

Машинное обучение. Вводная лекция.

Московский физико-технический институт, МФТИ

Москва

Этапы развития области искусственный интеллект

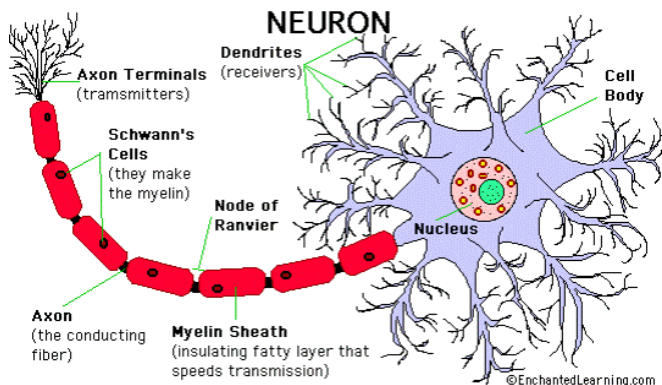
- Предыстория (мечты об ИИ и теоретические основы алгоритмистики);
- Появление самообучающихся машин конец 50х, начало 60х;
- Системы основанные на знаниях 70е;
- Первые нейронные сети (в основе алгоритм обратного распространения ошибки) начало 80х;
- Прорыв связанный с применением глубокого обучения в 2006г;

Причины бурного развития ИИ: глубокое обучение, вычислительные мощности, большие наборы качественных данных.

Модель МакКаллока-Питтса

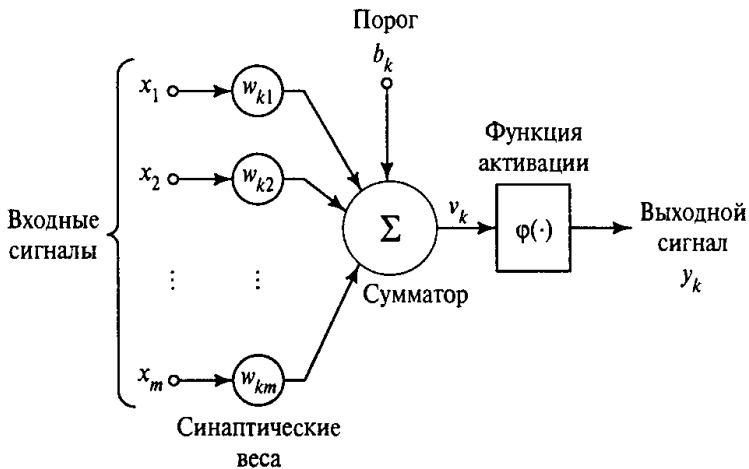
Линейный классификатор или персептрон является простейшей математической моделью нервной клетки- нейрона. Нейрон имеет множество разветвлённых отростков - дендритов, и одно длинное тонкое волокно - аксон, на конце которого находятся синапсы, примыкающие к дендритам других нервных клеток. Нервная клетка может находиться в двух состояниях: обычном и возбуждённом. Клетка возбуждается, когда в ней накапливается достаточное количество положительных зарядов. В возбуждённом состоянии клетка генерирует электрический импульс величиной около 100 мВ и длительностью около 1 мс, который проходит по аксону до синапсов.

Картинка нейрона



Особенности работы мозга: асинхронность, параллелизм, пластичность.

Мат модель нейрона



От нейрона к нейросетям

Тринадцатая проблема Гилберта решена Колмогоровым в 1957 году:

