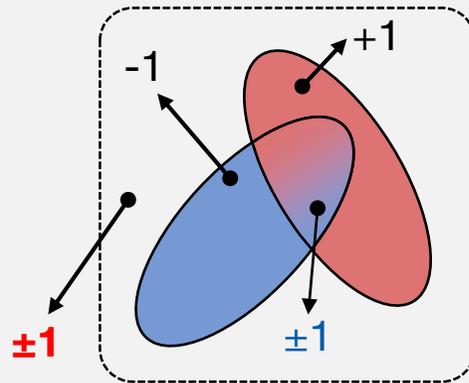
$\mathcal{P}$  – распределение «фона».

**Классический** классификатор  $s(x)$ :

$$s(x) = \begin{cases} +1, & g(x) > 0 \text{ и } x \in S \\ \pm 1, & g(x) = 0 \text{ и } x \in S \\ -1, & g(x) < 0 \text{ и } x \in S \\ \pm 1, & x \notin S \end{cases}$$

**Модифицированный** классификатор  $s_\alpha(x)$ :

$$s_\alpha(x) = \begin{cases} +1, & g_\alpha(x) > 0 \text{ и } x \in S \\ \pm 1, & g_\alpha(x) = 0 \text{ и } x \in S \\ -1, & g_\alpha(x) < 0 \text{ и } x \in S \\ \mathbf{0}, & x \notin S \end{cases}$$



# Аппроксимация байесовского классификатора

$D: (X, Y) = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$  – размеченный **обучающий** набор,  $X_i \in [0, 1]^d$ ,  $Y_i \in \{-1, +1\}$ ;

$c(X): [0, 1]^d \rightarrow R$  – равномерно непрерывная функция.

**Задача средне-квадратичной аппроксимации:**

$$E(c(X) - Y)^2 \rightarrow \min_{c(X)}$$

**Сводится к аппроксимации функции регрессии:**

$$E(c(X) - g(X))^2 \rightarrow \min_{g(X)}$$

**Универсальная теорема аппроксимации:** для любого  $\varepsilon > 0$  существуют такие  $k$  и  $L$ , что для любого  $x \in [0, 1]^d$ :

$$\sup_{x \in [0, 1]^d} |c(x) - g(x)| < \varepsilon$$

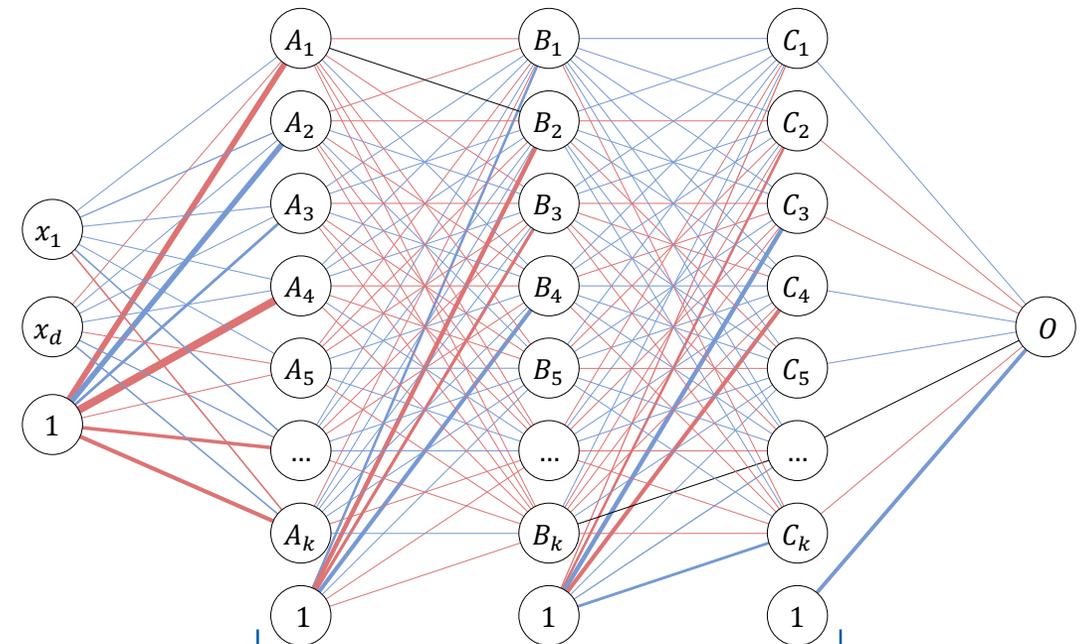
## Конфигурация персептрона:

01

$L$  **скрытых** слоёв по  $k$  нейронов с кусочно-линейной активацией

02

**Один** нейрон в **выходном** слое с линейной функцией активации



$L$  скрытых слоёв по  $k$  нейронов

# Функция нейросетевой регрессии

$\{(X_{n+1}, Y_{n+1}), \dots, (X_{2n}, Y_{2n})\}$  – сгенерированные **фоновые** данные,  $X_i \in [0, 1]^d$ ,  $Y_i = 0$ ;

$f(X)$  – равномерно непрерывная плотность выборки  $(X, Y)$  на  $S$ ;

$p(X)$  – плотность равномерно распределённых на  $[0, 1]^d$  фоновых векторов.

$C(L, k)$  – класс многослойных персептронов с активацией ***Abs*** и заданными  $L, k$ .

## Выборочная оценка решения задачи аппроксимации

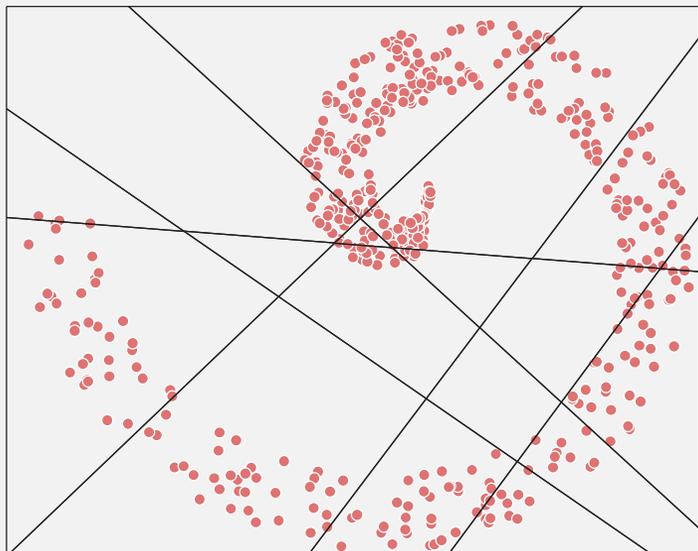
$$\sum_{i=1}^{2n} (c_n(X_i) - Y_i)^2 \rightarrow \min_{c_n(X) \in C(L, k)}$$

$c_n^*(X)$  – решение задачи, называемое **функцией нейросетевой регрессии**.

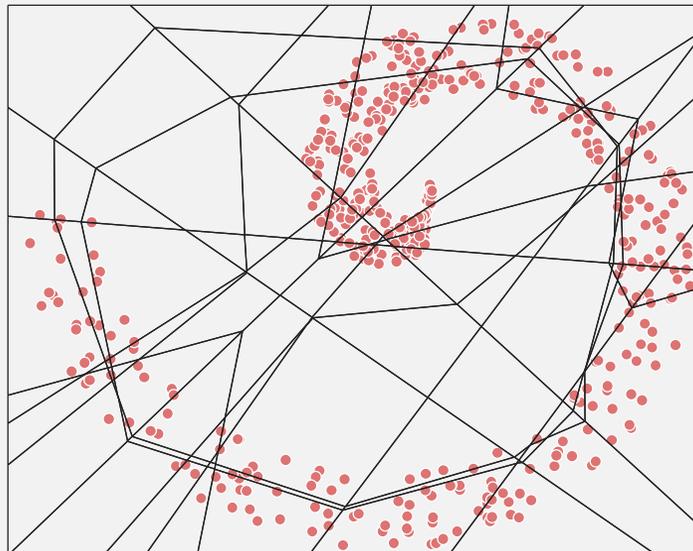
# Иерархическое разбиение на ячейки

$c_n^*(x)$  строит **иерархическое разбиение**  $[0, 1]^d$  на  $N$  непересекающихся ячеек:  $K = \{K_1, K_2, \dots, K_N\}$

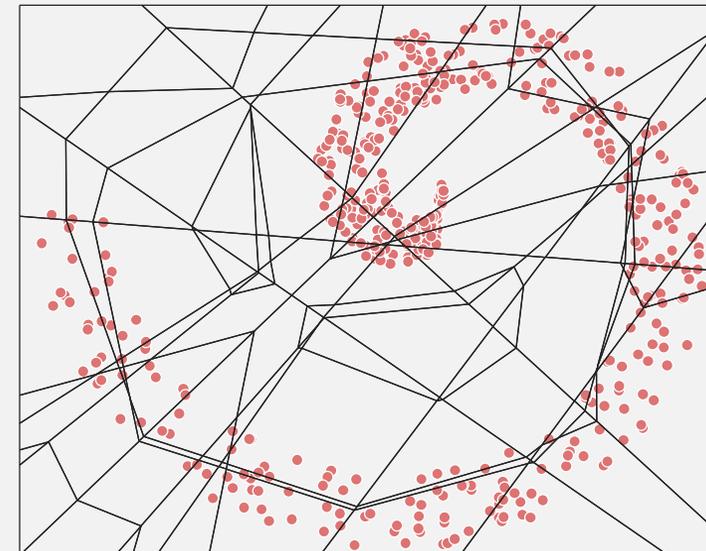
Пример разбиения некоторым MLP с  $L = 2$ ,  $k = 6$ :



Первый скрытый слой: 18 ячеек



Второй скрытый слой: 121 ячейка



Выходной слой: 148 ячеек

# Функция гистограммной регрессии

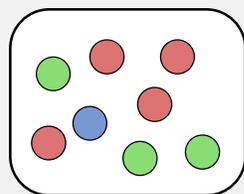
$h_n(X)$  – кусочно-постоянная функция гистограммной регрессии, определённая путём минимизации:

$$\sum_{i=1}^{2n} (h_n(X_i) - Y_i)^2 \rightarrow \min_{h_n(X)}$$

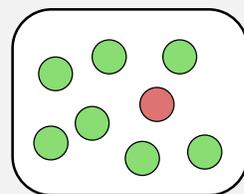
Пусть  $X \in K_r \rightarrow$  для каждой ячейки  $K_r$  оптимальное решение имеет вид:

$$h_n^*(X) = \frac{n_1(X) - n_{-1}(X)}{n_{-1}(X) + n_0(X) + n_1(X)}$$

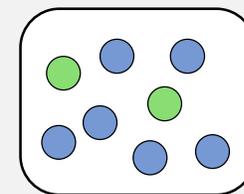
$n_j(X)$  – количество наблюдений с меткой  $j$  в ячейке.



$$n_{-1} = 1, n_1 = 4, n_0 = 3$$
$$h_n^*(X) = \frac{4 - 1}{1 + 3 + 4} = 0.375$$



$$n_{-1} = 0, n_1 = 1, n_0 = 7$$
$$h_n^*(X) = \frac{1 - 0}{0 + 7 + 1} = 0.125$$



$$n_{-1} = 6, n_1 = 0, n_0 = 2$$
$$h_n^*(X) = \frac{0 - 6}{6 + 2 + 0} = -0.75$$

# eXVTree: объясняющее дерево решений

## Идея:

Раскрытие модулей в каждом нейроне – это переход по ветвям двоичного дерева (по знаку неравенства).

**Узлы:** линейные неравенства нейронов.

**Листья:** линейная функция и гистограмма классов.

**Адрес ячейки (листа):** знаки нейронов до активации.

## Интерпретируемость и доверие:

**Геометрия** ячейки: выпуклый многогранник (система линейных неравенств).

Условия принятия **решения:** путь по дереву.

**Прецеденты:** объекты из той же или соседних ячеек.

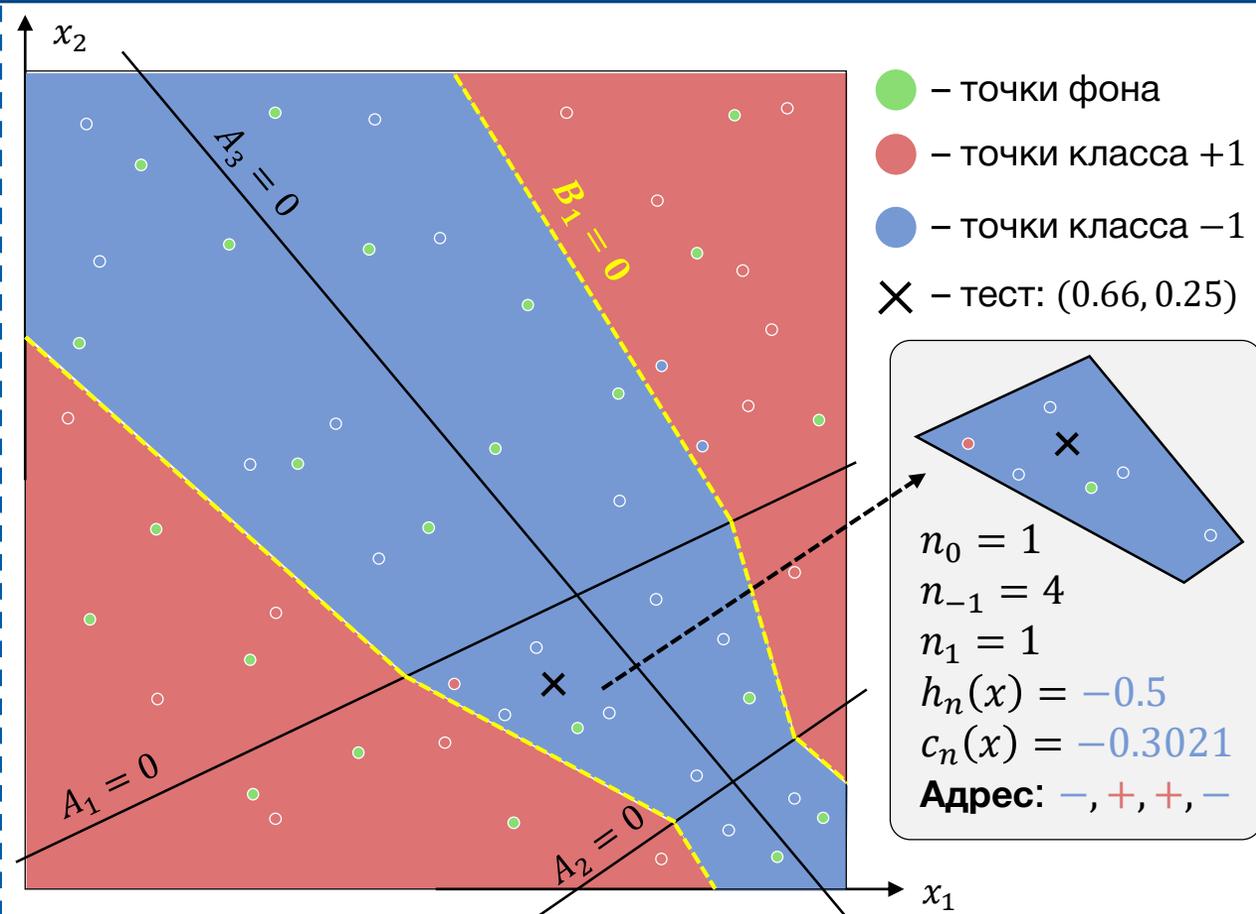
**Малое число точек:** адаптивное укрупнение области (подъём на уровень выше) или отказ от решения.

