

# Нейронные сети. Лекции.

Лабунец Л.В. & Паймин Антон

## Содержание

<b>Лекция 02.04</b>	<b>2</b>
1.1 Математическая модель многослойного персептрона . . . . .	2
1.2 Системный подход к синтезу алгоритмов обучения многослойного персептрона . . . . .	4
1.3 Алгоритм обучения нейронной сети (метод градиентного поиска) . . . . .	5
<b>Лекция 09.04 (дистанционная)</b>	<b>6</b>
2.1 Алгоритм обучения нейронной сети для случая набора обучающих примеров . . . . .	6
2.1.1 Ядерная модель Парзена-Розенблатта . . . . .	6
2.1.2 Интегрирование весовой функции среднего риска в конечном виде . . . . .	7
2.1.3 Пакетный, последовательный и комбинированный режимы обучения . . . . .	8
2.2 Алгоритм обучения персептронного нейрона . . . . .	9
2.2.1 Формулирование цели обучения . . . . .	10
2.2.2 Получение оценки вектора градиента . . . . .	11
2.2.3 Алгоритм градиентного поиска . . . . .	12
<b>Лекция 16.04</b>	<b>14</b>
3.1 Алгоритм обучения персептронного нейрона с квадратичной функцией потерь . . . . .	14
3.1.1 Формулирование цели обучения . . . . .	14
3.1.2 Вычисление оценки градиента функции ошибок . . . . .	15
3.1.3 Алгоритм градиентного поиска . . . . .	15
3.2 Аппроксимация многомерной скалярной функции . . . . .	16
3.2.1 Выбор модели функции стоимости . . . . .	16
3.2.2 Поиск градиента функции ошибок . . . . .	17
3.2.3 Применение алгоритма градиентного поиска . . . . .	17
3.3 Алгоритм обратного распространения ошибок обучения многослойного персептрона . . . . .	17
3.3.1 Этап прямого распространения . . . . .	18
3.3.2 Этап обратного распространения . . . . .	19
<b>Лекция 23.04 (пропустил)</b>	<b>21</b>
4.1 Практические рекомендации по улучшению сходимости алгоритма обратного распространения ошибок . . . . .	21
4.1.1 Выбор модели функции активации . . . . .	21
4.1.2 Выбор компромиссного режима обучения персептронных нейронов и сети в целом . . . . .	22
4.1.3 Максимизация информативности . . . . .	23
4.1.4 Инициализация. Выбор начального приближения для биаса и синаптических весов . . . . .	24
<b>Лекция 07.05</b>	<b>25</b>
5.1 Сеть главных компонент . . . . .	25
5.1.1 Метод главных компонент (МГК) . . . . .	25
5.1.2 Сеть главных компонент (введение) . . . . .	26
5.1.3 Гипотеза Хебба (принцип самоусиления (корреляции)) . . . . .	26
5.1.4 Принцип ковариации . . . . .	27
5.2 Метод настройки весов для сети главных компонент . . . . .	27
5.2.1 Правила обучения линейного нейрона . . . . .	27
5.2.2 Матричная формулировка алгоритма Ойа . . . . .	28
<b>Лекция 14.05 (не проводилась)</b>	<b>29</b>
<b>Лекция 21.05</b>	<b>30</b>
7.1 Архитектура сети главных компонент — обобщённый алгоритм Хебба . . . . .	30
7.1.1 Последовательное обучение . . . . .	30
7.1.2 Параллельное обучение . . . . .	31

## Лекция 02.04

## 1.1 Математическая модель многослойного персептрона

Рассмотрим на примере 3-слойной архитектуры (в количество слоёв входят только слои, реализующие нелинейное отображение). Рассмотрим сигналы, поступающие на вход нейронов, но для избежания путаницы, будем визуализировать связи только для первых нейронов соответствующих слоёв. Рассмотрим эти связи и соотв. им параметры.

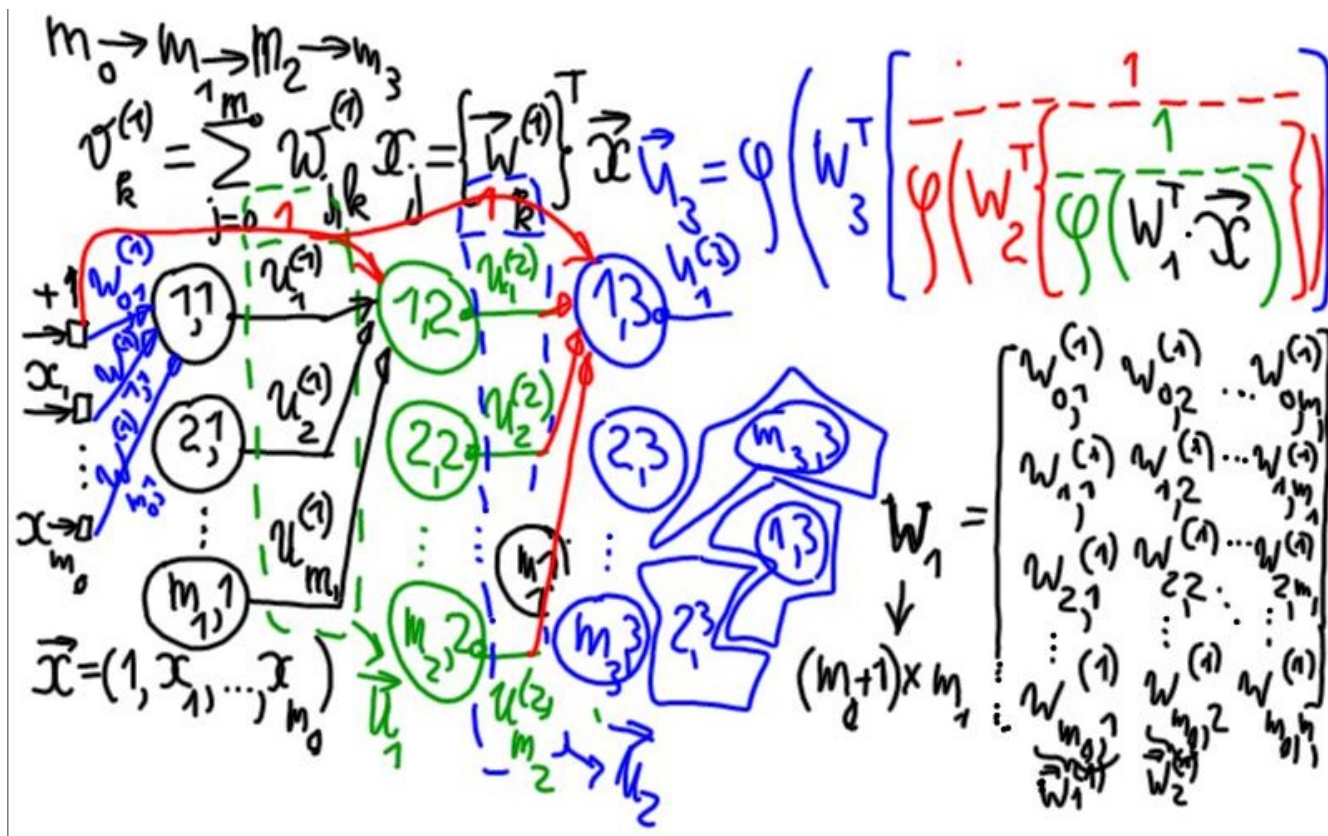


Рис. 1: Структура многослойного персептрона

3-слойный персептрон содержит 0-й слой - слой источников (черные квадратики). Имеется  $m_0$  входов, на которые поступают признаки  $\vec{x} = (1, x_0, \dots, x_{m_0})$ , 3 скрытых слоя.

Примем следующие правила индексации скрытых нейронов и нейронов выходного слоя: первый эл-т индекса - номер нейрона, второй - номер слоя:  $m_{i,j}$  -  $i$ -й нейрон  $j$ -го скрытого слоя.