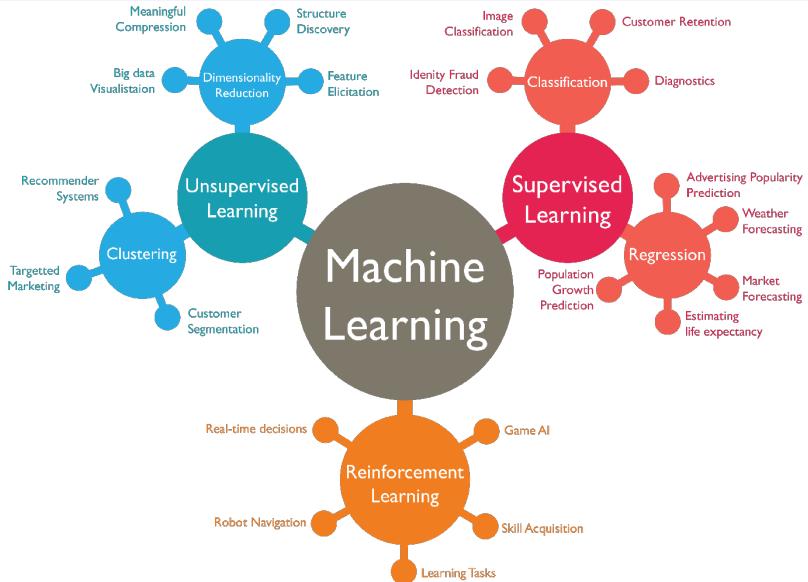
Теорема

Любая непрерывная n аргументов на единичном кубе $[0, 1]$ представима в виде суперпозиции непрерывных функций одного аргумента и операции сложения:

$$f(x^1, x^2, \dots, x^n) = \sum_{k=1}^{2n+1} h_k \left(\sum_{i=1}^n \varphi_{ik}(x^i) \right), \text{ где } h_k, \varphi_{ik}$$

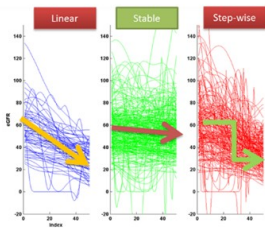
непрерывные функции, причем φ_{ik} не зависят от выбора f .

Нетрудно видеть, что записанное здесь выражение имеет структуру нейронной сети с одним скрытым слоем из $2n+1$ нейронов. Таким образом, двух слоёв уже достаточно, чтобы вычислять произвольные непрерывные функции, и не приближённо, а точно.

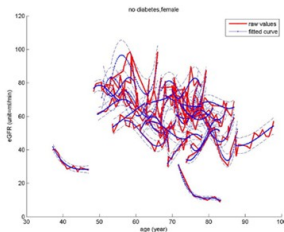


Решаемые задачи

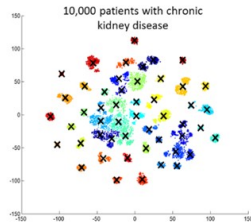
Classification



Regression



Clustering



Классификация

- байсовский метод (например наивный байсовский классификатор);
- метод опорных векторов (SVM);
- kNN (k -ближайших соседей);
- деревья решений;
- Random forest;

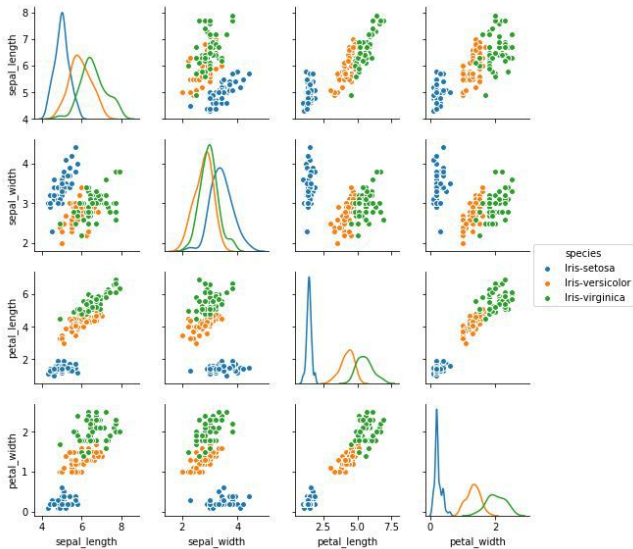
Параметры цветков

KNN Iris Classification



1. Длина наружной доли околоцветника (англ. sepal length);
2. Ширина наружной доли околоцветника (англ. sepal width);
3. Длина внутренней доли околоцветника (англ. petal length);
4. Ширина внутренней доли околоцветника (англ. petal width).

