

Лекция 11. Композиционные алгоритмы распознавания образов. Бэггинг и Бустинг. Нейронные сети

Композиционные (комитетные) алгоритмы (КА) основаны на объединении ансамбля из нескольких классификаторов (распознающих алгоритмов), каждый из которых дает свою оценку ситуации, пригодную для принятия решения. На основе совокупности результатов, полученных от всех классификаторов, выдается итоговый общий результат в виде окончательного решения.

Композиционные алгоритмы настраиваются на основе обучающих данных (прецедентов). Основная идея – ансамбль классификаторов может обеспечить существенно более высокую достоверность распознавания, чем каждый из них в отдельности.

Общая структура композиционных алгоритмов

Пусть есть классы образов $\omega_1, \dots, \omega_M$ и связанная с ними смешанная проиндексированная обучающая выборка

$$X^N = \{x^{(1)}, \dots, x^{(N)}\}, D^N = \{d^{(1)}, \dots, d^{(N)}\}, d^{(*)} \in D = \{1, 2, \dots, M\}, \\ x^{(*)} \in \omega_j, d^{(*)} = j, \quad \{x^{(i)}, d^{(i)}, i = \overline{1, N}\}.$$

Вводится вспомогательное множество оценок R , называемое пространством оценок. Рассматриваются алгоритмы, имеющие вид суперпозиции

$$g'(x) = g'[g(x)],$$

где $g: X \rightarrow R$ называется базовым алгоритмом (БА), а функция $g': R \rightarrow D$ – решающим правилом.

Композицией БА $g_k(x), k = \overline{1, L}$ называется решающее правило вида, основанное на использовании суперпозиции

$$G'(x) = G'[G(g_1(x), \dots, g_L(x))], G: R^L \rightarrow R, G': R \rightarrow D,$$

где $G: R^L \rightarrow R$ называется корректирующей операцией.

Замечание: введение корректирующей операции для получения сначала комбинированной оценки на множестве R вместо того, чтобы использовать совокупность собственных решений БА (взять в качестве суперпозиции $G: D^L \rightarrow D$) дает возможность расширить перечень рассматриваемых вариантов объединения БА.

Рассмотрим два класса ω_1, ω_2 образов. Пусть имеется ансамбль БА $\{g_k(x), k = \overline{1, L}\}$, каждый из которых обучен формировать значения оценки. Эти значения далее могут быть использованы для принятия решений вида

$$g'_k(x) = \text{sign}[g_k(x)] = \begin{cases} +1 \rightarrow x \in \omega_1, \\ -1 \rightarrow x \in \omega_2. \end{cases}$$

Для случая двух классов с использованием понятия вероятности ошибки p_g вводятся следующие определения для базовых классификаторов:

классификатор называется плохим, если $p_g \geq 1/2$;

классификатор называется слабым, если $p_g = 1/2 - \varepsilon$;

классификатор называется сильным, если $p_g = \varepsilon$.

Различают линейные композиции (комбинации) БА и смеси БА нелинейного вида. Стандартный вариант **линейного** ансамбля БА:

$$G(x) = G\left[\sum_{k=1}^L \alpha_k g_k(x)\right], \quad G'(x) = \text{sign}[G(x)], \quad \sum_{k=1}^L \alpha_k = 1,$$

где $\alpha_k \geq 0, k = \overline{1, L}$ – весовые коэффициенты, учитывающие вклад каждого базового алгоритма. Если $\alpha_k = 1/L, k = \overline{1, L}$, то такое объединение называется простым голосованием.

Простейший вариант: использование суммы бинарных пороговых классификаторов

$$G'(x) = \text{sign}[G(x)] = \text{sign}\left[\sum_{k=1}^L \text{sign}(g_k(x))\right].$$

Такой классификатор обеспечивает принятие решений путем «голосования» по большинству решений, принимаемых членами комитета в пользу того или иного класса.

Наиболее часто встречающиеся подходы к реализации КА: **бэггинг и бустинг**.

Бэггинг (bagging) – сокращение от bootstrap aggregating – агрегированный бутстреп. Bootstrap – подход к увеличению репрезентативности обучающей выборки, основанный на случайном извлечении (с последующим возвращением) из исходной обучающей выборки нескольких подмножеств примеров с целью получения устойчивых оценок. При выполнении бэггинга композиция составляется на основе обучения нескольких классификаторов по бутстреп - множествам, полученным из общей обучающей выборки, и агрегирования результатов их работы путем простого голосования. **Общая идея бэггинга** – снижение зависимости «экспертов» – базовых классификаторов ансамбля друг от друга.

Бустинг (boosting – улучшение) – процедура итеративного последовательного построения композиций базовых алгоритмов. Каждая следующая композиция включает большее количество БА, стремясь при этом компенсировать недостатки предыдущей композиции.

На каждой итерации осуществляется перевзвешивание наблюдений, добавление и коррекция весов базовых алгоритмов. В итоге происходит «усиление» слабых классификаторов и повышение результирующей эффективности распознавания. **Общая идея бустинга** – «эксперты» (базовые классификаторы) учатся на ранее допущенных ошибках других классификаторов.

1. Бэггинг: деревья решений и композиции «случайный лес»

Для иллюстрации идей бэггинга первоначально рассмотрим алгоритмы принятия решений, известные под названием «деревья решений». Идея, лежащая в основе деревьев решений (ДР), состоит в разбиении множества возможных значений вектора признаков на непересекающиеся множества и настройке простой модели решений для каждого такого множества.