Выпишем веса нейронов  $i$ -го слоя в виде матрицы  $W_i$ , столбец которой  $\vec{w}_i^{(k)}$  является вектором весов  $i$ -го нейрона  $k$ -го слоя. Количество строк матрицы - количество входов слоя, „+1“ (на вход поступает расширенный вектор), количество столбцов - количество нейронов слоя. Элемент матрицы  $w_{i,j}^{(k)}$  - вес  $i$ -го входа  $j$ -го нейрона  $k$ -го слоя. У первого скрытого слоя имеется  $m_0$  входов (слой источников) и он содержит  $m_1$  нейронов.

$$W_1 = \begin{pmatrix} w_{0,1}^{(1)} & w_{0,2}^{(1)} & \dots & w_{0,m_1}^{(1)} \\ w_{1,1}^{(1)} & w_{1,2}^{(1)} & \dots & w_{1,m_1}^{(1)} \\ w_{2,1}^{(1)} & w_{2,2}^{(1)} & \dots & w_{2,m_1}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m_0,1}^{(1)} & w_{m_0,2}^{(1)} & \dots & w_{m_0,m_1}^{(1)} \end{pmatrix} \quad (1.1.1)$$

Рассмотрим первый слой. На вход каждого нейрона слоя поступает сигнал смещения  $+1$ , взвешиваемый *биасом*  $w_{0,i}^{(1)}$ . Также на вход этого нейрона поступают входные признаки  $x_i$ , взвешиваемые *синаптическими весами*  $w_{1,j}^{(1)}$  до  $w_{m_0,j}^{(1)}$  - веса  $i$ -го признака  $j$ -го нейрона в первом слое. Для экономии места введём вектор синаптических весов конкретного нейрона:

$$\vec{w}_j^{(k)} = \begin{pmatrix} w_{0,j}^{(k)} & w_{1,j}^{(k)} & \dots & w_{m_0,j}^{(k)} \end{pmatrix}^T \quad (1.1.2)$$

Индекс в скобках - номер слоя, нижний индекс - номер нейрона. Т.е.,  $\vec{w}_1^{(1)}$  - вектор синаптических весов 1-го нейрона 1-го слоя.

На выходе  $i$ -го нейрона  $k$ -го слоя формируется сигнал  $u_i^{(k)}$ . Выпишем формулу для вычисления сигнала на выходе 1-го нейрона 1-го слоя  $u_1^{(1)}$ . Опять, в скобках - номер слоя, индекс - номер нейрона. Для этого необходимо посчитать сигнал на выходе сумматора  $k$ -го нейрона (*индуцированного локального поля*):

$$v_k^{(1)} = \sum_{j=0}^{m_0} w_{j,k}^{(1)} \cdot x_j = [\vec{w}_k^{(1)}]^T \cdot \vec{x} \quad (1.1.3)$$

$$v_2^{(1)} = [\vec{w}_2^{(1)}]^T \cdot \vec{x}$$

где  $j$  - номер нейрона,  $k$  - номер слоя,  $w_{j,k}^{(1)}$  - компоненты ...,  $\vec{x}$  - расширенный вектор.

Из приведённых выражений можем получить общее:  $W_k^T \cdot \vec{x}$  - совокупность элементов на выходе сумматоров всех нейронов  $k$ -го слоя.

Полученные на выходе сумматора значения подвергаются нелинейному преобразованию с помощью функции активации  $\varphi(W_k^T \cdot \vec{x})$ . Получаем совокупность сигналов на выходе функций активации нейронов. Для первого слоя:  $\varphi(W_1^T \cdot \vec{x})$ , для второго:  $\varphi(W_2^T \cdot \vec{x})$  и т.д.

Для получения из матрицы  $W_1$  матрицы  $W_2$  рисуем таблицу:

Количество входов и выходов скрытых слоёв (без учёта "расширенности" векторов):

	1 слой	2 слой	3 слой
Количество входов	$m_0$	$m_1$	$m_2$
Количество выходов	$m_1$	$m_2$	$m_3$

Выпишем матрицу  $W_2$ :

$$W_2 = \begin{pmatrix} w_{0,1}^{(2)} & w_{0,2}^{(2)} & \dots & w_{0,m_2}^{(2)} \\ w_{1,1}^{(2)} & w_{1,2}^{(2)} & \dots & w_{1,m_2}^{(2)} \\ w_{2,1}^{(2)} & w_{2,2}^{(2)} & \dots & w_{2,m_2}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m_1,1}^{(2)} & w_{m_1,2}^{(2)} & \dots & w_{m_1,m_2}^{(2)} \end{pmatrix} \quad (1.1.4)$$

И получим выходы сумматоров нейронов второго слоя:

$$\vec{u}_2 = W_2^T \cdot \left[ \frac{1}{\varphi(W_1^T \cdot \vec{x})} \right] \quad (1.1.5)$$

В формуле 1 в якобы числителе - добавленный первый элемент, чтобы получить расширенный вектор-столбец.

$$\vec{u}_3 = \varphi(W_3^T \left[ \frac{1}{\varphi(W_2^T \cdot \left[ \frac{1}{\varphi(W_1^T \cdot \vec{x})} \right])} \right]) \quad (1.1.6)$$

Увеличение слоёв нецелесообразно из-за рекурсивной структуры модели многослойного персептрона (последовательность вложенных блоков), что приводит к уменьшению скорости вычисления, увеличению занимаемой памяти и накопления ошибок. Также по мере увеличения параметров модели, возникает *эффект переобучения*, а по мере увеличения количества слоёв количество параметров увеличивается экспоненциально.

Возможно применение альтернативных архитектур - например, сети радиальных базисных функций. Также необходимо применение методов регуляризации, например, метод *хирургии мозга* — исключения из каждого слоя неинформативных нейронов или целых слоёв. Существуют методы обучения, подразумевающие цикличность.

Также существует теорема, доказанная Колмогоровым, утверждающая, что многомерная скалярная функция может воспроизводить любую закономерность сколь угодно большой сложности. Теорема об универсальной аппроксимации многослойным персептроном утверждает, что архитектура многослойного персептрона *теоретически* позволяет с любой заданной точностью синтезировать многомерную по количеству входов векторную функцию, что позволяет строить модели нелинейных дискриминантных границ в пространстве признаков и, соответственно, определять принадлежность данных кластерам.

## 1.2 Системный подход к синтезу алгоритмов обучения многослойного персептрона

Согласно системному подходу к синтезу алгоритмов обучения многослойного персептрона, необходимо выполнить 3 шага:

1. Выбор цели обучения;
2. Оценка градиента средних потерь целевой функции по параметрам нейронной сети;
3. Использование алгоритмом поиска оптимальных параметров.

Предполагается наличие функционала качества работы системы, оптимальное значение которого и ищется.

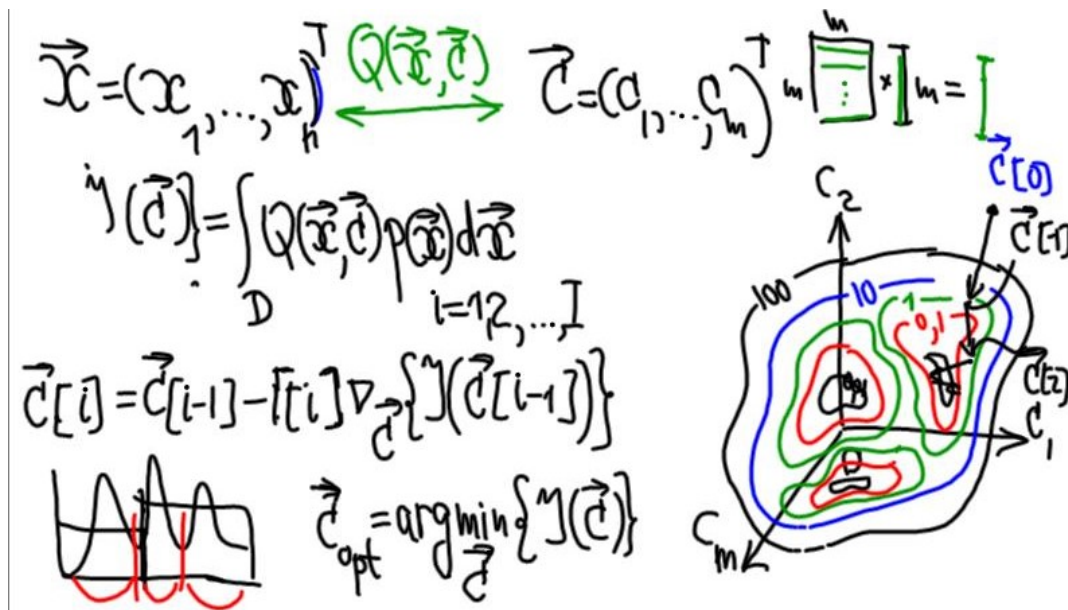


Рис. 2: Обучение нейронной сети

Будем рассматривать 3 объекта анализа

1. Пространство входных воздействий  $\vec{x}$  - сигналы, поступающие на вход системы (обучающие и контрольные примеры), и соответствующее пространство входных воздействий, в котором существует область  $D$  все возможных входных воздействий;
2. Вектор параметров системы  $\vec{c}$  и пространство, ему соответствующее, в котором существуют эквипотенциальные гиперповерхности, все точки которых соотв. определенным значениям функционала качества. Необходимо найти область минимальных/максимальных значений этого функционала качества;
3. Функция стоимости (потерь, приобретений)  $Q(\vec{x}, \vec{c})$ , определяющая потери/приобретения, которые можно измерить, если на вход системы поступает воздействие  $\vec{x}$ , а система находится в состоянии  $\vec{c}$ .

Набор этих 3 объектов позволяет математически строго сформулировать выбор цели обучения в виде функционала средних потерь (среднего риска). Сделать это несложно, необходимо протестировать область  $D$  всех входных воздействий: зафиксировать  $\vec{x}$  и выбрать параметры системы  $\vec{C}$ .

