

# Статистический анализ экспрессии генов: систематизация и классификация

Зиннатулина Белла Раифовна

Санкт-Петербургский Государственный Университет  
Математико-механический факультет  
Кафедра статистического моделирования

Научный руководитель — к.ф.-м.н. доцент Н.П. Алексеева  
Рецензент — д.ф.-м.н. проф. М.С. Ермаков



Санкт-Петербург  
2015г.

- Данные — экспрессии генов стволовых клеток (Северо-Западный федеральный медицинский исследовательский центр).
- Выборка ( $N = 42$ ) больных сердечной недостаточностью (**HF**).
- Признаки ( $p = 36$ ) — экспрессии генов, характеризующие **HF**.

**Проблема:** Безошибочная классификация данных по наличию диабета и ожирения по полному набору признаков линейным дискриминантным анализом. Возможно наличие неинформативных признаков и переобучение.

**Цель:** Выявить экспрессии генов, влияющих на **HF** опосредованно через факторы диабета (**D**) и ожирения (**O**).

- Проверка действия классификаторов
  - Логистическая регрессия
  - Случайный лес
  - $k$  ближайших соседей
  - Линейный дискриминантный анализ
- Отбор информативных признаков методами
  - Дискриминантный анализ с разрежением
  - Специальный метод
- Систематизация признаков при помощи кластерного анализа с учетом влияния на факторы **D** и **O**.

# Постановка задачи классификации

- $X \subset \mathbb{R}^p$  — множество наблюдений (индивидов).
- $Y = \{0, 1\}$  — метки классов в случае бинарной классификации.
- $(x_1, y_1), \dots, (x_n, y_n), y_i \in Y, x_i \in X$  — обучающая выборка.
- Целевая функция  $f : X \rightarrow Y, f(x_i) = y_i, i = \overline{1 : n}$ .
- Задача: аппроксимация  $f$  на все пространство  $X$ .

- $f(z) = \frac{1}{1 + e^{-z}}$  — логистическая (сигмоидная) функция.
- $\mathbb{P}(y = 1 \mid x) = f(\theta^T x) = f(\theta_1 x_1 + \dots + \theta_p x_p), x \in X,$   
 $\mathbb{P}(y = 0 \mid x) = 1 - f(\theta_1 x_1 + \dots + \theta_p x_p).$
- $\theta = (\theta_1, \theta_2, \dots, \theta_p)$  — вектор значений вещественных параметров.  
Подбирается по методу максимального правдоподобия.

**Применение метода в задаче классификации:** Класс  $y = 1$  или  $y = 0$  объекта  $x$  определяется как  $\mathbb{P}(y = 1 \mid x) > 0.5$  или  $\mathbb{P}(y = 1 \mid x) \leq 0.5$  соответственно.

# Метод: случайный лес ( $RF$ )

**Случайный лес** основан на построении ансамбля деревьев принятия решений. Параметр метода — количество деревьев  $N$ .

## Дерево принятия решений

- $X \subset \mathbb{R}^p$ .
- Узлы помечены предикатом  $\Pi_i : X \rightarrow \{\text{True}, \text{False}\}, i = 0, 1, \dots$ , который является критерием разбиения данных.
- Разбиение данных происходит в каждом узле по одному признаку.
- Листья —  $y_i \in Y = \{0, 1\}$ .

Строится  $N$  деревьев принятия решений по случайным подвыборкам  $X_i, i = \overline{1 : N}$ . Принадлежность элемента классу определяется путем голосования деревьев.

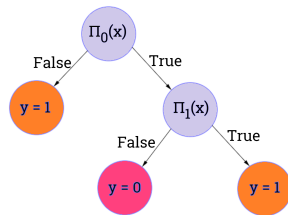


Рис.: Абстрактное дерево принятий решений.

