



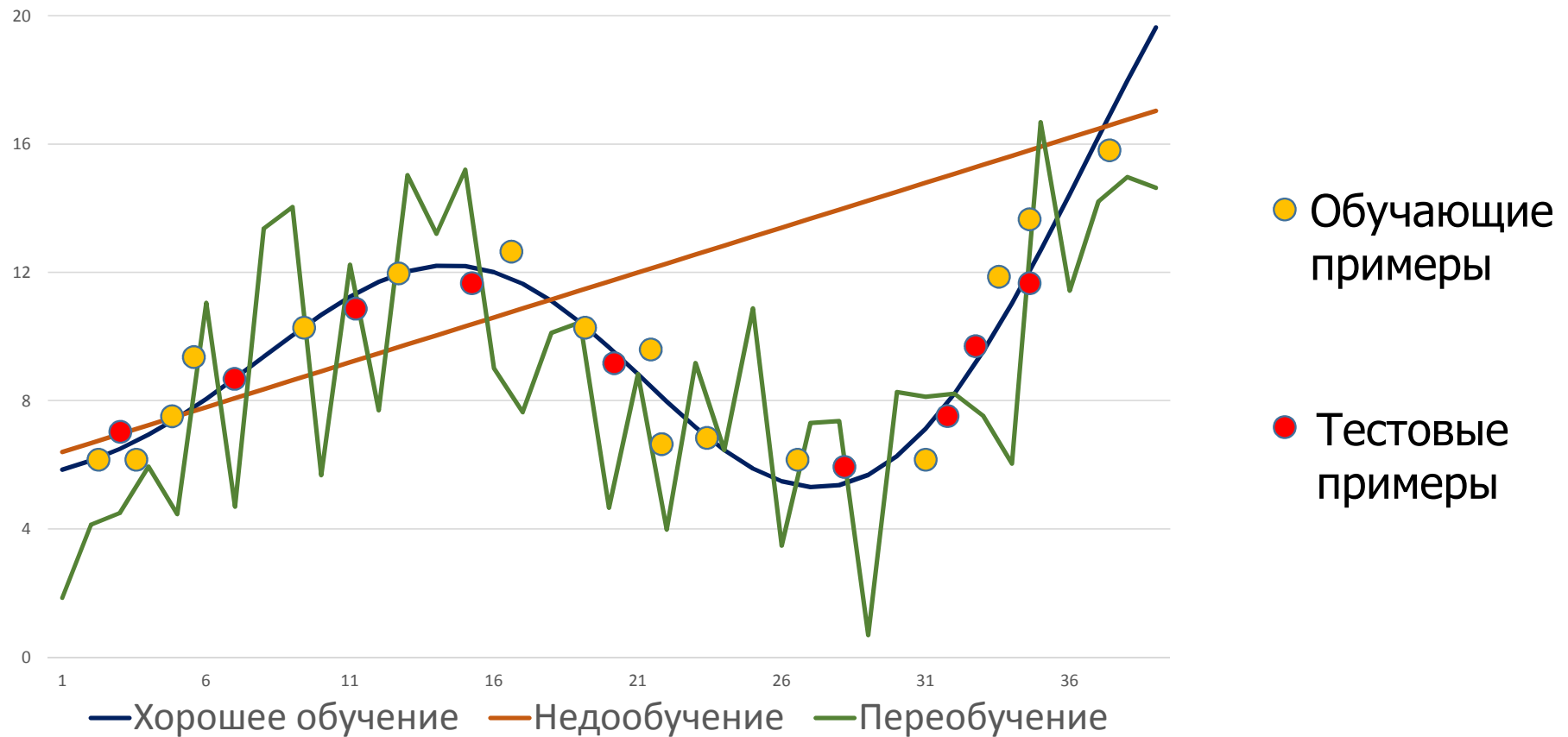
Национальный
исследовательский
Томский
государственный
университет



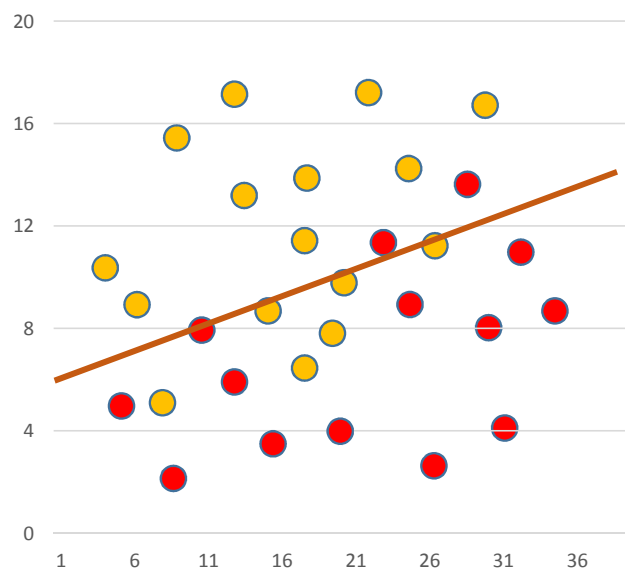
Подбор переменных и алгоритмов для модели