Для каждого зафиксированного выбора параметров системы необходимо вычислить потери по выбранной для каждого входного воздействия и просуммировать их:

$$J(\vec{C}) = \int_D Q(\vec{x}, \vec{C}) p(\vec{x}) d\vec{x} \quad (1.2.1)$$

Полученное  $J(\vec{C})$  - *средний риск (средние потери)*,  $p(\vec{x})$  — весовая функция, сглаживающая аномальные значения (плотность распределения вероятностей входных воздействий).

На практике возможны 2 ситуации:

1. Нетипичный случай полной априорной определённости - полное знание входных воздействий (весовой функции  $p(\vec{x})$ , Типа распределения и тд.);
2. Типичный случай - функция задана набором обучающих примеров

Во втором случае необходимо выполнить предварительную выборочную оценку весовой функции  $p(\vec{x})$  по набору обучающих примеров.

Однако, будем считать весовую функцию  $p(\vec{x})$  известной, тогда возможно посчитать функционал ошибок  $J(\vec{C})$ .

### 1.3 Алгоритм обучения нейронной сети (метод градиентного поиска)

1. Выбрать начальное приближение параметров нейронной сети  $\vec{C}[0]$  (0 - номер итерации), и оценить потери;
2. Выбрать оптимальное направление поиска (вектор антиградиента  $-\nabla \vec{C}$ , указывающий направление, в котором потери убывают с максимальной скоростью);
3. Выбрать шаг поиска (используется процедура стохастической аппроксимации Роббенса-Монро, в соответствии с которой с каждой итерацией длина шага уменьшается обратно пропорционально номеру итерации);

В полученной точке выполняется подтверждение сходимости (сравнение итераций), вычисление нового направления и шага.

Процедура описывается следующим выражением:

$$\vec{C}[i] = \vec{C}[i-1] - \Gamma[i] \nabla \vec{C} \{J(\vec{C}[i-1])\} \quad (1.3.1)$$

Здесь  $\Gamma[i]$  — матрица усиления, управляющая выбором величины шага поиска. Если она является диагональной, то по каждому признаку шаги будут независимы (характерно для случая гиперболоидальной топологии гиперповерхностей равных значений функции стоимости).

Цель обучения:

$$\vec{C}_{opt} = \arg \min_{\vec{C}} \{J(\vec{C})\} \quad (1.3.2)$$

Существуют также методы, не требующие вычисления частной производной, — метод деформируемого многогранника Нелдера-Мида (симплекс-поиск).

*Прикладное нелинейное программирование, Химмельблау*

Обучение сводится к решению задачи оптимизации. Существуют ограничения на значения параметров системы в виде равенств или неравенств.

Результаты выполнения алгоритма зависят от выбора начального приближения, что приводит к необходимости выбора нескольких начальных приближений для обнаружения всех локальных минимумов. Существенное преимущество в этом случае обеспечивают генетические алгоритмы.

## Лекция 09.04 (дистанционная)

### 2.1 Алгоритм обучения нейронной сети для случая набора обучающих примеров

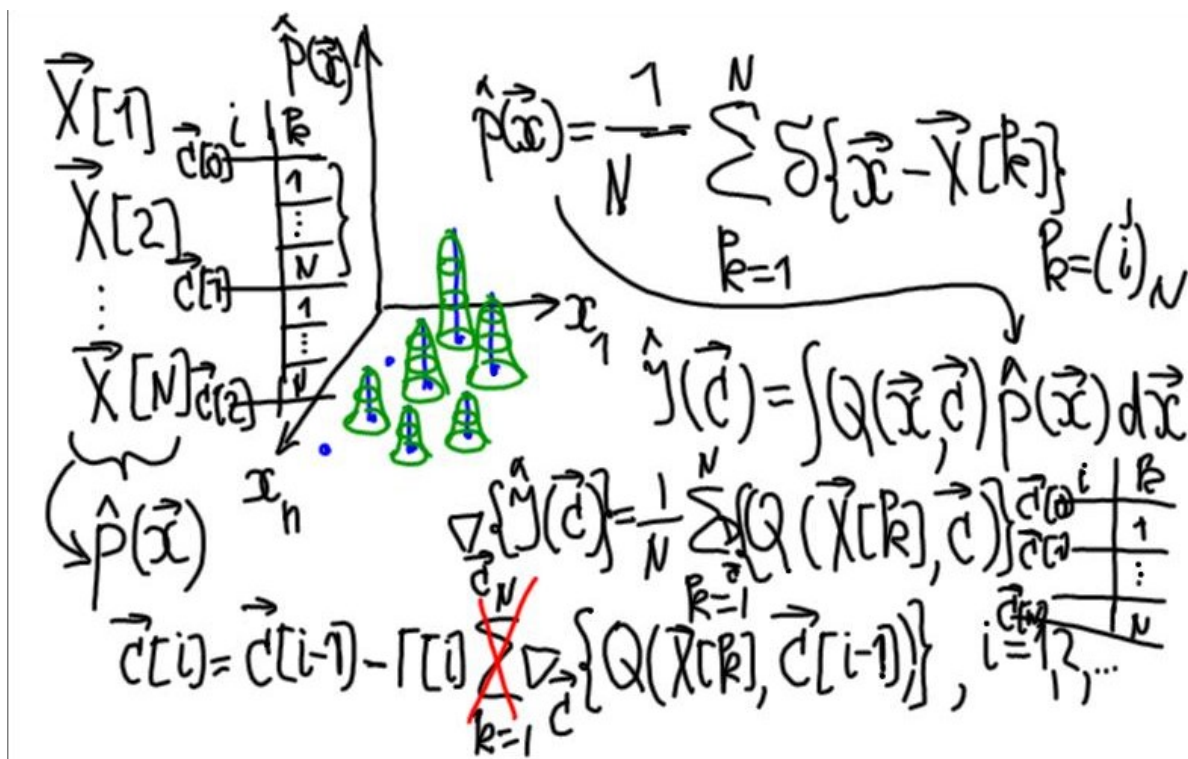


Рис. 3: Системный подход к обучению многослойного персептрона

В подавляющем большинстве практических приложений весовая функция  $p(\vec{x})$ , входящая в выражение функционала среднего риска 1.2.1 задана неявным образом в виде набора обучающих примеров  $\hat{p}(\vec{x})$ .

Здесь и далее обучающие примеры, ранее обозначавшиеся  $\vec{x}$ , будем обозначать  $\vec{X}$ . Объём выборки -  $N$  обучающих примеров (см. левый столбец на 3), размерность пространства признаков  $\{x_1, \dots, x_n\} - n$ .

Таким образом в  $x$ -пространстве имеется набор обучающих примеров  $\vec{X}[1], \dots, \vec{X}[N]$ , визуализированных диаграммой рассеяния (график на 3). Точка задаётся радиус-вектором, набор которых задаёт набор обучающих примеров.

Для приближения функционала среднего риска к реальности, дополним его ещё одной осью  $\hat{p}(\vec{x})$ , вдоль которой будем откладывать значения оценки весовой функции, где крышечка показывает, что это оценка.

В теории искусственного интеллекта существует огромное количество непараметрических, параметрических и полупараметрических описаний распределений, в т.ч. и вероятностных. Например, примером параметрического описания многомерного распределения является модель конечной смеси стандартных распределений, для идентификации параметров которой служат ЕМ-алгоритмы. Пример полупараметрической оценки — ядерная оценка Парзена-Розенблатта, одномерная модель ПРВ в виде гистограммы, сглаженной сдвигом (с обобщением до многомерной модели).

Ссылка на книгу David Scott, *Multivariate Density Estimation*

#### 2.1.1 Ядерная модель Парзена-Розенблатта

Идея заключается в ассоциировании функции ядра каждой точкой в подпространстве признаков с каждым наблюдением (обучающим примером) — „гиперхолмы“ на рисунке 3. Рассмотрим простейший случай, когда оценка обладает единственным параметром сглаживания для всех наблюдений — неадаптивная ядерная оценка Парзена-Розенблатта.

Такие функции обладают важными свойствами. Гиперобъём под поверхностью функции ядра всегда



равен 1 (функция ядра интегрируема с 1). Также, как и у любой оценки многомерного ПРВ, присутствует параметр сглаживания (в одномерном случае — ширина разрядного интервала, в многомерном случае — гиперобъём гиперпараллелепипеда), регулирующий степень убывания высоты холма по мере удаления текущего радиуса-вектора (аргумента ПРВ) от центра холма (выбранного обучающего примера).

