

Машинное обучение (Machine Learning)

Введение. Основные понятия

Уткин Л.В.



- 1 Что такое машинное обучение?
- 2 Постановки задач:
 - Обучение по прецедентам
 - Обучение без учителя
- 3 Примеры практических задач
- 4 О курсе

Презентация является компиляцией и заимствованием материалов из замечательных курсов и презентаций по машинному обучению:

К.В. Воронцова, А.Г. Дьяконова, Н.Ю. Золотых, С.И. Николенко, Andrew Moore, Lior Rokach, Rong Jin, Jessica Lin, Luis F. Teixeira, Alexander Statnikov и других.

Основные понятия

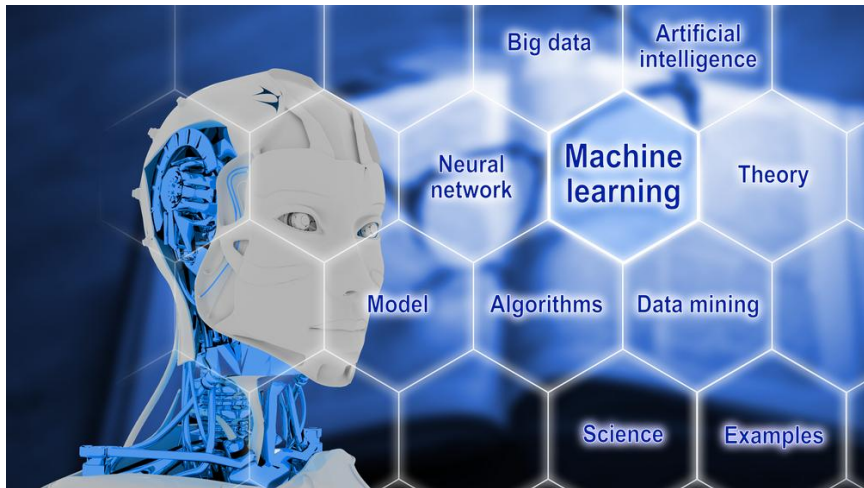
Что такое машинное обучение (machine learning)?

Машинное обучение – это подраздел ИИ, включающий методы построения алгоритмов, способных обучаться.

Машинное обучение – подраздел ИИ, математическая дисциплина, использующая разделы математической статистики, численных методов оптимизации, теории вероятностей, дискретного анализа, выделяющая знания из данных. (из Википедии)

Машинное обучение изучает методы построения алгоритмов, которые могут обучаться из данных и делать прогноз на данных.

Поток понятий



Что такое машинное обучение (machine learning)?

*Говорят, что компьютерная программа обучается на основе опыта E по отношению к некоторому классу задач T и меры качества P , если качество решения задач из T , измеренное на основе P , улучшается с приобретением опыта E . - Т.М. Mitchell
Machine Learning. McGraw-Hill, 1997.*

Дедуктивное и индуктивное методы обучения

Способы обучения и в компьютерных системах:

- 1 **Дедуктивное**, или аналитическое, обучение (экспертные системы). Имеются знания, сформулированные экспертом и как-то формализованные. Программа выводит из этих правил конкретные факты и новые правила.
- 2 **Индуктивное** обучение (статистическое обучение). На основе эмпирических данных программа строит общее правило. Эмпирические данные могут быть получены самой программой в предыдущие сеансы ее работы или просто предъявлены ей.
- 3 **Комбинированное** обучение.

“It is a capital mistake to theorize before one has data.”
- Arthur Conan Doyle



От данных к знаниям

Сферы приложения

- 1 Компьютерное зрение (computer vision)
- 2 Распознавание речи (speech recognition)
- 3 Компьютерная лингвистика и обработка естественных языков (natural language processing)
- 4 Медицинская диагностика
- 5 Биоинформатика
- 6 Техническая диагностика
- 7 Финансовые приложения
- 8 Рубрикация, аннотирование и упрощение текстов
- 9 Информационный поиск
- 10 . . .

Смежные и близкие области

- Pattern Recognition (распознавание образов)
- Data Mining (интеллектуальный анализ данных, включая Big Data)
- Artificial Intelligence (искусственный интеллект)

Разделы математики, используемые в машинном обучении

- Линейная алгебра
- Теория вероятностей и математическая статистика
- Методы оптимизации
- Численные методы
- Математический анализ
- Дискретная математика
- и др.

Классификация задач индуктивного обучения

- Обучение с учителем, или обучение по прецедентам (supervised learning): **классификация; восстановление регрессии; сегментация**
- Обучение без учителя (unsupervised learning): **кластеризация; визуализация данных; понижение размерности;**
- Обучение с подкреплением (reinforcement learning)

*Что определяет тип задачи машинного обучения: **Данные** (обучающая выборка) и **Цель***

Обучение с учителем - классификация

Мужчины



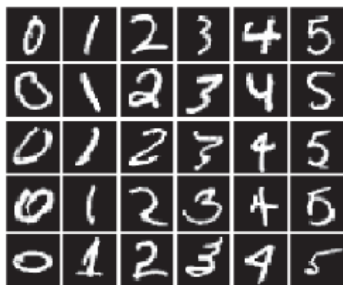
Женщины



Кто это?



Обучение с учителем - Классификация



Обучение с учителем - Классификация

Брак



Не брак



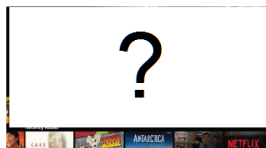
Брак или не брак?



Рекомендательные системы

Петроградский
Гражданка
Купчино

Каменный остров

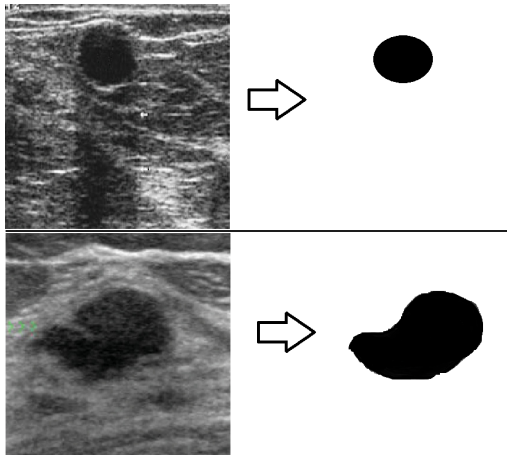


Обучение с учителем - пример регрессии

Котировки акций



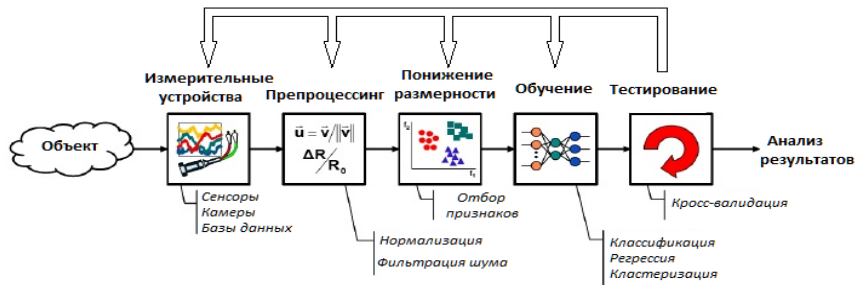
Обучение с учителем - сегментация



Обучение с учителем - сегментация



Схема всего процесса машинного обучения



Обучение по прецедентам или с учителем

Множество X — объекты, примеры, образцы (samples)