Сергей Аксёнов, к.т.н., доцент кафедры
Теоретических основ информатики ТГУ

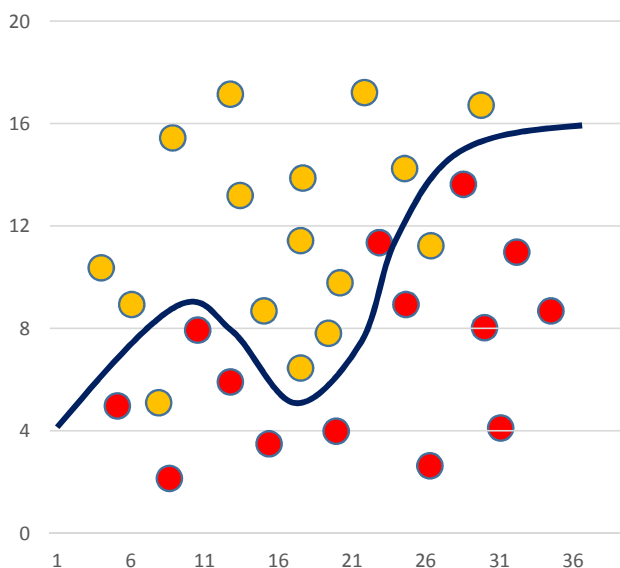
Плохое и хорошее обучение (регрессия)



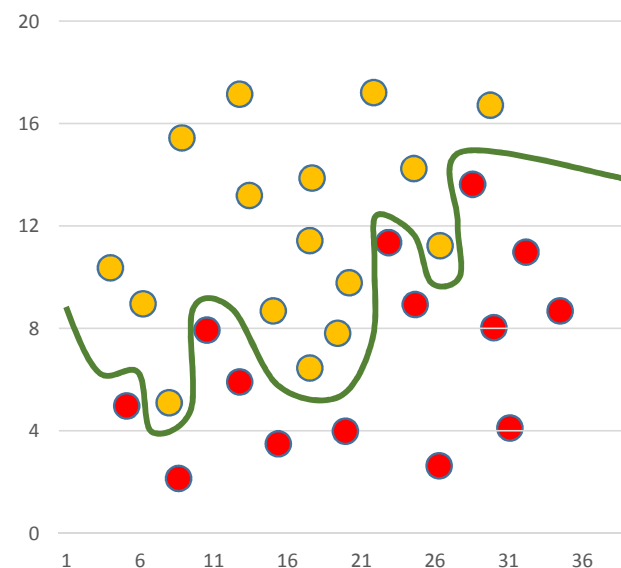
Плохое и хорошее обучение (Классификация)



Недообучение



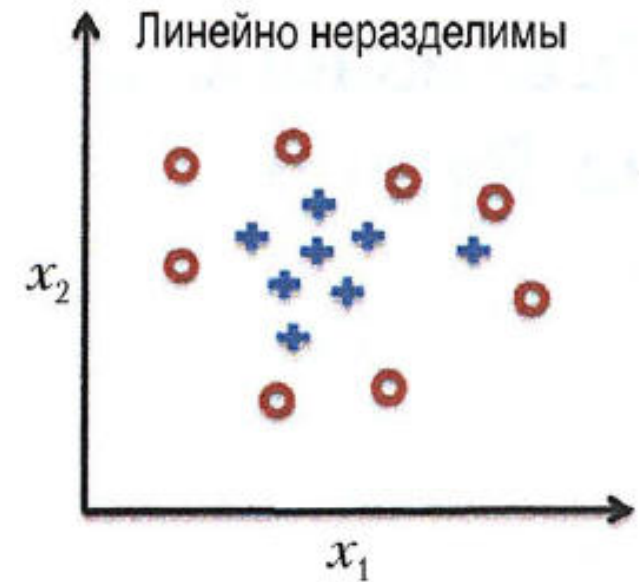
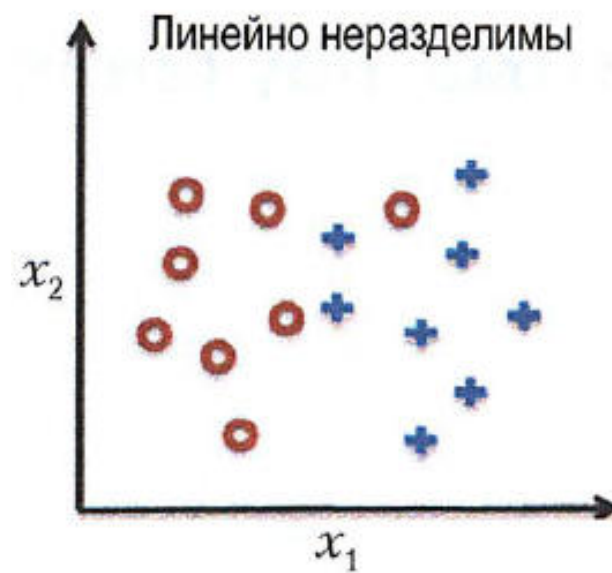
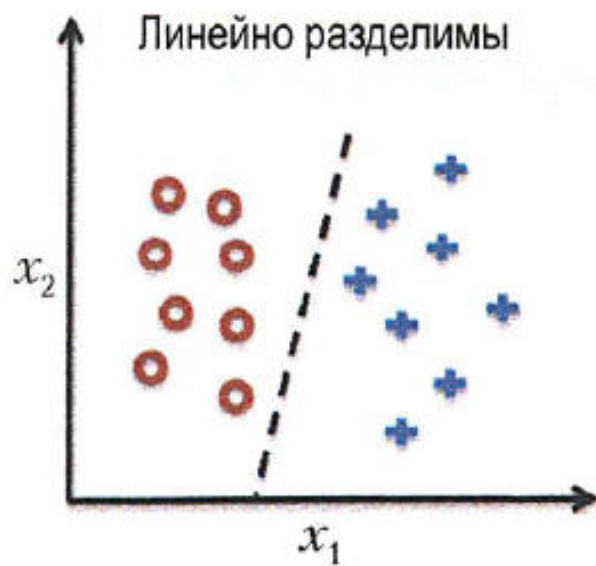
Хорошее обучение