ДР представляет собой ориентированный (заданы направления соединения вершин ребрами графа) связный граф без циклов (без обратных связей).

Под корневым деревом понимается дерево, в котором одна вершина выделена и называется корнем. В качестве ДР рассматриваются только ориентированные корневые деревья, в которых дуги (ориентированные ребра) направлены от корня к вершинам. Они удовлетворяют следующим условиям:

существует только одна вершина, называемая корнем, в которую не ведет ни одна дуга;

в каждую вершину (исключая корень) ведет только одна дуга; существует единственный путь от корня к любой вершине.

Если (v, w) – некоторая дуга, то вершина v называется родителем w , а вершина w – потомком вершины v . Вершина, не имеющая потомков, называется терминальной вершиной или листом. Дерево называется бинарным, если каждая его вершина (за исключением терминальных вершин) имеет ровно двух потомков.

Алгоритм формирования дерева решений для распознавания образов реализуется на основе смешанной проиндексированной обучающей выборки $XD^N = \{(x^{(i)}, d^{(i)}), i = \overline{1, N}\}$. Пусть $X \subseteq \mathbf{R}^n$ пространство образов. Деревом решений называется граф, в котором для каждой вершины t заданы:

- подмножество в признаковом пространстве $X_t \subset X$, при этом с корневой вершиной связывается все пространство X ;
- подвыборка обучающей выборки $XD_t \subset XD^N$, такая, что $XD_t = \{(x, d) \in XD^N, x \in X_t\}$, при этом с корневой вершиной связывается вся XD^N ;
- правило (функция) $h_t: X_t \rightarrow \{0, 1, \dots, k_{t-1}\}$, где $k_t \geq 2$ – количество потомков вершины t , определяющее разбиение множества X_t на k_t непересекающиеся подмножества, при этом для терминальных вершин правило не вводится.

Обозначим $t_{i(t)}$, $i = 0, 1, \dots, k_{t-1}$ вершину, являющуюся i -м потомком вершины t . Подмножество X_t и правило h_t определяют подмножества $X_{t_{i(t)}}$, возникающие после разбиения, следующим образом:

$$X_{t_{i(t)}} = X_t \cap \{x \in X : h_t(x) = i\}.$$

Цель построения дерева решений состоит в распознавании нового вектора x . Процесс принятия решений начинается с корневой вершины и состоит в последовательном применении правил, связанных с вершинами. Результатом этого процесса является определение терминальной вершины t такой, что $x \in X_t$. Вектор x относят к классу, который наиболее часто встречается в подвыборке, связанной с этой терминальной вершиной.

Классический алгоритм CART (Classification And Regression Trees), реализующий данный подход, основан на идее рекурсивного разбиения обучающей выборки на две более однородные подвыборки с помощью одного из признаков. Для реализации этой идеи необходимо определить понятие однородности, которая обычно рассчитывается на основе показателей, характеризующей степень загрязненности подвыборок при выполнении разбиения.

Показатель загрязненности (impurity). Общая идея состоит в следующем. Пусть t – вершина дерева решений, $XD_t \subset XD^N$ – подвыборка, связанная с этой вершиной, и $J(t)$ – загрязненность вершины. Загрязненность вершины равна 0, если XD_t содержит прецеденты только одного класса и максимальна в случае, если XD_t содержит одинаковое число прецедентов каждого класса.

Используются различные показатели загрязненности. Индекс Джини представляет собой частоту ошибочной классификации при случайном назначении меток классов образам подвыборки XD_t .

$$J(t) = 1 - \sum_{i=1}^M N_t^2(\omega_i)$$

Расщепление деревьев (split). Правило разбиения множества X , связанное с каждой вершиной дерева решений, называется расщеплением.

Бинарное расщепление вершины t можно рассматривать как функцию $h_t: X_t \rightarrow \{0,1\}$, $x \in X_t$, где в случае $h_t(x)=0$ вектор относится к первому (левому) потомку, а в случае $h_t(x)=1$ – ко второму (правому).

Расщепление состоит в выборе порогового значения x_{*0} , которое минимизирует используемый показатель загрязненности. В этом случае

$$h_t(x)=0, x_* < x_{*0}, h_t(x)=1, x_* \geq x_{*0}.$$

Области решений будут представлять собой многомерные параллелепипеды, при этом любая часть границы области решений, соответствующая данному расщеплению, будет представлять собой часть гиперплоскости, параллельной соответствующей координатной оси в пространстве признаков.

Критерии остановки расщепления. Процесс расщепления вершин имеет естественный предел, когда становится невозможно уменьшить загрязненность очередной вершины. Такие вершины объявляются терминальными, а соответствующее дерево называется полным.

В полном ДР каждая терминальная вершина содержит примеры только одного класса, а само дерево содержит большое количество вершин.

Полное дерево, как правило, обладает низкой достоверностью классификации. Это связано с уже упомянутой проблемой переобучения, заключающейся в том, что построенная деревом слишком точ-

ная модель обучающих данных фактически описывает только эту выборку и непригодна в качестве модели для другой выборки.

Использование критериев остановки расщепления ДР позволяет решить проблему построения слишком больших и детальных деревьев играют. Наиболее простой способ состоит в задании минимального числа для количества наблюдений в подвыборках, соответствующих терминальным вершинам (или минимальной доли наблюдений обучающей выборки).

Усечение деревьев (pruning). Недостатком подхода, основанного на способах остановки расщепления вершин для предотвращения построения больших деревьев, является то, что решение об остановке принимается без учета ситуации, которая могла бы сложиться при продолжении расщепления.