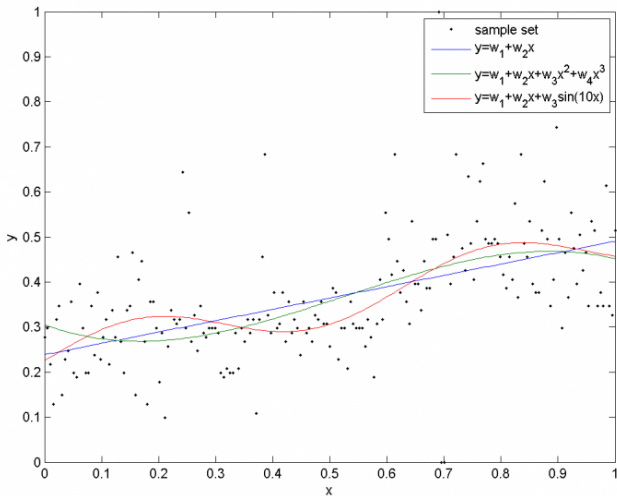
Иллюстрация классификации цветков ириса



Регрессия

Регрессионная модель есть функция независимой переменной и параметров с добавленной случайной переменной. Параметры модели настраиваются таким образом, что модель наилучшим образом приближает данные. Критерием качества приближения (целевой функцией) обычно является среднеквадратичная ошибка.

Регрессия. Иллюстрация.



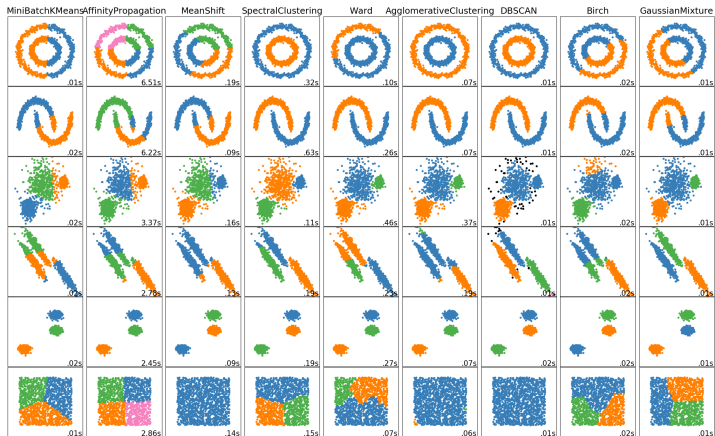
Кластеризация

Задача кластеризации (или обучение без учителя) заключается в следующем. Имеется обучающая выборка $X' = x_1, \dots, x_l \subset X$ и функция расстояния между объектами $\rho(x, x')$. Требуется разбить выборку на непересекающиеся подмножества, называемые кластерами, так, чтобы каждый кластер состоял из объектов, близких по метрике ρ , а объекты разных классов существенно различались. При этом каждому объекту $x_i \in X'$ приписывается метка (номер) кластера y_i .

Кластеризация

- Сгущения: внутрикластерные расстояния, как правило меньше межкластерных;
- Ленты: для любого объекта найдется близкий ему объект того же кластера, в то же время существуют не близкие друг другу объекты кластера;
- Кластеры с центром
- Кластеры могут соединяться перемычками;
- На фоне кластеров могут присутствовать лишние объекты;
- Перекрывание кластеров;
- Кластеры образованные по специфическим признакам;
- Отсутствие кластеров;

Кластеризация. Иллюстрация.