Множество Y — ответы, отклики, «метки», классы (responses)

Имеется некоторая зависимость $g : X \rightarrow Y$, позволяющая по $x \in X$ предсказать (или оценить вероятность появления) $y \in Y$.

Зависимость известна только на объектах из обучающей выборки:

$$T = (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Пара $(x_i, y_i) \in X \times Y$ - прецедент.

Задача обучения по прецедентам: научиться по новым объектам $x \in X$ предсказывать ответы $y \in Y$.

Пример обучающей выборки (классификация)

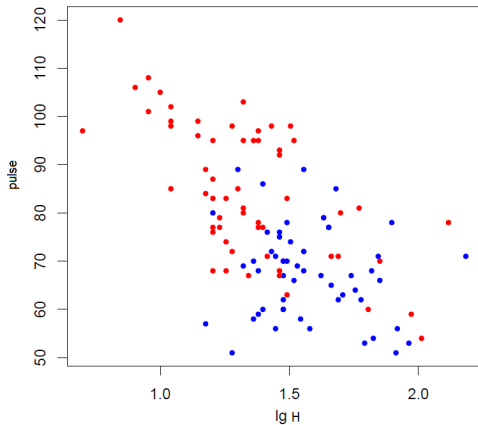
	пульс	гемоглобин	диагноз
x_1	70	140	здоров ($y = -1$)
x_2	60	160	здоров ($y = -1$)
x_3	94	120	миокардит ($y = 1$)
...
x_{114}	86	98	миокардит ($y = 1$)

Обучающая выборка:

$((70, 140), -1), (60, 160), -1), (94, 120), 1) \dots, (86, 98), 1))$

Задача обучения: новый пациент $x = (75, 128)$, $y = ?$

Графическое представление обучающей выборки



Другой пример обучающей выборки (классификация)

	вес	рост	возраст	ср.дл.волос	пол
x_1	96	170	42	0	м ($y = -1$)
x_2	60	180	25	8	м ($y = -1$)
x_3	54	165	30	21	ж ($y = 1$)
x_4	83	178	47	18	ж ($y = 1$)
...
x_{100}	108	193	32	40	ж ($y = 1$)

Задача обучения: $x = (75, 184, 28, 10)$, $y = ?$

Обучающая выборка с категориальными данными

	вес	рост	возраст	ср. дл. волос	пол
x_1	96	170	42	короткие	м ($y = -1$)
x_2	60	180	25	короткие	м ($y = -1$)
x_3	54	165	30	длинные	ж ($y = 1$)
x_4	83	178	47	короткие	ж ($y = 1$)
...
x_{100}	108	193	32	длинные	ж ($y = 1$)

Задача обучения: $x = (75, 184, 28, \text{"короткие"})$, $y = ?$

Пример пропущенных данных (missing data)

	вес	рост	возраст	ср.дл.волос	пол
x_1	96	170	42	короткие	м ($y = -1$)
x_2	60	180	25	короткие	-
x_3	54	165	-	длинные	ж ($y = 1$)
x_4	-	178	47	короткие	ж ($y = 1$)
...
x_{100}	108	193	32	длинные	ж ($y = 1$)

Задача обучения: $x = (75, 184, 28, \text{"короткие"})$, $y = ?$

Пример ненужного признака

	вес	рост	возраст	ср. дл. волос	оценка по маш.обуч.	пол
x_1	96	170	42	короткие	5	м
x_2	60	180	25	короткие	3	-
x_3	54	165	-	длинные	5	ж
x_4	-	178	47	короткие	4	ж
...
x_{100}	108	193	32	длинные	3	ж

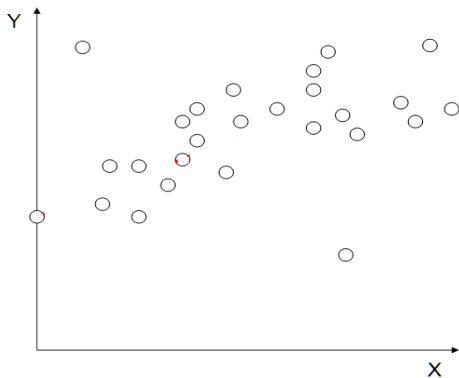
Задача обучения: $x = (75, 184, 28, \text{"короткие"}, 5)$, $y = ?$

Пример регрессионных данных

	вес	рост	ср.дл. волос	пол	возраст (y)
x_1	96	170	короткие	м	42
x_2	60	180	короткие	м	25
x_3	54	165	длинные	ж	30
x_4	83	178	короткие	ж	47
...
x_{100}	108	193	длинные	ж	32

Задача обучения: определить возраст
 $x = (75, 184, \text{"короткие"}, \text{"м"})$, $y = ?$

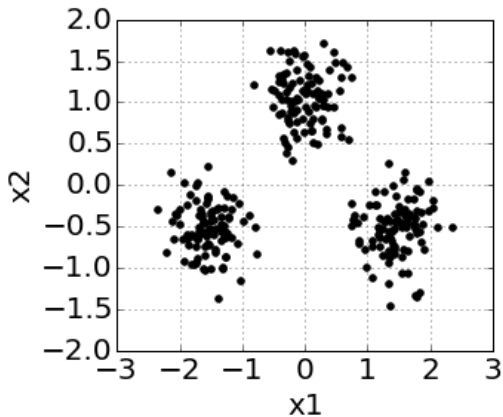
Графическое представление данных для регрессии



Обучение без учителя

- В этом случае нет “учителя” и “обучающая выборка” состоит только из объектов, т.е. Y отсутствует.
- **Задача кластеризации:** разбить объекты на группы (кластеры), так, чтобы в одном кластере оказались близкие друг к другу объекты, а в разных кластерах объекты были существенно различные.
- **Кластер** можно охарактеризовать как группу объектов, имеющих общие свойства.