Переобучение

● Класс А ● Класс В

Линейно разделимые и линейно неразделимые классы



Бинарный линейный классификатор

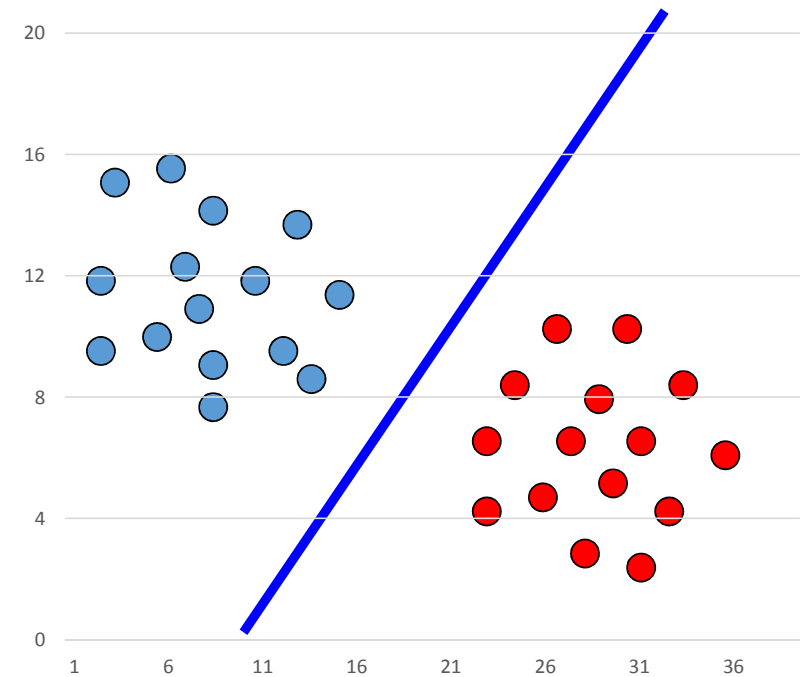
Результат обучения: входной вектор относится либо к положительному ($\hat{y}=+1$), либо отрицательному ($\hat{y}=-1$) классу

Вектор признаков:

$$x = (x_1, x_2, x_3, \dots, x_N)$$

Выход модели:

$$\hat{y} = \hat{y}(x, w) = \text{sign}\left(w_0 + \sum_i^N w_i x_i\right) = \text{sign}(w^T x)$$



Логистическая регрессия

Прогнозируют вероятность p_+ отнесения примера x к классу $+1$

$$z = \sum_i^N w_i x_i$$

Функция стоимости

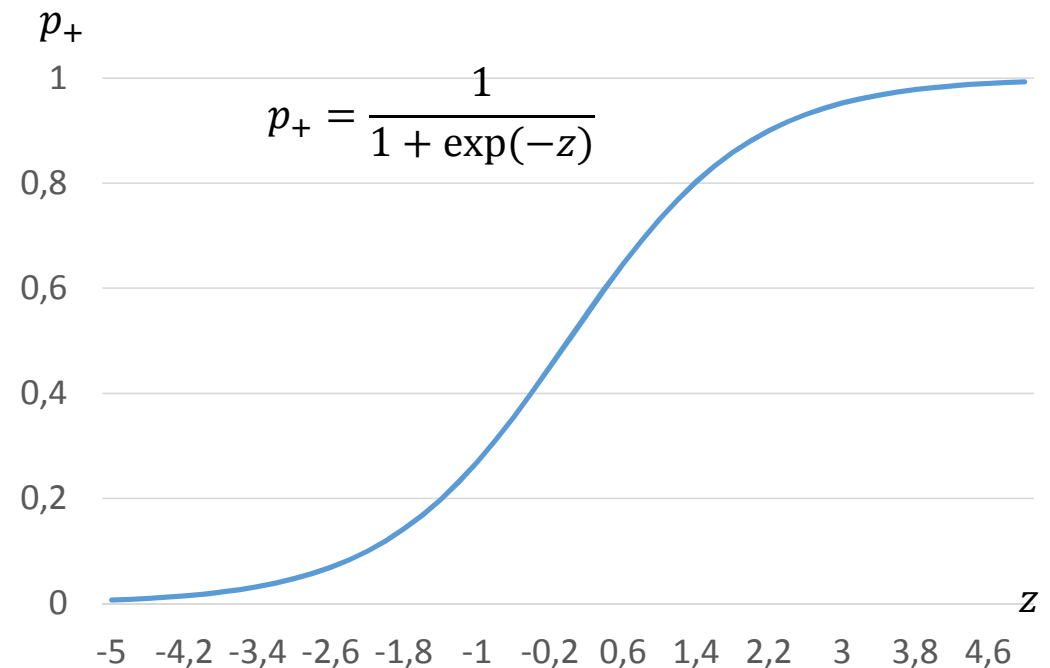
$$J(w) = \sum_i \frac{1}{2} (\phi(z^{(i)}) - y^{(i)})^2$$

Функция правдоподобия

$$L(w) = P(y | x; w) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; w) = \prod_{i=1}^n (\phi(z^{(i)}))^{y^{(i)}} (1 - \phi(z^{(i)}))^{1-y^{(i)}}$$