Альтернативный подход – построение полных деревьев и их дальнейшее усечение.

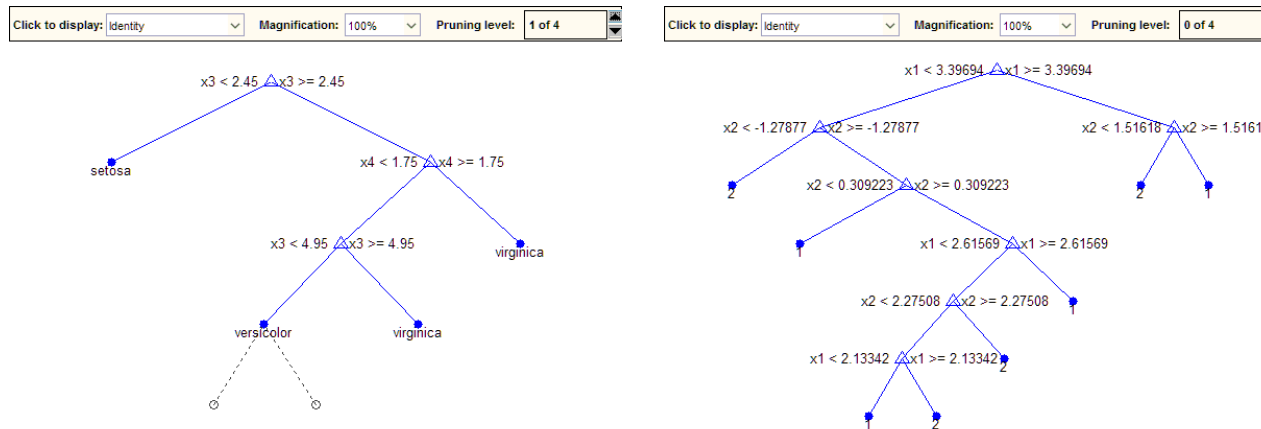
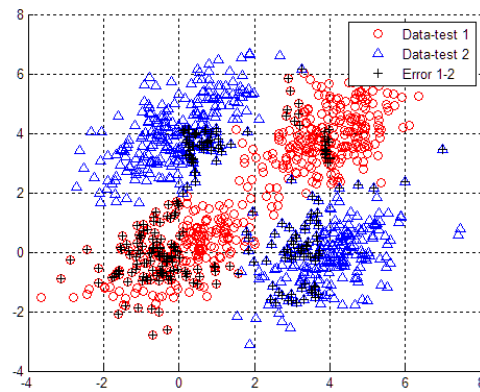
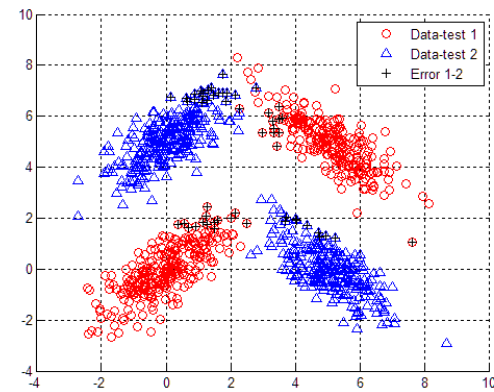


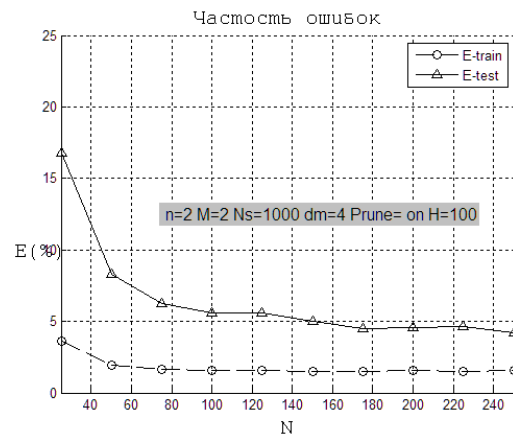
Рис.1. Исходное и усеченное на один уровень деревья решений в примере fisheriris (а), дерево решений в примере «запутанные восьмерки» (б)



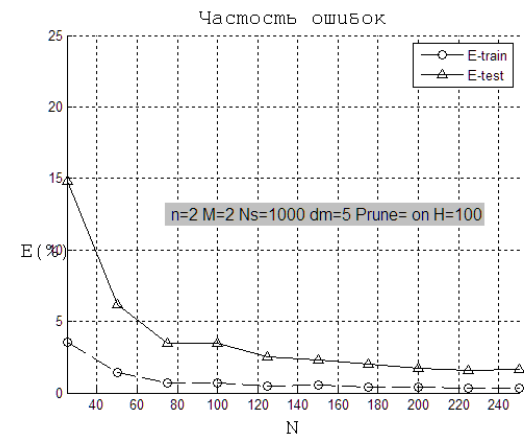
а)



б)



в)



г)

Рис.2. Результаты тестирования алгоритма ДР на основе дерева решений для примера «запутанные восьмерки»

Усечение означает процедуру замены в построенном полном дереве некоторой вершины и связанного с ней поддерева терминальной вершиной.

Большинство методов усечения основано на оценке чувствительности поддеревьев по отношению к некоторой мере и удалению поддеревьев, которые оказывают минимальное влияние на эту меру.