Области применения кластеризации

- Биология и биоинформатика: группировка сложных геномных последовательностей, выделение пространственных и временных сообществ;
- Маркетинг: выделение типичных групп покупателей, создание персональных предложений;
- Интернет: группировка веб-сайтов по смысловым значениям поисковых запросов;
- Компьютерные науки: определение границ - распознавание объектов;

Математические обозначения

Пусть задано множество объектов X и множество допустимых ответов Y . Исследуемый процесс описывается целевой функцией $y^* : X \rightarrow Y$, задано конечное множество объектов $\{x_1, \dots, x_l\} \subset X$. Пары объектов $X^l = (x_i, y_i)_{i=1}^l$ представляют собой обучающую выборку.

Основная задача состоит в восстановлении зависимости y^* , а именно в построении **решающей функции** $a : X \rightarrow Y$, причем на всем множестве X .

Решаются следующие вопросы

- как задаются объекты и какими могут быть ответы (представление данных);
- в каком смысле $a(x)$ приближает $y(x)$ (выбор нормы);
- правило построения $a(x)$;

Пространство признаков

- $\mathbb{Y} = \{0, 1\}$ бинарная классификация;
- $\mathbb{Y} = \{1, \dots, K\}$ - многоклассовая классификация;
- $\mathbb{Y} = \{0, 1\}^K$ - многоклассовая классификация с пересекающимися классами;

Методика обучения

В качестве **модели алгоритма** рассмотрим функцию $g(x, \theta) : X \times \Theta \rightarrow Y$, где Θ множество допустимых значений параметра.

$$g(x, \theta) = \sum_{j=1}^n \theta_j f_j(x), \quad Y = \mathbb{R}$$

$$g(x, \theta) = \text{sign} \sum_{j=1}^n \theta_j f_j(x), \quad Y = \{-1, +1\}$$

Метод обучения - это отображение $\mu : (X \times Y)^I \rightarrow A$, которое произвольной коеченой выборке ставит некоторый алгоритм.

Функционал качества

Функционал качества

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

Разновидности функции потерь (loss function):

- $\mathcal{L}(a, x_i) = [a(x) \neq y^*(x)]$ - индикатор ошибки,
- $\mathcal{L}(a, x_i) = |a(x) - y^*(x)|$ - отклонение от правильного ответа;
- $\mathcal{L}(a, x_i) = (a(x) - y^*(x))^2$ - квадратичная функция потерь;

Иными словами задача восстановления регрессии есть не что иное, как метод наименьших квадратов

$$\mu(X^I) = \arg \min_{\theta} \sum_{i=1}^I (g(x_i, \theta) - y_i)^2,$$

Принцип максимума правдоподобия

Вместо модели алгоритма $g(x, \theta)$, аппроксимирующей $y^*(x)$ зададим модель совместной плотности распределения объектов и ответов $\varphi(x, y, \theta)$.

В случае независимой последовательности наблюдений $p(X^I) = p(x_1, y_1) \dots p(x_I, y_I)$, получаем функцию правдоподобия

$$L(\theta, X^I) = \prod_{i=1}^I \varphi(x_i, y_i, \theta).$$

Задача найти θ при котором функция правдоподобия максимальна. [Пример с нормальным распределением и средней дисперсией]

Оценки обучающей способности

- Эмпирический риск на тестовых данных

$$HO(\mu, X^l, X^k) = Q(\mu(X^l), X^k) \rightarrow \min,$$

- Кросс-проверка, $L = l + k$, $X^L = X_n^l \sqcup X_n^k$:

$$CV(\mu, X^L) = \frac{1}{|N|} \sum_{n \in N} Q(\mu(X_n^l), X_n^k) \rightarrow \min,$$

- Эмпирическая оценка вероятности переобучения:

$$Q_\varepsilon(\mu, X^L) = \frac{1}{|N|} \sum_{n \in N} \left[Q(\mu(X_n^l), X_n^k) - Q(\mu(X_n^l), X_n^l) \geq \varepsilon \right] \rightarrow \min$$