Логарифмическая функция правдоподобия

$$l(w) = \log L(w) = \sum_{i=1}^n \left[y^{(i)} \log(\phi(z^{(i)})) + (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \right]$$



Регуляризация в логистической регрессии

Функция, используемая для поиска параметров модели

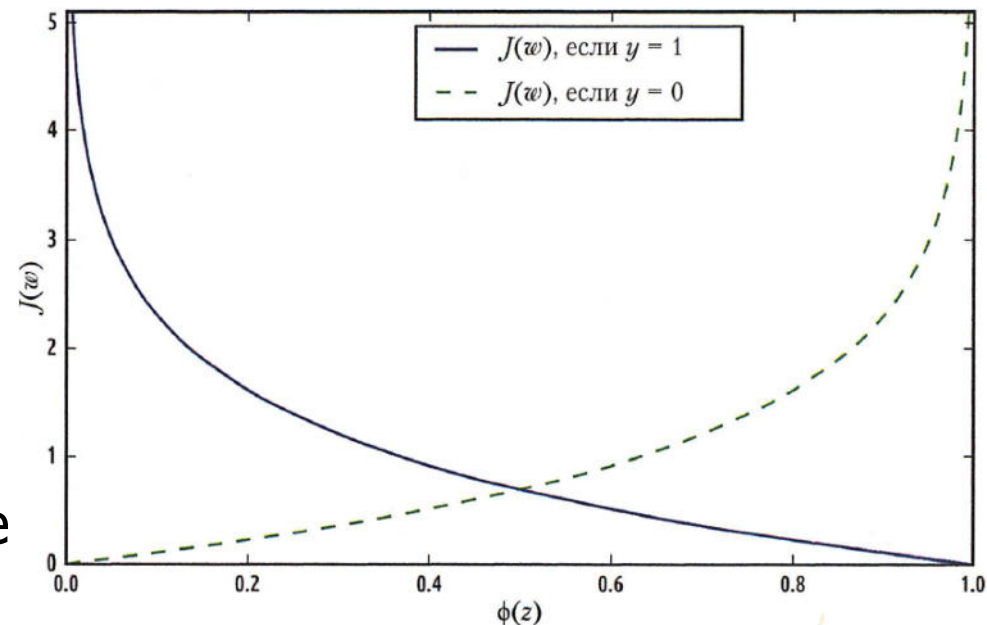
$$J(\mathbf{w}) = \sum_{i=1}^n \left[-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \right]$$

Регуляция переобучения выполняется при помощи регуляризации.
Наложение штрафов на экстремальные значения параметров.

$$\frac{\lambda}{2} \|\mathbf{w}\|^2 = \frac{\lambda}{2} \sum_{j=1}^m w_j^2 \quad \text{L2 - регуляризация}$$

Новая функция,
учитывающая штрафы

$$J(\mathbf{w}) = \sum_{i=1}^n \left[-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \right] + \frac{\lambda}{2} \|\mathbf{w}\|^2$$



Дерево решений

Разбиение данных на подмножества, приводящему к самому большому приросту информации (получению однородных регионов решения)

Функция прироста информации:

$$IG(D_p, f) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

Меры неоднородности:

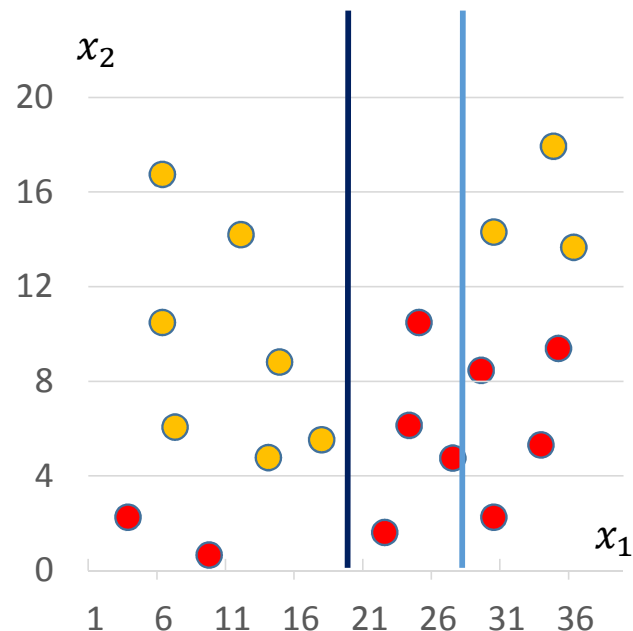
Энтропия: $I_G(t) = 1 - \sum_{i=1}^c p(i|t)^2$

Мера неопределенности Джини: $I_H(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$

Ошибка классификации: $I_E(t) = 1 - \max(p(i|t))$

$p(i|t)$ -доля образцов, принадлежащая классу i для узла t

Построение деревьев решений. Пример-1



$$IG(D_p, f) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

В качестве критерия взята ошибка классификации:

