

# Методы достижения интерпретируемости алгоритмов машинного обучения

*Сенько Олег Валентинович, Кунецова Анна Викторовна*

ФИЦ ИУ РАН, ИБХФ им. Н.М. Эмануэля РАН

При использовании большинства методов машинного обучения процесс вычисления компьютерных решений оказывается недоступным пользователю.

## Непрозрачность

- не позволяет понять причины выбора компьютерной программой предлагаемого решения
- снижает убедительность выдаваемых результатов
- затрудняет выработку оптимальных действий, связанных с предложенным решением

Непрозрачность или неинтерпретируемость свойственна следующим известным методам машинного обучения:

- **Нейросетевые методы**
- **Метод опорных векторов**

В значительной степени интерпретируемыми являются решающие или регрессионные деревья, а также ассоциативные правила.

Обобщающая способность отдельных решающих деревьев часто оказывается слабой по сравнению с упомянутыми выше технологиями.

Обобщающая способность существенно возрастает при использовании ансамблей решающих деревьев, получаемых с использованием технологий

- бэггинга
- адаптивного бустинга
- градиентного бустинга

Прозрачность и интерпретируемость также в значительной мере теряются при переходе к ансамблям

Интерпретируемость может быть достигнута через аппроксимацию работы обученного алгоритма с помощью набора моделей, связывающих целевую переменную  $Y$  с подмножествами признаков из набора  $X_1, \dots, X_n$ :

- модели должны достаточно просты для того, чтобы быть наглядно представлены в понятной для пользователей интерпретируемой форме
- модели должны быть статистически достоверны
- набор моделей должен достаточно полно описывать работу алгоритма

## Способы достижения интерпретируемости. Двумерные закономерности.

Для достижения наглядности могут быть использованы, например,

- модели, связывающие  $Y$  с отдельными переменными из  $X_1, \dots, X_n$
- Бокс-плоты, демонстрирующие различия средних значений переменных из множества  $X_1, \dots, X_n$  в группах, задаваемых категориальной целевой переменной  $Y$

Более точная аппроксимация работы обученного алгоритма может быть достигнута с помощью **двумерных моделей**, связывающих  $Y$  с парами переменных из набора  $X_1, \dots, X_n$ .

Использование **двумерных моделей** позволяет **сохранить высокую степень наглядности**.

- При верификации более сложных двумерных моделей возникает дополнительная **проблема избыточности**
- Проблема избыточности возникает, если достоверность более сложной модели оценивать только **в целом с помощью одного  $p$ -значения**
- Двумерную модель, связывающую переменную  $Y$  с переменными  $X_i$  и  $X_j$  назовём избыточной, если статистическая достоверность связана с реальным существованием регрессионной связи только с одной из двух переменных
- Например, реальная связь существует с переменной  $X_i$ , а переменная  $X_j$  на самом деле является шумовой
- Число избыточных закономерностей может оказаться очень велико. Компьютерный анализ должен включать средства их отсева.

При высокой исходной размерности данных существует высокая вероятность чисто случайного возникновения закономерностей, которые могут оказаться статистически достоверными при формальном применении статистических тестов.

- **Компьютерный анализ должен учитывать множественное тестирование.** Иными словами, должна проводиться дополнительная коррекция целевых уровней значимости с целью достижения истинной достоверности.



В случае, если  $Y$  является непрерывной величиной, для поиска двумерных моделей, удовлетворяющих перечисленным требованиям, потенциально могли бы использоваться

- двумерные регрессионные модели
- дисперсионный анализ (Two way ANOVA)

При этом достоверными могут считаться только такие двумерные модели, у которых регрессионные коэффициенты перед обоими регрессорами оказываются значимыми. Подход является математически обоснованным при справедливости предположения о действительном существовании линейной связи с добавленной нормально распределённой шумовой составляющей. Практика показывает, что существенное отклонение от нормальности нередко приводит к ложным выводам о значимости закономерности.

Метод **Two way ANOVA** может быть использован только, если обе независимые переменные являются категориальными. Кроме того, предположение о **нормально распределённой шумовой составляющей** также необходимо. В методе **Two way ANOVA** не имеет прямого решения проблема избыточности.

Общим способом решения **проблемы избыточности** могло бы являться использование информационных критериев **Akaike** или **BIC**.