Обобщающая способность ДР не так хороша, как у алгоритма SVM: наблюдается существенное отличие числа ошибок, допущенных на обучающей и на тестирующей выборках (рис.1,2). В примере «запутанные восьмерки» степень «запутанности» классов определяется параметром dm , задающим относительный сдвиг центров областей локализации данных в каждом классе (рис.2а,б). Если $dm=0$, то классы сливаются, если dm становится достаточно большим, то каждый класс состоит из двух не пересекающихся областей эллипсоидальной формы.

Случайный лес (Random Forest) на основе бэггинга. Случайные леса реализуют подходы к построению ансамбля классификаторов в виде ДР, основанные на манипулировании с обучающими данными и манипулировании с признаками с включением в эти процессы элементов случайности.

«Случайный лес», основан на формировании бутстреп - обучающей выборки для каждого классификатора ансамбля путем случайной выборки с возвращением из исходной обучающей выборки. При этом каждый раз получается подвыборка того же объема, что и исходная обучающая выборка.

Каждая бутстреп - выборка содержит в среднем примерно 63% наблюдений исходной обучающей выборки: поскольку выборка с возвращением, то некоторые наблюдения в нее не попадают, а некоторые попадают несколько раз).

В качестве ансамбля БА в случайном лесе рассматривается ансамбль ДР, каждое из которых строится на основе бутстреп-подвыборки из исходной обучающей, причем для расщепления вершин используется только доля случайно отбираемых признаков. При объединении, или агрегирования решений отдельных классификаторов используется метод голосования.

Информация, касающаяся деталей алгоритма, обзор и примеры могут быть найдена по адресу:

<http://www.stat.berkeley.edu/users/breiman/RandomForests>

Стандартный алгоритм построения случайного леса. На этапе обучения реализуется индукция леса. В цикле для $k = \overline{1, L}$ деревьев ансамбля выполнить следующие действия:

- по исходной обучающей выборке $XD^N = \{(x^{(i)}, d^{(i)}), i = \overline{1, N}\}$ сформировать бутстреп выборку XD_k^N (случайную выборку того же объема с возвращением элементов);
- по бутстреп выборке XD_k^N индуцировать не усеченное дерево решений T_k с минимальным количеством наблюдений в терминальных вершинах, равным n_{\min} , рекурсивно следуя следующей процедуре:
 - а) из исходного набора n признаков случайно выбрать p признаков;
 - б) из p признаков выбрать признак, который обеспечивает наилучшее расщепление;
 - в) расщепить выборку, соответствующую обрабатываемой вершине, на две подвыборки;

В результате получаем ансамбль деревьев решений $\{T_k, k = \overline{1, L}\}$.

Распознавание новых наблюдений: выбирается класс образов $x \in \omega_j$, который максимальное число раз выбран базовыми классификаторами:

$$\tilde{d}(x) = i = \arg \max_{j=1, \overline{L}} \left\{ \sum_{k=1}^L I_k(x \in \omega_j) \right\}, \quad I_k(x \in \omega_j) = \begin{cases} 1, x \in \omega_j, \\ 0, x \notin \omega_j. \end{cases}$$

Одним из достоинств случайных лесов является то, что для оценки вероятности ошибочной классификации нет необходимости использовать тестовую выборку. Оценка вероятности ошибочной классификации случайного леса может осуществляться внутренним образом на основе метода **"Out-Of-Bag" (OOB)**.

При OOB выполняется распознавание каждого вектора $x^{(*)} \in X^N$ с использованием при голосовании только тех деревьев, которые строились по бутстреп - выборкам, не содержащим $x^{(*)}$. Частота ошибочно классифицированных векторов обучающей выборки при таком способе представляет собой оценку вероятности ошибочной классификации случайного леса методом OOB.

Случайные леса обладают рядом привлекательных свойств:

- повышение достоверности распознавания за счет слабой зависимости деревьев вследствие двойной инъекции случайности – посредством бэггинга и использования случайного набора признаков при расщеплении каждой вершины;
- сложная задача усечения полного дерева решений снимается, поскольку деревья в случайном лесу не усекаются (это также приводит к высокой вычислительной эффективности).
- проблема переобучения не стоит так остро, поскольку ансамбль ДР уже не является единственным детально настроенным алгоритмом;
- простота настройки алгоритма: единственными параметрами являются количество деревьев и количество признаков, случайно отбираемых для расщепления в каждой вершине.

Алгоритм «случайный лес» показывает в целом лучшие результаты по сравнению алгоритмом, основанным на использовании одного дерева решений (рис.3).

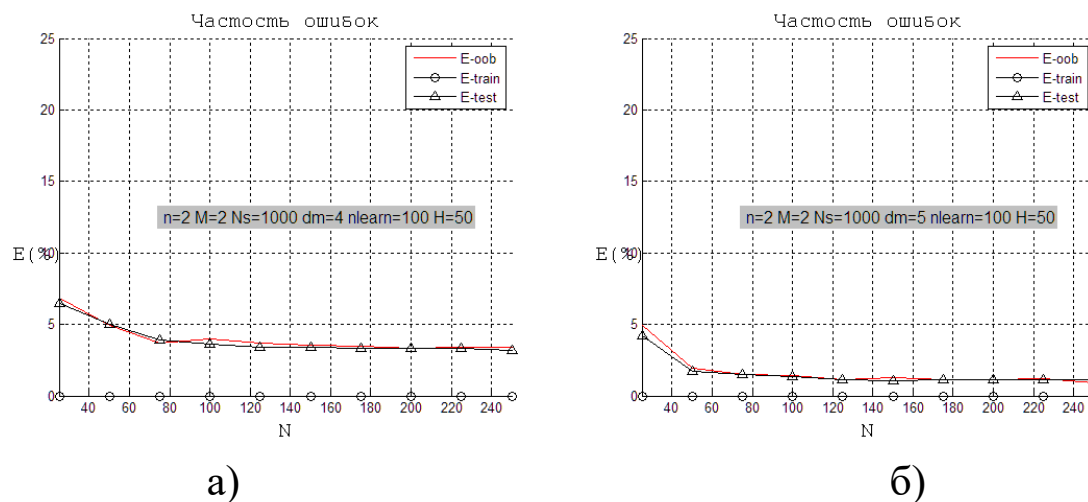


Рис.3. Результаты тестирования алгоритма «случайный лес» для примера «запутанные восьмерки»