Если текущее значение совпадает с выбранным примером, то высота максимальна. По мере удаления от *центра опорной области* (центра холма) высота асимптотически стремится к 0. Аналогично случаю гистограммной оценки, при стремлении параметра сглаживания к 0 (масштаб функции ядра нулевой) высота будет неограниченно возрастать, а опорная область асимптотически стягиваться к 0 и приближаться к дельта-функции Дирака — набору „иглоков“ единичного объёма в гиперпространстве. Такие „иглоки“ не могут иметь общие области и, таким образом, ядерная оценка при таком выборе параметров превращается в „лес иглоков“ и (что плохо) сглаживание экспериментальных производных не будет — *эмпирическая оценка весовой функции*.

$$\hat{p}(\vec{x}) = \frac{1}{N} \sum_{k=1}^N \delta(\vec{x} - \vec{X}[k]) \quad (2.1.1)$$

Для эмпирической оценки суммируем дельта-функции, ассоциированные с обучающими примерами и масштабируем значение в зависимости от числа обучающих примеров.

В ядерной оценке Парзена-Розенблатта увеличивается параметр сглаживания и опорная область расширяется по мере увеличения параметра сглаживания, а высота „холмов“ начинает уменьшаться и ближайшие обучающие примеры начинают перекрываться, происходит постепенное сглаживание совокупности обучающих примеров и получаем „гиперрельеф“.

До некоторого момента при определённом числовом диапазоне параметров сглаживания гиперрельеф сохраняет свою морфологию — наблюдается определённое количество вершин, что свидетельствует о наличии кластерной структуры в данных и характеризует статистическую устойчивость оценки ПРВ при выборе параметра сглаживания из этого диапазона.

При дальнейшем увеличении параметра сглаживания опорные области функций-ядер начинают чрезмерно увеличиваться и кластерная структура гиперрельефа начинает исчезать и в асимптотике стремится к плато — равномерному распределению по всем признакам.

Существует набор методов, позволяющих выбрать оптимальное значение параметров сглаживания, который позволяет обнаружить кластерную структуру, т.е. *многомодовый характер* гиперрельефа в  $n$ -мерном пространстве признаков.

Один из таких методов — *метод скользящей проверки* выбора оптимального параметра сглаживания. Существуют также адаптивные ядерные оценки, в которых для каждого наблюдения выбирается свой оптимальный параметр сглаживания и каждый „гиперхолм“ будет обладать своей оптимальной опорной областью (похоже на выбор характеристик масштаба модели конечной смеси стандартных распределений). В ещё более сложных оценках индивидуальными являются не только параметры сглаживания, но и модель функции ядра, например, *адаптивная фильтрационная ядерная оценка плотности распределения вероятности*. Это отдельное направление, связанное с построением оптимальных оценок одномерных ПРВ.

Есть статья на [labnet.ru](http://labnet.ru) - Рандомизация многомерных распределений в метрике Махаланобиса

### 2.1.2 Интегрирование весовой функции среднего риска в конечном виде

Конечная цель — получить возможность проинтегрировать 1.2.1 в конечном виде, чтобы получить в конечном виде оценку среднего риска, поскольку необходимо сформировать оценки компонентов вектора градиента и реализовать этот алгоритм обучения для практических приложений.

В этом случае поможет эмпирическая оценка — плата за попытку получить выражение 1.2.1 является грубая эмпирическая оценка весовой функции, которая преднамеренно игнорирует процедуру сглаживания обучающих данных.

Преимущество — фильтрующее свойство дельта-функции. При подстановке 2.1.1 в 1.2.1:

$$\hat{J}(\vec{C}) = \int_D Q(\vec{x}, \vec{C}) \hat{p}(\vec{x}) d\vec{x} = \frac{1}{N} \int_D Q(\vec{x}, \vec{C}) \sum_{k=1}^N \delta(\vec{x} - \vec{X}[k]) d\vec{x} \quad (2.1.2)$$

при смене порядка суммирования и интегрирования дельта-функция „профильтрует“ значения весовой функции для каждого обучающего примера:

$$\hat{J}(\vec{C}) = \frac{1}{N} \sum_{k=1}^N \int_D Q(\vec{x}, \vec{C}) \cdot \delta(\vec{x} - \vec{X}[k]) d\vec{x} \quad (2.1.3)$$

Произведение функции стоимости и дельта-функции, ассоциированной с конкретным обучающим примером  $Q(\vec{x}, \vec{C}) \cdot \delta(\vec{x} - \vec{X}[k])$  даст значение функции стоимости для конкретного обучающего примера.

Таким образом, эмпирическая оценка функционала среднего риска представляет собой сумму функции стоимости по всем обучающим примерам для выбранного вектора параметров системы, нормированную по количеству обучающих примеров, т.е. выборочное среднее мат. ожидания функции стоимости:

$$\hat{J}(\vec{C}) = \frac{1}{N} \sum_{k=1}^N Q(\vec{X}[k], \vec{C}) \quad (2.1.4)$$

Применение выборочного среднего означает неробастность функции, т.е. неустойчивость к аномалиям. Финальный алгоритм обучения будет обладать значимой шумовой компонентой и игнорирование влияния аномалий тоже скажется на качестве текущих оценок параметров обучения системы в процессе их настройки.

Поскольку вектор градиента есть набор частных производных, а операции суммирования и дифференцирования перестановочны, то запишем выражение для эмпирической оценки вектора градиента:

$$\nabla_{\vec{C}} \hat{J}(\vec{C}) = \frac{1}{N} \sum_{k=1}^N \nabla_{\vec{C}} Q(\vec{X}[k], \vec{C}) \quad (2.1.5)$$

### 2.1.3 Пакетный, последовательный и комбинированный режимы обучения

Подставим полученную в конечном виде оценку градиента в синтезированный ранее алгоритм, получим адаптированный к случаю обучения системы по набору параметров алгоритм:

$$\vec{C}[i] = \vec{C}[i-1] - \Gamma[i] \sum_{k=1}^N \nabla_{\vec{C}} \{Q(\vec{X}[k], \vec{C}[i-1])\} \quad (2.1.6)$$

Пока операцию суммирования (усреднения по всем возможным значениям компонент вектора градиента) сохраняем.

Проанализируем обучение системы при наличии этого усреднения. В этом случае процедура обучения реализует *пакетный режим* обучения (см. 1).

В распоряжении имеется набор из  $N$  обучающих примеров, происходит инициализация алгоритма обучения начальными значениями параметров системы. Прежде чем откорректировать начальное приближение параметров системы в соответствии с существованием этой суммы, необходимо перебрать все обучающие примеры, вычислить значения функции стоимости и частные производные для всех обучающих примеров, получить в  $C$ -пространстве набор направлений оптимальных поисков (свой для каждого примера), просуммировать эти вектора, взять усреднённое направление поиска с отрицательным знаком, гарантирующее приемлемую скорость убывания, выбрать величину шага и только после этого от  $\vec{C}[0]$  перейти к  $\vec{C}[1]$ . На следующей итерации шаги повторяются.

Пакетный режим обучения требует максимально возможной памяти.

Таблица 1: Описание пакетного режима обучения

Вектор параметров системы	$i$ — Номер итерации обучения	$k$ — Номер обучающего примера
$\vec{C}[0]$	1	1
		$\vdots$
		$N$
$\vec{C}[1]$	2	1
		$\vdots$
		$N$
$\vdots$	$\vdots$	$\vdots$



Альтернативой пакетному режиму обучения является *последовательный режим обучения* (см. 2). Для этого опускаем операцию суммирования всех возможных направлений и после каждого обучающего примера выполняем коррекцию:

$$\vec{C}[i] = \vec{C}[i-1] - \Gamma[i] \nabla \vec{C}\{Q(\vec{X}[k], \vec{C}[i-1])\} \quad (2.1.7)$$

при последовательном режиме обучения эффект шумовой компоненты существенно усиливается, поэтому на практике выбирают „золотую середину“, что согласуется с процедурами сглаживания временных рядов.

Таблица 2: Описание последовательного режима обучения

Вектор параметров системы	$i$ — Номер итерации обучения	$k$ — Номер обучающего примера
$\vec{C}[0]$	1	1
$\vec{C}[1]$	2	2
$\vdots$	$\vdots$	$\vdots$
$\vec{C}[N]$	N	N

Выберем некоторую компоненту в параметрах системы — временной ряд в терминологии номера итерации обучения. Этот ВР будет иметь интересующую трендовую составляющую, шумовую и, возможно, циклические компоненты.

Пользуясь подходами выделения трендовых компонент ВР, используем цифровую фильтрацию: для текущего обучающего примера и нескольких предыдущих примеров (текущую и несколько предыдущих итераций), количество которых определяет интервал сглаживания. Обладая этой информацией, можно вычислить скользящую среднюю (или экспоненциальную скользящую среднюю) и реализовать подавление шумовой компоненты, выделив трендовые составляющие.

Таким образом, *комбинированный режим обучения* предполагает учёт текущего обучающего примера и нескольких предыдущих, что обеспечивает подавление шумовой составляющей и циклической компоненты.