Недостатком обоих подходов является оценивание ими только адекватности общей сложности моделей. Информационные критерии не оценивают значимость вклада отдельных независимых переменных.

В случаях, когда целевая переменная является категориальной, поиск оптимальных двумерных моделей мог бы производиться с помощью **методов распознавания с визуализацией границ между классами**.

Однако методы распознавания не включают статистических оценок достоверности вкладов отдельных переменных.

В случае, если категориальными являются одновременно целевая переменная и независимые переменные, для оценки значимости могут быть использованы статистические критерии:

- критерий  $\chi^2$
- G-тест

Недостатки - отсутствуют прямые способы статистической оценки вкладов каждой из независимых переменных.

Для учёта множественного тестирования могут быть использованы известные методы

- Бонферрони
- Бонферрони-Холма

Однако данные методы являются чрезмерно жёсткими, отвергая на самом деле значимые модели.

# Способы достижения интерпретируемости. Методы поиска двумерных закономерностей

Одним из возможных способов аппроксимации является использование двумерных моделей, выделяющих внутри совместной области допустимых значений пары переменных  $(X_i, X_j)$  подобластей, где преобладают определённые значения переменной  $Y$ .

Поиск оптимальных системы двумерных моделей может производиться с помощью

## Метода оптимальных достоверных разбиений (МОДР)

МОДР предназначен для поиска в данных статистически достоверных закономерностей, связанных с разбиениями интервалов допустимых значений отдельных переменных или совместных областей допустимых значений для пар переменных.

Метод ОДР удовлетворяет перечисленным выше требованиям:

- Метод обладает высокой наглядностью.
- Вследствие использования непараметрических перестановочных тестов при верификации не требуется предположений о типе распределений шумовых составляющих. Также отсутствуют ограничения на размер выборок.
- Метод ОДР позволяет для более сложных закономерностей оценивать значимость вклада каждого из элементов с помощью соответствующих  $p$ -значений. Например, для двумерных закономерностей метод оценивает значимость каждой из двух независимых переменных
- Метод включает эффективную современную процедуру учёта множественного тестирования, основанную опять же на перестановочных тестах

Метод ОДР является методом интеллектуального анализа данных (ИАД). Однако в отличие от входящих в группу технологий ИАД методов кластерного анализа (включая метод k-средних, иерархическая кластеризация SOTA и др.), ОДР предназначен для задач с выделенной целевой переменной.

## Метод оптимальных достоверных разбиений. Переход к многомерности.

Основная цель метода ОДР - наглядно представить на экране компьютера по-возможности все достоверные эффекты существующие в данных. Предполагается, что таким способом можно значительно улучшить интерпретируемость и прозрачность, решений, полученных с помощью технологий машинного обучения.

Вместе с тем существует многомерная технология машинного обучения, являющаяся фактически производной от метода ОДР - метод Статистически взвешенных синдромов (СВС). Метод СВС основан на взвешенном голосовании по квадрантам разбиений, в которые попадают вновь распознаваемые объекты.

Метод СВС был успешно использован на ряде биомедицинских и экономических задач, включая задачу прогнозирования отзыва банковских лицензий.



## Метод оптимальных достоверных разбиений

В МОДР оптимальные разбиения ищутся внутри фиксированных семейств различного фиксированного уровня сложности.

Использовались 4 типа семейств, представленных на рисунке

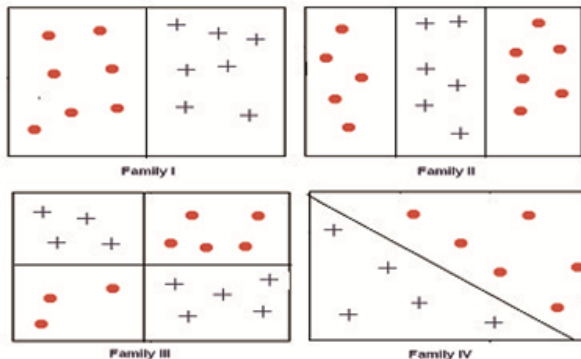


Рис.: 4 типа семейств разбиений

- Семейство I включает все разбиения интервалов значений одиночных признаков с помощью одной граничной точки
- Семейство II включает все разбиения интервалов значений одиночных признаков с помощью двух граничных точек
- Семейство III включает все разбиения совместной области допустимых значений пары признаков на четыре подобласти с помощью двух границ, параллельных координатным осям
- Семейство IV включает все разбиения совместной области допустимых значений пар признаков на две подобласти с помощью прямой, произвольно ориентированной относительно координатных осей