$$I_E(t) = 1 - \max(p(i|t))$$

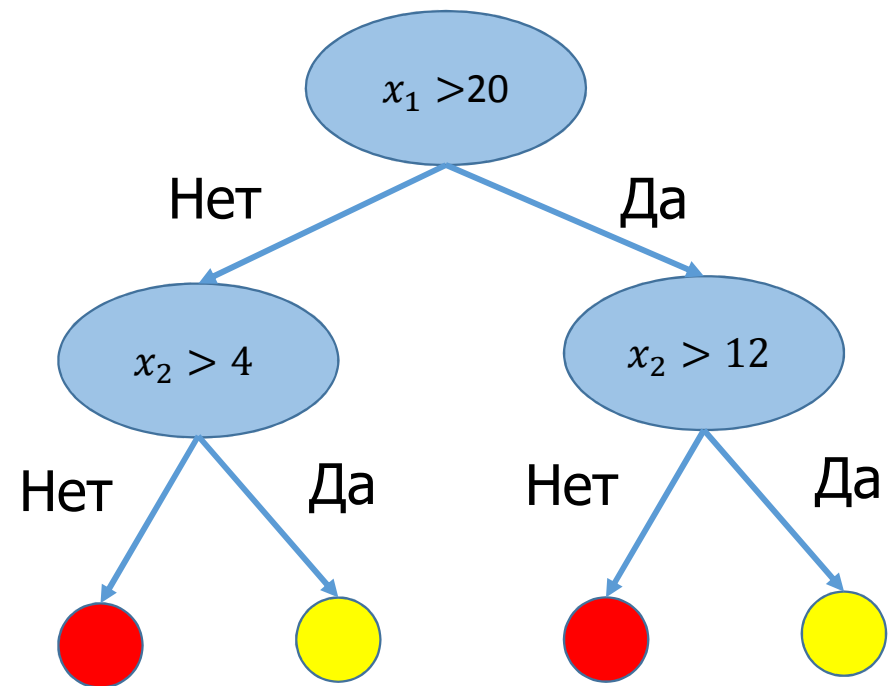
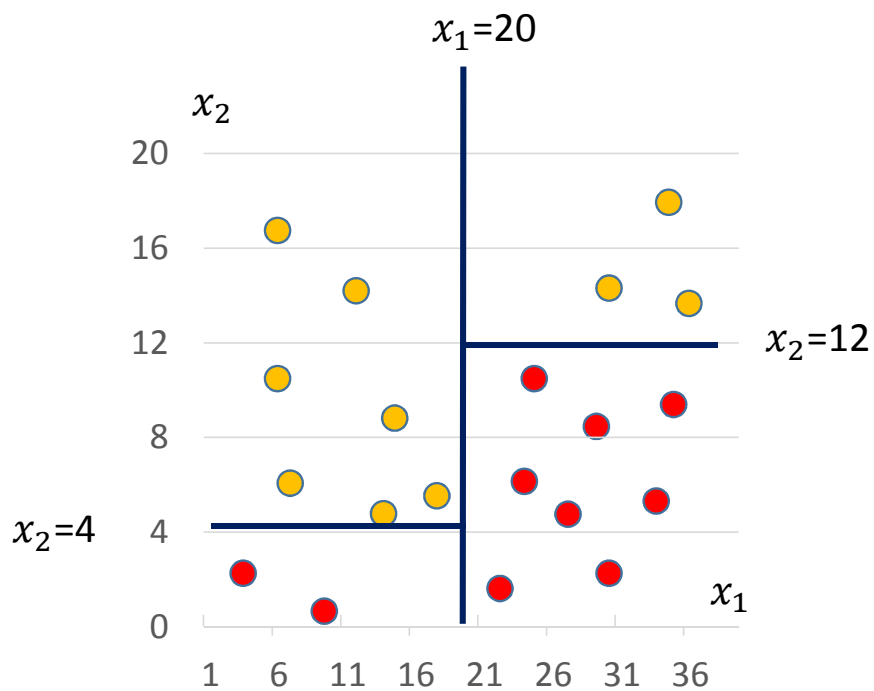
Неоднородность корневого узла:

$$I(D_0) = 1 - \max\left(\frac{10}{20}, \frac{10}{20}\right) = 1 - 0.5 = 0.5$$

Для расщепления $x_1 = 20$: $IG(D_0, x_1 = 20) = 0.5 - \frac{9}{20} \left(1 - \frac{7}{9}\right) - \frac{11}{20} \left(1 - \frac{8}{11}\right) = 0.25$

Для расщепления $x_1 = 28$: $IG(D_0, x_1 = 28) = 0.5 - \frac{13}{20} \left(1 - \frac{7}{13}\right) - \frac{7}{20} \left(1 - \frac{4}{7}\right) = 0.05$

Построение деревьев решений. Пример-2



Примеры классов из Scikit-learn. Параметры

```
class sklearn.tree.DecisionTreeClassifier(criterion='gini', splitter='best',
max_depth=None, min_samples_split=2, min_samples_leaf=1,
min_weight_fraction_leaf=0.0, max_features=None, random_state=None,
max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None,
class_weight=None, presort=False)
```

```
class sklearn.linear_model.LogisticRegression(penalty='l2', *, dual=False,
tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None,
random_state=None, solver='lbfgs', max_iter=100, multi_class='auto',
verbose=0, warm_start=False, n_jobs=None, l1_ratio=None)
```

Хорошая модель должна использовать не параметры алгоритмов по умолчанию, а исследование результатов алгоритмов с разными параметрами!!!