## 2.2 Алгоритм обучения персептронного нейрона

Прежде чем переходить к алгоритму обучения многослойного персептрона, рассмотрим синтез алгоритмов обучения персептронного нейрона — основного структурного элемента многослойного персептрона.

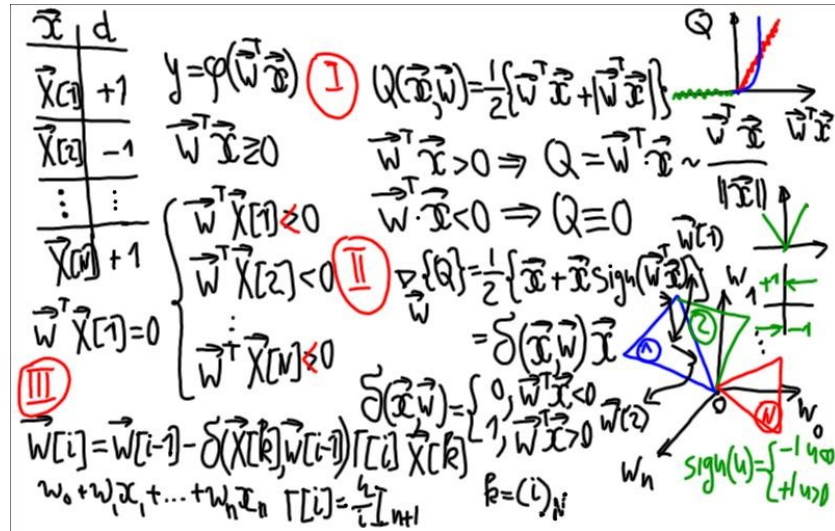


Рис. 4: Алгоритм обучения персептронного нейрона

Рассмотрим подробнее задачу *дихотомии* (разделение обучающих примеров на 2 класса). Имеется метка +1 и метка -1 принадлежности обучающих примеров к 2 альтернативным классам. Рассмотрим синтез алгоритма обучения персептронного нейрона в рамках задачи дихотомии. Реализуется принцип обучения с

поощрением, т.е. для настройки параметров персептронного нейрона (биас и синаптические веса) применяются обучающие пары (представлены в 3): обучающие примеры  $\vec{x}$  на входе персептронного нейрона и желаемое указание (или метка класса)  $d$ , поскольку персептронный нейрон решает задачу *дихотомии*.

$\vec{x}$ — Обучающие примеры на входе	$d$ — желаемое указание или метка класса
$\vec{X}[1]$	+1
$\vec{X}[2]$	-1
$\vdots$	$\vdots$
$\vec{X}[N]$	+1

Таблица 3: Обучающие пары

Необходимо в пространстве признаков сформировать модель линейной дискриминантной границы для того, чтобы реализовать классификацию. Как было сказано выше, существуют 2 случая: линейно разделимых и линейно неразделимых обучающих примеров. Случай линейно разделимых обучающих примеров достаточно редко встречается на практике, поэтому для снижения количества неизбежных ошибочных классификаций для линейно неразделимых векторов образов была придумана *система опорных векторов*.

Рассмотрим сначала простейший случай линейно разделимых векторов образов. Для этого понадобится математическая модель персептронного нейрона:

$$y = \varphi(\vec{w}^T \cdot \vec{x}), \quad (2.2.1)$$

причём в случае задач дихотомии функция активации  $\varphi$  анализирует знак сигнала на выходе сумматора нейрона (линейной дискриминантной функции) и принимает вид знаковой функции („ступенька“):  $\varphi(x) = -1$  для  $x < 0$  и  $\varphi(x) = +1$  для  $x \geq 0$ . Геометрически представляет собой гиперплоскость.

Таким образом, задачу дихотомии удобно сформулировать в терминах решения системы линейных неравенств. С этой целью для каждой обучающей пары выпишем соответствующее неравенство:

$$\begin{cases} \vec{w}^T \cdot \vec{X}[1] > 0 \\ \vec{w}^T \cdot \vec{X}[2] < 0 \\ \vdots \\ \vec{w}^T \cdot \vec{X}[N] > 0 \end{cases} \quad (2.2.2)$$

Стандартизуем систему линейных неравенств — приведём все неравенства к знаку  $<$ . Для этого откорректируем все обучающие примеры для тех пар, которые обладают меткой +1 путём домножения их на  $-1$ . Получим:

$$\begin{cases} \vec{w}^T \cdot \vec{X}[1] < 0 \\ \vec{w}^T \cdot \vec{X}[2] < 0 \\ \vdots \\ \vec{w}^T \cdot \vec{X}[N] < 0 \end{cases} \quad (2.2.3)$$

Воспользуемся 3 этапами системного подхода Я.З. Цыпкина к синтезу алгоритма обучения персептронного нейрона:

1. Сформулировать цель обучения;
2. Найти оценки вектора градиента функции стоимости;
3. Выписать алгоритм градиентного поиска в  $C$ -пространстве параметров системы (состоит из биаса  $w_0$  и синаптических весов  $w_i$ ).

### 2.2.1 Формулирование цели обучения

Выберем подходящую модель функции стоимости.

*Ссылка на книгу Ту, Гонсалес — Принципы распознавания образов.*

В качестве цели обучения выберем среднее арифметическое линейной дискриминантной функции и её модуля:

$$Q(\vec{x}, \vec{w}) = \frac{1}{2} \{ \vec{w}^T \cdot \vec{x} + |\vec{w}^T \cdot \vec{x}| \} \quad (2.2.4)$$

Для понимания смысла выбора модели проанализируем 2 случая:

1. Неудачный выбор параметров персептронного нейрона;
2. Удачный выбор параметров модели.

При неудачном выборе параметров персептронного нейрона знак неравенства из 2.2.3 становится  $> 0$ :

$$\vec{w}^T \cdot \vec{x} > 0 \Rightarrow Q = \vec{w}^T \cdot \vec{x} \sim \frac{\vec{w}^T \cdot \vec{x}}{\|\vec{x}\|} \quad (2.2.5)$$

Т.е., для положительных значений линейной дискриминантной функции потери будут возрастать линейно — неудачный выбор будет „наказан“.

В случае удачного выбора параметров персептронного нейрона:

$$\vec{w}^T \vec{x} > 0 \Rightarrow Q \equiv 0 \quad (2.2.6)$$

Это позволит реализовать 2 режима: *принцип подкрепления* (при удачном выборе параметров персептронного нейрона) и *принцип наказания* (при неудачном выборе параметров персептронного нейрона).

В данном случае важно рассматривать выражения не в  $x$ -пространстве, а в  $w$ -пространстве: при попадании в область неудачных выборов в  $w$ -пространстве, необходимо перейти из неё в область удачных выборов. Чем больше расстояние от текущей точки в  $w$ -пространстве от отрицательной области, тем сильнее „наказание“.

До этого „числами“ являлись параметры персептронного нейрона  $w_i$ , а „буквами“ — обучающие признаки  $x_i$ . В  $w$ -пространстве всё наоборот, поскольку мы обладаем набором обучающих примеров и  $w_i$  — „буквы“ (аргументы, решение системы линейных неравенств), а  $x_i$  — „числа“.

## 2.2.2 Получение оценки вектора градиента

Выпишем в явном виде скалярное произведение  $\vec{w}^T \vec{x}$ :

$$\vec{w}^T \vec{x} = w_0 + w_1 \cdot x_1 + \dots + w_n \cdot x_n \quad (2.2.7)$$

Рассмотрим частные производные этого выражения по компонентам вектора  $\vec{w}$ :

$$\begin{aligned} \frac{\partial (\vec{w}^T \vec{x})}{\partial w_0} &= 1 \\ \frac{\partial (\vec{w}^T \vec{x})}{\partial w_1} &= x_1 \\ &\vdots \\ \frac{\partial (\vec{w}^T \vec{x})}{\partial w_n} &= x_n \end{aligned}$$

Получим, что градиент произведения  $\vec{w}^T \vec{x}$  по параметрам персептронного нейрона — это обучающий пример  $\vec{x}$ :

$$\nabla_{\vec{w}} \{ \vec{w}^T \vec{x} \} = \vec{x} \quad (2.2.8)$$

При дифференцировании модуля  $|x|$  получим:

$$\frac{d|x|}{dx} = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \end{cases} \Rightarrow \frac{d|x|}{dx} = \text{sign}(x) \quad (2.2.9)$$

То есть:

$$\frac{\partial |\vec{w}^T \vec{x}|}{\partial w_i} = x_i \cdot \text{sign}(w_i \cdot x_i) \Rightarrow \nabla_{\vec{w}} \{|\vec{w}^T \vec{x}|\} = \vec{x} \cdot \text{sign} \cdot (\vec{w}^T \vec{x}) \quad (2.2.10)$$

Получим градиент функции стоимости 2.2.4:

$$\nabla_{\vec{w}} \{Q\} = \nabla_{\vec{w}} \left[ \frac{1}{2} \{ \vec{w}^T \cdot \vec{x} + |\vec{w}^T \cdot \vec{x}| \} \right] = \frac{1}{2} \nabla_{\vec{w}} [\vec{w}^T \vec{x}] + \nabla_{\vec{w}} [|\vec{w}^T \cdot \vec{x}|] = \frac{1}{2} (\vec{x} + \vec{x} \cdot \text{sign}(\vec{w}^T \cdot \vec{x})) \quad (2.2.11)$$

Введём обозначение:

$$\delta(\vec{x}, \vec{w}) = \frac{1}{2} (1 + \text{sign}(\vec{w}^T \cdot \vec{x})) \quad (2.2.12)$$

Заметим также, что:

$$\delta(\vec{x}, \vec{w}) = \begin{cases} 0, & \vec{w}^T \cdot \vec{x} < 0 \quad \text{— удачный выбор параметров персептронного нейрона} \\ 1, & \vec{w}^T \cdot \vec{x} > 0 \quad \text{— неудачный выбор параметров} \end{cases} \quad (2.2.13)$$

Содержательный смысл функции  $\delta(\vec{x}, \vec{w})$  — счётчик ошибок, накапливающий количество ошибок при неверном выборе параметров персептронного нейрона.