## Сложность

**Построение** (полное):  $O(k^{dL})$ .

**Построение** (по выборке):  $O(n)$ .

**Применение** (1 наблюдение):  $O(Lkd) = O(1)$ .



$$A_1 = -0.4x_1 + 0.8x_2$$

$$A_2 = -1.2x_1 + 1.8x_2 + 0.8$$

$$A_3 = -1.5x_1 - 1.3x_2 + 1.5$$

$$B_1 = 0.6|A_1| - 0.5|A_2| + 2.1|A_3| - 0.5$$

**Линейная функция ячейки:**  $-2.31x_1 - 4.11x_2 + 2.25$

$$A_1(x) = -0.064 < 0 \rightarrow -$$

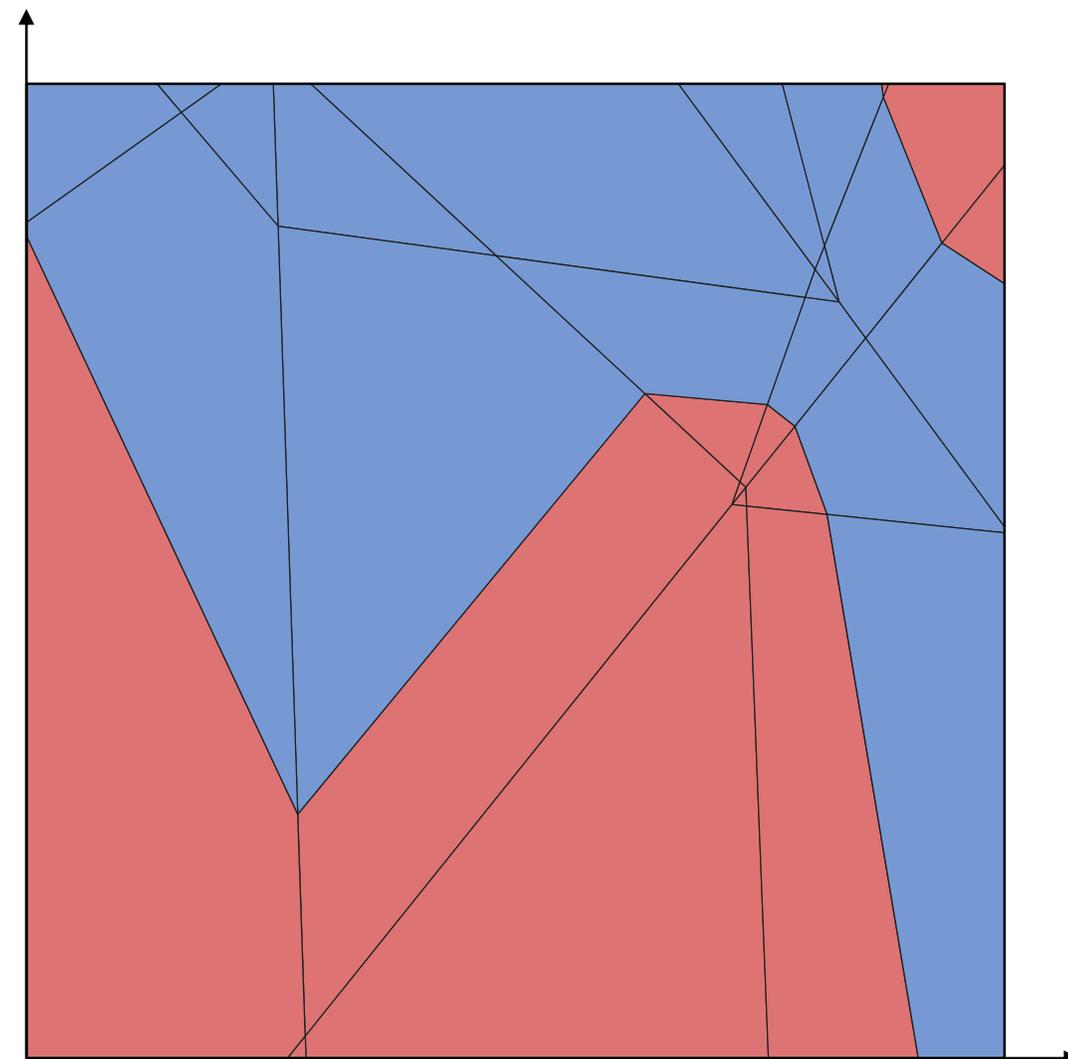
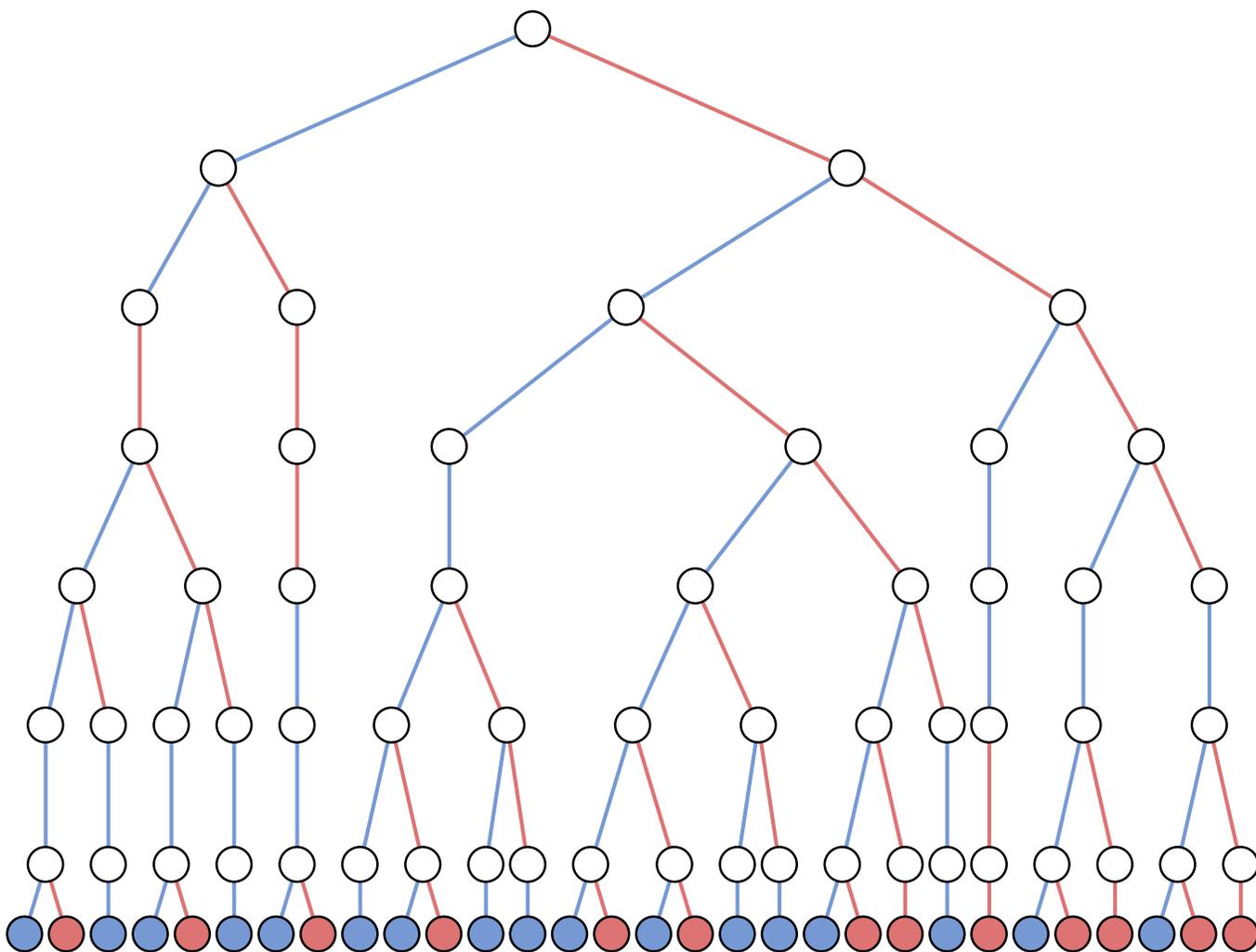
$$A_2(x) = 0.458 > 0 \rightarrow +$$

$$A_3(x) = 0.185 > 0 \rightarrow +$$

$$B_1(x) = -0.3021 < 0 \rightarrow -$$

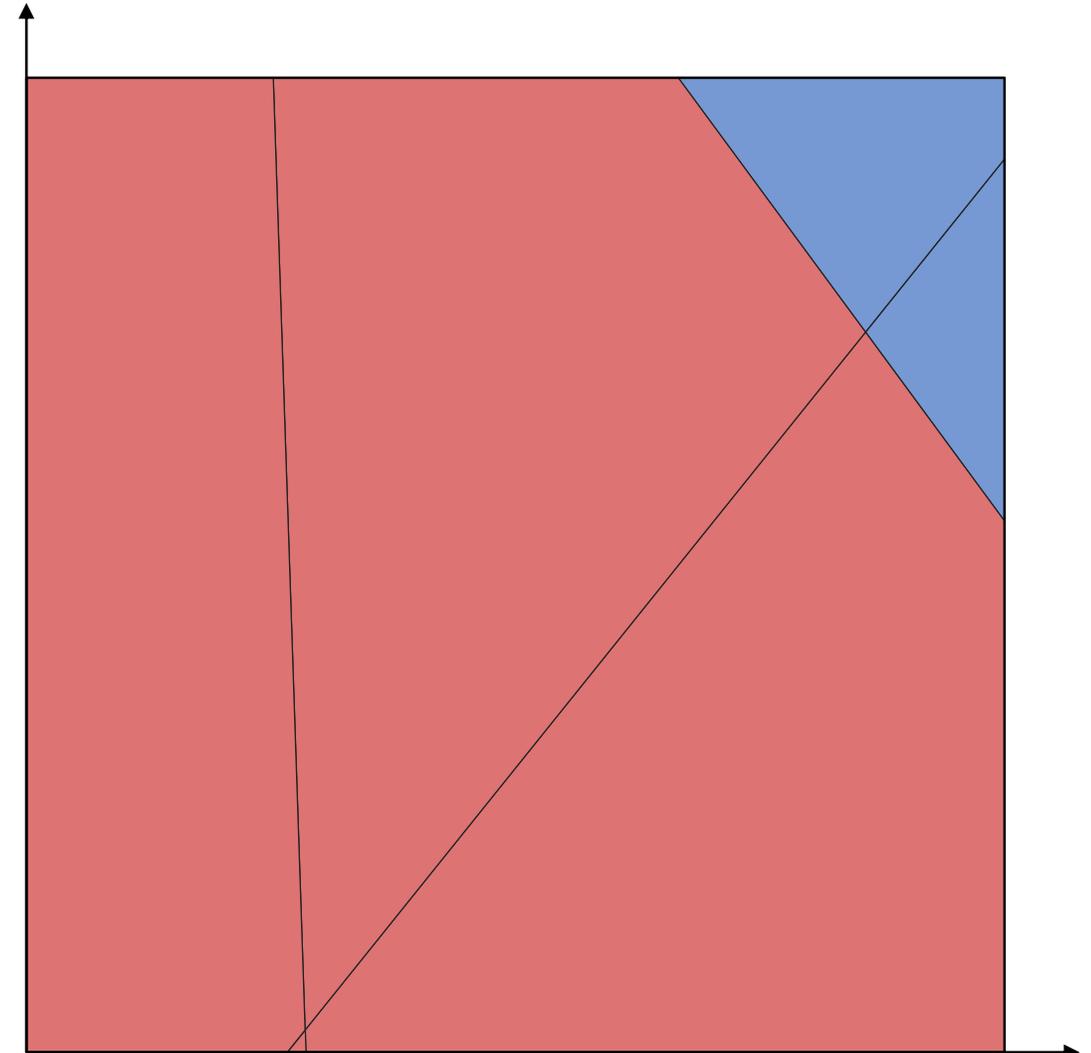
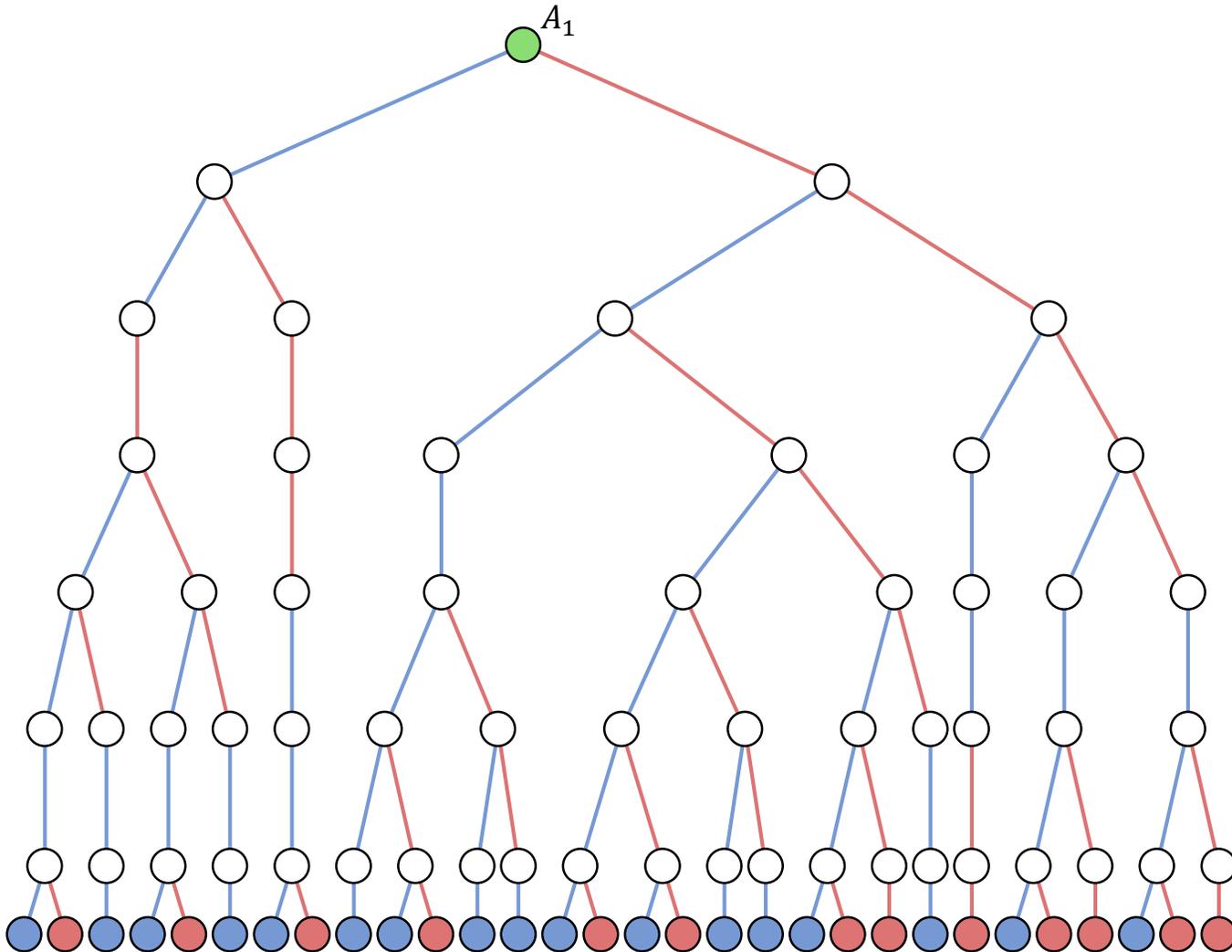
# eXVTree: подробный пример

Модель с двумя скрытыми слоями по 3 нейрона: полное дерево



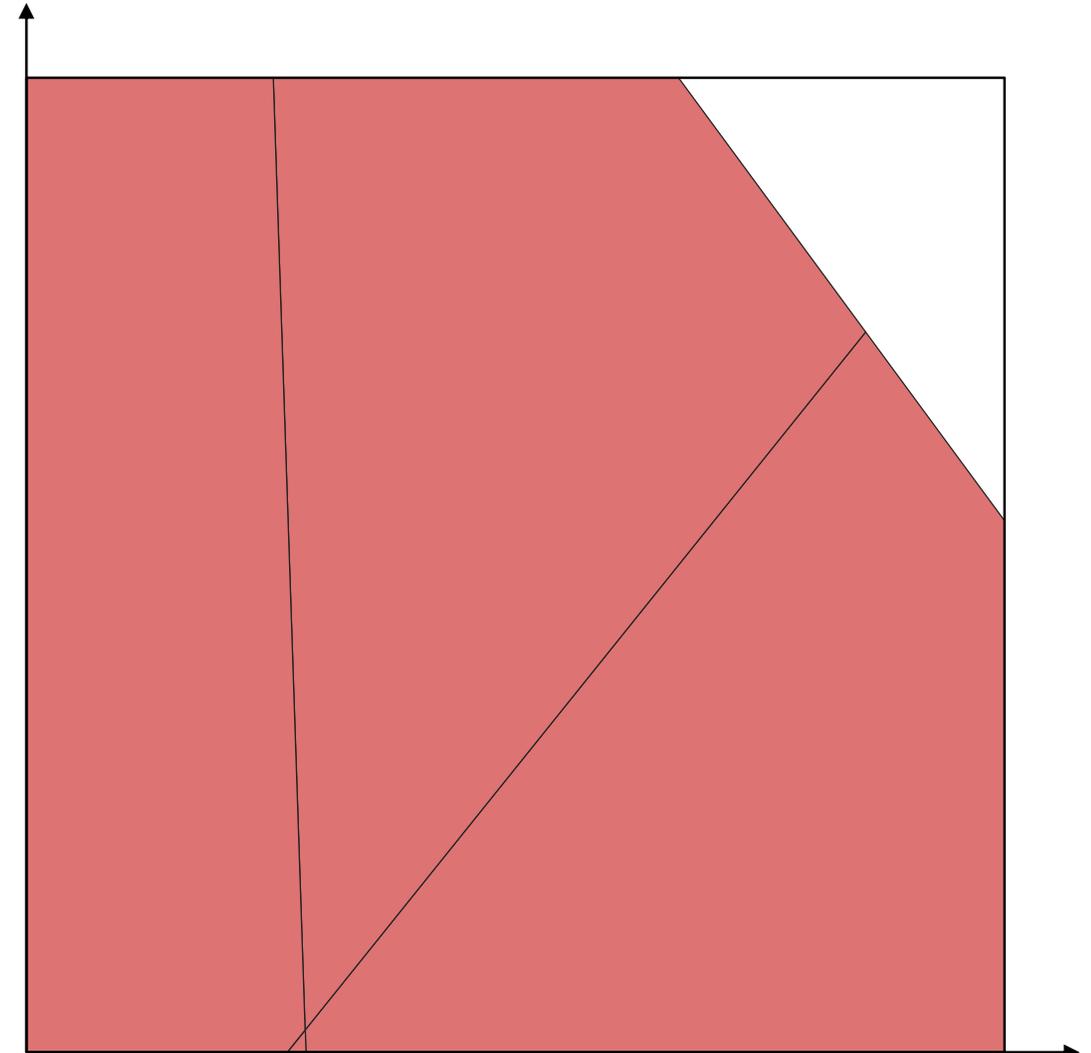
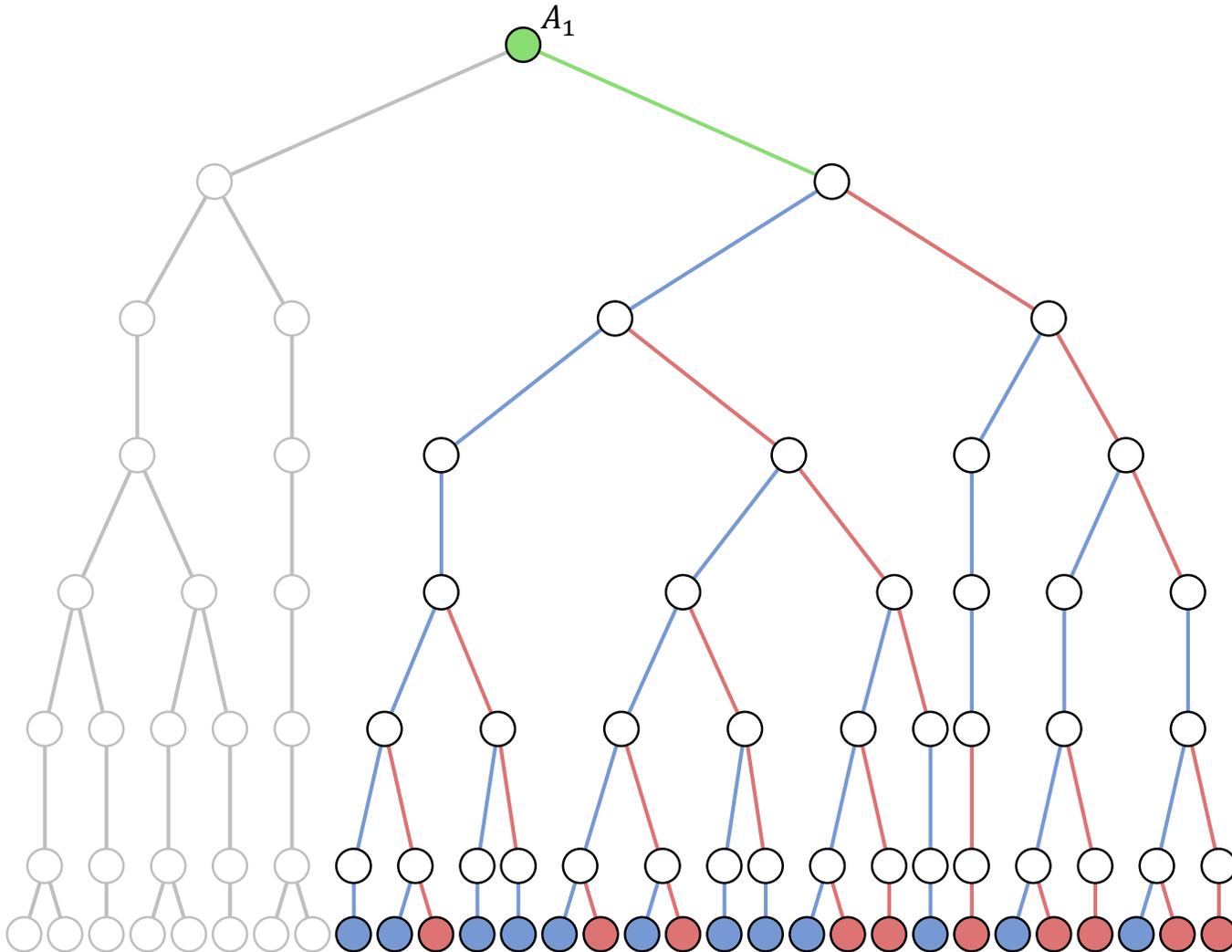
# eXTree: ячейки первого скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $A_1$