Случайный лес. Выборки для обучения

#	A1	A2	A3	A4	A5
1					
2					
3					
4					
5					

Исходный
набор:
признаки

Выборка 1

#	A2	A4
1		
3		
4		

Выборка 2

#	A1	A3	A4
3			
5			

Выборка 3

#	A1	A5
1		
2		
3		

Выборка 4

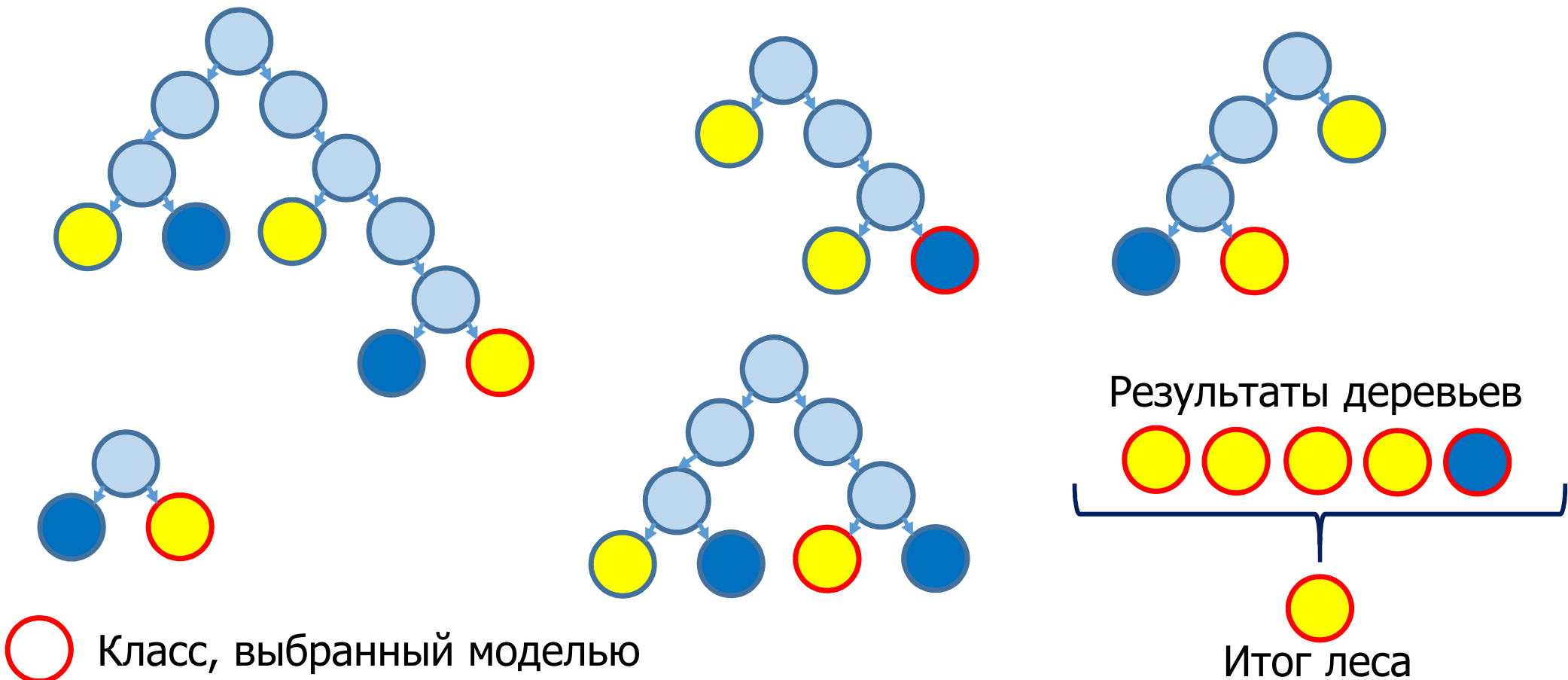
#	A2	A3	A5
2			
4			
5			

#	A1	A2
1		
2		

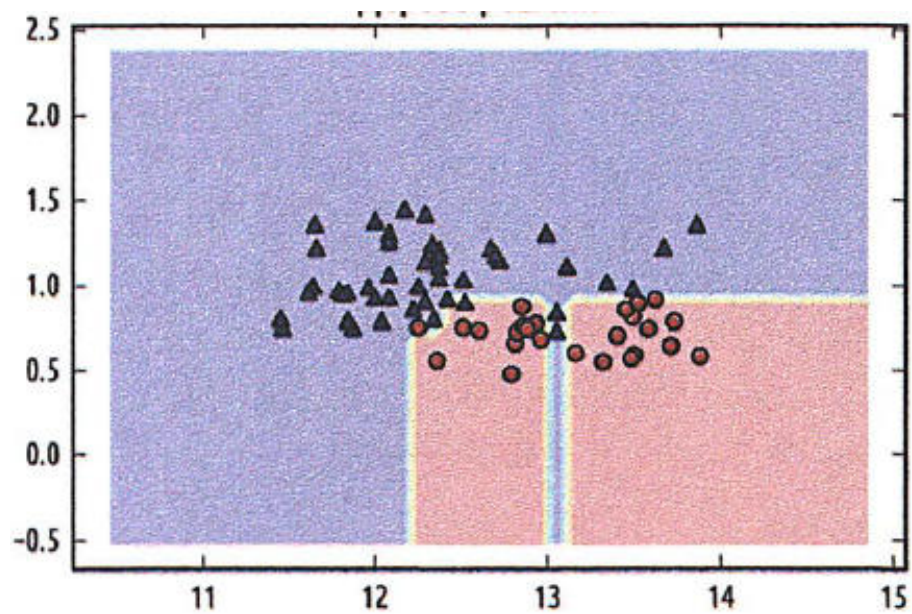
Выборка 5

Случайный лес

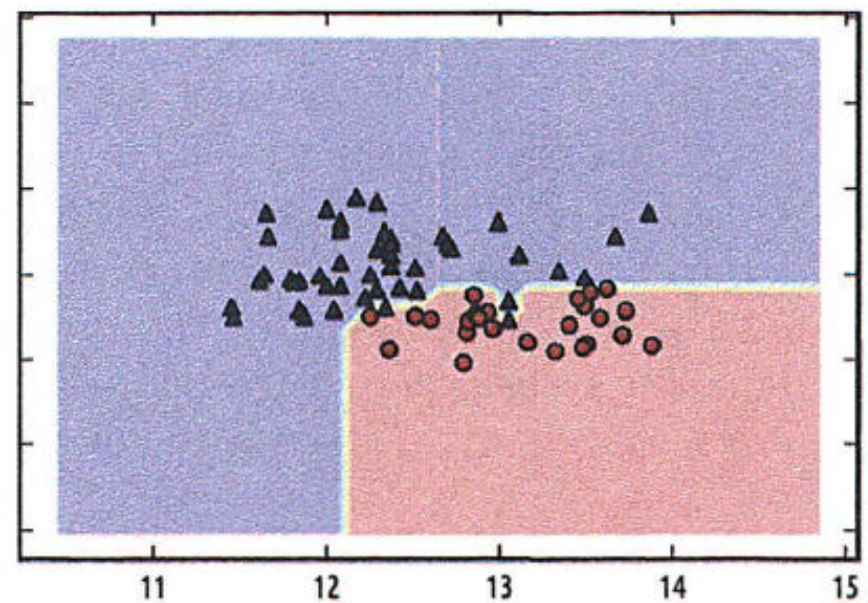
Объединение работы нескольких деревьев.
Мажоритарное голосование



Случайный лес



Дерево решений



Случайный лес

Определение важности признаков на случайном лесе

Основано на критерии прироста информации
для всех деревьев в случайном лесе

$$IG(D_p, f)$$

Пример: Оценка важности признаков для задачи оценки состояния тяжести заболевания по результатам общего анализа крови (из собств. практики)

Министерство здравоохранения
и социального развития РФ
Наименование учреждения
Лаборатория

Код формы по ОКУД
Код учреждения по ОКУД
Медицинская документация
Формы № 2380