Предположим, что разбиение области допустимых значений на  $r$  подобластей индуцирует разбиение обучающей выборки  $\tilde{S}_t$  на подвыборки  $\tilde{s}_1, \dots, \tilde{s}_r$ . Предположим, что

- $\hat{y}_i$  - среднее значение целевой переменной  $Y$  на подвыборке  $\tilde{s}_i$ .
- $\hat{y}_0$  - среднее значение целевой переменной  $Y$  на всей выборке  $\tilde{S}_t$ .
- $m_i$  - размер выборки  $\tilde{s}_i$

Оптимальными считаются разбиения, для которых достигается максимума интегральный функционал

$$Q_I(\tilde{S}_t) = \sum_{i=1}^r \rho(\hat{y}_i, \hat{y}_0) m_i,$$

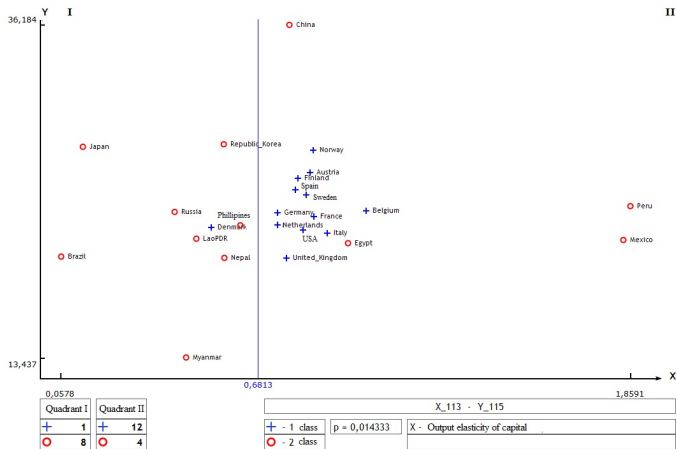
суммирующий отклонения от средних значений для полной выборки  $\tilde{S}_t$  по всем подобластям разбиения.

Наряду с интегральным функционалом может использоваться также альтернативный локальный функционал

$$Q_L(\tilde{S}_t) = \max_{i=1, \dots, r} [\rho(\hat{y}_i, \hat{y}_0) m_i],$$

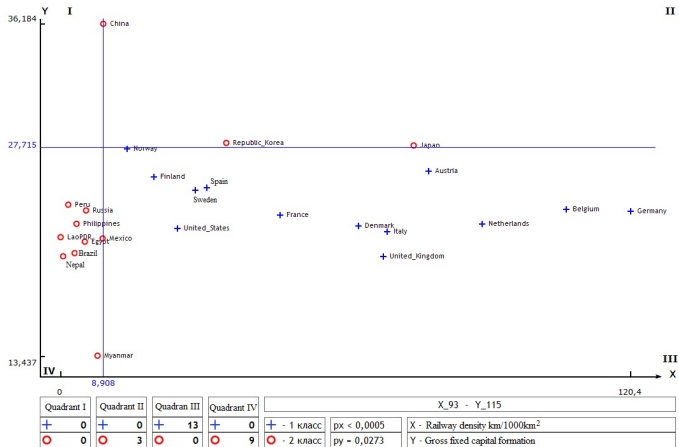
выделяющий подобласть с максимальным отклонением от средних значений по выборке.

# МОДР, примеры



**Рис.:** Связь типа институциональной матрицы с эластичностью по капиталу в Модели Кобба-Дугласа

# МОДР, примеры



**Рис.:** Связь типа институциональной матрицы с густотой железнодорожной сети и валовым накоплением основного капитала

Предыдущий слайд относится к исследованиям, связанным с верификацией теории С.Г.Кирдиной о возможности выделения двух базовых типов стран

- страны типа X с редистрибутивной экономикой, унитарной политической системой и доминирующим коммунитарным сознанием
- страны типа Y с рыночной экономикой, федералистской политической системой и доминирующей индивидуалистической компонентой в общественном сознании

Стандартные статистические тесты, основанные на сравнении средних значений целевой переменной  $Y$  в группах, индуцируемых разбиением, правомерно использовать если

оптимальное разбиение и статистика критерия вычисляются по двум независимым выборкам