Окончательно получим:

$$\nabla_{\vec{w}} \{Q\} = \delta(\vec{x}, \vec{w}) \cdot \vec{x} \quad (2.2.14)$$

### 2.2.3 Алгоритм градиентного поиска

Выпишем для режима последовательного обучения:

$$\vec{w}[i] = \vec{w}[i-1] - \delta(\vec{X}[k], \vec{w}[i-1]) \cdot \Gamma[i] \cdot \vec{X}[k] \quad (2.2.15)$$

Если  $\delta(\vec{X}[k], \vec{w}[i-1]) = 0$ , то параметры на предыдущей итерации были выбраны удачно и ничего менять не надо — *фаза подкрепления*, т.е.  $\vec{w}[i] = \vec{w}[i-1]$ , коррекция отсутствует.

Если для вновь поступившего обучающего примера оценка является неудачной, т.е.  $\delta(\vec{X}[k], \vec{w}[i-1]) = 0$ , то необходимо изменить параметры нейрона — *фаза наказания*, осуществляется движение в направлении антиградиента так, чтобы попасть в отрицательную зону.

В  $w$ -пространстве граничная гиперплоскость, соответствующая первому неравенству из 2.2.3, делящая пространство на случай удачного и неудачного выбора, проходит через начало координат, поскольку при  $\forall w_i = 0 \quad \vec{w}^T \cdot \vec{X}[0] = 0$ . Поворот гиперплоскости (управляющие косинусы ей нормали) определяются компонентами обучающего примера:

$$\frac{\vec{X}[i]}{\|\vec{X}[i]\|}$$

Вся совокупность граничных гиперплоскостей для всех обучающих примеров проходят через начало координат, но ориентированы различным образом (см. правый нижний угол 4, синий — гиперплоскость 1-го обучающего примера, зелёный — гиперплоскость 2-го обучающего примера, красный — гиперплоскость  $N$ -го обучающего примера).

В случае линейно разделимых векторов образов эти плоскости („нерегулярно“) будут расположены более-менее регулярно, а в идеальном случае — образуют многогранную поверхность в виде разрезанного гиперконуса (гиперпирамиды), которая делит  $w$ -пространство на 2 зоны (положительных и отрицательных значений линейной дискриминантной функции). В этом случае существует минимум одно решение системы линейных неравенств 2.2.3.

Отталкиваясь от начального предложения, необходимо сформировать посегментную траекторию градиентного поиска так, чтобы, даже находясь в положительной зоне рано или поздно, попасть в отрицательную зону  $w$ -пространства.

В случае линейно неразделимых обучающих примеров „фацеты“ будут расположены нерегулярно (хаотически), так, что нельзя будет найти точку в  $w$ -пространстве, удовлетворяющей всем неравенствам 2.2.3.

Как было сказано выше, эту задачу решает *машина опорных векторов*, обеспечивающая выбор параметров персептронного нейрона с минимальным количеством нарушений.

В случае процедуры стохастической аппроксимации Роббенса-Монро шаг поиска выбирают, отталкиваясь от начального приближения в направлении антиградиента пропорционально величине  $\eta$  — *параметру скорости обучения*.

$$\Gamma[i] = \frac{\eta}{i} I_{n+1} \quad (2.2.16)$$

где  $I_{n+1}$  — единичная матрица размерности  $n + 1$ , поскольку помимо  $n$  синаптических весов присутствует биас  $w_0$ .

Уменьшение шага поиска пропорциональное номеру шага обучения гарантирует сходимость метода как минимум в случае линейно разделимых векторов образов. Номер текущего обучающего примера  $k$  выбирается, как значение целочисленной переменной  $i$  по модулю  $N$  — остаток от деления нацело номера итерации на количество обучающих примеров:

$$k = (i)_N$$

## Лекция 16.04

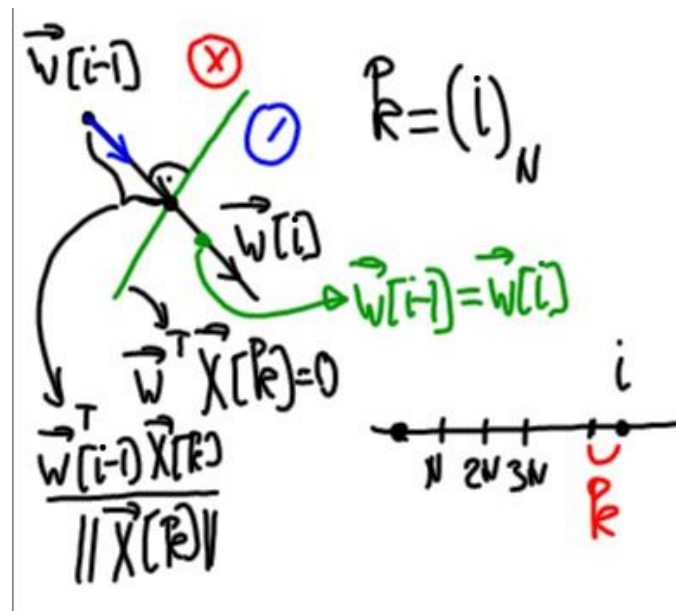


Рис. 5: Аппроксимация многомерной скалярной функции (hyperflat\_7\_1)

**Геометрическая интерпретация** Параметр скорости обучения  $\eta$  выбирают, как .. Переходя от одного к другому обучающему примеру, приближаемся к поверхности разрезанной пирамиды.

### 3.1 Алгоритм обучения персептронного нейрона с квадратичной функцией потерь

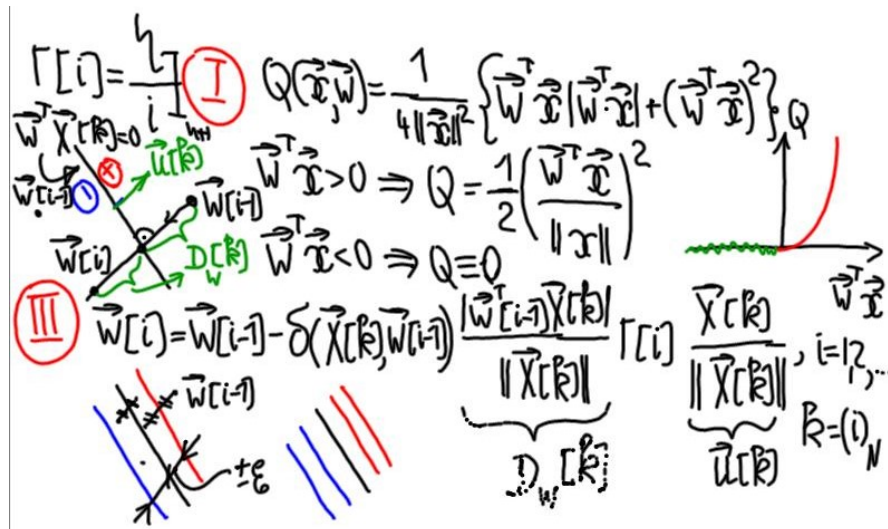


Рис. 6: Алгоритм обучения персептронного нейрона с квадратичной функцией потерь

#### 3.1.1 Формулирование цели обучения

Рассмотрим новую модель функции потерь:

$$Q(\vec{x}, \vec{w}) = \frac{1}{4\|\vec{x}\|^2} \left[ \vec{w}^T \vec{x} |\vec{w}^T \vec{x}| + (\vec{w}^T \vec{x})^2 \right] \quad (3.1.1)$$

Для неё получим:

$$Q(\vec{x}, \vec{w}) = \begin{cases} \frac{1}{2} \left( \frac{\vec{w}^T \vec{x}}{\|\vec{x}\|} \right)^2, & \vec{w}^T \vec{x} > 0 \\ 0, & \vec{w}^T \vec{x} < 0 \end{cases} \quad (3.1.2)$$

То есть, при неудачном выборе потери возрастают пропорционально квадрату расстояния текущего значения линейной дискриминантной функции от 0 — функция потерь квадратичная.