Ул. Минеральная СССР 64-19.88
№ 1030

АНАЛИЗ КРОВИ №
02.05.00
дата взятия биоматериала

Фамилия, И.О. **Л-в В.В.**
Возраст **41**
Учреждение
Участок медицинская карта № **2284**

	Результат	Единицы СИ	Единицы, подпадающие в норму	Норма
Гемоглобин М.Ж.	121	г/л	130-160	г/л
Эритроциты М.Ж.	3.84	$\cdot 10^{12}/л$	4.0-5.0	$\cdot 10^{12}/л$
Цветовой показатель	0.85		0.85-1.05	
Среднее содержание гемоглобина (в 1 эритроците)	31.51	пг	30-35	пг
Ретикулы	1	%	0-10	%
Тромбоциты	195	$\cdot 10^9/л$	180-320	$\cdot 10^9/л$
Лейкоциты	18.8	$\cdot 10^9/л$	4.0-9.0	$\cdot 10^9/л$
Нейтрофилы	+	%	50-70	%
Лимфоциты	1	%	20-40	%
Моноциты	21	%	2-10	%
Эозинофилы	52	%	1-5	%
Базофилы	0	%	0-1	%
Скорость оседания эритроцитов (СОЭ)	61	мм/ч	2-15	мм/ч

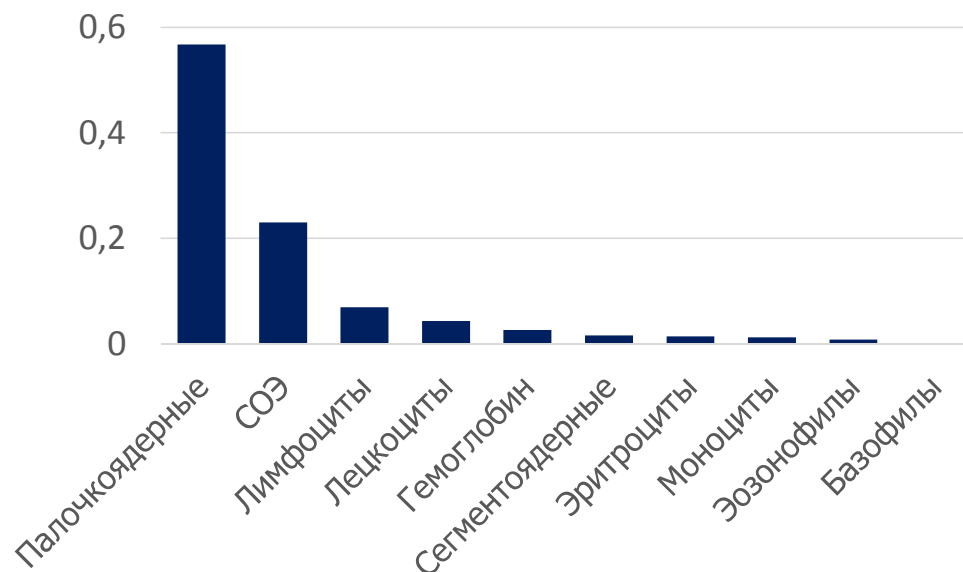
Министерство здравоохранения
и социального развития РФ
Наименование учреждения
Лаборатория

Код формы по ОКУД
Код учреждения по ОКУД
Медицинская документация
Формы № 2380
Ул. Минеральная СССР 64-19.88
№ 1030