# Метод: $k$ ближайших соседей ( $KNN$ )

- Параметр метода —  $k$ .
- Метод основан на оценивании сходства объектов.
- $\rho : X \times X \rightarrow \mathbb{N}$  — функция расстояния, евклидова метрика .
- $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), x_i \in X, y_i \in Y$  — обучающая выборка.
- $\tilde{x}$  — элемент тестовой выборки:  
$$\rho(\tilde{x}, x^{[1]}) \leq \rho(\tilde{x}, x^{[2]}) \leq \dots \leq \rho(\tilde{x}, x^{[n]}).$$
- Выбор класса  $\tilde{x}$  определяется  $\tilde{y} = \lceil (y^{[1]} + \dots + y^{[k]})/k - 0.5 \rceil$ .

# Метод: дискриминантный анализ с разрежением (sparseLDA)

- $X \subset \mathbb{R}^p$  — множество наблюдений, принадлежащих  $K$  классам.
- $H = \{H_{ij} = \mathbb{1}_{\{x_i \rightarrow y_j\}}, i \in \overline{1:n}, j \in \overline{1:k}\}$ .
- $\Omega$  — положительно определенная матрица.
- $\theta_k$  — корректирующий вектор весов.
- $\lambda$  и  $\gamma$  — неотрицательные настраиваемые параметры.
- $\beta_k \in \mathbb{R}^p$  —  $k$ -й дискриминантный вектор.
- $(\theta_k, \beta_k)$  является  $k$ -м решением задачи

$$\min_{\beta_k, \theta_k} \{ \| H\theta_k - X\beta_k \|^2 + \gamma \beta_k^T \Omega \beta_k + \lambda \| \beta_k \|_1 \},$$
$$\frac{1}{n} \theta_k^T H^T H \theta_k = 1, \quad \theta_k^T H^T H \theta_l = 0, \quad \forall l < k.$$

- Классификация применением стандартного линейного дискриминантного анализа к матрице  $(X\beta_1 \dots X\beta_q)$ ,  $q < K$ .



Таблица: Доля правильной классификации выборки

Метод	Тренировочная выборка		Тестовая выборка	
	Диабет	Ожирение	Диабет	Ожирение
Random Forest	0.7	0.72	0.6	0.66
KNN	0.78	0.76	0.75	0.75
sparseLDA	0.83	0.75	0.8	0.7
Logistic Regression	1	1	1	1
LDA	1	1	1	1

# Выделение информативных признаков: этап 1

- Матрица наблюдений  $X(n, p) = [X_1 \dotscdot X_p]$ , где  $n$  — число индивидов,  $p$  — число признаков.  $Y = \{\mathbf{O}, \mathbf{D}\}$ .
- Усеченная матрица наблюдений  $X_\tau(n, \theta) = [X_{\tau_1} \dotscdot X_{\tau_\theta}]$ , где  $\tau = (\tau_1, \dots, \tau_\theta) \subset (1, 2, \dots, p) = \mathcal{N}_p$ .
- Множество  $\Theta(p, \theta)$ ,  $|\Theta(p, \theta)| = C_p^\theta$  всех  $\theta$ -подмножеств из  $\mathcal{N}_p$ .
- Классификатор  $f^Y : X_\tau \rightarrow [0, 1]$ .
- Информативные  $\theta$ -наборы по классификаторам  $f_r^Y$ ,  $r \in \overline{1:m}$  с уровнем вероятности  $P$

$$\Sigma_r = \Sigma_r(\theta, p, f_r^Y, P, Y) = \{i | i \in \Theta(p, \theta), f_r^Y(X_\tau) \geq P\}.$$

# Выделение информативных признаков: этап2

- Значения линейных дискриминантных функций

$$DF_j(X_\tau, Y) = \alpha_0 + X_\tau \alpha^T, \quad \tau \in \bigcup_{r=1}^m \Sigma_r, \quad j = \overline{1:n}.$$

- Персональные наборы  $\tilde{\tau}_j = \arg \max_{\tau} |DF_j(X_\tau)|$ .

- Финальный  $\theta$ -набор — мода  $(\tilde{\tau}_1, \dots, \tilde{\tau}_n)$ .