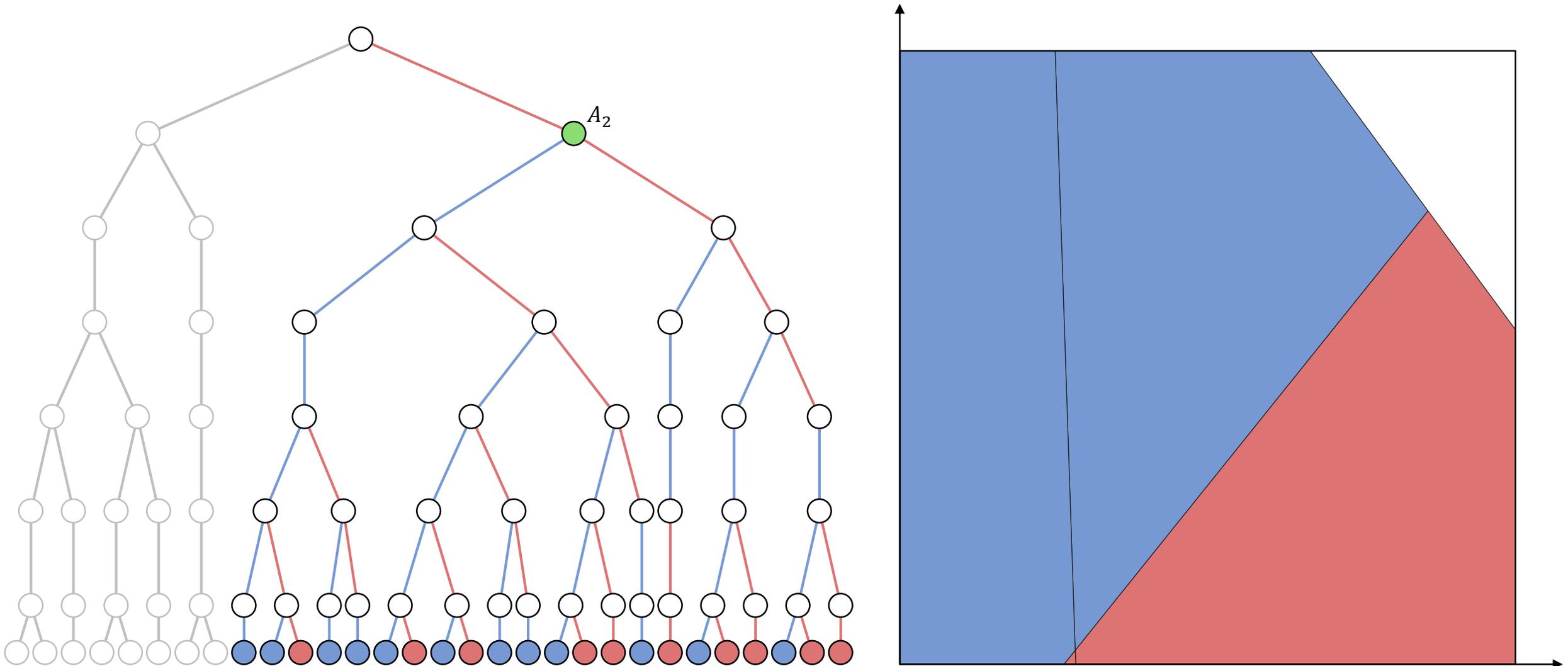
# eXVTree: ячейки первого скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $A_1 > 0$



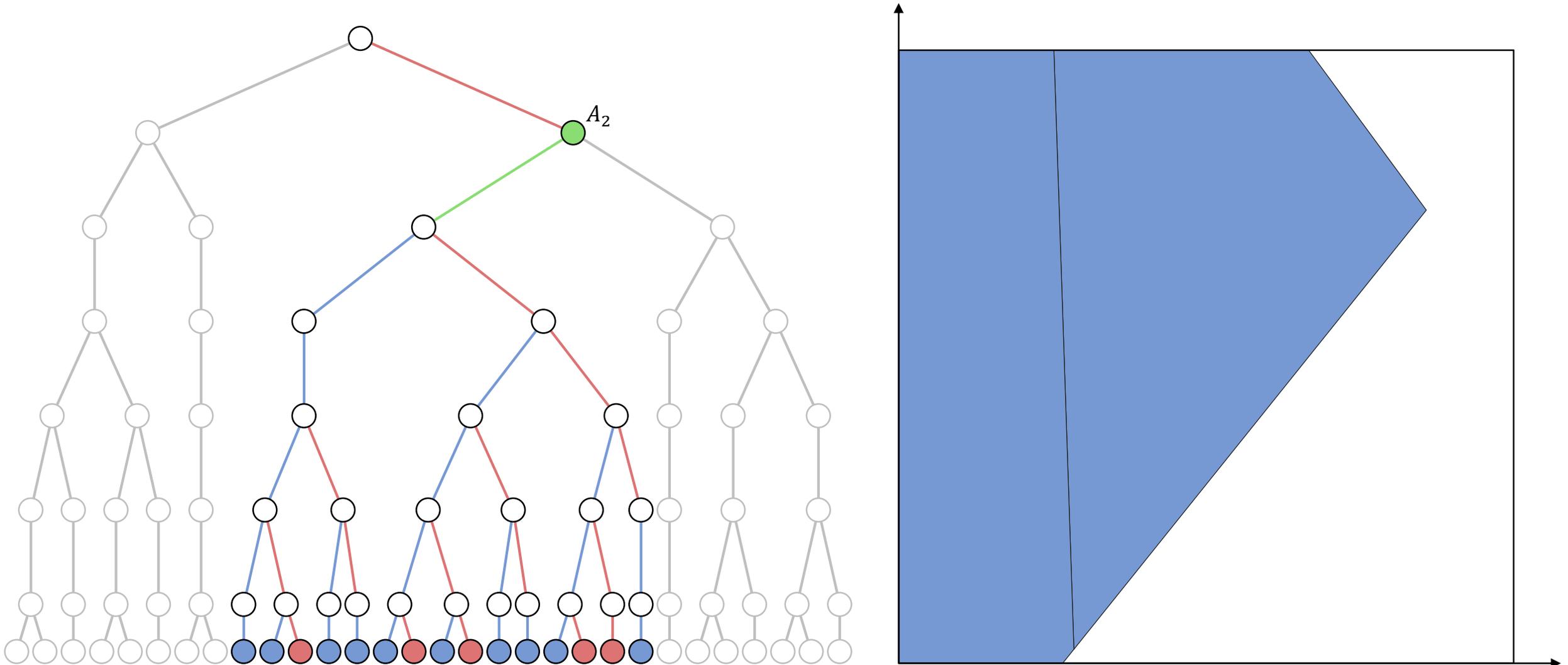
# eXVTree: ячейки первого скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $A_2$



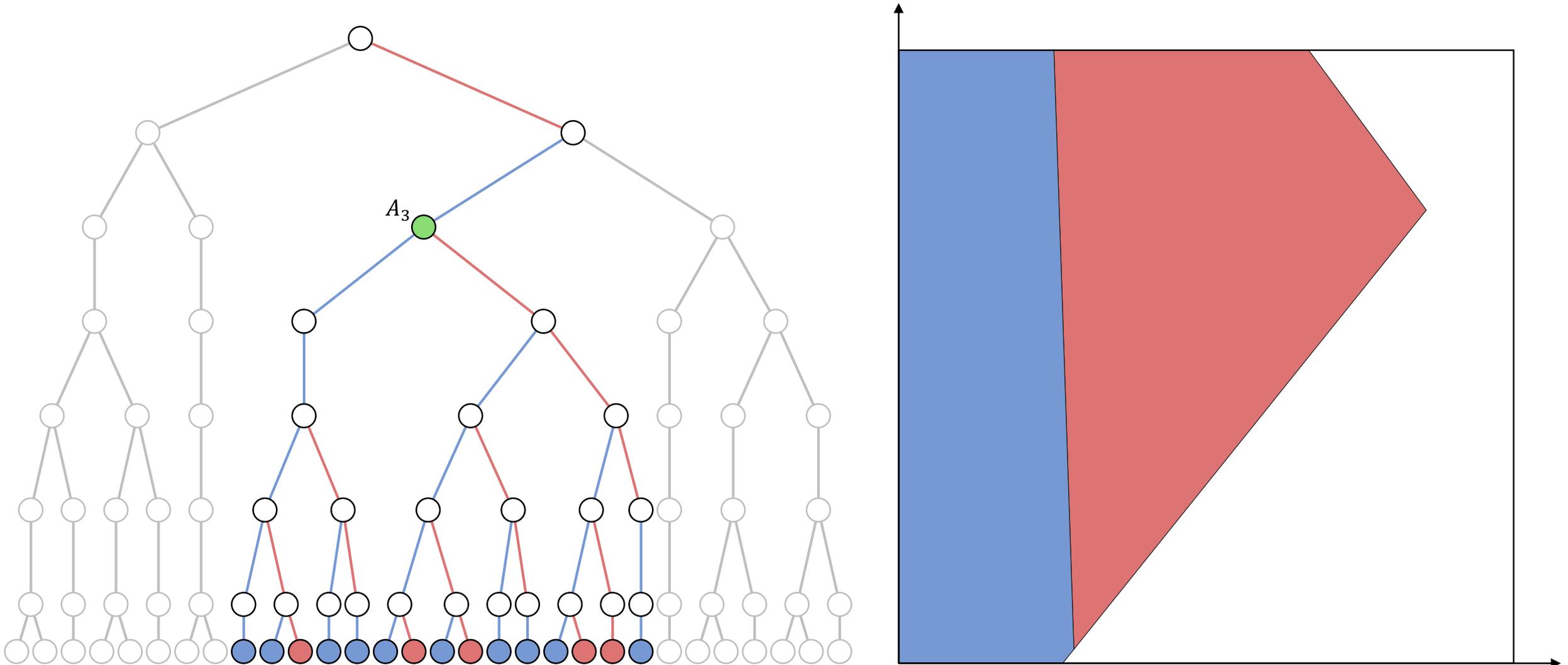
# eXVTree: ячейки первого скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $A_2 < 0$



# eXVTree: ячейки первого скрытого слоя

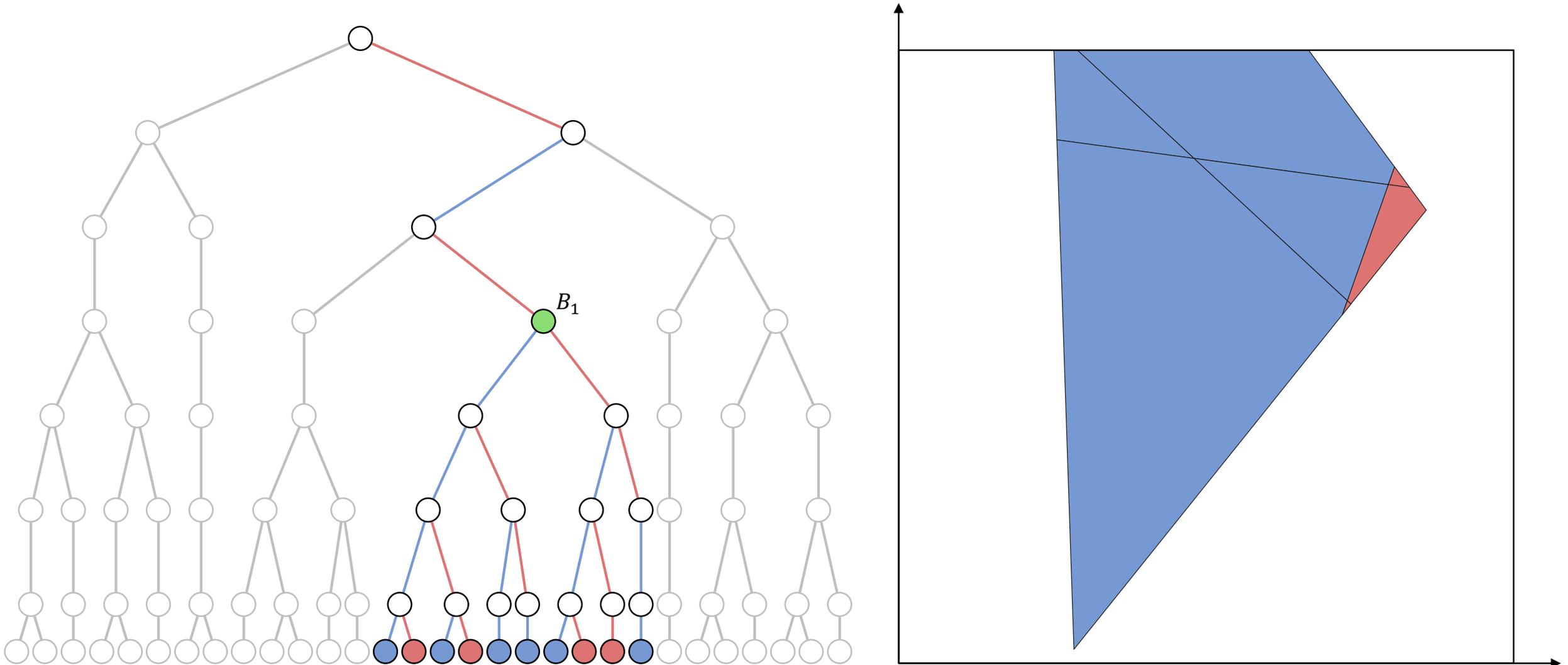
Проследим путь (+ - + - - + -) по дереву: нейрон  $A_3$