2. Композиции базовых алгоритмов на основе бустинга

Композиции, формируемые на основе бустинга, обладают следующими особенностями:

реализация последовательного процесса обучения за несколько итераций;

учет предыдущих результатов классификации образов из обучающей выборки при выполнении каждой итераций в плане допущенных на них ошибок;

учет результатов работы ранее используемых базовых классификаторов в плане допущенных ими ошибок.

Цель подобной обработки состоит в том, чтобы **усилить слабые классификаторы**.

Описание стандартного алгоритма AdaBoost (adaptive boosting). Пусть на входе имеется $XD^N = \{(x^{(i)}, d^{(i)}), i = \overline{1, N}\}$ обучающая выборка для двух классов образов ω_1, ω_2 $d^{(i)} \in \{-1, +1\}$. При формировании композиции выполняется L итераций. В качестве БА рассматриваются алгоритмы

$$g_k(x) = \begin{cases} +1 \rightarrow x \in \omega_1, \\ -1 \rightarrow x \in \omega_2. \end{cases}$$

В качестве итоговой композиции рассматривается алгоритм

$$G(x) = \text{sign} \left[\sum_{k=1}^L \alpha_k g_k(x) \right].$$

Для произвольного БА вводится взвешенное число допущенных ошибок

$$E(g_k, w) = \sum_{i=1}^N w_i I(g_k(x^{(i)}) \neq d^{(i)}), \quad \sum_{i=1}^N w_i = 1$$

Общий функционал качества распознавания, который следует минимизировать, определяется как

$$Q(g, \alpha) = \sum_{i=1}^N [M(x^{(i)}) < 0] = \sum_{i=1}^N \left[d^{(i)} \sum_{k=1}^L \alpha_k g_k(x^{(i)}) < 0 \right] \rightarrow \min.$$

1. Назначить одинаковые начальные веса для элементов обучающей выборки $w_i = 1/N$, $i = \overline{1, N}$ (имеют смысл априорных вероятностей).
2. Для каждого шага $k = 1, 2, \dots, L$

- обучить базовый алгоритм $g_k(x)$ на XD^N , который минимизирует взвешенную ошибку классификации $E_k = \min E(g_k, w)$;
- если $E_k = 0$, то $G(x) = g_k(x)$ и переход к п. 3.
- если $E_k > 1/2$, то переход к п. 4.
- вычислить коэффициент $\alpha_k = \frac{1}{2} \ln((1 - E_k)/E_k)$ и зафиксировать его как вес соответствующего БА;
- пересчитать веса обучающих примеров следующим образом $w'_i = w_i \exp(-\alpha_k d^{(i)} g_k(x^{(i)}))$, $i = \overline{1, N}$
- выполнить нормировку весовых коэффициентов образов $\sum_{i=1}^N w'_i = 1$.

3. Зафиксировать композицию $G(x) = \text{sign} \left[\sum_{k=1}^L \alpha_k g_k(x) \right]$.

4. Окончание.

В алгоритме при перевзвешивании наблюдений усиливается роль тех из них, на которых допущены ошибки. Большие веса получают те объекты, которые «плохо» классифицировались на предыдущих шагах.

Весовые коэффициенты БА учитывают допущенные ими ошибки как напрямую, так и косвенно, через веса обучающих образов. Таким образом, на каждом последующем шаге новый классификатор приспособляется исправить ранее допущенные ошибки предыдущих классификаторов.

При выполнении алгоритма по схеме AdaBoost на каждом шаг добавляется новый классификатор, при этом значения коэффициентов БА, полученные на предыдущих шагах не изменяются.

Использование приведенных соотношений для перевзвешивания весов образов эквивалентно замене ступов на экспоненциальную функцию

$$[M(x^{(i)}) < 0] \leq \exp[-M(x^{(i)})]$$

и минимизации функционала вида

$$\bar{Q}(g, X) = \sum_{i=1}^N \exp[-M(x^{(i)})] \geq Q(g, X), H(M) = \exp(-M)$$

В качестве БА в алгоритмах типа AdaBoost чаще всего используются **деревья решений малой высоты** и простые правила порогового типа.

Среда MATLAB предоставляет широкие возможности по применению следующих вариантов организации бустинга и, прежде всего, следующие методы:

- 'AdaBoostM1' – классический AdaBoost для распознавания двух классов;
- 'AdaBoostM2' – классический AdaBoost для распознавания многих классов;
- 'LogitBoost' – бустинг с использованием в качестве $H(M)$ логарифмической функции;
- 'GentleBoost' – бустинг с использованием в качестве $H(M)$ квадратичной функции;
- 'RobustBoost' – бустинг с использованием приемов, направленных на исключение чрезмерного повышения весов отдельных наблюдений и последующего игнорирования большей частью обучающих данных;
- 'LPBoost' – бустинг с использованием для оптимизации коэффициентов методов линейного программирования.