При условии невозможности генерации двух независимых выборок достаточного объёма целесообразно использовать

Перестановочные тесты



Предположим, что требуется найти оптимальное разбиение интервала допустимых значений признака  $X_i$  внутри Семейства I по выборке

$$\tilde{S}_t = \{(y_1, x_{1i}), \dots, (y_m, x_{mi})\}$$

- На первом шаге найдём пороговое значение для признака  $X_i$ , при котором достигается максимальное значение функционала  $Q_L(\tilde{S}_t)$ , которое обозначим как  $Q_0$ .
- Сгенерируем  $N$  случайных независимых перестановок  $\mathbf{f}^1, \dots, \mathbf{f}^N$  натуральных чисел из множества  $\{1, \dots, m\}$ , где  $\mathbf{f}^j = (f_1^j, \dots, f_m^j)$

- Зафиксируем равным 0 значение счётчика  $ng$ .

Повторим следующие шаги при значениях  $j = 1, \dots, N$ .

- Получим из  $\tilde{S}_t$  по перестановке  $\mathbf{f}^j$  выборку  $\tilde{S}_j^p = \{(y_{f_1^j}, x_{1i}), \dots, (y_{f_m^j}, x_{mi})\}$
- Найдём пороговое значение для признака  $X_i$ , при котором достигается максимальное значение функционала  $Q_L(\tilde{S}_j^p)$ , которое обозначим как  $Q_j^p$
- В случае выполнения неравенства  $Q_j^p \geq Q_0$  значение счётчика  $ng$  увеличивается на 1

В качестве оценки  $p$ -значения используется отношение  $ng$  к общему числу перестановок:

$$p = \frac{ng}{N}$$

Описанный вариант перестановочного теста проверяет нулевую гипотезу о независимости переменных  $Y$  и  $X$  и позволяет **эффективно выявлять простейшие закономерности с одной граничной точкой**.

Применение данного варианта для верификации более сложных двумерных моделей приводит к ошибочному выявлению частично ложных закономерностей.

Мы будем называть закономерность, характеризующую связь переменной  $Y$  с переменными  $X_i$  и  $X_j$ , частично ложной, если её статистическая достоверность целиком определяется связью  $Y$  только с одной из двух переменных:  $X_i$  или  $X_j$ .

Было предложено два подхода к верификации, направленных на исключение частично ложных закономерностей. Оба подхода предусматривают вычисление двух  $p$ -значений, характеризующих отдельно вклады переменных  $X_i$  и  $X_j$ . Подход I основан на сравнении функционала качества для двумерной закономерности, характеризующей совместную связь  $Y$  с переменными  $X_i$  и  $X_j$ , с функционалами качества для двух одномерных закономерностей, характеризующих связь  $Y$  с каждой из переменных. Подход II основан на попытке опровержения с помощью двумерной модели возможности исчерпывающего описания данных простыми моделями.

Предположим, что требуется оценить достоверность двумерной закономерности, связанной с оптимальным разбиением совместной области допустимых значений признаков  $X_{i'}$  и  $X_{i''}$  внутри Семейства II по выборке

$$\tilde{S}_t = \{(y_1, x_{1i'}, x_{1i''}), \dots, (y_m, x_{mi'}, x_{mi''})\}$$

- На первом этапе по  $\tilde{S}_t$  найдём оптимальную совместную пару пороговых значений для признаков  $X_{i'}$  и  $X_{i''}$ , вычислим соответствующее значение  $Q_L(X_{i'}, X_{i''}, \tilde{S}_t)$
- Найдём оптимальные пороговые значения для признаков  $X_{i'}$  и  $X_{i''}$  по отдельности, вычислим соответствующие значения  $Q_L(X_{i'}, \tilde{S}_t), Q_L(X_{i''}, \tilde{S}_t)$

- Вычислим величины  $D'_0 = Q_L(X_{i'}, X_{i''}, \tilde{S}_t) - Q_L(X_{i'}, \tilde{S}_t)$   
 $D''_0 = Q_L(X_{i'}, X_{i''}, \tilde{S}_t) - Q_L(X_{i''}, \tilde{S}_t)$
- Зафиксируем равным 0 значение счётчиков  $ng'$  и  $ng''$
- Сгенерируем  $N$  случайных независимых перестановок  $\mathbf{f}^1, \dots, \mathbf{f}^N$  натуральных чисел из множества  $\{1, \dots, m\}$ , где  $\mathbf{f}^j = (f_1^j, \dots, f_m^j)$

Повторим следующие шаги при значениях  $j = 1, \dots, N$ .