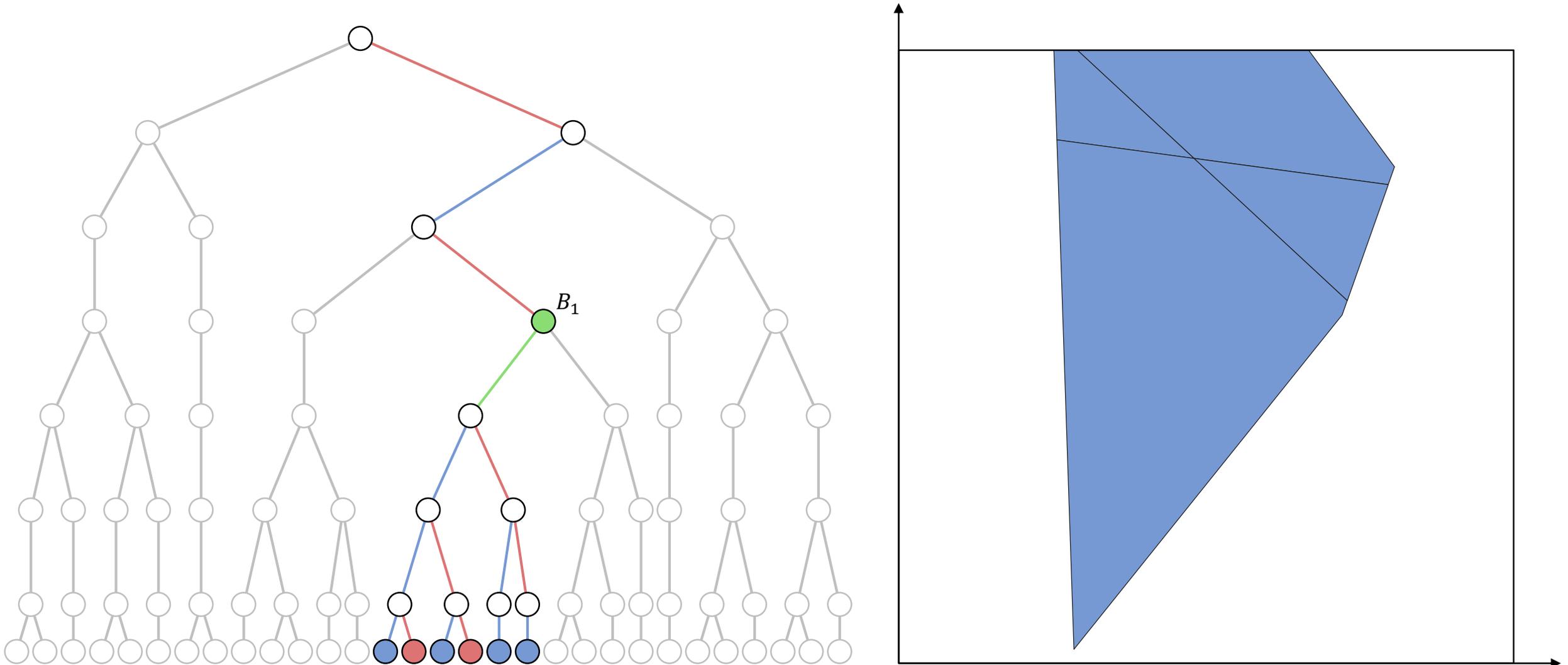
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_1$



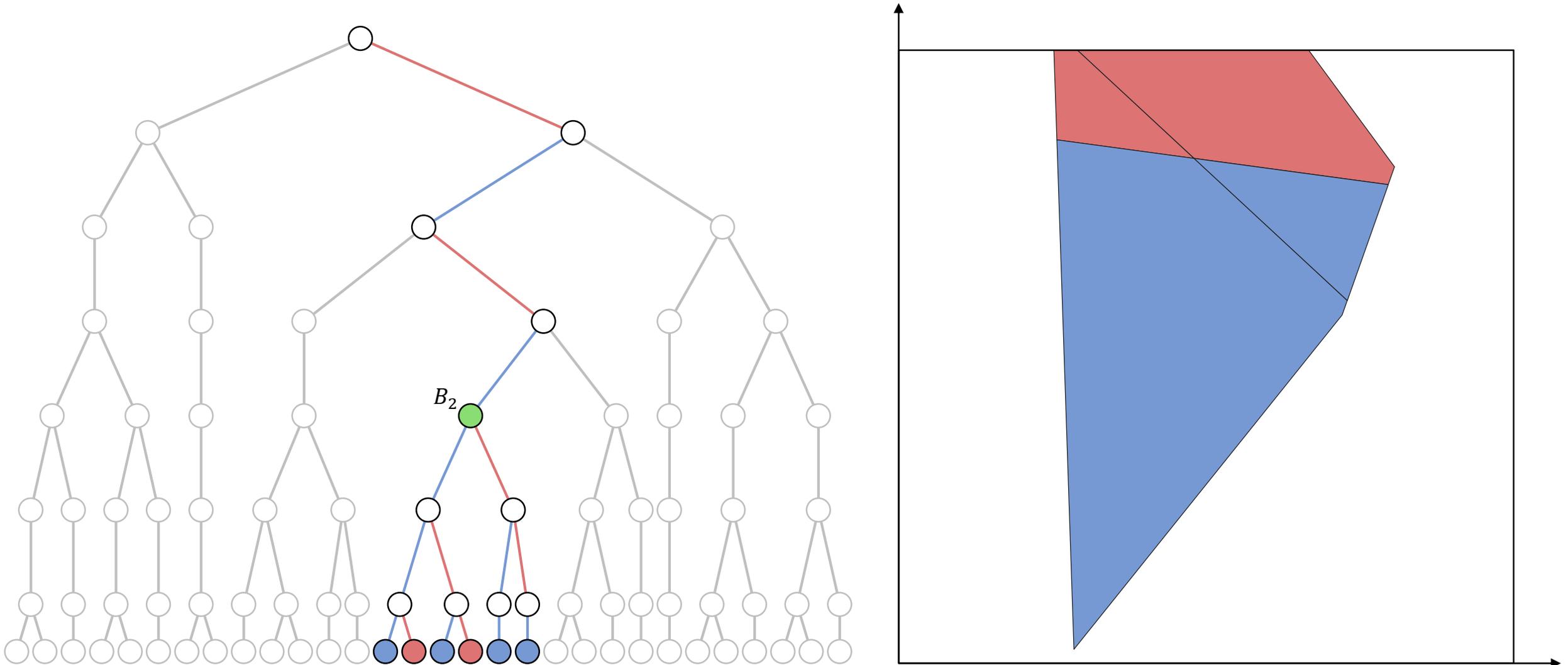
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_1 < 0$



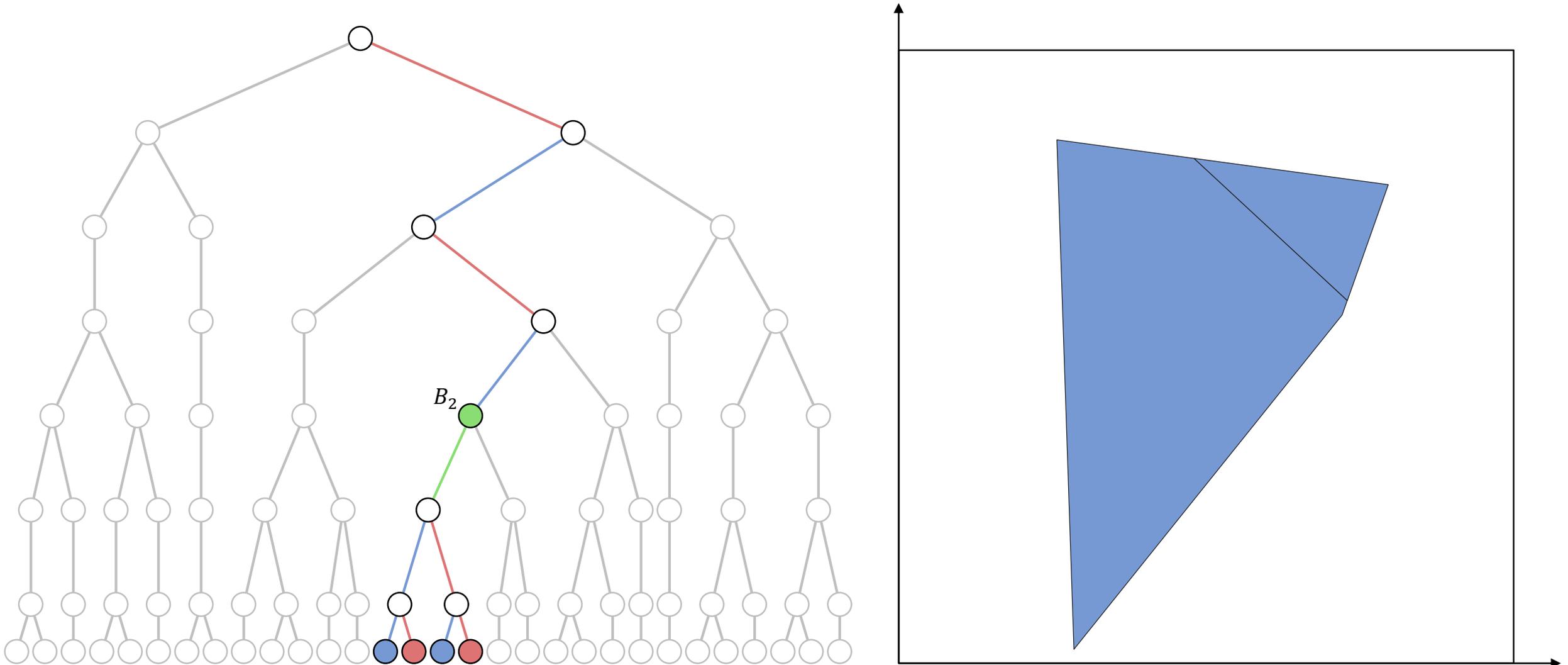
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_2$



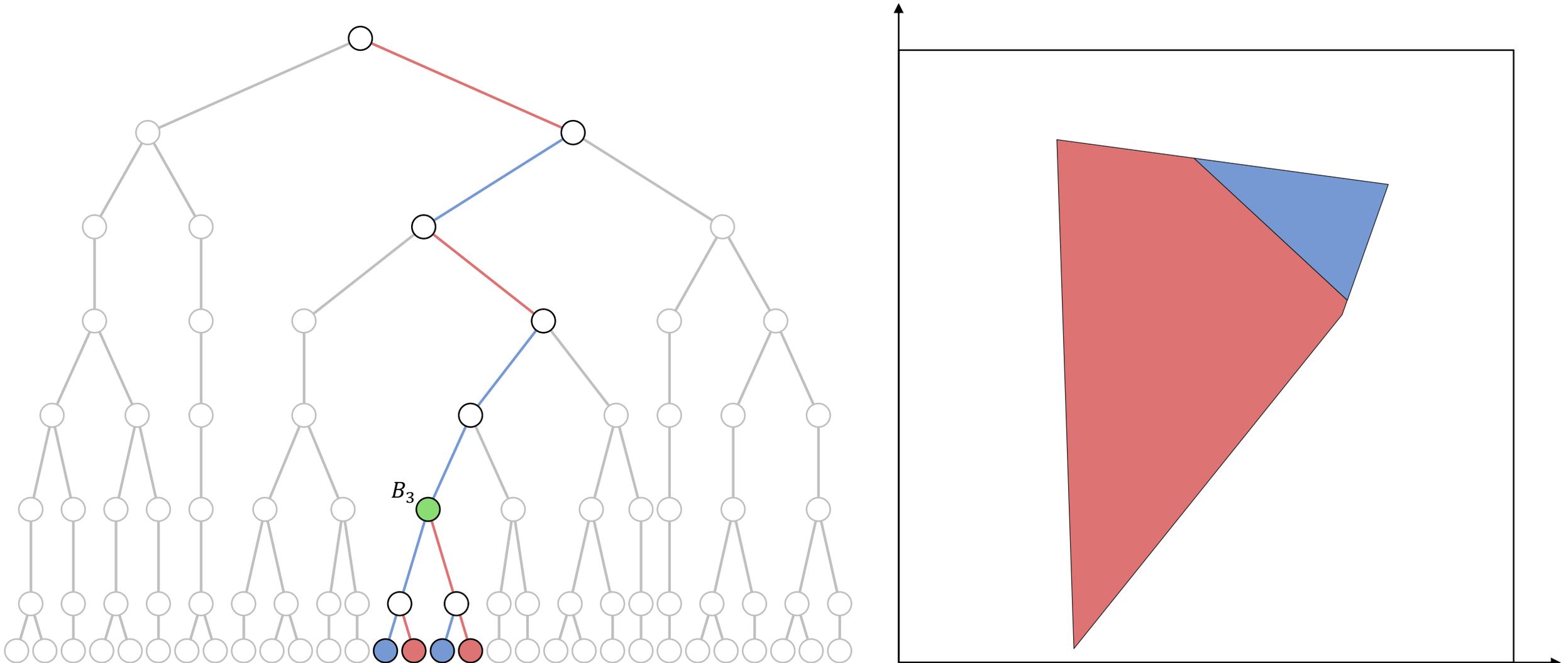
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_2 < 0$



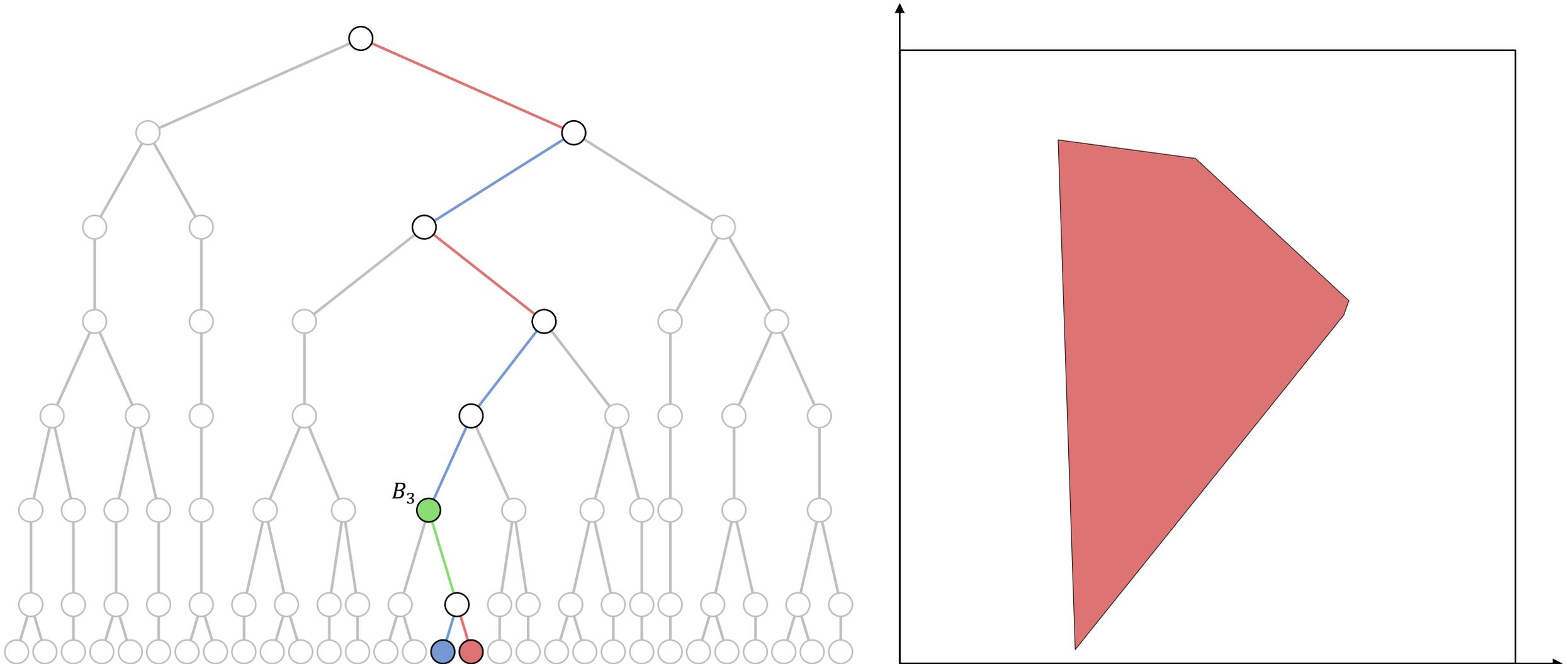
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_3$