Результаты тестирования различных вариантов алгоритмов, основанных на бустинге (AdaBoostM1 и GentleBoost) для примера «запутанные восьмерки» при различных конфигурациях областей локализации классов представлены на рис.4 (стандартная, рис.2а) и рис. 5 (отличающаяся от стандартной, рис.5а). Графики отражают зависимости частоты ошибок от объема обучающих данных.

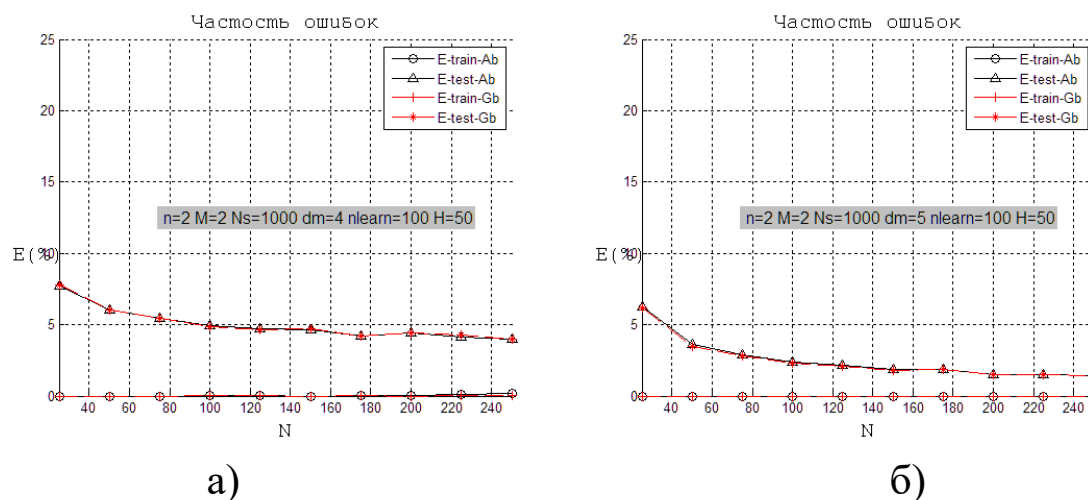
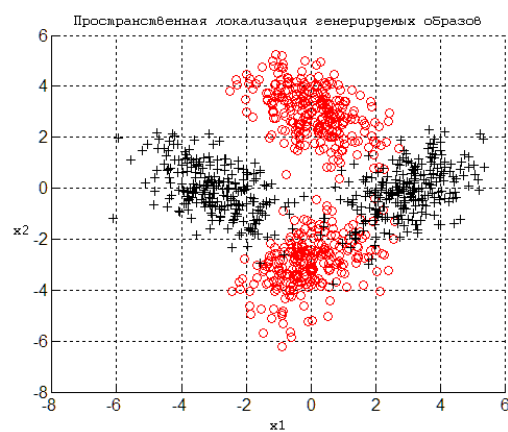
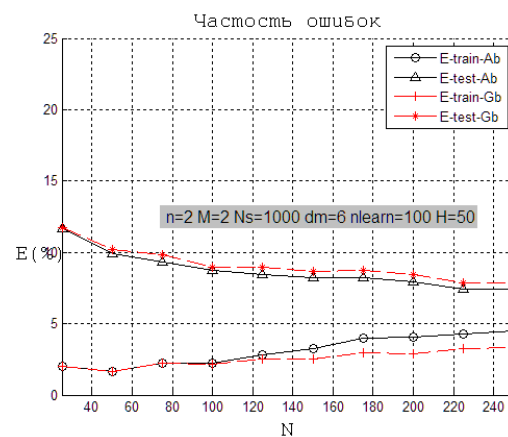


Рис.4. Результаты тестирования алгоритмов AdaBoostM1 и GentleBoost для примера «запутанные восьмерки» при стандартной конфигурации рис.2а



а)



б)

Рис.5. Результаты тестирования алгоритмов AdaBoostM1 и GentleBoost для примера «запутанные восьмерки» для конфигурации рис.5а

3. Нейронные сети

Одним из альтернативных подходов, который может быть использован для построения алгоритмов анализа данных в задачах оценивания, регрессии, распознавания является нейросетевой подход, т.е. подход, основанный на использовании аппарата искусственных нейронных сетей (ИНС). Его возможности позволяют реализовать функциональные модели любых процессов как «черных ящиков», наделенных однозначным описанием входа и выхода и осуществляющих преобразование входной информации в выходную в соответствии с требуемым результатом. При использовании нейросетевой модели преобразования данных обычно в явном виде не воспроизводится логика выполняемого алгоритма.

Для построения алгоритма обработки информации формируется некоторая достаточно **универсальная вычислительная среда**, которая на определенных примерах обучается выдаче нужных выходных реакций при поступлении тех или иных воздействий на вход.

Во многом нейросетевые алгоритмы обработки информации близки к композиционным, так как используют ансамбли базовых процессорных элементов (БПЭ) или нейронов осуществляющих относительно независимую обработку входных данных с объединением полученных в каждом элементе информации в конечный результат.. В тоже время, есть существенные различия, которые далее будут вполне очевидны. Главное отличие состоит в том, что в отличие от композиционных алгоритмов, в ИНС ансамбль БВЭ организуется послойно. При этом результат (сигнал), получаемый на выходе нейронов одного слоя, поступает на входы нейронов другого слоя или на один из выходов сети.

В зависимости от количества слоев, различают не глубокие и глубокие нейронные сети.