Графическое представление данных для кластеризации



Пример задачи без учителя

	вес	рост	возраст	ср. дл. волос
x_1	96	170	42	короткие
x_2	60	180	25	короткие
x_3	54	165	30	длинные
x_4	83	178	47	короткие
...
x_{100}	108	193	32	длинные

Задача обучения: “отгадать” пол всех людей из обучающей выборки

Признаковые описания

Каждый объект характеризуется набором признаков (свойств, атрибутов, features) $f_j : X \rightarrow D_j, j = 1, \dots, n$

Типы признаков:

- $D_j = \{0, 1\}$ бинарный признак;
- $D_j = \{1, 2, 3, \dots, s\}$ номинальный (категориальный) признак (красный, зеленый, синий);
- D_j упорядочено - порядковый признак, например, вес:(малый, средний, большой).
- $D_j = \mathbb{R}$ количественный признак

Вектор $(f_1(x), f_2(x), \dots, f_n(x))$ - признаковое описание объекта x .

Признаки в примерах определения пола

- **вес:** количественный
- **рост:** количественный
- **возраст:** количественный
- **ср.дл. волос:** бинарный или упорядочено - порядковый или количественный
- **оценка по маш.обуч.:** упорядочено - порядковый или категориальный

Описание меток классов

В зависимости от множества Y выделяют разные типы задачи обучения:

- 1 **Задачи классификации (classification):**
 $Y = \{-1, +1\}$ классификация на 2 класса.
 $Y = \{1, \dots, M\}$ на M непересекающихся классов.
 $Y = \{0, 1\}^M$ на M классов, которые могут пересекаться.
- 2 **Задачи восстановления регрессии (regression):**
 $Y = \mathbb{R}$.
- 3 **Задачи ранжирования (ranking, learning to rank):** Y - конечное упорядоченное множество.

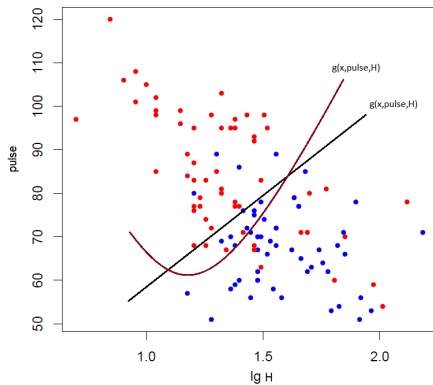
Модель алгоритма

Решить задачу машинного обучения означает разработать алгоритм или модель алгоритма, зависящего от параметров и позволяющих определить значение метки класса (Y) для нового объекта (x).

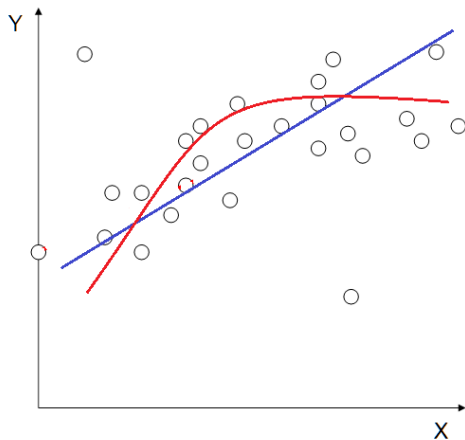
Модель алгоритма

- **Моделью алгоритма** a называется параметрическое семейство функций $g : X \rightarrow Y$ или $g(x, \theta)$, где $\theta \in \Theta$ параметры в пространстве параметров.
- **Пример:** В задачах с m признаками $f_j(x)$, $j = 1, \dots, m$ используются линейные модели с $\theta = (\theta_1, \dots, \theta_m)$:
$$g(x, \theta) = \sum_{j=1}^m \theta_j f_j(x)$$
- Процесс подбора оптимальной функции g и оптимального параметра θ по обучающей выборке называют **настройкой** (fitting, tuning) или **обучением** (training) алгоритма a .

Модели алгоритмов класификации



Модели алгоритмов регрессии



“Essentially, all models are wrong, but some are useful”
- George E. P. Box



Функционал качества

- **Функционал качества** может определяться как средняя ошибка ответов.
- **Функционал риска** или качества алгоритма a обучения есть

$$Q(a, X) = \int (\mathcal{L}(a, x) \cdot P(X, y)) dXdy$$

- **Функция потерь** (loss function) - это неотрицательная функция $L(a, x)$, характеризующая величину ошибки алгоритма a на объекте x . Если $\mathcal{L}(a, x) = 0$, то ответ $a(x)$ называется **корректным**.
- $P(X, y)$ - совместная плотность вероятностей

Функции потерь

- Функции потерь для классификации:
 - $\mathcal{L}(a, x) = [a(x) \neq y(x)]$ - индикатор ошибки
 - $\mathcal{L}(a, x) = \max(0, 1 - y_i a(x))$ - петлевая функция (hinge-loss function)
- Функции потерь для регрессии:
 - $\mathcal{L}(a, x) = |a(x) - y(x)|$ - абсолютное значение ошибки
 - $\mathcal{L}(a, x) = (a(x) - y(x))^2$ - квадратичная ошибка
 - $\mathcal{L}(a, x) = \begin{cases} (y - a)^2/2, & \text{если } |y - a| \leq \delta \\ \delta(|y - a|) - \delta/2, & \text{если } |y - a| > \delta \end{cases}$ - функция потерь Хьюбера
- Функции потерь для кластеризации:
$$\mathcal{L}(a, x) = \sum_{i=1}^n \min_c \|x_i - a_c\|^2$$

Эмпирический функционал качества

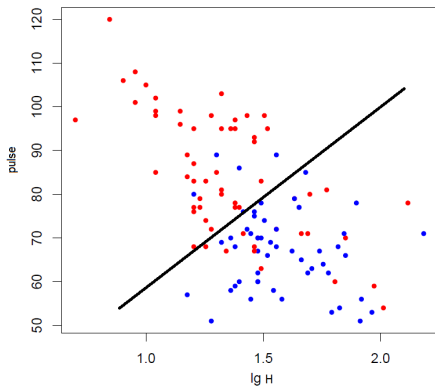
$$Q(a, X) = \int (\mathcal{L}(a, x) \cdot P(X, y)) dXdy$$

- **Эмпирический функционал риска** или качества алгоритма a на выборке X есть

$$Q(a, X) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(a, x_i)$$

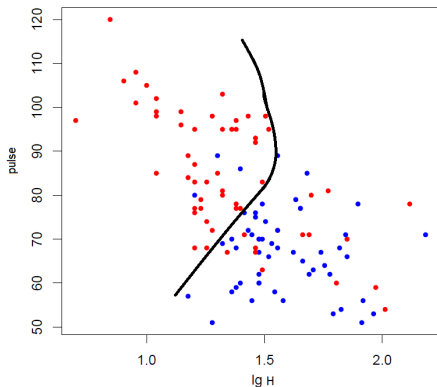
- Плотность $P(X, y)$ в функционале риска заменена на эмпирическое распределение (равномерное распределение) на элементах обучающей выборки.
- **Задача выбора “наилучшего” метода обучения** - это минимизация функционала риска по множеству A или по множеству параметров Θ .