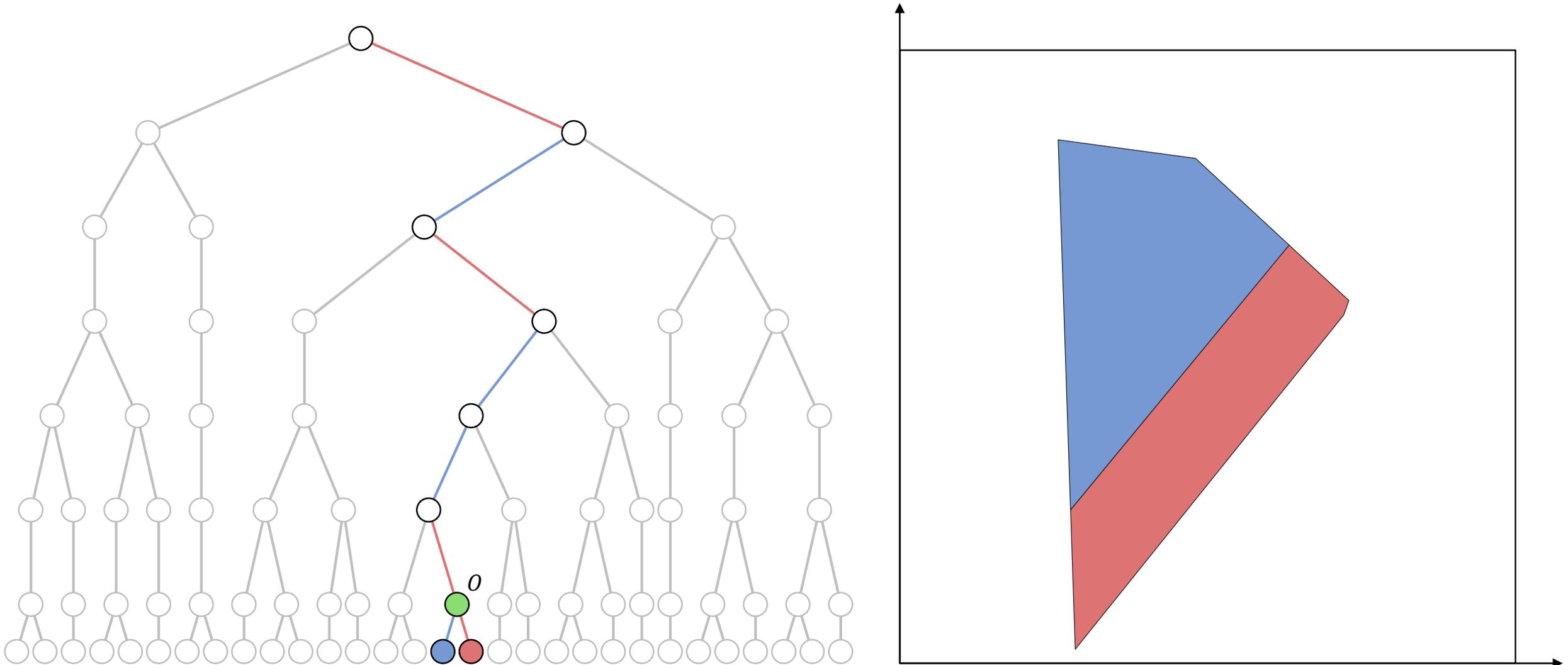
# eXVTree: ячейки второго скрытого слоя

Проследим путь (+ - + - - + -) по дереву: нейрон  $B_3 > 0$



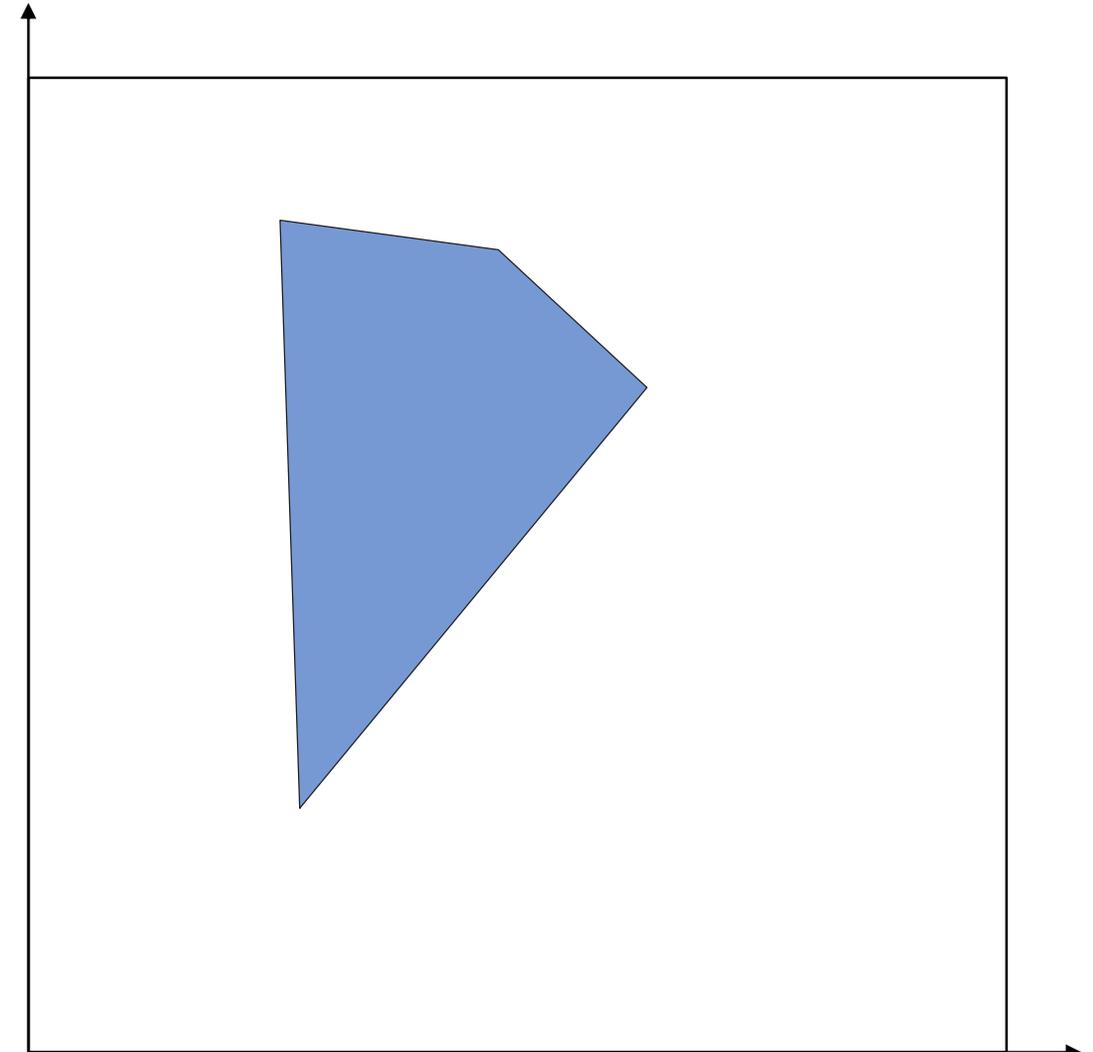
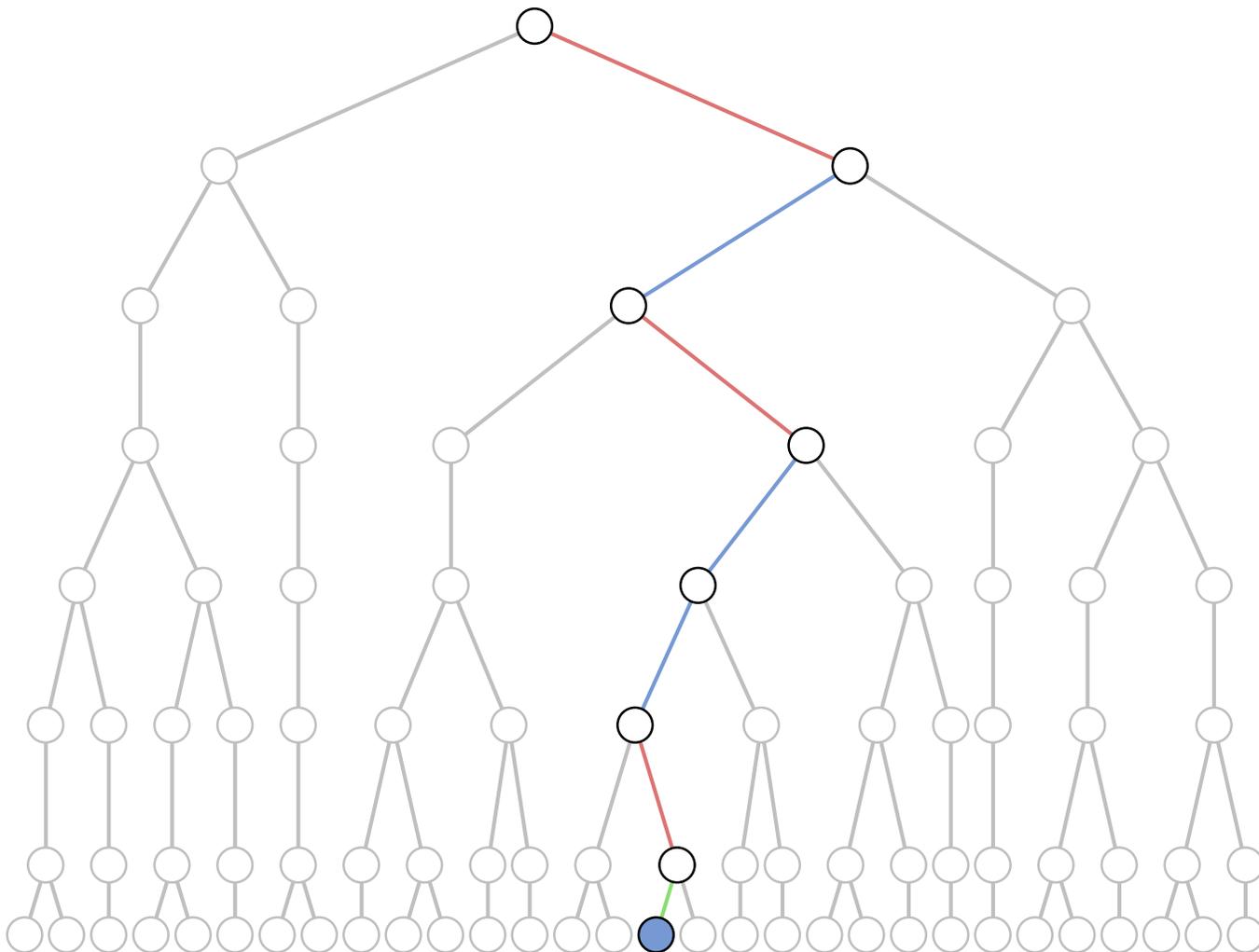
# eXVTree: ячейки выходного слоя

Проследим путь (+ - + - - + -) по дереву: выходной нейрон



# eXVTree: лист дерева

Проследим путь (+ - + - - + -) по дереву: итоговая ячейка



# Связь нейросетевой и гистограммной регрессий

## Асимптотические свойства:

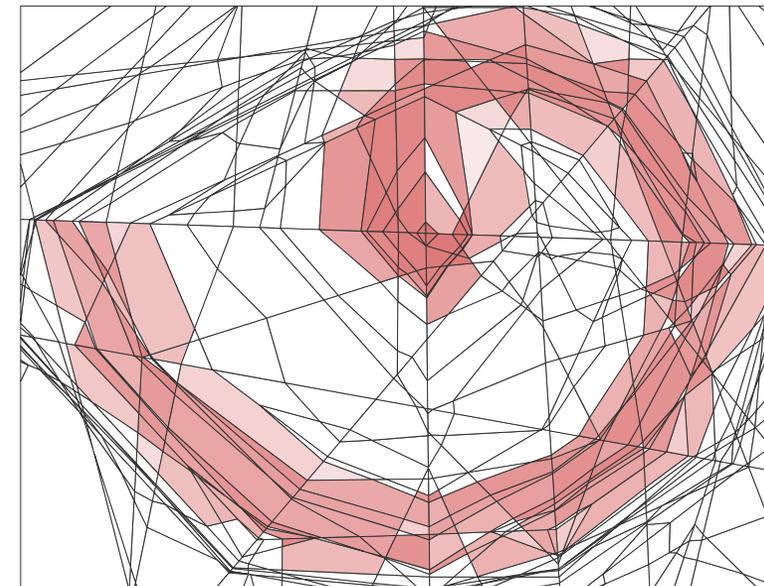
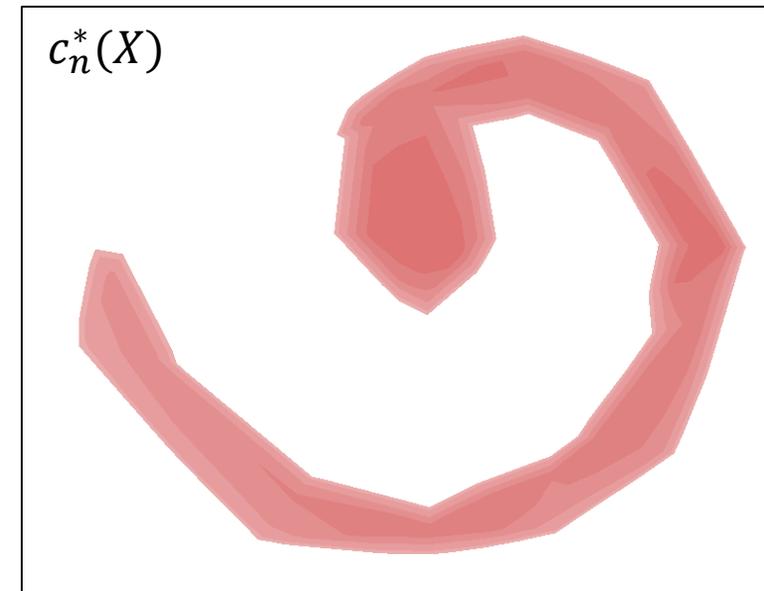
При  $n \rightarrow \infty$  и  $kL \rightarrow \infty$ :

$$\mathbb{E}\{h_n^*(X) - c_n^*(X)\} \rightarrow 0$$

Возможна замена гистограммной регрессии на нейросетевую, и при этом MLP вычислительно **эффективнее** и **не хранит** обучающую выборку.

## Практическое применение:

- Принятие решения с **порогом доверия**  $\beta$ :  
 $|c_n^*(X)| > \beta$ .
- Контроль уверенности классификатора.
- Отказ от ответа при низкой апостериорной вероятности.



$h_n^*(X)$

# Состоятельность классификатора

$c_n(x)$  – выборочная оценка апостериорной вероятности  $\mathbb{P}(Y = 1|X = x)$ .

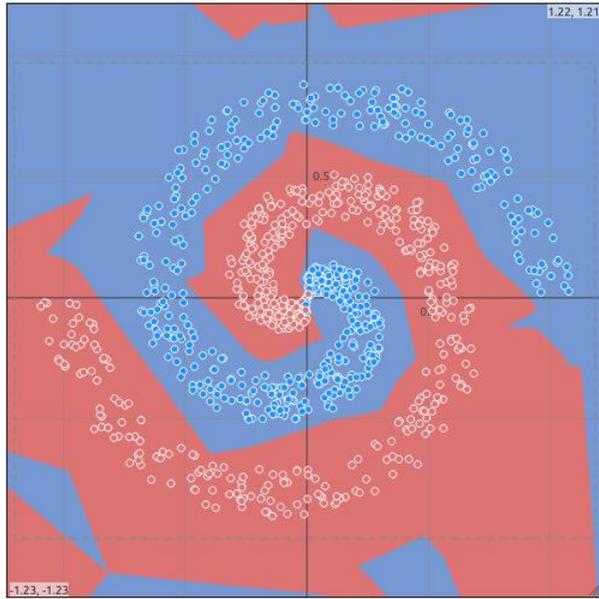
$g(x)$  – истинная байесовская апостериорная вероятность.

Классификатор на основе  $c_n(x)$  **состоятельный**, если при  $n \rightarrow \infty$  его ошибка стремится к минимально возможной (ошибке Байеса):

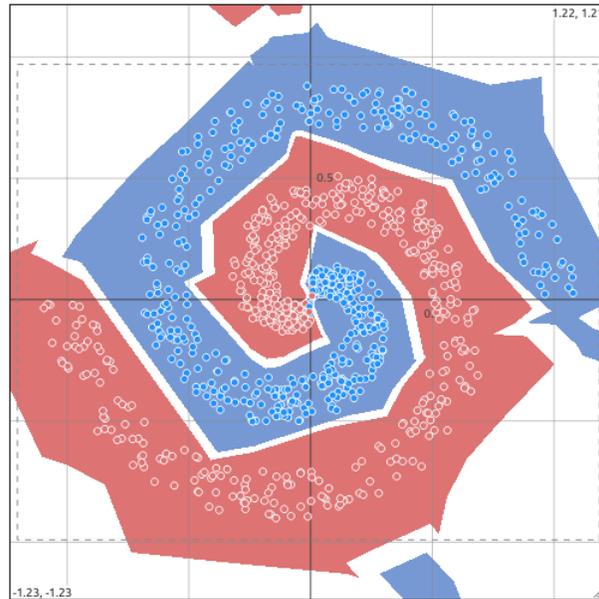
$$c_n(x) \xrightarrow{\mathbb{P}} g(x)$$

Гарантирует, что с ростом данных классификатор приближается к **оптимальному** решению.