Каждая искусственная нейронная сеть состоит из множества БПЭ – нейронов (рис. 6). В качестве исторической модели БПЭ известна модель нейрона Мак-Каллока и Питса (рис.6).

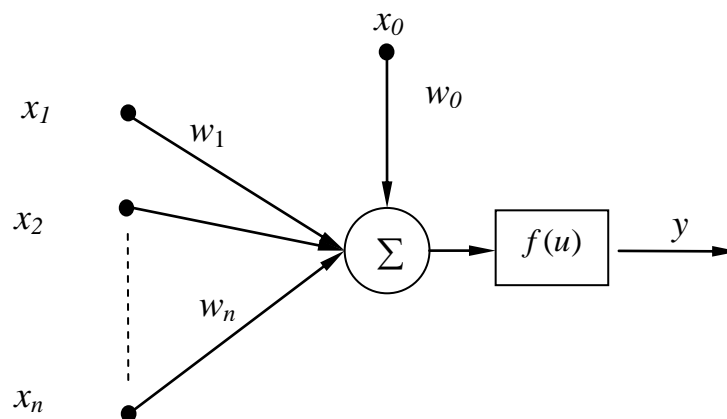


Рис.6. Модель искусственного нейрона Мак-Каллока и Питса

Каждый такой нейрон выполняет стандартный набор операций. Сначала производится суммирование совокупности входных элементарных сигналов – компонентов входного вектора $x = (x_1, \dots, x_n)^T$ с весовыми коэффициентами $w = (w_0, w_1, \dots, w_n)^T$, то есть вычисление величины

$$u = \sum_{i=1}^n w_i x_i + w_0 = \sum_{i=0}^n w_i x_i .$$

Здесь формально $x_0 \equiv 1$, а w_0 играет роль постоянного смещения при взвешенном суммировании входных сигналов нейрона.

Веса $w = (w_0, w_1, \dots, w_n)^T$ называются весами нейронных связей; они определяют «силу» влияния каждого входного сигнала на выходную реакцию нейрона. После взвешенного суммирования в общем случае осуществляется нелинейное преобразование и формирование выходного сигнала нейрона – в данном случае скалярной величины y как функции $y = f(u)$

$$y = f(u) = f\left(\sum_{i=0}^n w_i x_i\right).$$

Функция $f(u)$, осуществляющая это преобразование, называется функцией активации. Она может иметь различный вид, в том числе и линейный. В модели Мак-Каллока и Питса функция активации была пороговой. В качестве пороговой функции используется как униполярная, так и биполярная функции:

$$f(u) = \begin{cases} 1 & \text{для } u > 0, \\ 0 & \text{для } u \leq 0, \end{cases} \quad f(u) = \text{sign}(u) = \begin{cases} 1 & \text{для } u > 0, \\ -1 & \text{для } u \leq 0. \end{cases}$$

В других реализациях БПЭ наиболее часто используются активационные функции сигмоидального вида: униполярная функция «логический сигмоид» и биполярная функция «гиперболический тангенс». Их выражения представляются следующим образом:

$$f(u) = \frac{1}{1 + e^{-\beta u}}, \quad 0 < f(u) < 1, \quad f(u) = \frac{e^{\beta u} - e^{-\beta u}}{e^{\beta u} + e^{-\beta u}}, \quad -1 < f(u) < 1.$$

Указанные функции в диапазонах значений своего аргумента монотонно возрастают и хороши тем, что являются дифференцируемыми, что важно с точки зрения применения стандартных способов обучения.

Архитектура ИНС с точки зрения организации связей между нейронами может быть самой разнообразной. Подобные связи могут осуществляться: в прямом (от входа к выходу) и обратном направлениях (с обратной связью); при отсутствии задержек по времени (безынерционная сеть) при передаче сигналов от одного нейрона к другому и при наличии таких задержек; при наличии связей, реализующих конкурентное взаимодействие между нейронами, и при отсутствии таких взаимодействий. Важное значение для алгоритма обучения и функционирования ИНС имеет вид используемых активационных функций.

В качестве примера можно рассмотреть типовую архитектуру многослойной ИНС прямого распространения, часто используемой в задачах распознавания, а также оценивания и регрессии. Сеть обучается с учителем на основе совокупности примеров, организованной в виде индексированной обучающей выборки.