Эмпирический функционал качества



$$Q(a, X) = \frac{1}{n} \sum_{i=1}^n [a(x_i) \neq y_i] = \frac{1}{114} (5 + 15) = 0.175$$

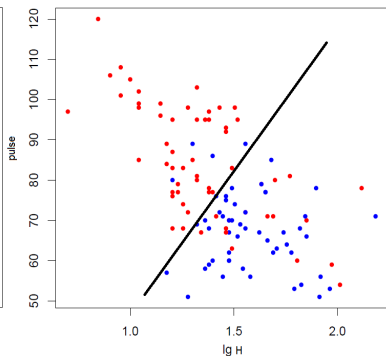
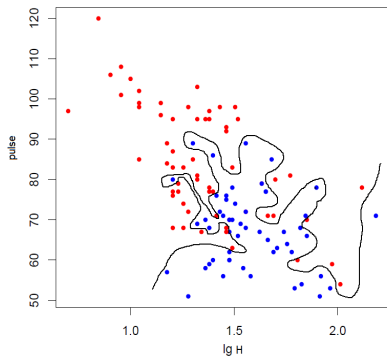
Эмпирический функционал качества



$$Q(a, X) = \frac{1}{n} \sum_{i=1}^n [a(x_i) \neq y_i] = \frac{1}{114} (3 + 14) = 0.149$$

Оценка качества обучения

Переобучение и недообучение в классификации

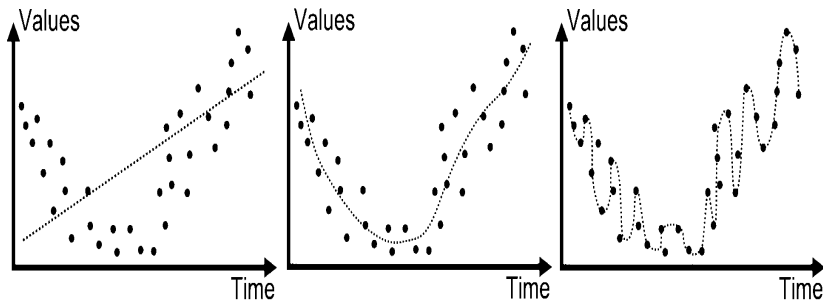


Проблема переобучения и недообучения

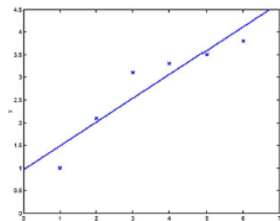
Переобучение (overfitting) — нежелательное явление, возникающее при решении задач обучения по прецедентам, когда вероятность ошибки обученного алгоритма на объектах тестовой выборки оказывается существенно выше, чем средняя ошибка на обучающей выборке. Переобучение возникает при использовании избыточно сложных моделей.

Недообучение — нежелательное явление, возникающее при решении задач обучения по прецедентам, когда алгоритм обучения не обеспечивает достаточно малой величины средней ошибки на обучающей выборке. Недообучение возникает при использовании недостаточно сложных моделей.

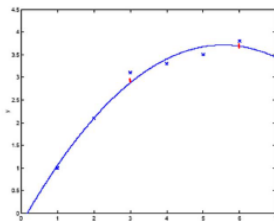
Переобучение и недообучение в регрессии



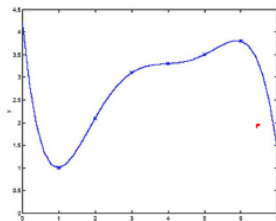
Переобучение и недообучение в регрессии



$$y = \theta_0 + \theta_1 x$$



$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$y = \sum_{j=0}^5 \theta_j x^j$$

Этапы решения задачи обучения

В задачах обучения по прецедентам всегда есть два этапа:

- 1 **Этап обучения (training)**: по выборке X строится алгоритм a и определяется функция $g(x, \theta)$ с учетом функционала риска алгоритма a
- 2 **Этап применения или тестирования (testing)**: насколько правильные или неправильные ответы $a(x)$ выдает алгоритм a для новых объектов x .

Тестовые данные

“A model for data, no matter how elegant or correctly derived, must be discarded or revised if it does not fit the data or when new or better data are found and it fails to fit them.”

- Paul Velleman



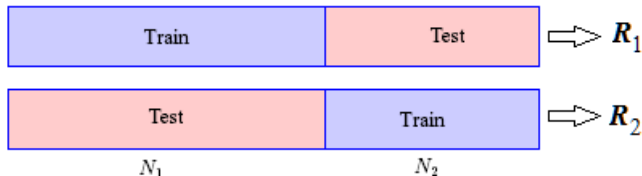
Обучающая и тестовая выборки

Случайно разделим все имеющиеся данные на:

- **обучающую** (train) выборку, которая используется для построения моделей
- **тестовую** (test) выборку, которая используется для оценки как модель ведет себя на новых данных



Метод перекрестного (скользящего) контроля

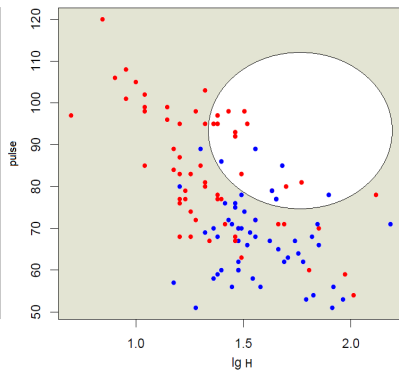
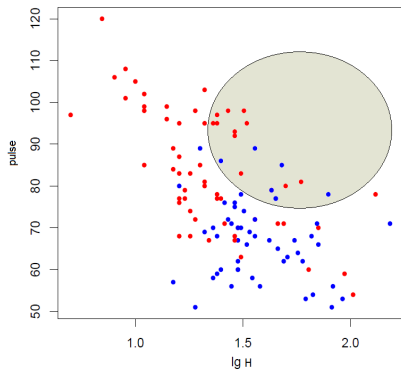


- Модель обучим на обучающей (train) выборке, а оценку ошибки R произведем на тестовой (test) выборке. Получим оценку R_1 .
- Поменяем выборки ролями. Получим оценку R_2 .
- Итоговая оценка качества - среднее взвешенное оценок R_1 и R_2 .