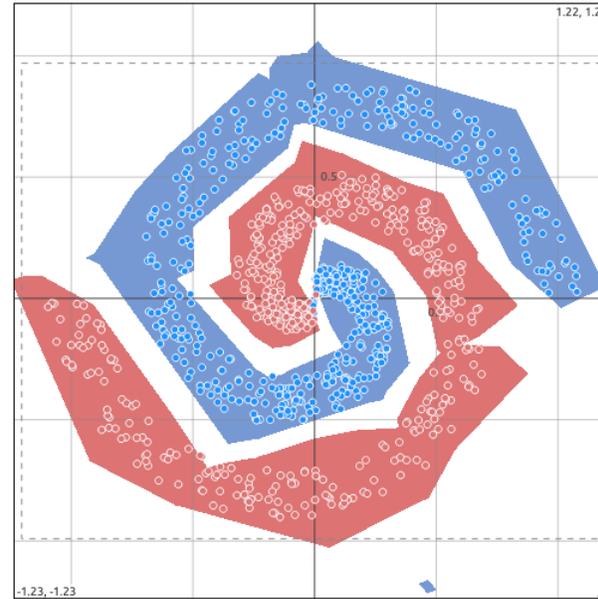
# Влияние порога доверия



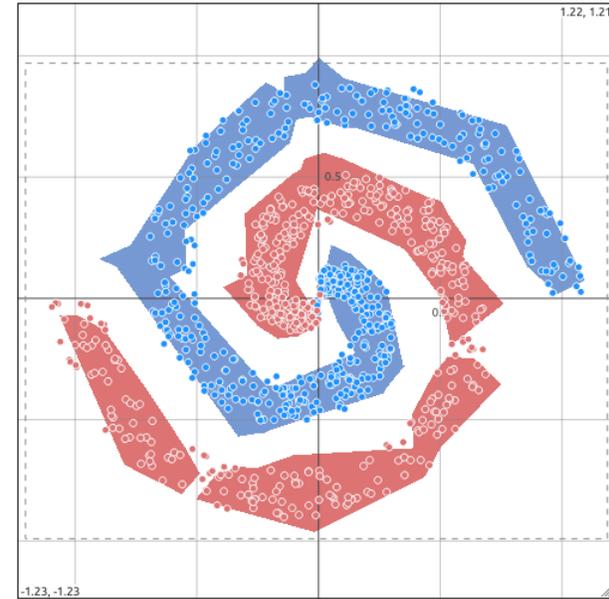
$$\beta = 0$$



$$\beta = 0.1$$



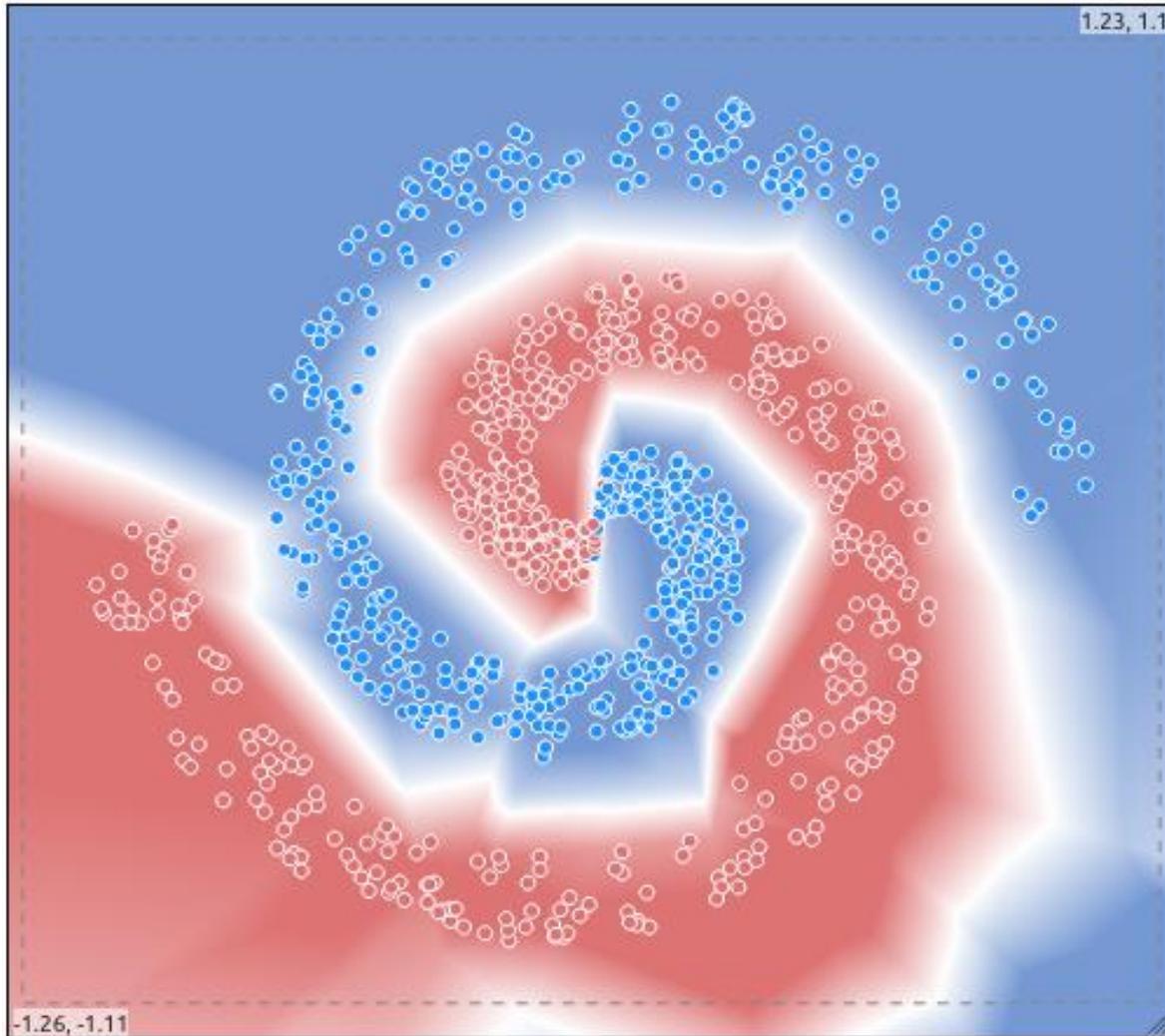
$$\beta = 0.3$$



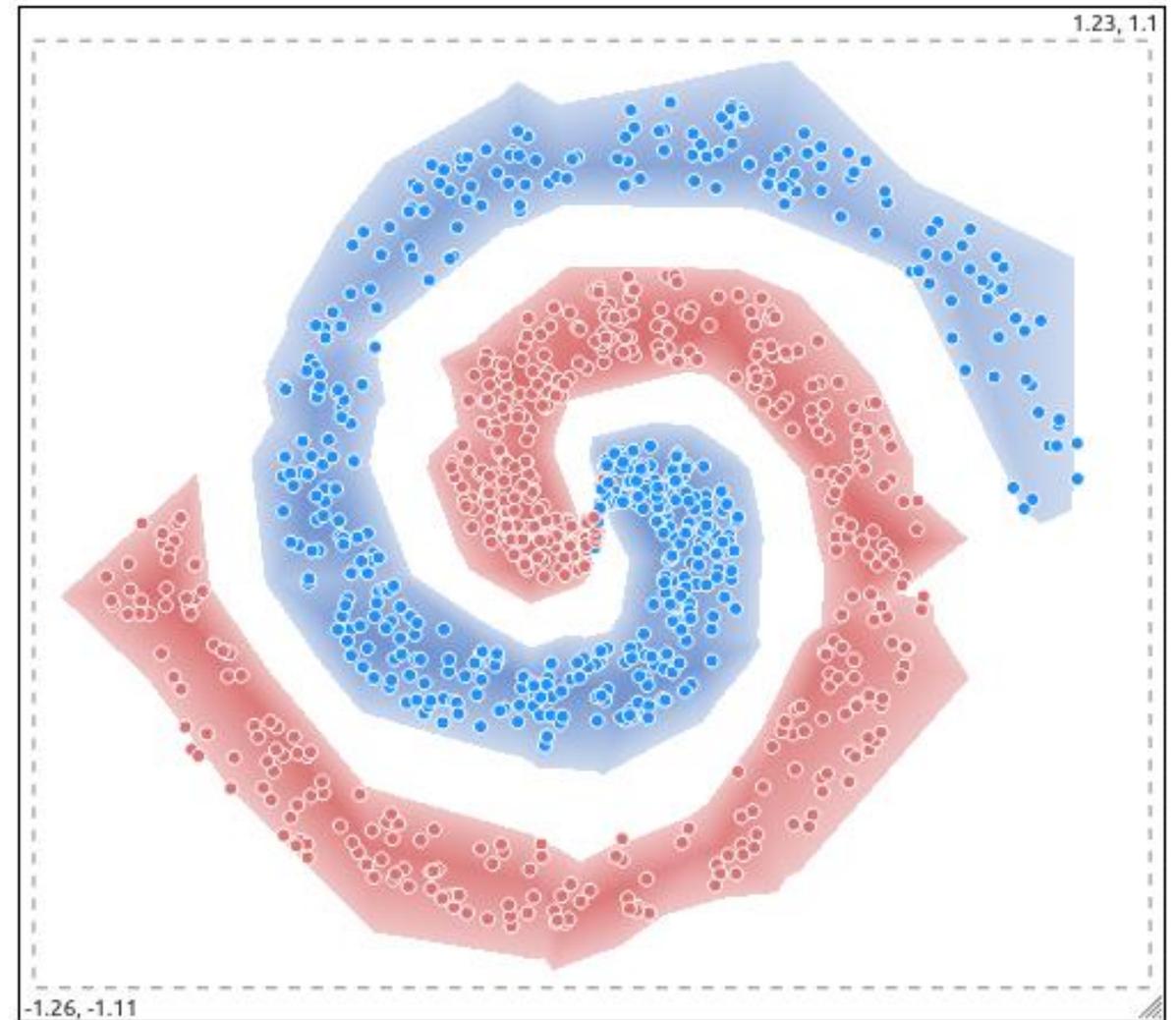
$$\beta = 0.5$$

При увеличении  $\beta$  область отказов расширяется (обозначена белым цветом), что соответствует желаемому поведению системы в условиях ограниченной уверенности модели.

# Поведение вне носителя



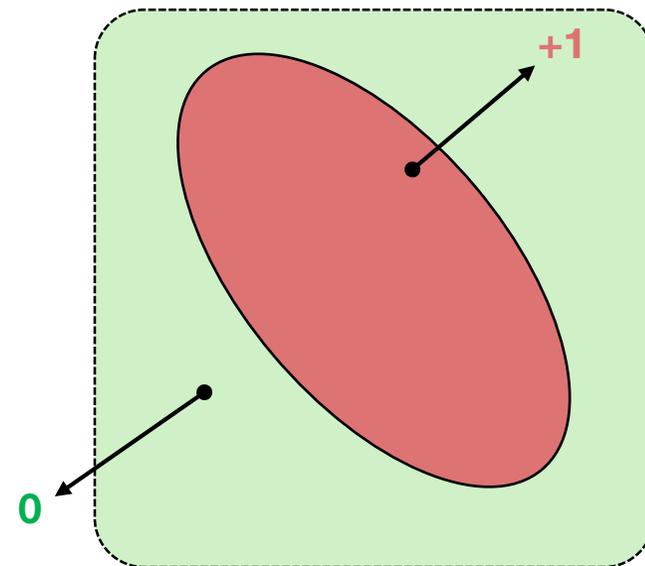
**Классический классификатор**  
(даёт решение в любой точке)



**Модифицированный классификатор**  
(отказывается от классификации)

# Метод унарной классификации

- 1 Оставим только наблюдения класса +1 и фоновые.
- 2 На каждый класс **свой унарный** классификатор  $c^{(i)}(X)$ .
- 3 Классификатор является многослойным **персептроном**.
- 4 Отказ, если предсказание меньше порога  $\beta$ :  $c^{(i)}(X) < \beta$ .



## Преимущества:

### 01 Устойчивость к дисбалансу

- Обучение на **положительных** объектах класса и фоне в равных пропорциях.
- Исключается **влияние** других классов.
- Устойчивость к **неоднородности** в данных.

### 02 Векторное предсказание

- Каждый классификатор выдаёт  $c^{(i)}(X)$  – оценку принадлежности к классу  $i$ .
- Класс с максимальным значением или отказ при

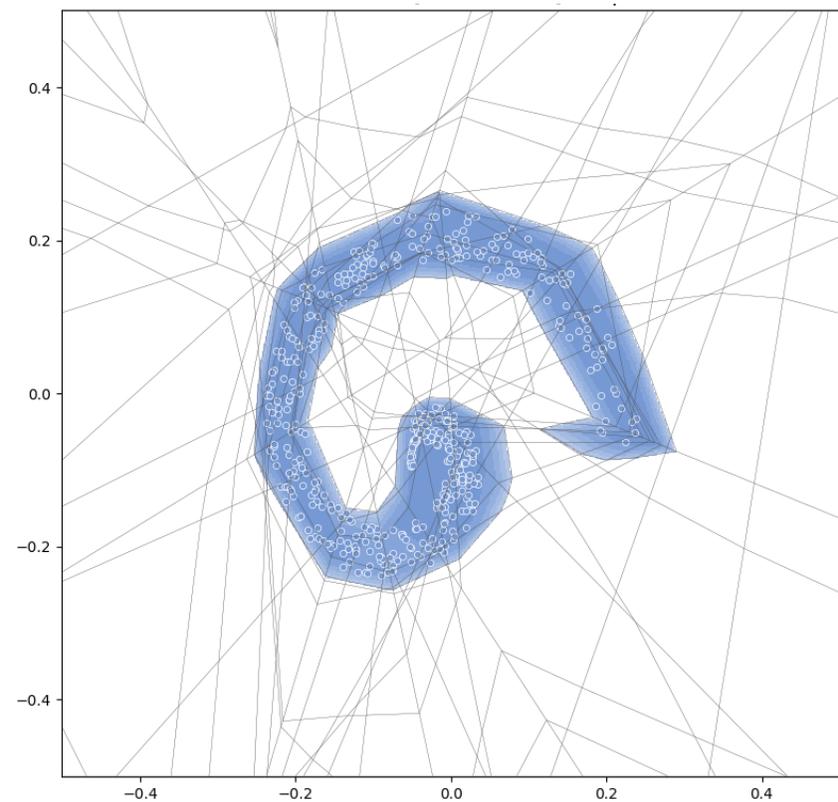
$$\max_{i=1\dots C} c^{(i)}(x) \leq \beta$$

### 03 Модульность архитектуры

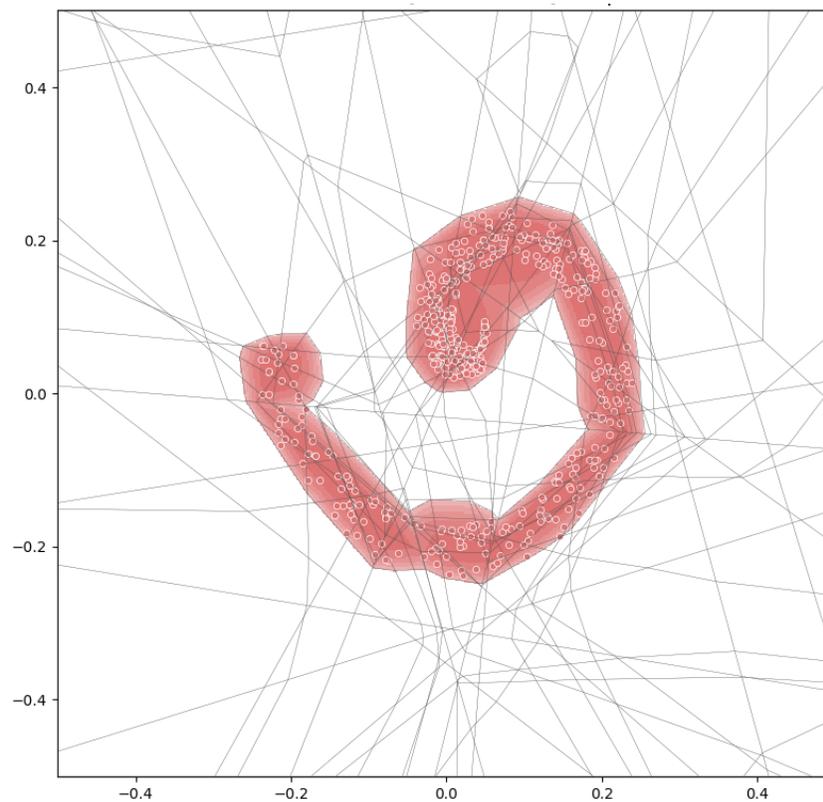
- Классификаторы **независимы**.
- Возможны **разные архитектуры** и глубина для разных классов.
- Адаптация под **особенности** классов.