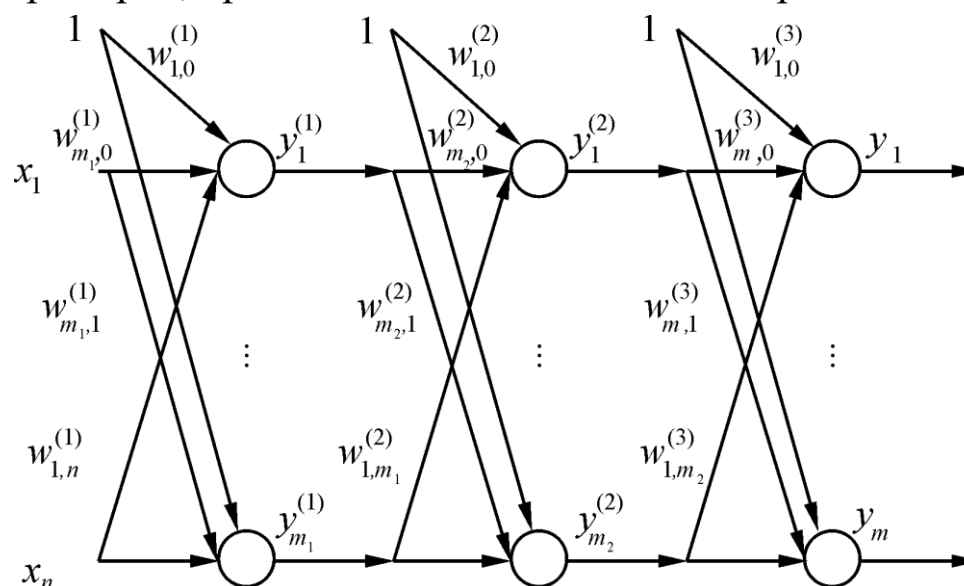


Рис.7. Типовая архитектура многослойной нейронной сети распространения

Рассматриваемая сеть называется нейронной сетью персептронного типа или многослойным персептроном (MLP). В MLP соединение нейронов осуществляются каскадным образом, по слоям; при этом находящиеся в одном слое нейроны не обмениваются сигналами. Задержки при передаче сигналов отсутствуют, используются активационные функции сигмоидального вида и линейного вида. Как правило, устанавливается общий вид активационной функции для БПЭ одного слоя.

На входе сети размещается совокупность контактов, являющихся распределителями входной информации. Слои нейронов разделяются на скрытые (внутренние) и выходной. Скрытым слоем является любой слой, который в прямом направлении соединен с другим слоем. Выходной слой формирует общий выход сети. Совокупность входных контактов иногда определяют как еще один, нулевой слой.

Представленная, как пример, на рис.7 нейронная сеть прямого распространения имеет $K=3$ слоя с m_i , $1 \leq i \leq K$ нейронов в каждом слое.

Векторы входов $x^{(i)}$ и выходов $y^{(i)}$ для каждого слоя имеют, соответственно, длину m_{i-1} и m_i , причем $y^{(i)} = x^{(i+1)}$, $1 \leq i \leq K-1$, то есть выходы нейронов i -го скрытого слоя являются одновременно входами нейронов $i+1$ -го слоя. Для скрытого первого слоя вектор входов является вектором входов всей сети и имеет длину $m_0 = n$, то есть $x^{(0)} = x = (x_1, \dots, x_n)^T$. Аналогично вектор выходов K -го слоя является вектором выходов всей сети y и для него удобно определить $m_K = m$ и, соответственно, $y^{(K)} = y = (y_1, \dots, y_m)^T$. Для каждого слоя определяются матрицы весов нейронных связей $W^{(i)} = \|w_{k,r}^{(i)}\|$ размером $m_{i-1} \times m_i$, а также матрицы (векторы – строки) весов постоянных смещений $W_0^{(i)} = (w_{1,0}^{(i)}, \dots, w_{m_i,0}^{(i)})$. Значения индексов k, r для каждого элемента матриц весов

$W^{(i)}$ и $W_0^{(i)}$ определяют, соответственно, номер нейрона, к которому направляется данная связь, и номер нейрона, от которого она исходит.

При фиксированных значениях весов нейронных связей можно записать результирующее выражение, определяющее преобразование входного вектора $x = (x_1, \dots, x_n)^T$ в выходной вектор $y = (y_1, \dots, y_m)^T$ НС

$$y_k = f\left(\sum_{r=0}^{mK-1} w_{k,r}^{(K)} f\left(\sum_{t=0}^{mK-2} w_{r,t}^{(K-1)} \dots f\left(\sum_{h=0}^n w_{t,h}^{(1)} x_h\right)\right)\right), \quad k = \overline{1, m}, \quad (1)$$

$$y = (y_1, \dots, y_m)^T = \Psi(x_1, \dots, x_n) = \Psi(x).$$

Таким образом, нейронные сети выполняют достаточно простой набор стандартных операций, вычисляя суперпозиции простых функций одного переменного и их линейных комбинаций.

При обучении ИНС прямого распространения давать «правильный ответ» на выходе ограничиваются, как правило, настройкой весовых коэффициентов связей между БПЭ и весов смещений $W = \{W^{(i)}, W_0^{(i)}, i = \overline{1, K}\}$. Технология обучения состоит в данном в следующем. Пусть имеется смешанная обучающая выборка $X^N = \{x^{(1)}, \dots, x^{(N)}\}$, $Y^N = \{y^{(1)}, \dots, y^{(N)}\}$ $y^{(p)} \in \mathbf{R}^m$, $i = \overline{1, N}$, где Y^N требуемые (желаемые) «ответы» ИНС. Далее проводится собственно обучение, состоящее в адаптации весовых коэффициентов сети так, чтобы минимизировать среднюю квадратичную ошибку представления целевых векторов на выходе

$$E(W) = \frac{1}{2} \sum_{p=1}^N (y^{(p)} - \tilde{y}^{(p)})^T (y^{(p)} - \tilde{y}^{(p)}) \xrightarrow{W} \min, \quad (2)$$

где $\tilde{y}^{(*)} = \tilde{y}(x^{(*)}) \neq y^{(*)}$ – реальная реакция сети на входное воздействие $x^{(*)}$ в процессе обучения, определяемая при заданной архитектуре сети и фиксированных значениях весов нейронных связей выражением (1).

Минимизация (2) проводится путем настройки матриц весовых коэффициентов W в (1) с использованием алгоритма обратного распространения (алгоритма ВР, back propagation) ошибки и его многочисленных модификаций реализующих градиентный подход к поиску минимума функционала ошибки относительно всей совокупности весовых коэффициентов. В результате обучения ИНС формирует искомое преобразование информации в универсальной форме.

Неглубокие и глубокие нейронные сети: много слоев и много кое-чего другого