АНАЛИЗ КРОВИ №
03.05.00
дата взятия биоматериала

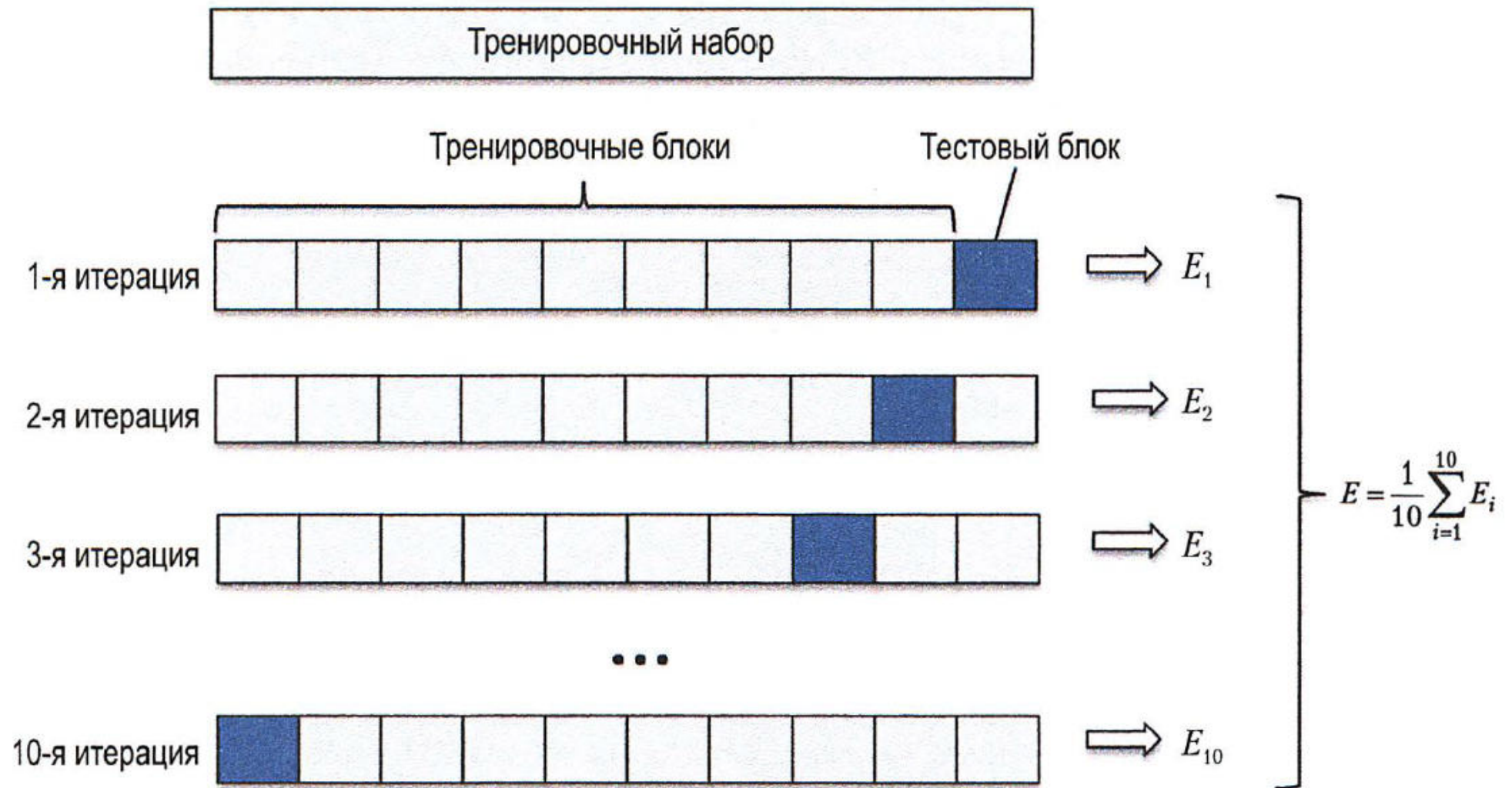
Фамилия, И.О. **Л-в В.В.**
Возраст **41**
Учреждение
Участок **6 ОТД** медицинская карта № **2284**

	Результат	Единицы СИ	Единицы, подпадающие в норму	Норма
Гемоглобин М.Ж.	110	г/л	130-160	г/л
Эритроциты М.Ж.	3.75	$\cdot 10^{12}/л$	4.0-5.0	$\cdot 10^{12}/л$
Цветовой показатель	0.88		0.85-1.05	
Среднее содержание гемоглобина (в 1 эритроците)	29.33	пг	30-35	пг
Ретикулы	10	%	0-10	%
Тромбоциты	210	$\cdot 10^9/л$	180-320	$\cdot 10^9/л$
Лейкоциты	13	$\cdot 10^9/л$	4.0-9.0	$\cdot 10^9/л$
Нейтрофилы	0	%	50-70	%
Лимфоциты	33	%	20-40	%
Моноциты	52	%	2-10	%
Эозинофилы	1	%	1-5	%
Базофилы	0	%	0-1	%
Скорость оседания эритроцитов (СОЭ)	32	мм/ч	2-15	мм/ч



Как изменилось состояние пациента?

К-блочная перекрестная проверка



Мажоритарное голосование в ансамбле

Номер модели	Пример 1	Пример 2	Пример 3	Пример 4	Пример 5	Количество ошибок
Модель 1	+	-	+	-	+	2
Модель 2	+	+	-	+	-	2
Модель 3	-	-	+	+	-	2
Модель 4	+	-	-	+	+	2
Модель 5	-	+	+	-	+	2
Мажорит. голосование	+	-	+	+	+	1

"+" – Правильно классифицированные примеры

"-" – Ошибочно классифицированные примеры

Пример мажоритарного голосования