**Ограничения:** размерность пространства признаков невелика ( $d \leq 30$ )

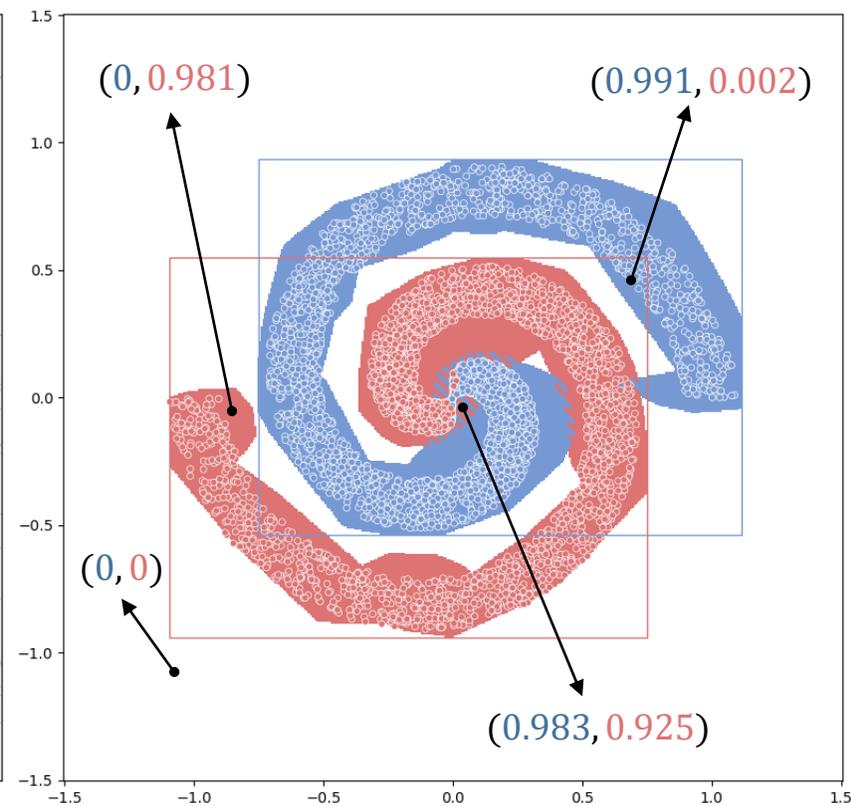
# Унарная бинарная классификация



модель для **первого** класса  $c^{(1)}(x)$



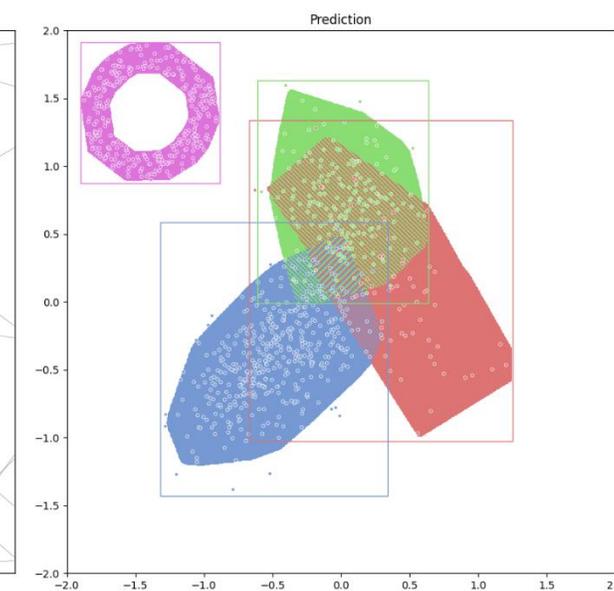
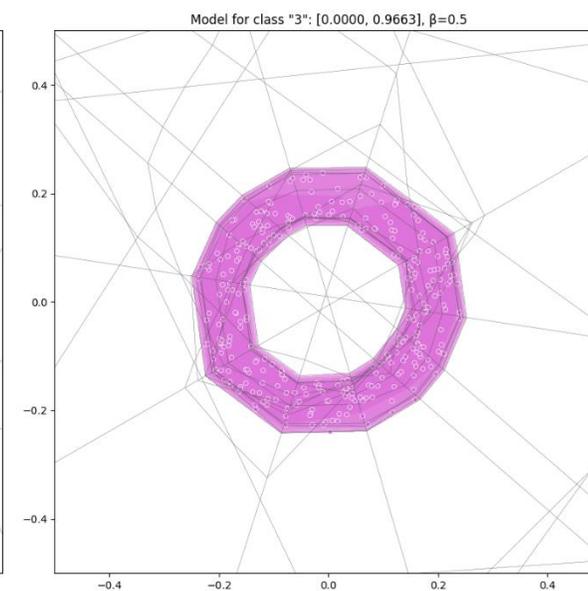
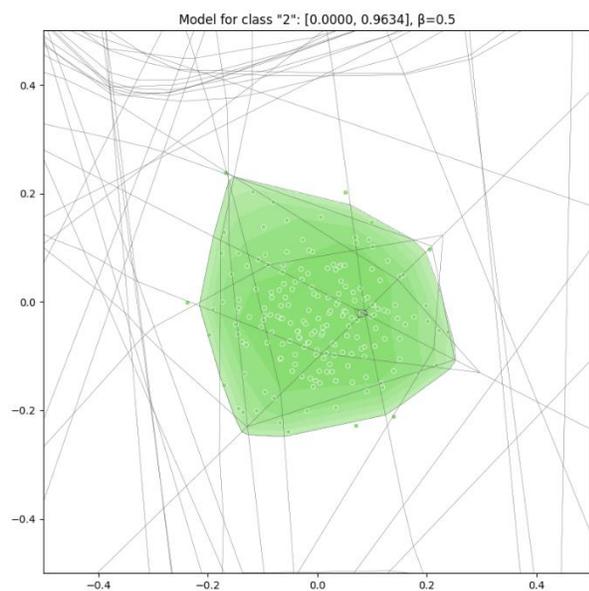
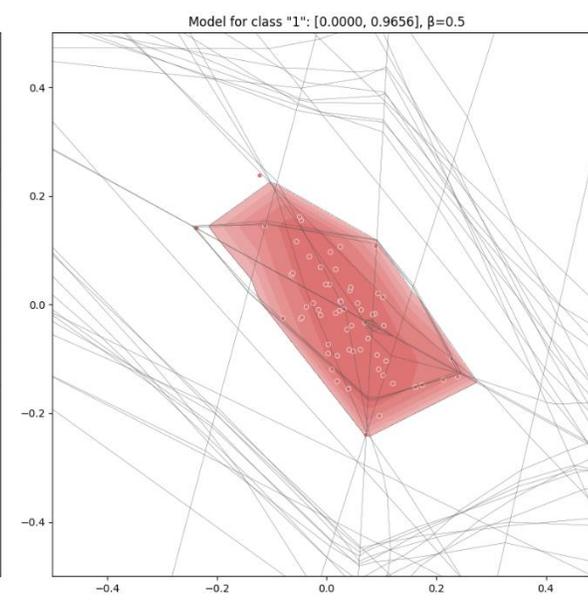
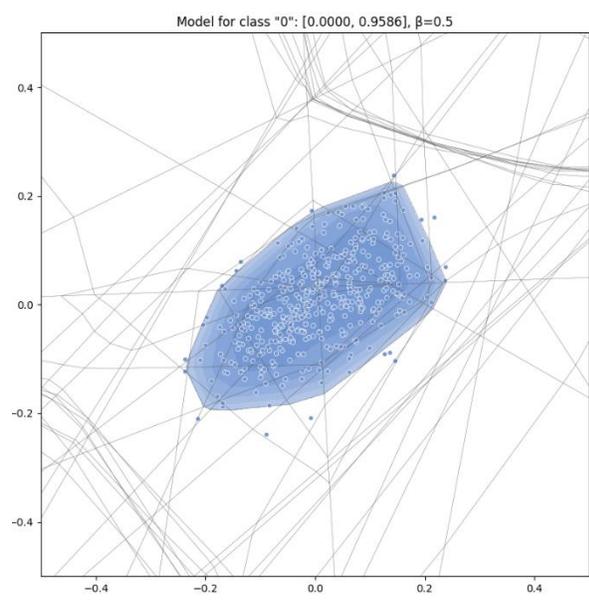
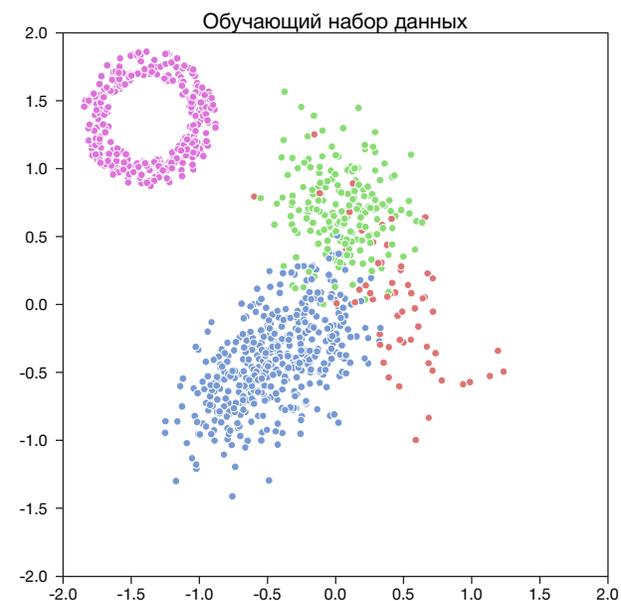
модель для **второго** класса  $c^{(2)}(x)$



**ИТОГОВЫЙ** классификатор  $c(x)$

- 1 Для каждого класса **своя область** принятия решений и **модель**.
- 2 **Вне области** происходит **отказ** от распознавания.
- 3 В смеси распределений требуется **дополнительный анализ**.

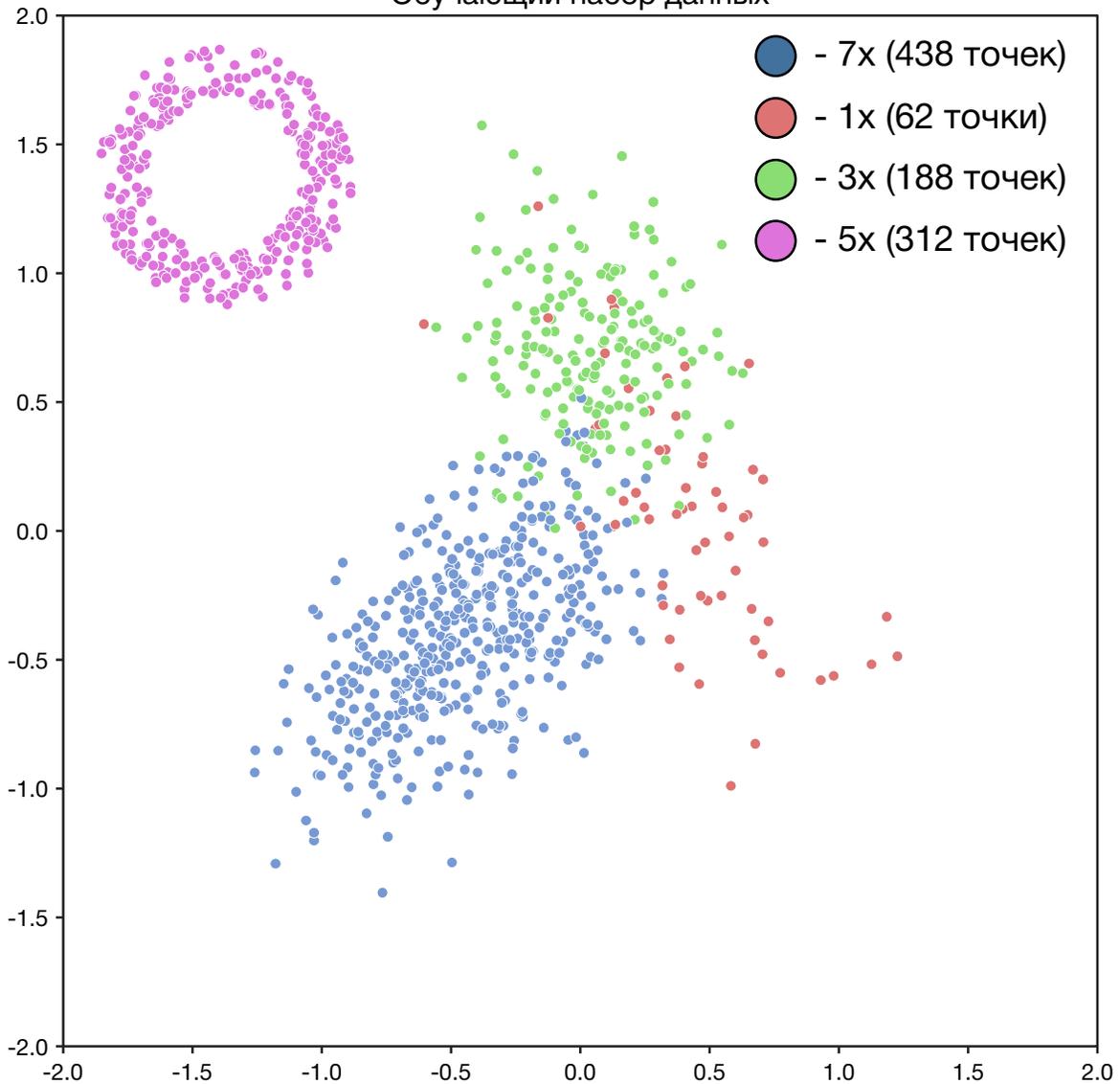
# Унарная классификация: устойчивость к дисбалансу



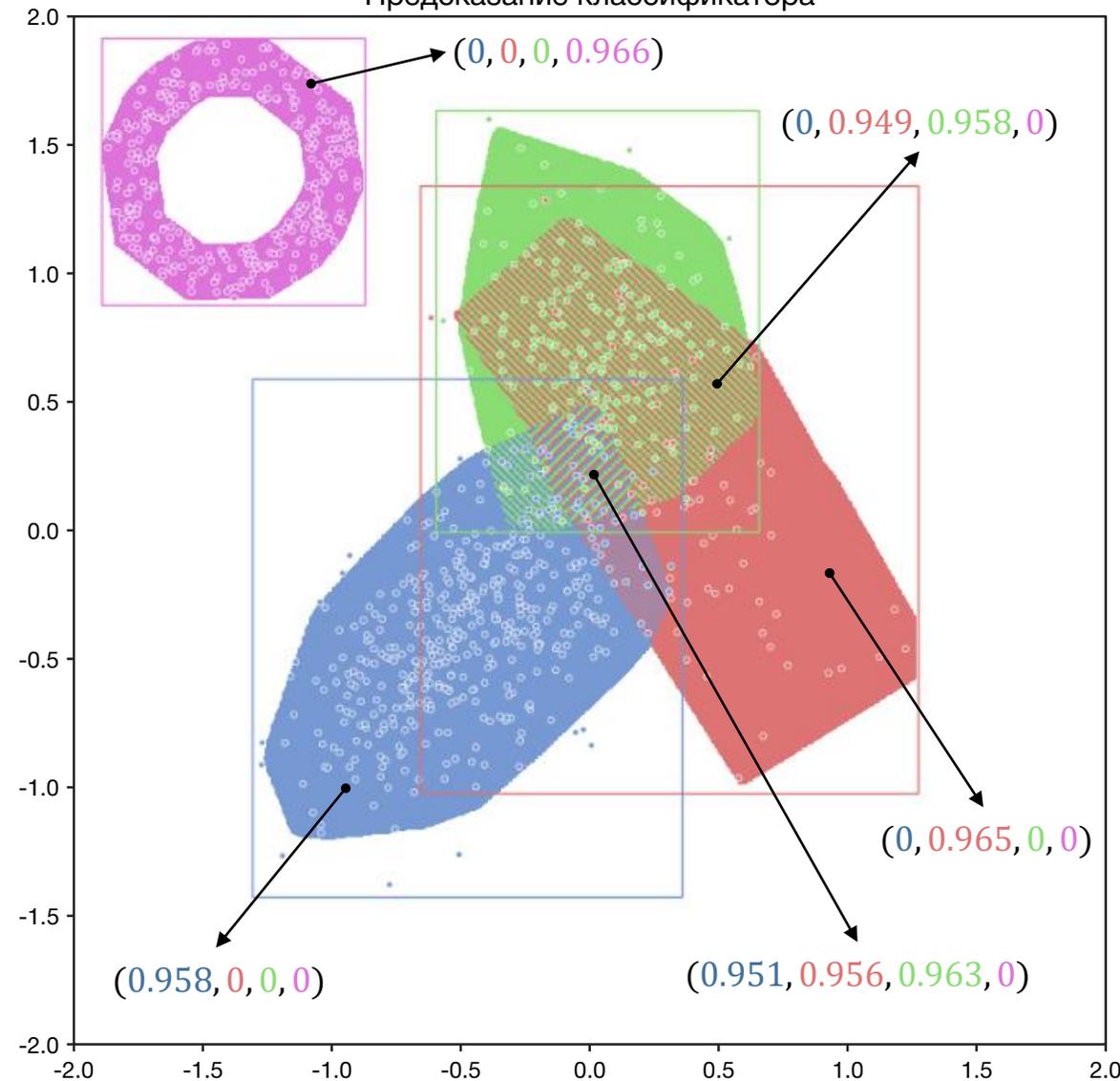
- - 7x (438 точек)
- - 1x (62 точки)
- - 3x (188 точек)
- - 5x (312 точек)

# Унарная классификация: устойчивость к дисбалансу

Обучающий набор данных



Предсказание классификатора



# Унарная классификация: проблемы

**Проклятие размерности** ограничивает применимость для больших размерностей пространства.

При увеличении размерности пространства число ячеек  $N$  увеличивается существенно быстрее, чем растёт объём обучающей выборки  $n$ , доминируют «одноэлементные» и «пустые» ячейки, понятие **частоты** просто **исчезает**.

При  $d = 10$  необходимы **миллионы точек фона**, чтобы в достаточной мере заполнить компакт.

Можно ли как-то побороть эту проблему?

# SLAP – Simple Linear Attack for Perceptron

Простая линейная атака на персептрон



- Cat – 2%
- Dog – 98%

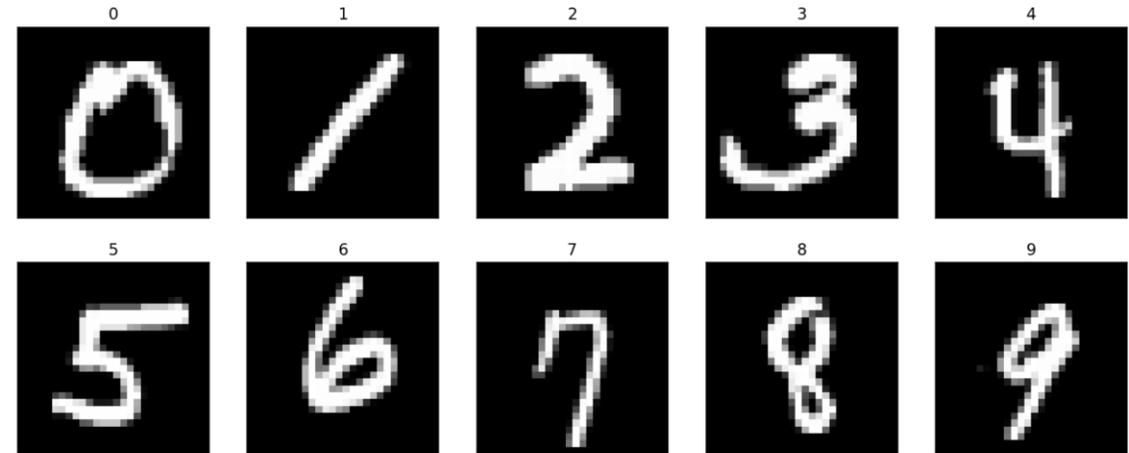
# Постановка задачи

## Что имеем?

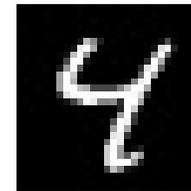
- полносвязная сеть для классификации изображений с **кусочно-линейными** функциями активации (*ReLU*, *Abs*, ...);
- пара изображений – целевое ( $x_t$ ) и атакуемое ( $x_a$ ) (в виде векторов);
- магия линейной алгебры.

## Что хотим получить?

- изображение ( $x$ ), имеющее тот же выход сети (или тот же класс) как и целевое, но похожее на атакуемое.