### 3.1.2 Вычисление оценки градиента функции ошибок

Необходимо вычислить вектор градиента квадратичной функции ошибок:

$$\nabla_{\vec{w}} [Q(\vec{x}, \vec{w})] = \nabla_{\vec{w}} \left[ \frac{1}{4\|\vec{x}\|^2} \left[ \vec{w}^T \vec{x} |\vec{w}^T \vec{x}| + (\vec{w}^T \vec{x})^2 \right] \right] = \frac{1}{4\|\vec{x}\|^2} \cdot \left( \nabla_{\vec{w}} [\vec{w}^T \vec{x} |\vec{w}^T \vec{x}|] + \nabla_{\vec{w}} [(\vec{w}^T \vec{x})^2] \right) \quad (3.1.3)$$

### 3.1.3 Алгоритм градиентного поиска

Имея текущий обучающий пример  $x_k$  и текущие параметры нейрона  $\vec{W}[i-1]$ , ..

*ВЫписать выражение для градиентного поиска после слов*

В  $W$ -пространстве нормаль  $\vec{u}[k]$  вычисляется как:

$$\frac{\vec{X}_k}{\|\vec{X}_k\|} \quad (3.1.4)$$

Положим, что для  $k$ -го обучающего примера параметры неудачны. В этом случае опускаем перпендикуляр к гиперплоскости и вычисляем расстояние до неё:

$$D_W[k] = \frac{|\vec{W}^T[i-1] \vec{X}[k]|}{\|\vec{X}[k]\|} \quad (3.1.5)$$

Долю масштаба определяют компоненты матрицы  $\Gamma[i] = \frac{\eta}{i} I_{n \times m}$ , где  $I_{n \times m}$  - единичная матрица.

*Харман, Математические основы математической вычислительной томографии, глава "РВТ"*

*Также в книге присутствует алгоритм, в котором для гиперплоскости каждого обучающего примера выбирается некая погрешность и гиперплоскость обволакивается полосой шириной  $\pm \varepsilon$  в связи с присутствием погрешностей в "реальных" примерах.*

Для простоты анализа здесь был выбран режим последовательного обучения, однако ранее были обсуждены также алгоритмы смешанного обучения.

*Рекомендовано выписать этот алгоритм для простой и экспоненциальной (задача со звёздочкой) скользящей средней*

Здесь рассмотрен случай линейно разделимых разделимых векторов образов, исходя из предположения, что система линейных неравенств имеет минимум одно решение. Иначе будет минимум одна зона в  $W$ -пространстве, характеризуемая количеством ошибочных классификаций и необходим алгоритм, строящий оптимальную гиперплоскость.

Для оптимальной гиперплоскости количество неизбежных ошибочных классификаций минимально. Такой метод был реализован Владимиром Вапником в методе опорных векторов.

Булавский (НГУ) предложил симплекс-метод для решения несовместных систем линейных неравенств — поиска компромиссного решения. Метод применим для случая линейно неразделимых векторов образов



и является альтернативой методу опорных векторов. Метод не получил широкого распространения, имеется только препринт.

Препринт есть у ЛЛВ

### 3.2 Аппроксимация многомерной скалярной функции

Handwritten notes showing a table of input-output pairs and associated equations for a neural network model. The table has columns  $\vec{x}$ ,  $d$ ,  $y$ , and  $e$ . The rows are labeled  $\vec{x}[1]$ ,  $\vec{x}[2]$ , ...,  $\vec{x}[N]$ . To the right of the table, the following equations are written:

- $Q(\vec{x}, \vec{w}) = e^2(\vec{x}, \vec{w})$
- $E(\vec{w}) = \sum_{k=1}^N e^2(\vec{x}[k], \vec{w})$
- $\nabla_{\vec{w}} E(\vec{w}) = \sum_{k=1}^N e(\vec{x}[k], \vec{w}) (-1) \varphi'(\vec{w}^T \vec{x}[k])$
- $y = \varphi(\vec{w}^T \vec{x})$
- $\delta(\vec{x}, \vec{w}) = \varphi'(\vec{w}^T \vec{x}) e(\vec{x}, \vec{w})$
- $e = d - y$
- $\vec{w}[k] = \vec{w}[k-1] + \delta(\vec{x}[k], \vec{w}[k-1]) \vec{x}[k]$

Рис. 7: Аппроксимация многомерной скалярной функции

Рассмотрим альтернативную дихотомии задачу - аппроксимацию многомерной скалярной функции, т.е. синтез закономерности, скрытой в данных, с помощью многослойного персептрона. В данном случае стоит воспринимать персептрон как устройство воспроизведения математической модели.

$\vec{x}$	$d$	$y$	$e$
$\vec{X}[1]$	$d[1]$	$y[1]$	$e[1]$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\vec{X}[N]$	$d[N]$	$y[N]$	$e[N]$

Таблица 4: Табличная форма задания многомерной скалярной функции

Имеется табличная форма задания многомерной скалярной функции (3.2), набор обучающих примеров в виде совокупности факторов (входных аргументов)  $\vec{x}[i]$ ,  $d$  — желаемое выходной значение скрытой закономерности,  $y = \varphi(\vec{W}^T \vec{x})$  — фактическое значение (реакция) персептрона и ошибка  $e = d - y$ .

Необходимо синтезировать математическую модель с помощью персептронного нейрона.  $\vec{x}$  — обобщенный вектор образов с +1 в первой позиции,  $\vec{W}$  — обобщенный вектор синаптических весов с биасом в первой позиции.

Будем выполнять задачу согласно этапам системного подхода.

#### 3.2.1 Выбор модели функции стоимости

Выберем квадратичную функцию ошибок:

$$Q(\vec{x}, \vec{W}) = e^2(\vec{x}, \vec{W}) \quad (3.2.1)$$

Тогда средние потери будут равны:

$$E(\vec{W}) = \sum_{k=1}^T e^2(\vec{X}[k], \vec{W}) \quad (3.2.2)$$

### 3.2.2 Поиск градиента функции ошибок

Необходимо найти в  $W$ -пространстве такую точку, которая минимизировала бы суммарную энергию ошибок  $E(\vec{W})$  3.2.2.

$$\nabla_{\vec{w}}\{E(\vec{W})\} = 2 \sum_{k=1}^N e(\vec{X}[k], \vec{W}) \cdot (-1) \cdot \varphi'(\vec{W}^T \cdot \vec{X}[k]) \vec{X}[k] \quad (3.2.3)$$

В первую очередь интересно произведение первой производной функции активации и функции ошибки, которую обозначим  $\delta(\vec{x}, \vec{W})$ :

$$\delta(\vec{x}, \vec{W}) = \varphi'(\vec{W}^T \cdot \vec{x}) \cdot e(\vec{x}, \vec{W}) \quad (3.2.4)$$

и будем называть *локальным градиентом*. Он введён, поскольку при выполнении поиска выгоднее находиться в области ненулевых значений функции активации (иначе итерация "пропадает").

### 3.2.3 Применение алгоритма градиентного поиска

Распишем алгоритм градиентного поиска (пакетный режим):

$$\vec{W}[i] = \vec{W}[i-1] + \sum_{k=1}^N \delta(\vec{X}[k], \vec{W}[i-1]) \cdot \Gamma[i] \cdot \vec{X}[k] \quad (3.2.5)$$

Рекомендовано выписывать алгоритм слева направо - скаляр, матрица, вектор.

Для усиления шумового эффекта игнорируется усреднение по всем примерам (удаление знака суммы) и выполняется пересчёт после каждого нового значения - последовательное обучение:

$$\vec{W}[i] = \vec{W}[i-1] + \delta(\vec{X}[k], \vec{W}[i-1]) \cdot \Gamma[i] \cdot \vec{X}[k] \quad (3.2.6)$$

Здесь  $\delta(\vec{X}[k], \vec{W}[i-1])$  – локальный градиент 3.2.4. Режим последовательного обучения реализует *оптимизирующий алгоритм Уидроу-Хоффа* 3.2.6, иначе называемый *дельта-правилом*.

При введении сглаживания с помощью модели ЕМА получим *обобщённое дельта-правило*.

## 3.3 Алгоритм обратного распространения ошибок обучения многослойного персептрона

Lec\_neuro\_08.docx

Основным параметром алгоритма обучения Уидроу-Хоффа является ошибка реакции, которая недоступна при обучении нейронов скрытых слоёв. Эту проблему решает алгоритм обратного распространения ошибок.