Метод перекрестного (скользящего) контроля

Обобщение этой процедуры называется методом перекрестного (скользящего) контроля (cross-validation).

Метод перекрестного (скользящего) контроля



Метод перекрестного контроля в общем виде

1. Случайным образом разобьем исходную выборку на M непересекающихся примерно равных по размеру частей.
2. Последовательно каждую из этих частей рассмотрим в качестве тестовой выборки, а объединение остальных частей — в качестве обучающей выборки.
3. Таким образом построим M моделей и соответственно M оценок для ошибки предсказания.
4. В качестве окончательной оценки ошибки возьмем их среднее взвешенное.

Метод перекрестного (скользящего) контроля

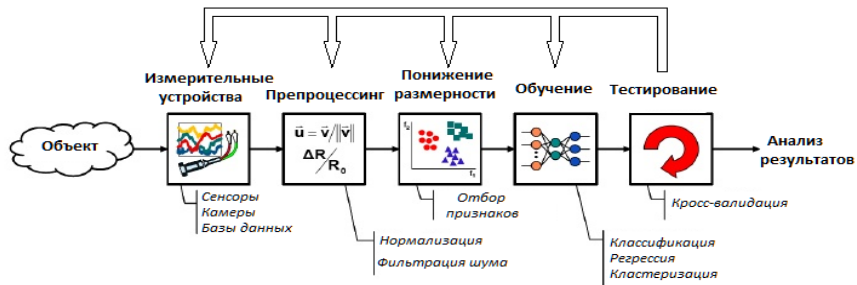
Test	Train	Train	Train	Train
Train	Test	Train	Train	Train
Train	Train	Test	Train	Train
Train	Train	Train	Test	Train
Train	Train	Train	Train	Test

Частный случай - один отделяемый элемент

- $M = N$ - метод перекрестного контроля с **одним отделяемым элементом** или число частей равно числу элементов выборки
- (leave-one-out cross-validation, **LOO**)
- LOO — самый точный, но требует много времени

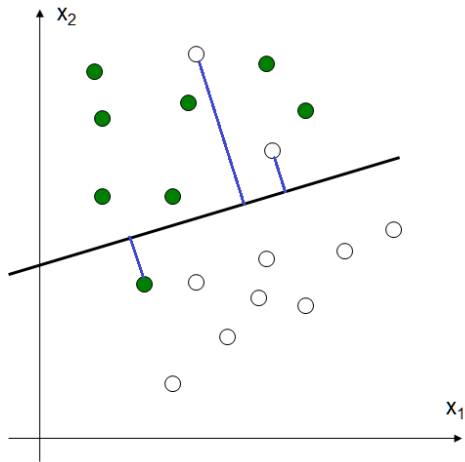
Этапы задачи обучения

Схема всего процесса машинного обучения

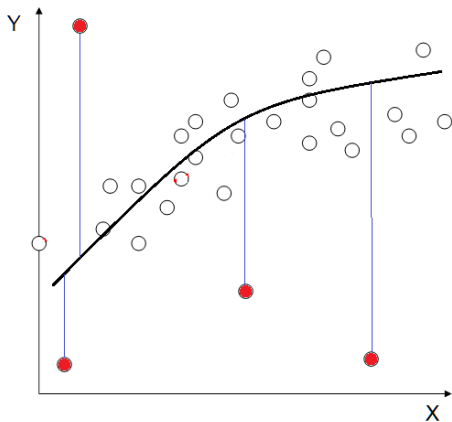


- 1 **Задача фильтрации выбросов (outliers detection)** — обнаружение в обучающей выборке небольшого числа нетипичных объектов.
 - В некоторых приложениях их поиск является самоцелью (например, обнаружение мошенничества).
 - Следствие ошибок в данных или неточности модели, то есть шум.
 - Используются робастные методы и одноклассовая классификация.
- 2 **Задача заполнения пропущенных значений (missing values)** — замена недостающих значений признаков их прогнозными значениями.

Фильтрации выбросов в классификации



Фильтрация выбросов в регрессии



Задача фильтрации выбросов

- Пусть $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ - обучающая выборка, $\mathcal{Y} = \{1, 2, \dots, c\}$ - множество классов
- **Отступ:** $M(x_i, y_i) = g_{y_i}(x_i) - \max_{y \in \mathcal{Y} \setminus \{y_i\}} g_y(x_i)$
 - отступ отрицательный означает, что объект x_i был неправильно классифицирован
 - величина отступа показывает, насколько классификатор уверен, что объект x_i может быть отнесен к истинному классу y_i

Удаление шумов

Шумы - это объекты, сильно выбивающиеся из закономерности, определяемой алгоритмом обучения, т.е. их можно определить как

$$\{x_i : M(x_i, y_i) < -\delta\}$$

для достаточно большого $\delta > 0$.

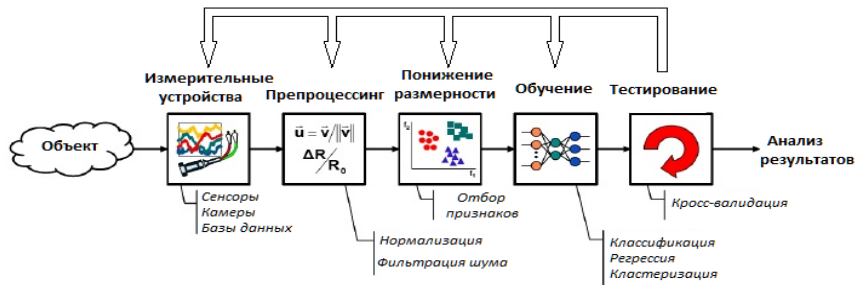
Алгоритм фильтрации шумов

- 1 для каждого (x_i, y_i) в обучающей выборке T вычислить $M(x_i, y_i)$
- 2 вернуть отфильтрованную обучающую выборку $T^* = \{(x_i, y_i) : M(x_i, y_i) \geq -\delta\}$

Задача заполнения пропущенных значений

	вес	рост	возраст	ср. дл. волос	пол
x_1	96	170	42	короткие	м ($y = -1$)
x_2	60	180	25	короткие	-
x_3	54	165	-	длинные	ж ($y = 1$)
x_4	-	178	47	короткие	ж ($y = 1$)
...
x_{100}	108	193	32	длинные	ж ($y = 1$)

Схема всего процесса машинного обучения



Отбор признаков и сокращение размерности (1)

- **Задача сокращения размерности** (dimensionality reduction) заключается в том, чтобы по исходным признакам с помощью некоторых функций преобразования перейти к наименьшему числу новых признаков, не потеряв при этом никакой существенной информации об объектах выборки.
- В классе линейных преобразований наиболее известным примером является **метод главных компонент**.