- Получим из  $\tilde{S}_t$  по перестановке  $\mathbf{f}^j$  выборку  $\tilde{S}_j^p = \{(y_{f_1^j}, x_{1i'}, x_{1i''}), \dots, (y_{f_m^j}, x_{mi}, x_{mi''})\}$
- По  $\tilde{S}_j^p$  найдём оптимальную совместную пару пороговых значений для признаков  $X_{i'}$  и  $X_{i''}$ , вычислим соответствующее значение  $Q_L(X_{i'}, X_{i''}, \tilde{S}_j^p)$

- Найдём оптимальные пороговые значения для признаков  $X_{i'}$  и  $X_{i''}$  по отдельности, вычислим соответствующие значения  $Q_L(X_{i'}, \tilde{S}_j^p)$ ,  $Q_L(X_{i''}, \tilde{S}_j^p)$
- Вычислим величины  $D'_j = Q_L(X_{i'}, X_{i''}, \tilde{S}_j^p) - Q_L(X_{i'}, \tilde{S}_j^p)$   
 $D''_j = Q_L(X_{i'}, X_{i''}, \tilde{S}_j^p) - Q_L(X_{i''}, \tilde{S}_j^p)$
- В случае выполнения неравенства  $D'_j \geq D'_0$  значение счётчика  $ng'$  увеличивается на 1
- В случае выполнения неравенства  $D''_j \geq D''_0$  значение счётчика  $ng''$  увеличивается на 1

Статистическая значимость вклада переменной  $X_{i'}$  оценивается с помощью  $p$ -значения, равного отношению  $p' = \frac{ng''}{N}$ .

Статистическая значимость вклада переменной  $X_{i''}$  оценивается с помощью  $p$ -значения, равного отношению  $p'' = \frac{ng'}{N}$ .

Значимость закономерности будем оценивать как максимальное из двух рассчитанных  $p$ -значений:

$$p = \max(p', p'')$$



Подход II основан на попытке опровержения с использованием двумерного оптимального разбиения нулевой гипотезы об исчерпывающем описании данных с помощью более простой одномерной модели.

Предположим, что оптимальное разбиение интервала значений признака  $X_{i''}$  включает два подинтервала  $q_l$  и  $q_r$ .

При справедливости нулевой гипотезы целевая переменная  $Y$  очевидно не зависит от  $X_{i'}$  и  $X_{i''}$  внутри каждого из подинтервалов  $q_l$  и  $q_r$ . В этом случае любые перестановки  $Y$  относительно фиксированных позиций переменных  $X_{i'}$  и  $X_{i''}$  внутри каждого из подинтервалов оказываются равновероятными.

Высокое значение функционала  $Q_l(X_{i'}, X_{i''}, \tilde{S}_t)$  должно свидетельствовать зависимости  $Y$  не только от  $X_{i''}$  но и от  $X_{i'}$

На первом этапе по  $\tilde{S}_t$  найдём оптимальную совместную пару пороговых значений для признаков  $X_{i'}$  и  $X_{i''}$ , вычислим соответствующее значение  $Q_L(X_{i'}, X_{i''}, \tilde{S}_t)$ , которое обозначим  $Q_0$ . Сначала оценим значимость вклада в закономерность признака  $X_{i'}$ . С этой целью по граничной точке  $\delta_{i''}$  для найденного двумерного разбиения сформируем две выборки:

- выборку  $\tilde{S}_l$ , содержащую объекты из  $\tilde{S}_t$ , удовлетворяющие условию  $x_{ji''} < \delta_{i''}$
- выборку  $\tilde{S}_r$ , содержащую объекты из  $\tilde{S}_t$ , удовлетворяющие условию  $x_{ji''} > \delta_{i''}$

Пусть

- $m_l$  - число объектов в выборке  $\tilde{S}_l$
- $m_r$  - число объектов в выборке  $\tilde{S}_r$

- Сгенерируем  $N$  случайных независимых перестановок  $\mathbf{f}^1(l), \dots, \mathbf{f}^N(l)$  натуральных чисел из множества  $\{1, \dots, m\}$ , где  $\mathbf{f}^j = [f_1^j(l), \dots, f_m^j(l)]$
- Сгенерируем  $N$  случайных независимых перестановок  $\mathbf{f}^1(r), \dots, \mathbf{f}^N(r)$  натуральных чисел из множества  $\{1, \dots, m\}$ , где  $\mathbf{f}^j(r) = [f_1^j(r), \dots, f_m^j(r)]$

Повторим следующие шаги при значениях  $j = 1, \dots, N$ .