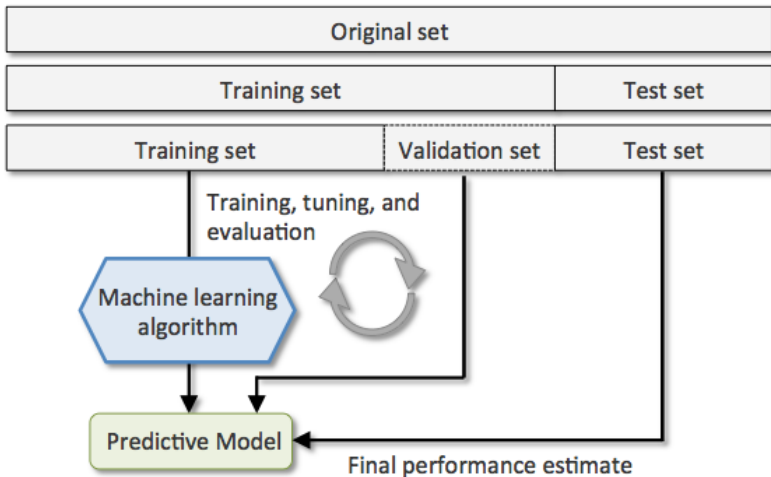
Кросс-проверка

Выборку $X^L = (x_i, y_i)_{i=1}^l$ разобьем N различными способами.

Для каждого разбиение $n = 1..N$ построим алгоритм $a_n = \mu(X_n^l)$ и вычислим среднее значение ошибки $Q_n = Q(a_n, X_n^k)$. Среднее арифметическое невязок всех указанных алгоритмов называется **скользящим контролем**

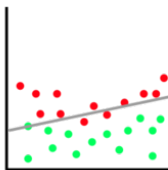
$$CV(\mu, X^L) = \frac{1}{N} \sum_{n=1}^l Q(\mu(X_n^l), X_n^l).$$

Скользящий контроль. Иллюстрация.

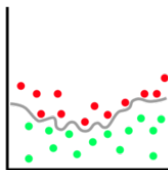


Переобучение

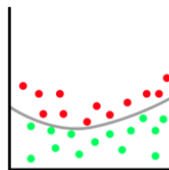
learning & regularization



Underfitting



Overfitting



Balanced

Еще раз регрессия

Задача линейной регрессии $g(x, \alpha) = \sum_{j=1}^n \alpha_j x_j$ сводится к следующей системе уравнений

$$F^T F \alpha = F^T y,$$

так что решение $\alpha^* = F^+ y$ ищется через псевдообратную матрицу $F^+ = (F^T F)^{-1} F^T$.

Мультиколлинеарность

В случае когда $\Sigma = F^T F$ имеет полный ранг, но близка к матрице с неполным рангом. Объекты выборки сосредоточены вблизи линейного подпространства меньшей размерности. Или малеенькие собств значения. Определим число обусловленности

$$\mu(\Sigma) = \frac{\max_{\|u\|=1} \|\Sigma u\|}{\min_{\|u\|=1} \|\Sigma u\|} = \frac{\lambda_{max}}{\lambda_{min}},$$

происходит увеличение погрешности при умножении на обратную матрицу

$$\frac{\|\partial z\|}{\|z\|} \leq \frac{\|\partial u\|}{\|u\|}.$$

Сглаживающий функционал (гребневая регрессия, L2 регуляризация), эффективная размерность, выбор константы регуляризации, L1 - регуляризация (лассо тибширани).

ДЗЗ изучения растительного покрова

Особое место среди попиксельной классификации занимают методы машинного обучения: классификация с использованием самоорганизующихся нейронных сетей, деревьев решений, опорных векторов. **Нейронные сети** более всего подходят для объектов с нечеткими границами.