Пример ненужного признака

	вес	рост	возраст	ср. дл. волос	оценка по маш. обуч.	пол
x_1	96	170	42	короткие	5	м
x_2	60	180	25	короткие	3	-
x_3	54	165	-	длинные	5	ж
x_4	-	178	47	короткие	4	ж
...
x_{100}	108	193	32	длинные	3	ж

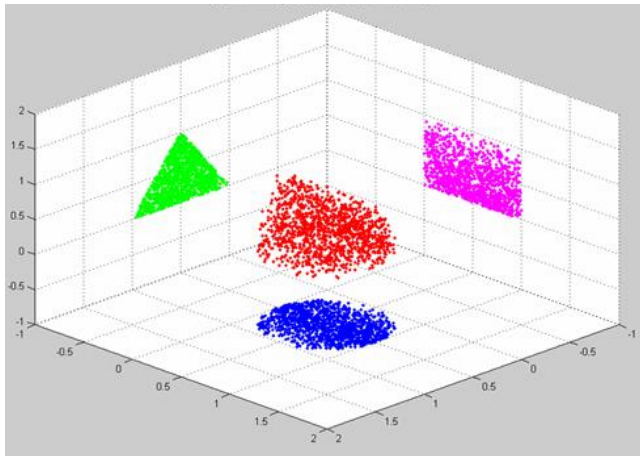
Отбор признаков и сокращение размерности (2)



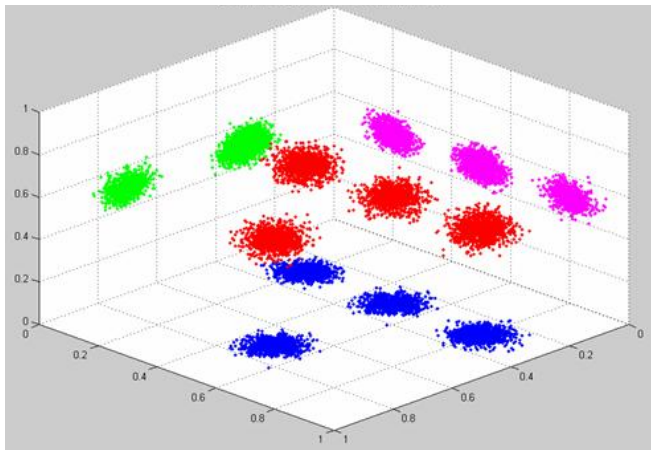
Отбор признаков и сокращение размерности (3)



Отбор признаков и сокращение размерности (4)



Отбор признаков и сокращение размерности (5)



Примеры прикладных задач

Задача о сортировке рыбы

Объект - морская рыба.

Классы: лосось и морской окунь.

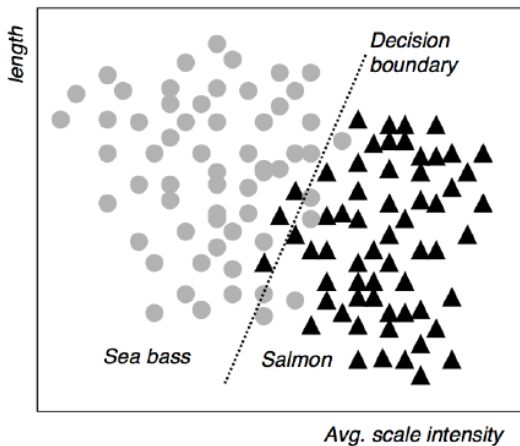
Примеры признаков:

- *количественные:* длина (чаще окунь длинее, чем лосось), освещенность (чаще лосось светлее, чем окунь)

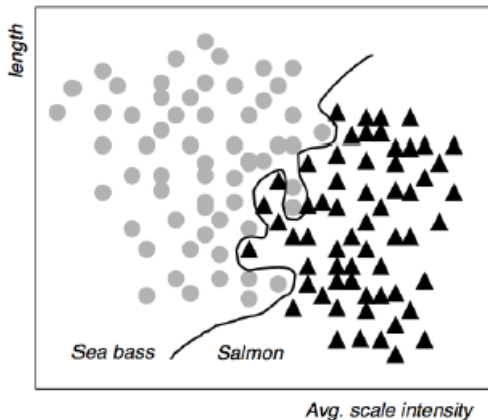
Задача о сортировке рыбы



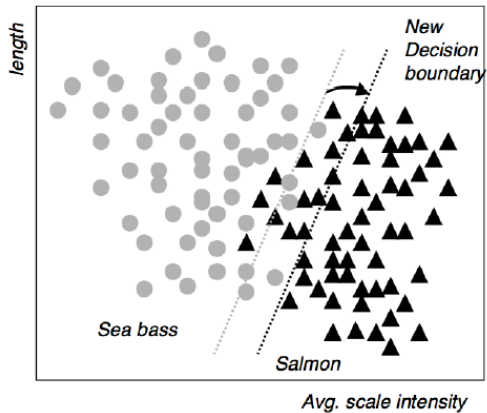
Задача о сортировке рыбы



Задача о сортировке рыбы (переобучение)



Задача о сортировке рыбы (веса классов)



Задачи медицинской диагностики

Объект - пациент в определенный момент времени.

Классы: диагноз или способ лечения или исход заболевания.

Примеры признаков:

- *бинарные*: пол, головная боль, слабость и т.д.
- *порядковые*: тяжесть состояния, желтушность и т.д.
- *количественные*: возраст, пульс, артериальное давление, содержание гемоглобина в крови и т.д.