- Получим из  $\tilde{S}_l$  по перестановке  $\mathbf{f}^j(l)$  выборку  $\tilde{S}_j^p(l) = \{(y_{f_1^j(l)}, x_{1i'}, x_{1i''}), \dots, (y_{f_m^j(l)}, x_{mi}, x_{mi''})\}$
- Получим из  $\tilde{S}_r$  по перестановке  $\mathbf{f}^j(r)$  выборку  $\tilde{S}_j^p(r) = \{(y_{f_1^j(r)}, x_{1i'}, x_{1i''}), \dots, (y_{f_m^j(r)}, x_{mi}, x_{mi''})\}$

- По объединению  $\tilde{S}_j^p = \tilde{S}_j^p(l) \cup \tilde{S}_j^p(r)$  найдём оптимальную совместную пару пороговых значений для признаков  $X_{i'}$  и  $X_{i''}$ , вычислим соответствующее значение  $Q_L(X_{i'}, X_{i''}, \tilde{S}_j^p)$ , которое обозначим  $Q_j^p$
- В случае выполнения неравенства  $Q_j^p \geq Q_0$  значение счётчика  $ng'$  увеличивается на 1

Статистическая значимость вклада переменной  $X_{i'}$  оценивается с помощью  $p$ -значения, равного отношению  $\frac{ng'}{N}$ .

Аналогичным образом оценим вклад признака  $X_{i''}$

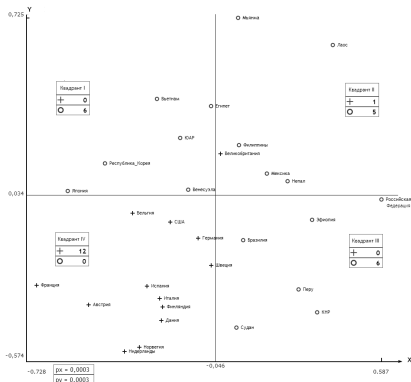
По граничной точки  $\delta_{i'}$  для найденного двумерного разбиения сформируем две выборки:

- выборку  $\tilde{S}_l$ , содержащую объекты из  $\tilde{S}_t$ , удовлетворяющие условию  $x_{ji'} < \delta_{i'}$
- выборку  $\tilde{S}_r$ , содержащую объекты из  $\tilde{S}_t$ , удовлетворяющие условию  $x_{ji'} > \delta_{i'}$

Далее повторяем те же самые операции, которые использовались при верификации вклада признака  $X_{i'}$ . То есть генерируем по  $\tilde{S}_l$  и  $\tilde{S}_r$  случайную выборку  $\tilde{S}_j^p$ , по которой вычисляем значение функционала  $Q_j^p$ . При выполнении неравенства  $Q_j^p \geq Q_0$  значение счётчика  $ng'$  увеличивается на 1.

Статистическую значимость вклада переменной  $X_{i'}$  оценивается с помощью  $p$ -значения, равного отношению  $\frac{ng''}{N}$

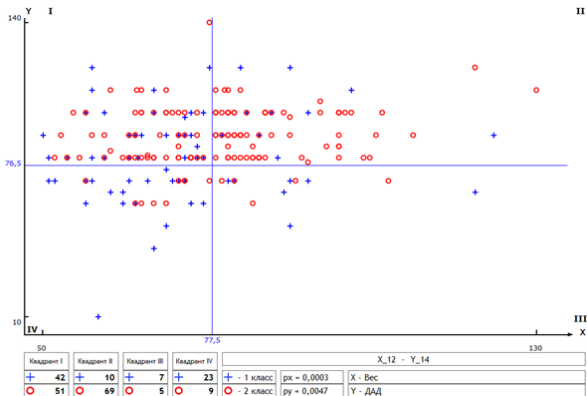
# Примеры двумерных закономерностей. Связь корреляций роста ВВП с внутренним долгом и государственными тратами с типом институциональных матриц