- Рейтинги в  $P$ -значимой  $f_r^Y$ -классификации

$$Z_{ir}^Y = |\{\tau \in \Sigma_r | i \in \tau\}|, \quad i \in \overline{1:p}, \quad r \in \overline{1:m}.$$

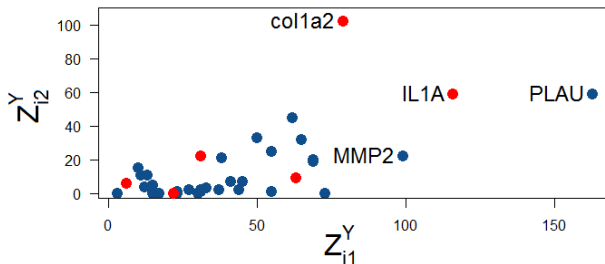
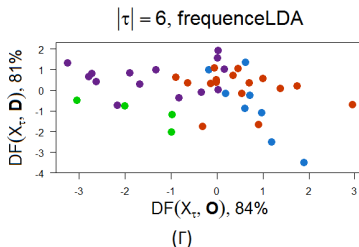
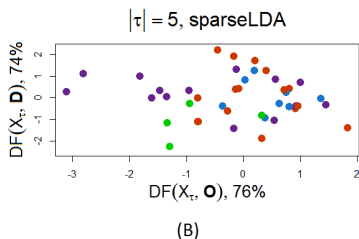
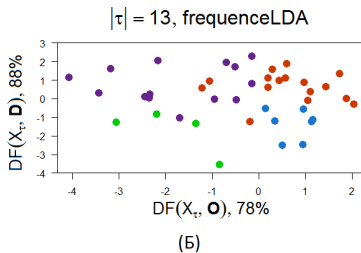
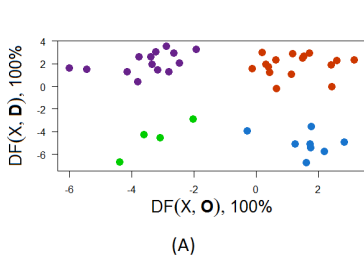


Рис.: Рейтинги признаков  $f_1^Y = KNN$ ,  $f_2^Y = RF$ ,  $\theta = 4$ ,  $Y = \mathbf{O}$ .

# Результаты: двумерная классификация



●  $-\mathbf{O}, -\mathbf{D}$  ●  $+\mathbf{O}, -\mathbf{D}$  ●  $-\mathbf{O}, +\mathbf{D}$  ●  $+\mathbf{O}, +\mathbf{D}$

Рис.: Визуализация данных по значениям дискриминантных функций.

# Систематизация. Кластерный анализ признаков

- Метрика — информационное разнообразие.
- Стратегия — информационный выигрыш от объединения групп.
- Объект кластеризации — рейтинги  $N_{jr}^Y$  в порядковой шкале признаков  $X_\tau$  классификатора  $f_\tau$  и разделителя  $Y$ ,

$$j \in \overline{1:n}, Y = \{\mathbf{O}, \mathbf{D}\}, r \in \overline{1:m}, \tau \in \bigcup_{r=1}^m \Sigma_r.$$

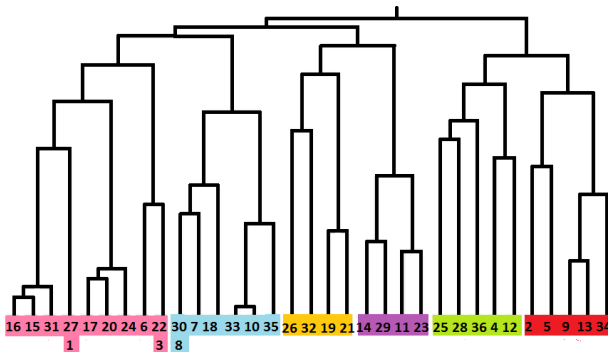


Рис.: Дерево кластеризации признаков.

## Результаты:

- Алгоритмы классификации и отбора признаков реализованы на языке  $R$ .
- Выделены группы признаков, дающих наилучшую классификацию по наличию диабета и ожирения с наименьшими потерями точности. Удалось сократить пространство признаков с 36 до 6.
- Признаки разделены на 6 групп по схожести роли, которую они играют в классификации.

## Планы:

- Исследовать структуру данных относительно разделителей и классификаторов.