Особенности задачи:

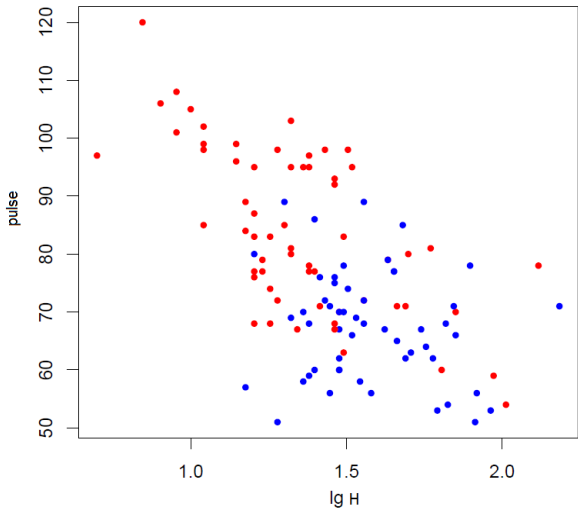
- обычно много “пропусков” в данных;
- нужен интерпретируемый алгоритм классификации;
- нужна оценка вероятности (риска | успеха | исхода).

Имеются данные о 114 лицах с заболеванием сердца: у 61 — проблемы, у 53 — нет проблем.

Для каждого пациента известны показатели:

- *pulse* — пульс,
- *H* — содержание гемоглобина в крови.

Можно ли научиться предсказывать (допуская небольшие ошибки) наличие проблем по *pulse* и *H* у новых пациентов?



Задача прогнозирования стоимости недвижимости

Объект - квартира в Санкт-Петербурге.

Примеры признаков:

- *бинарные*: наличие балкона, лифта, мусоропровода, охраны, и т. д.
- *номинальные*: район города, тип дома (кирпичный/панельный/блочный/монолит), и т. д.
- *количественные*: число комнат, жилая площадь, расстояние до центра, до метро, возраст дома, и т. д.

Особенности задачи: выборка неоднородна, стоимость меняется со временем; разнотипные признаки; для линейной модели нужны преобразования признаков.

Задача категоризации текстовых документов

Объект - текстовый документ.

Классы: рубрики иерархического тематического каталога.

Примеры признаков:

- *номинальные:* автор, издание, год, и т. д.
- *количественные:* для каждого термина частота в тексте, в заголовках, в аннотации, и т. д.

Особенности задачи: лишь небольшая часть документов имеют метки u_j ; документ может относиться к нескольким рубрикам.

Задача ранжирования поисковой выдачи

Объект - пара <запрос, документ>.

Классы: релевантен или не релевантен (разметка делается людьми ассессорами).

Примеры признаков:

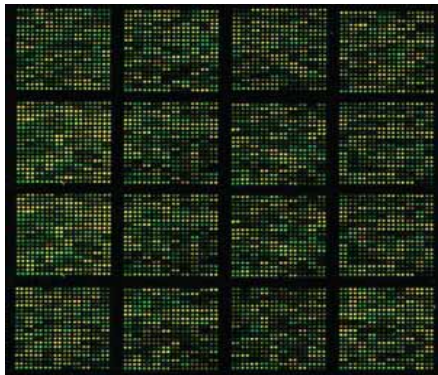
- *количественные:* частота слов запроса в документе, число ссылок на документ, число кликов на документ: всего, по данному запросу, и т. д.

Особенности задачи:

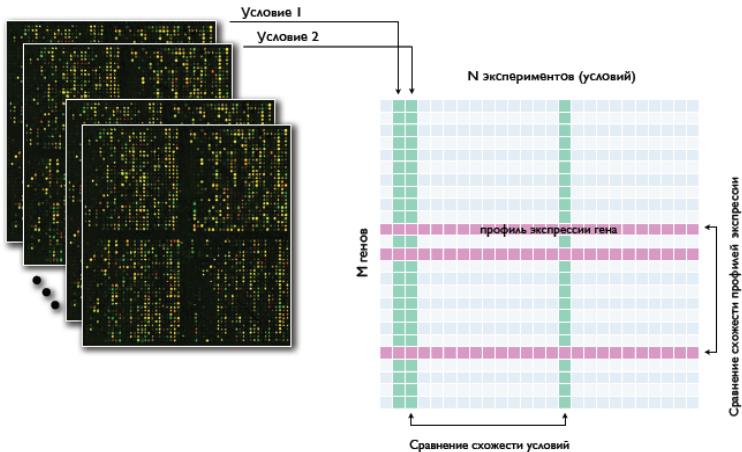
- оптимизируется не число ошибок, а качество ранжирования;
- сверхбольшие выборки;
- проблема конструирования признаков по сырым данным.

Анализ данных по экспрессии генов

ДНК-микрочипы - двумерный массив ДНК-зондов для тысяч нуклеотидных последовательностей, позволяющий измерять экспрессию генов при разных условиях















Анализ данных по экспрессии генов



Анализ данных по экспрессии генов

- **Кластеризация:** Группы генов выполняющие схожие функции имеют схожие профили экспрессии.
 - Задача: Поиск функциональных групп генов.
- **Классификация:** Клетка может находиться в разных состояниях (здоровая/раковая), различающихся уровнями экспрессии генов.
 - Задача: Определение состояния клетки на основе данных о профилях экспрессии генов.

Полногеномный поиск ассоциаций

												
	50	15	15	50	60	10	50	15	10	60	10	60
SNP-1	0	1	1	0	1	1	0	0	1	0	0	1
SNP-2	0	1	1	1	1	1	0	1	1	0	0	1
SNP-3	0	1	0	1	1	1	1	1	1	1	0	1
SNP-4	0	0	1	1	0	1	1	1	0	1	0	0
SNP-5	0	1	1	0	1	1	0	1	1	1	0	0
SNP-6	0	1	1	1	1	1	0	1	1	1	0	1
SNP-7	0	1	0	1	1	1	0	0	1	1	1	0
SNP-8	0	1	0	1	1	0	0	1	1	1	1	0
SNP-9	0	1	1	1	1	0	1	1	1	1	1	1
SNP-10	0	1	1	0	0	0	0	0	1	1	0	1

Полногеномный поиск ассоциаций



case



control

SNP-1	0	1	1	0	1	1	0	0	1	0	0	1
SNP-2	0	1	1	1	1	1	0	1	1	0	0	1
SNP-3	0	1	0	1	1	1	1	1	1	1	0	1
SNP-4	0	0	1	1	0	1	1	1	0	1	0	0
SNP-5	0	1	1	0	1	1	0	1	1	1	0	0
SNP-6	0	1	1	1	1	1	0	1	1	1	0	1
SNP-7	0	1	0	1	1	1	0	0	1	1	1	0
SNP-8	0	1	0	1	1	0	0	1	1	1	1	0
SNP-9	0	1	1	1	1	0	1	1	1	1	1	1
SNP-10	0	1	1	0	0	0	0	0	1	1	0	1

Распознавание рукописных символов (цифр)

Объект - рукописный символ (цифра).