**Целевое** изображение –  $x_t$   
Выход модели –  $y_t: x_1, x_2, x_3, \dots$



**Атакуемое** изображение –  $x_a$   
Выход модели –  $y_a: y_1, y_2, y_3, \dots$



**Атакующее** изображение –  $x$   
Выход модели –  $y: x_1, x_2, x_3, \dots$

# Атака: один слой

Рассматриваем первый слой сети, количество выходов которого меньше, чем входов:  $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$

$$w_{1,1}x_1 + w_{1,2}x_2 + \dots + w_{1,d}x_d + b_1 = y_1$$

$$w_{2,1}x_1 + w_{2,2}x_2 + \dots + w_{2,d}x_d + b_2 = y_2$$

...

$$w_{c,1}x_1 + w_{c,2}x_2 + \dots + w_{c,d}x_d + b_c = y_c$$

- $d$  – размер изображения,  $C$  – количество классов (и размер слоя)
- $\mathbf{W}$  – матрица весов ( $C \times d$ ),  $\mathbf{b}$  – вектор смещений ( $C \times 1$ )
- $\mathbf{x}_t$  – целевое изображение ( $d \times 1$ )
- $\mathbf{x}_a$  – атакуемое изображение ( $d \times 1$ )
- $\mathbf{x}$  – атакующее изображение ( $d \times 1$ )
- $\mathbf{y}_t$  – выход слоя ( $C \times 1$ )

# Атака: один слой

Решаем **СЛАУ** относительно коэффициентов первого слоя и ищем решения, удовлетворяющие следующим условиям:

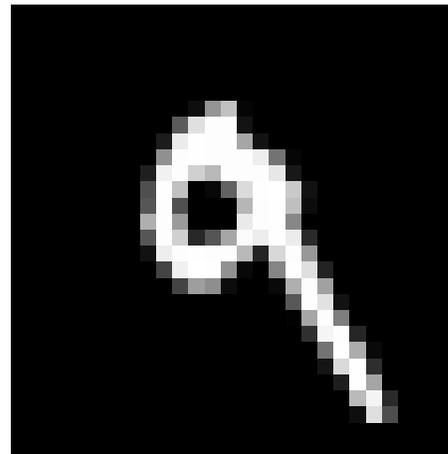
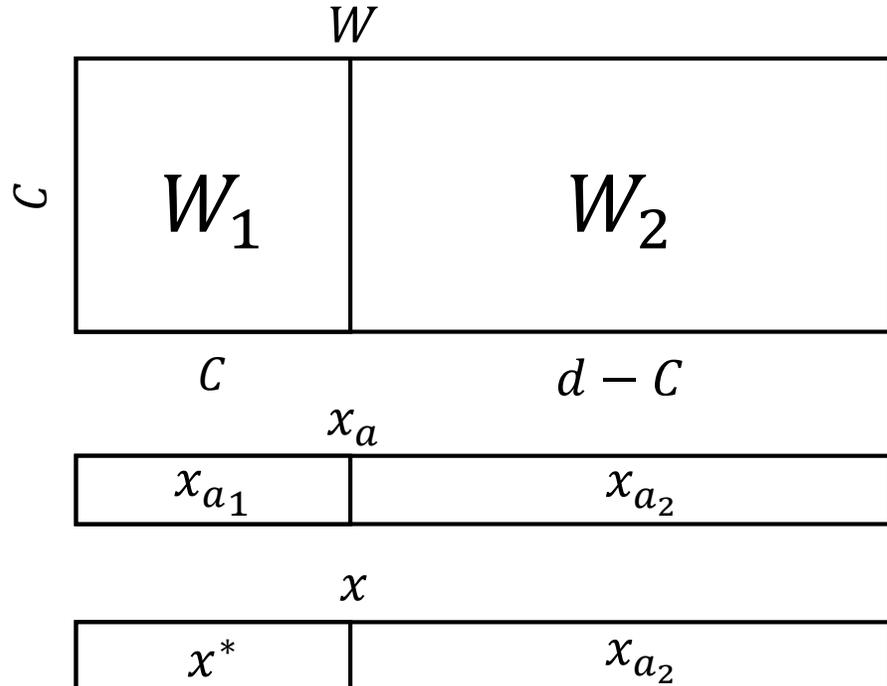
- $Wx + b = y_t$  (равенство выходов сети)
- $\|x_a - x\| \rightarrow \min$  (максимальная похожесть на атакуемое)
- $\|x_t - x\| > 0$  (**в идеале** не получить целевое)
- $0 \leq x_i \leq 1$  (**в идеале** остаться в диапазоне значений)

**Альтернативная цель:** равенство классов

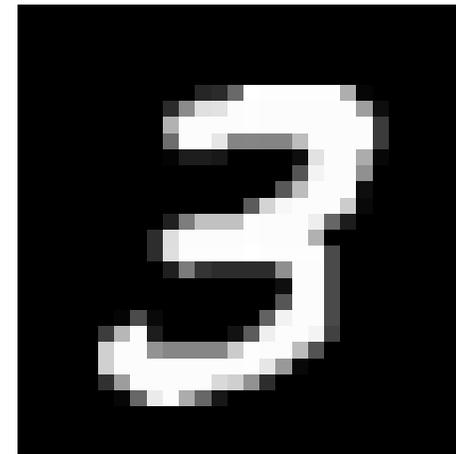
- $\operatorname{argmax}(Wx + b) = \operatorname{argmax}(y_t)$

# Атака: один слой, без учёта входного диапазона

- выберем подматрицу  $W_1$  размера  $C \times C$  (первые или случайные столбцы);
- оставшуюся подматрицу размера  $C \times (d - C)$  будем называть  $W_2$ ;
- аналогично разделим  $x_a$  на  $x_{a_1}$  и  $x_{a_2}$ ;
- первые  $C$  значений атаки получим как  $x^* = W_1^{-1} \cdot (b^T - W_2 \cdot x_{a_2})$ ;
- атакующее изображение получается **конкатенацией**  $x^*$  и  $x_{a_2}$ .



$x_t \in [0, 1]$



$x_a \in [0, 1]$



$x \in [-68, 82]$

# Атака: один слой, без учёта входного диапазона

## Достоинства:

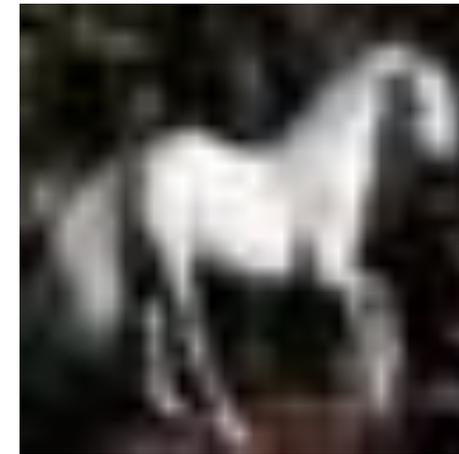
- Всегда работает (вероятность необратимости  $W_1$  крайне мала);
- Минимально заметна (работает исключительно в пределах тензора).

## Недостатки:

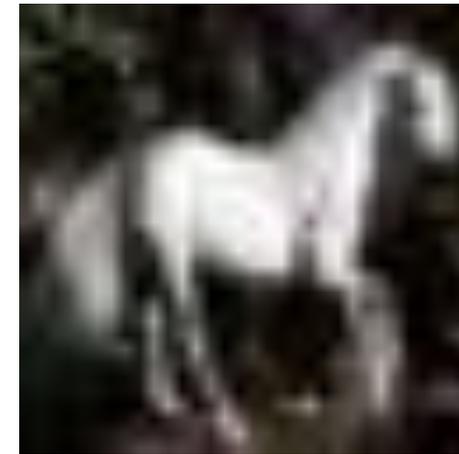
- Применима только при доступе к тензорам;
- Не учитывается диапазон входных значений;
- Не воспроизводима при сохранении  $x$  (может быть **достоинством**, поскольку визуально не будет видно заметных дефектов).



$$x_t \in [0, 1]$$



$$x_a \in [0, 1]$$



$$x \in [-938, 919]$$

# Атака: один слой – QP решатель

## QP task

$$\begin{cases} \frac{1}{2} x^T P x + q^T x \rightarrow \min \\ Ax = b \\ Gx \leq h \\ l_b \leq x_i \leq u_b \end{cases}$$

$$\begin{aligned} P &= E, q = -x_a \\ A &= W, b = y_t - b \\ G &= 0, h = 0 \\ l_b &= 0, u_b = 1 \end{aligned}$$

## QP attack

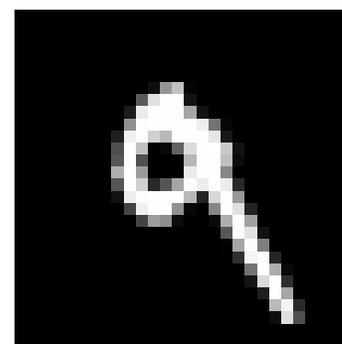
$$\begin{cases} \frac{1}{2} x^T x - x_a^T x \rightarrow \min \\ Wx = y - b \\ 0 \leq x_i \leq 1 \end{cases}$$

## Достоинства:

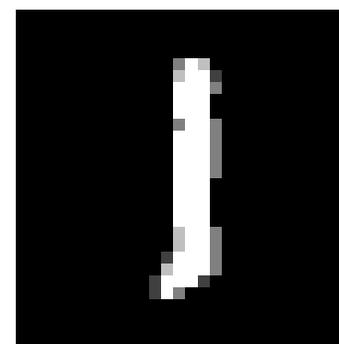
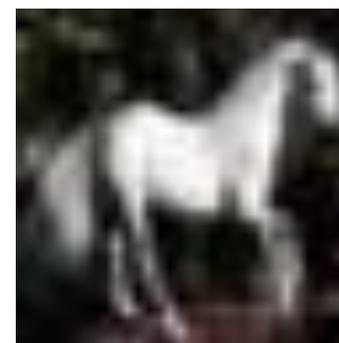
- Сохранение «не убивает» атаку;
- Учитывается входной диапазон.

## Недостатки:

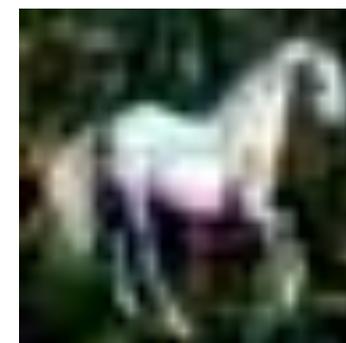
- Не всегда применима;
- Требуется вносить больше искажений.



$x_t \in [0, 1]$



$x_a \in [0, 1]$



$x \in [0, 1]$

# Атака: много слоёв

Пусть модель – трёхслойный персептрон с *Abs* активацией:

$$y = \sigma(W_3 \cdot \sigma(W_2 \cdot \sigma(W_1 x + b_1) + b_2) + b_3)$$

**Проблема:** присутствуют нелинейности, **не позволяющие** использовать QR решатель.

# Атака: много слоёв

Пусть модель – трёхслойный персептрон с *Abs* активацией:

$$y = \sigma(W_3 \cdot \sigma(W_2 \cdot \sigma(W_1 x + b_1) + b_2) + b_3)$$

**Проблема:** присутствуют нелинейности, **не позволяющие** использовать QR решатель.

**Решение:** используем **кусочную линейность** функций активации: фиксируем знаки и раскрываем нелинейности как в eXVTree.

# Атака: много слоёв – раскрываем первый слой

Пусть все значения первого слоя **неотрицательны**:

$$\begin{cases} W_1 x + b_1 \geq 0 \\ y = f(W_3 \cdot f(W_2 W_1 x + W_2 b_1 + b_2) + b_3) \end{cases}$$

# Атака: много слоёв – раскрываем второй слой

Пусть все значения второго слоя **неположительны**:

$$\left\{ \begin{array}{l} W_1 x + b_1 \geq 0 \\ W_2 W_1 x + W_2 b_1 + b_2 \leq 0 \\ y = f(-(W_3 W_2 W_1 x + W_3 W_2 b_1 + W_3 b_2) + b_3) \end{array} \right.$$

# Атака: много слоёв – упрощаем уравнения

Введём новые переменные:

- $W_{21} = W_2 W_1$
- $b_{21} = W_2 b_1 + b_2$
- $W_{321} = -W_3 W_{21}$
- $b_{321} = b_3 - W_3 b_{21}$

Имеем **однослойный случай** с системой неравенств:

$$\begin{cases} W_1 x + b_1 \geq 0 \\ W_{21} x + b_{21} \leq 0 \\ y = f(W_{321} x + b_{321}) \end{cases}$$



$x_t \in [0, 1]$



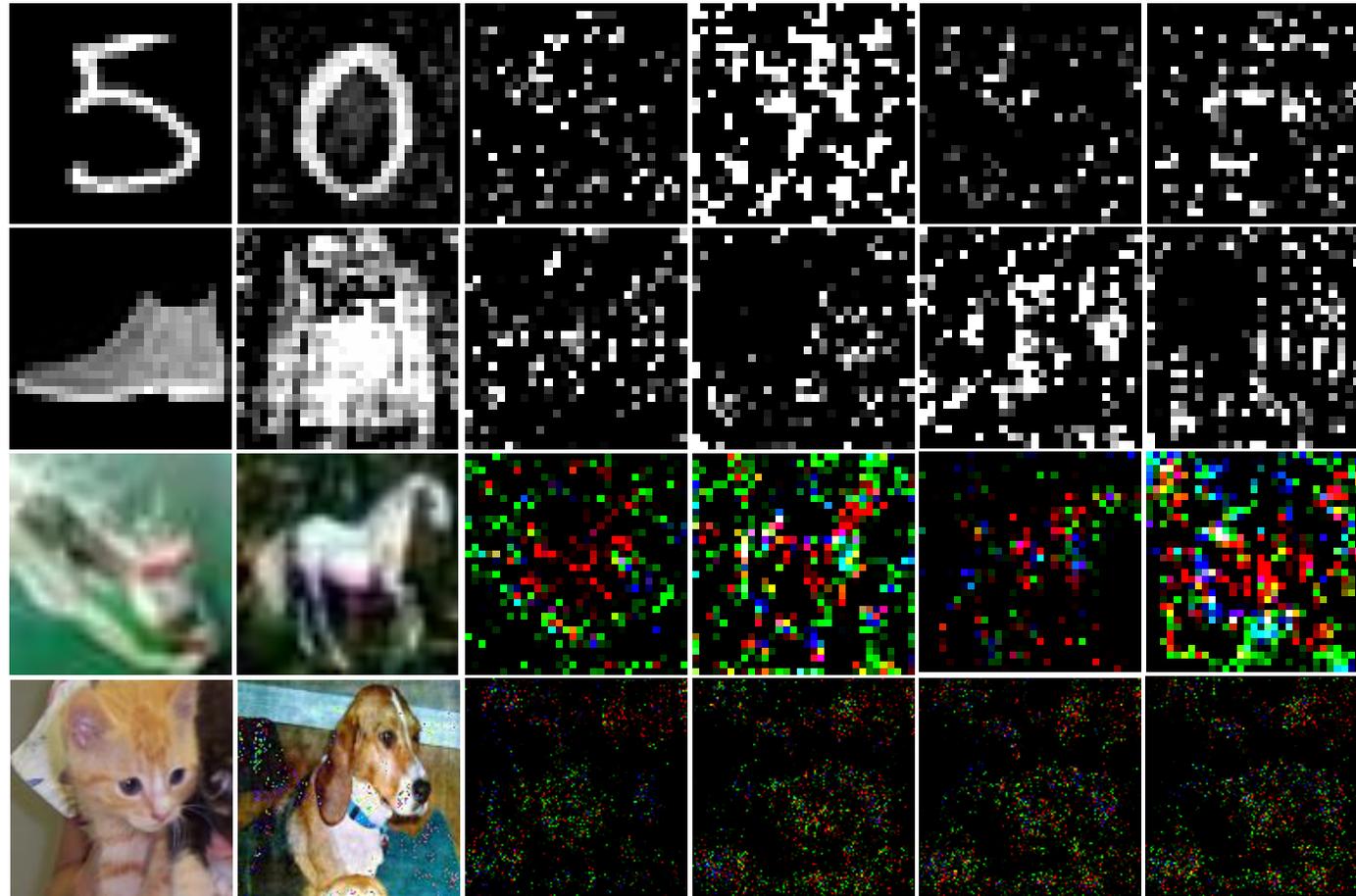
$x_a \in [0, 1]$



$x \in [0, 1]$

# Бонус: произвольная атака с выходом, равным $y_t$

- Заменяем  $P$  **случайной диагональной** матрицей;
- Заменяем  $q$  **случайным** вектором;
- Запустим QR решатель.



**все** эти изображения имеют **тот же выход модели**, что и  $x_t$

# SLAP атака: выводы

- Персептрон – такая **простая** модель, но такая **уязвимая** перед атаками.
- **Однослойные** модели **более уязвимы**, чем многослойные.
- **Одному** и тому же выходу сети почти всегда соответствует **множество** входов.
- Предлагаемая атака позволяет получить **все возможные** входы (с применением аналитических решателей).

# Что на счёт свёрточного слоя?

Свёрточный слой **сводится к полносвязному**, если аккуратно выписать все получаемые уравнения:

$$\begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix} \times \begin{bmatrix} f_1 & f_2 \\ f_3 & f_4 \end{bmatrix} = \begin{bmatrix} f_1 a_1 + f_2 a_2 + f_3 a_4 + f_4 a_5 \\ f_1 a_2 + f_2 a_3 + f_3 a_5 + f_4 a_6 \\ f_1 a_4 + f_2 a_5 + f_3 a_7 + f_4 a_8 \\ f_1 a_5 + f_2 a_6 + f_3 a_8 + f_4 a_9 \end{bmatrix}$$
$$\begin{bmatrix} f_1 & f_2 & 0 & f_3 & f_4 & 0 & 0 & 0 & 0 \\ 0 & f_1 & f_2 & 0 & f_3 & f_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & f_1 & f_2 & 0 & f_3 & f_4 & 0 \\ 0 & 0 & 0 & 0 & f_1 & f_2 & 0 & f_3 & f_4 \end{bmatrix} \times \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \end{bmatrix}$$

# Менее строгая цель

Использование **приблизительного сходства** выходов модели вместо точного соответствия (достигается с помощью матрицы  $G$  и вектора  $h$ ):

$$\begin{cases} Wx + b \geq y_t - \varepsilon \\ Wx + b \leq y_t + \varepsilon \end{cases} \longrightarrow \begin{cases} -Wx \leq b - y_t + \varepsilon \\ Wx \leq y_t - b + \varepsilon \end{cases}$$

$$G = \text{concat}(-Wx, Wx)$$

$$h = \text{concat}(b - y_t + \varepsilon, y_t - b + \varepsilon)$$

# Как следует выбирать знаки выходов слоёв?

- Перебор всех возможных комбинаций (**крайне непрактично**).
- Использование знаков, получаемых во время распространения по сети целевого и атакуемого изображений (или их комбинации с шумом).

# Вопросы?