**Рис.:** По оси X отложен коэффициент корреляции роста ВВП с Внутренним долгом с лагом 4 года. По оси Y отложен коэффициент корреляции Государственных трат с Внутренним долгом с лагом 4 года. + страны типа Y, o - страны типа X

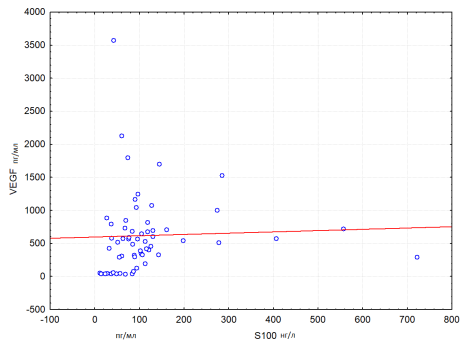
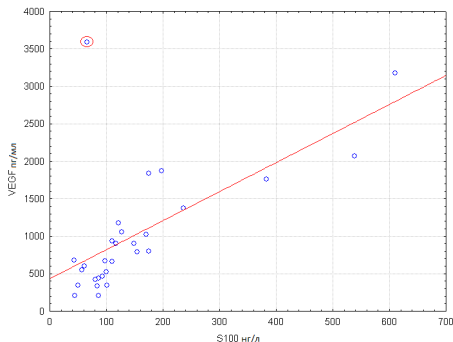


# Примеры закономерностей. Медицина. Связь летальности от сердечно-сосудистых заболеваний в учреждениях уголовно-исполнительной системы с весом и диастолическим давлением





# Примеры закономерностей. Медицина. Влияние гипоксии на связь фактора роста сосудов с уровнем протеинов S100



Для учёта **множественного тестирования в методе ОДР** разработана техника, являющаяся по сути модификацией технологии, предложенной в работе Westfall и Young в 1993 г.

- Предположим, что нам требуется найти скорректированную значимость закономерности, связанной с оптимальным разбиением области допустимых значений признаков  $X_{i'}$  и  $X_{i''}$  и оцениваемой  $p$ -значениями  $p'$  и  $p''$ . Предполагается, что  $p$ -значения были получены с помощью описанного ранее [Подхода I](#). Общее  $p$ -значение для закономерности оценивается тогда как  $p_o = \max(p', p'')$

- С помощью описанной ранее методики генерируется множество случайных выборок  $\tilde{S}_p^j = \{(y_{f_1^j}, \mathbf{x}_1), \dots, (y_{f_m^j}, \mathbf{x}_m)\}$ , где  $j = 1, \dots, N$
- С помощью опять же **Подхода I** для каждой пары признаков и для каждой выборки  $\tilde{S}_p^j$  вычисляем общие  $p$ -значения
- для каждой выборки  $\tilde{S}_p^j$  вычисляем минимальное общее  $p$ -значение  $p_m^j$
- В качестве общего скорректированного  $p$ -значения используется доля случайных выборок для которых  $p_m^j \leq p_o$

Использование перестановочного теста при верификации регрессионной модели допустимо при справедливости предположения о равновероятности различных перестановок целевой переменной  $Y$  относительно фиксированных позиций векторов независимых переменных.

Для того, чтобы такое предположение оказалось справедливым, необходимо включение в нулевую гипотезу наряду с предположением о независимости  $Y$  от независимых переменных также и дополнительного предположения о независимости между собой наблюдений

Подобное предположение чаще всего является несостоятельным для показателей, входящих в многомерные временные ряды. Например, гипотеза о взаимной независимости наблюдений очевидно является несостоятельной для в случаях когда компоненты многомерных рядов соответствуют процессу случайного блуждания.

# Применение методов верификации, основанных на ресэмплинге, при поиске закономерностей во временных рядах

Перестановочный тест является **процедурой верификации, основанной на ресэмплинге**, то есть

- Качество закономерностей, полученных на реальных данных, сравнивается с качеством закономерностей, полученных на данных, генерируемых случайным процессом.
- Для случайного процесса соблюдаются условия нулевой гипотезы.

Для верификации закономерностей, найденных по многомерным временным рядам, могут быть использованы технологии ресэмплинга, основанные, например, на **бутстрэпе**. Качество аппроксимации для регрессионных моделей, построенных по наблюдаемым временным рядам сравнивается с качеством аппроксимации для регрессионных моделей, построенных по временным рядам искусственно сгенерированных в соответствии с процессом случайного блуждания.