Классы: 0,1,...,9

Примеры признаков:

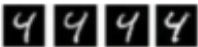



- *бинарные*: код (признаковое описание) - битовая матрица размера 32×32 .

1 — пиксел черный, 0 — пиксел белый.

Распознавание рукописных символов (цифр)

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Распознавание рукописных символов (цифр)

	Обучающая выборка	Новый пример
"четыре"		
"девять"		

Страховая компания (кластеризация)

- Информация об автомобилях и их владельцах:
марка автомобиля; стоимость автомобиля; возраст водителя; стаж водителя; возраст автомобиля
- Цель - разбиение автомобилей и их владельцев на классы, каждый из которых соответствует определенной рискованной группе.
- Наблюдения, попавшие в одну группу, характеризуются одинаковой вероятностью наступления страхового случая, которая впоследствии оценивается страховщиком.

No Free Lunch theorem

- D. H. Wolpert, The lack of a priori distinctions between learning algorithms and the existence of a priori distinctions between learning algorithms, Neural Computation 8 (1996)
- Теорема “No Free Lunch” утверждает, что не существует модели, которая работает лучше для всех задач
- Предположение о наилучшей модели для одних данных может быть не справедливым для других данных
- Необходимо пытаться использовать несколько моделей и выбирать из них лучшую для решения конкретной задачи

No Free Lunch theorem (постановка)

- $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ - обучающая выборка
- $T = \{(\mathbf{x}_{m+1}, y_{m+1}), \dots, (\mathbf{x}_n, y_n)\}$ - тестирующая выборка
- $A_S(\mathbf{x})$ - прогнозное значение модели для \mathbf{x}
- Ошибка на тестовой выборке

$$R(A) = \frac{1}{n - m} \sum_{i=m+1}^n [A_S(\mathbf{x}_i) \neq y_i]$$

No Free Lunch theorem (простейшая формулировка)

Теорема 1. Если предположить, что метки классов y_i вычисляются как $f(\mathbf{x}_i)$ для функции, выбранной случайным образом равномерно из всех возможных функций, то

$$\mathbb{E}_f(R(A|f)) = \frac{1}{2}.$$

- **Основные понятия машинного обучения:**
 - объект, ответ, метка признак, алгоритм, модель алгоритмов, метод обучения, эмпирический риск, переобучение.
- **Этапы решения задач машинного обучения:**
 - понимание задачи и данных;
 - предобработка данных и изобретение признаков;
 - построение модели;
 - сведение обучения к оптимизации;
 - решение проблем оптимизации и переобучения;
 - оценивание качества;
 - внедрение и эксплуатация.
- **Прикладные задачи машинного обучения:**
 - очень много, очень разных,
 - во всех областях бизнеса, науки, производства.

Различные алгоритмы и подходы к решению задач машинного обучения:

- Линейная регрессия
- Метод ближайших соседей
- Байесовский подход
- Машина опорных векторов
- Нейронные сети
- Деревья решений
- Бустинг (AdaBoost, Random Forest) и бэггинг
- Обучение без учителя, кластеризация

Ресурсы

- Wiki-портал <http://www.machinelearning.ru>
- Воронцов К.В. Машинное обучение (курс лекций)
 - см. <http://www.machinelearning.ru>,
 - видео-лекции
http://shad.yandex.ru/lectures/machine_learning.xml
- Ng A. Machine Learning Course (video, lecture notes, presentations, labs) <http://ml-class.org>
- Hastie T., Tibshirani R., Friedman J. The elements of statistical learning: Data Mining, Inference, and Prediction. Springer, 2009
- Загоруйко Н.Г. Прикладные методы анализа данных и знаний. Новосибирск: Изд-во Ин-та математики, 1999.

Ресурсы

- Вапник В.Н., Червоненкис А.Я. Теория распознавания образов. М.: Наука, 1974.
- Мерков А.Б. Распознавание образов. Введение в методы статистического обучения. М.: Едиториал УРСС, 2011.
- Барский А. Б. Нейронные сети: распознавание, управление, принятие решений. М.: Финансы и статистика, 2004.
- Флах П. Машинное обучение. Наука и искусство построения алгоритмов, которые извлекают знания из данных. М.: ДМК Пресс, 2015
- Домингос П. Верховный алгоритм: как машинное обучение изменит наш мир. М. : Манн, Иванов и Фербер, 2016.

- Система для статистических вычислений **R**
<http://www.r-project.org/>
- Библиотека алгоритмов для анализа данных Weka (Java) <http://www.cs.waikato.ac.nz/~ml/weka/>
- Пакет для решения задач машинного обучения и анализа данных Orange <http://orange.biolab.si/>
- Microsoft Azure Stack - a new hybrid cloud platform product (<https://azure.microsoft.com/ru-ru>)
- DL4J (Deeplearning4j) Deep Learning for Java (<http://deeplearning4j.org>)

Software еще

- MLlib – Apache's own machine learning library for Spark and Hadoop (<https://spark.apache.org/mllib/>)
- 0xdata's H2O's algorithms (<http://www.h2o.ai/>)
- Cloudera Oryx
(<https://code.google.com/archive/p/cuda-convnet2/>)
- ConvNetJS - deep learning
(<http://cs.stanford.edu/people/karpathy/convnetjs/>)
- WSO2 Machine Learner (<http://wso2.com/>)
- **Данные для экспериментов:** UCI Machine Learning Repository <http://archive.ics.uci.edu/ml/>

Литература по R

- Мастицкий С. Э., Шитиков В. К. (2014) Статистический анализ и визуализация данных с помощью R. - Электронная книга, 400 с.
- Савельев А. А. и др. (2007) Основные понятия языка R. Уч.-мет. пособие. Казань: КГУ, 29 с.
- Буховец А. Г. и др. (2010) Статистический анализ данных в системе R. Уч. пос. Воронеж: ВГАУ, 124 с.
- Зарядов И. С. (2010) Введение в статистический пакет R: типы переменных, структуры данных, чтение и запись информации, графика. М.: Изд-во РУДН, 207 с.
- Шипунов А. Б. и др. (2012) Наглядная статистика. Используем R! - М.: ДМК Пресс, 298 с.

Вопросы

?