Деревья решений здесь в отличие от "черного ящика" нейронных сетей можно указать, какие показатели использованы для разделения на классы (так называемые "ключи")

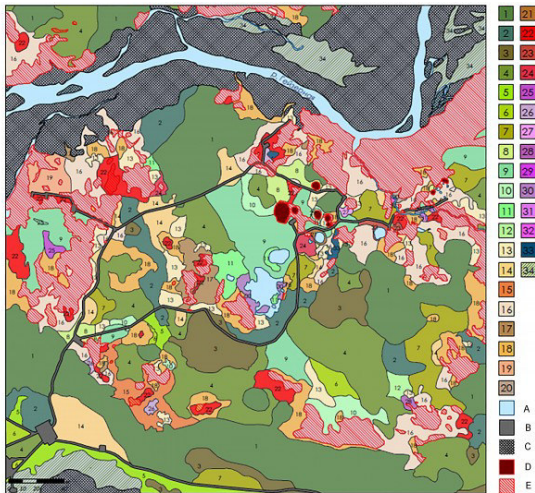


Рис. 1. Растительность Долины р. Гейзерной:

1 – каменистоберезняк разнотравный, 2 – ивняк шеломайниковый, 3 – сообщества ольхового стланика, 4 – крупнотравно-шеломайниковые луговые сообщества, 5 – высокотравные луговые сообщества с преобладанием лабазника и бодяка, 6 – высокотравные луговые сообщества с преобладанием вол-





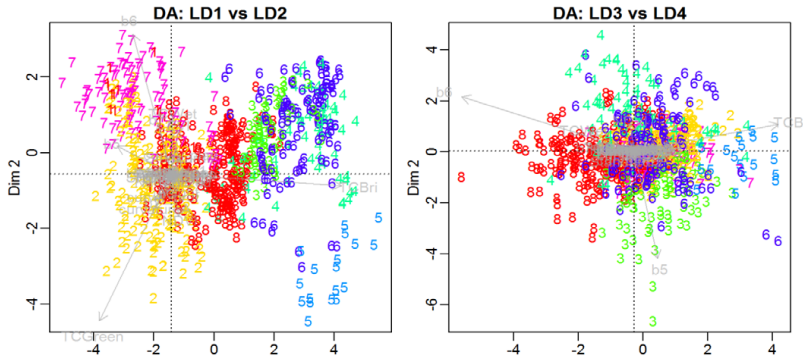


Рисунок 5. Варианты темныхвойных лесов хорошо разделяются в пространстве спектральных и морфометрических признаков методом дискриминантного анализа (цифрами обозначены синтаксоны в соответствии с таблицей 1)

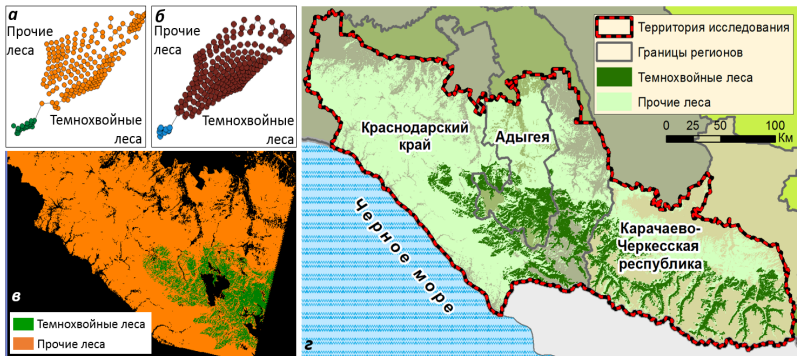


Рисунок 8. Картографирование темнохвойных лесов. а, б – структура нейронных сетей, построенных независимо по разным снимкам; в – результат классификации снимка нейронной сетью; г – конечный результат

Заключение

1. Анализ предметной области;
2. Коррекция данных, формализация признаков;
3. Построение прямой мат модели;
4. Подбор невязки, решение задачи оптимизации;
5. Поиск оптимального решения, регуляризация решения;
6. Оценки качества полученной модели;
7. Тестирование модельное и аппробация;