Упростим схему индексации нейронов (см. 8):

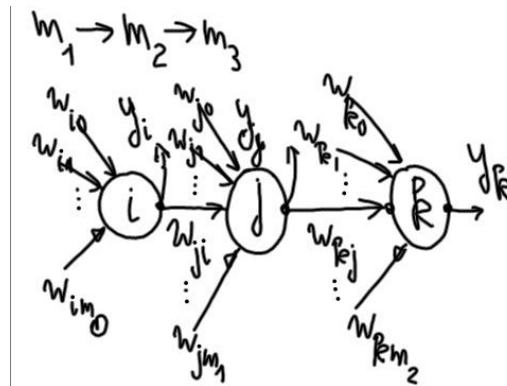


Рис. 8: Новая схема индексации нейронов

Для этого примем следующие соглашения:

1. Для обозначения 1-го скрытого слоя используем символ  $i$ , нейроны 2-го скрытого слоя обозначим символами  $j$ , выходные нейроны обозначим символами  $k$ ,  $l$  - символ линейных нейронов 0-го слоя.
2. Связь между нейронами младшего и старшего слоя (синаптический вес) —  $W_{ji}$ ,  $j$  — символ нейрона старшего слоя,  $i$  — символ нейрона младшего слоя, для биаса второй индекс 0.
3. Во входном слое  $m_0$  нейронов (функциональных входов), во 2-м слое  $m_1$  нейронов, в 3-м слое  $m_2$  нейронов.

Алгоритм обратного распространения содержит 2 вычислительных этапа

1. **Этап прямого распространения.** В соответствии с математической моделью многослойного персептрона можем распространить все сигналы в прямом направлении — посчитать фактические реакции нейронов скрытых слоёв и выходные реакции нейронов.

В финале этого этапа обладаем выходными реакциями нейронов, на основе которых можно вычислить сигналы ошибок и локальные градиенты, а значит, применить дельта-правило и скорректировать параметры нейронов выходного слоя.

2. **Этап обратного распространения.** На этом этапе распространяются сигналы ошибок для нейронов скрытых слоёв.

### 3.3.1 Этап прямого распространения

Этап начинается с расчёта фактических реакций на выходе нейронов 1-го скрытого слоя.

Рассчитаем сигнал на выходе сумматора текущего  $i$ -го нейрона на выходе 1-го скрытого слоя, на вход которого поступает сигнал смещения „+1“, взвешиваемый биасом, и входные признаки  $x_m$ , взвешиваемые синаптическим весом  $w_{im_0}$ .

На выходе сумматоров получим:

$$v_i(n) = \sum_{l=0}^{m_0} w_{il}(n)x_l(n) \quad (3.3.1)$$

$n$  - номер итерации.

Подвергаем его нелинейному преобразованию и получаем реакции на выходе нейронов 1-го скрытого слоя:

$$y_i(n) = \varphi(v_i(n)), \quad i = \overline{1, m_1} \quad (3.3.2)$$

Посчитаем сигналы на выходе сумматоров нейронов 2-го скрытого слоя:

$$v_j(n) = \sum_{i=0}^{m_1} w_{ji}(n)y_i(n) \quad (3.3.3)$$

На выходе нейронов 2-го скрытого слоя после нелинейного преобразования:

$$y_j(n) = \varphi(v_j(n)), \quad j = \overline{1, m_2} \quad (3.3.4)$$

Аналогично на выходе сумматоров нейронов 3-го скрытого слоя:

$$v_k(n) = \sum_{j=0}^{m_2} w_{kj}(n)y_j(n) \quad (3.3.5)$$

На выходе нейронов 3-го скрытого слоя после нелинейного преобразования:

$$y_k(n) = \varphi(v_k(n)), \quad k = \overline{1, m_3} \quad (3.3.6)$$

Зная фактические реакции выходных нейронов, вычислим набор ошибок:

$$e_k(n) = d_k(n) - y_k(n) \quad (3.3.7)$$

Вычисляем величину локального градиента выходных нейронов:

$$\delta_k(n) = \varphi'(v_k(n)) \cdot e_k(n) \quad (3.3.8)$$

Алгоритм коррекции  $k$ -го выходного нейрона:

$$w_{kj}(n) = w_{kj}(n-1) + \frac{\eta}{N} \cdot \delta_k(n) \cdot y_j(n), k = \overline{1, m_3} \quad (3.3.9)$$

Этап прямого распространения закончен.

### 3.3.2 Этап обратного распространения

В соответствии с методом потоковых графов, разработанным Станиславом Осовским и описанным в книге „Нейронные сети для обработки информации“, любую сеть можно представить в виде набора графов, имеющих узлы (сигналы) и рёбра (направления распространения сигналов и коэффициенты усиления). Если рёбра сходятся в одном узле, то это означает суммирование соответствующих сигналов.

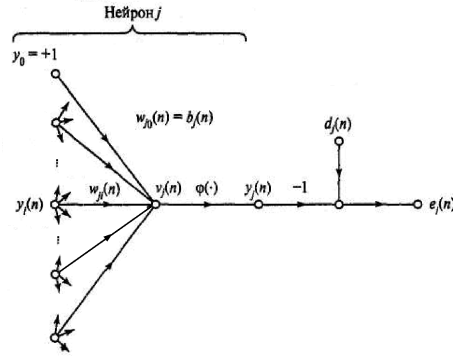


Рис. 9: Граф передачи сигнала в пределах  $j$ -го нейрона

Для нейронов выходного слоя имеем сигналы ошибки  $e_1(n), \dots, e_k(n), \dots, e_{m_3}$ , которые умножаются на вычисленные значения функции активации  $\varphi'(v_k(n))$  и получаем значения локальных градиентов выходных нейронов  $\delta_k(n)$ , а затем используем алгоритм обратного распространения ошибок.

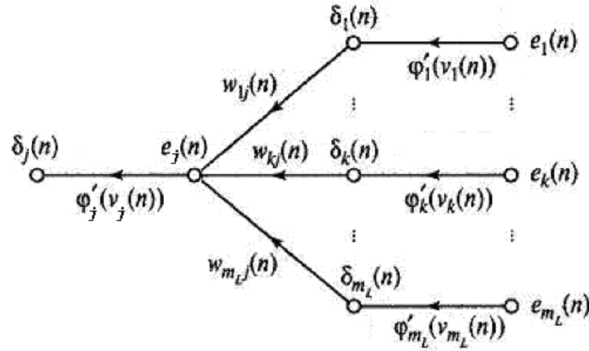


Рис. 10: Обратное распространение сигнала ошибки

Основной результат заключается в правиле вычисления ошибок для нейронов младшего слоя. Показано, что сигнал ошибки может быть вычислен в виде взвешенной суммы локальных градиентов нейронов  $n$ -го слоя:

$$e_j(n) = \sum_{k=1}^{m_3} w_{kj}(n) \delta_k(n) \quad (3.3.10)$$

*Возможно, суммирование должно быть от 0.*

Вычислив ошибку для нейронов второго слоя, вычислим для него локальный градиент, т.е. распространим ошибку на шаг назад:

$$\delta_j(n) = \varphi'(v_j(n)) \sum_{k=1}^{m_3} \delta_k(n) w_{kj}(n) = e_j(n) \cdot \varphi'(v_j(n)) \quad (3.3.11)$$

Значит, возможно произвести „редактирование“ синаптических весов нейронов второго слоя. Т.е.:

$$w_{ji}(n) = w_{ji}(n-1) + \frac{\eta}{N} \cdot \delta_j \cdot y_i(n) \quad (3.3.12)$$

Для нейронов первого скрытого слоя:

$$e_i(n) = \sum_{j=1}^{m_2} w_{ji}(n) \delta_j \quad (3.3.13)$$

*Возможно, суммирование должно быть от 0.*

$$\delta_i(n) = e_i(n) \cdot \varphi'(v_i(n)) \quad (3.3.14)$$

Для входного слоя получим:

$$w_{il}(n) = w_{il}(n-1) + \frac{\eta}{N} \cdot \delta_i \cdot x_l \quad (3.3.15)$$

Произошло распространение ошибок и локального градиента в обратном направлении, что позволило применить обычное дельта-правило.

## Лекция 23.04 (пропустил)

### 4.1 Практические рекомендации по улучшению сходимости алгоритма обратного распространения ошибок

#### 4.1.1 Выбор модели функции активации

Важную роль в вычислительных затратах играет вычисление производной функции активации  $\varphi'(x)$  каждого значения индуцированного локального поля  $v_j(n)$  (сигнала на выходе сумматора каждого нейрона).

Имеются рекомендации по выбору модели функции активации, обеспечивающей снижение вычислительных затрат на вычисление её производной. Наиболее популярными и широко применяющимися на практике моделями являются *логистическая функция* и *функция гиперболического тангенса*.

**Логистическая функция**  $S$ -образная кривая (*sigmoid*), нижняя горизонтальная асимптота которой совпадает с осью абсцисс, верхняя горизонтальная асимптота равна  $-1$ . Пересечение с осью ординат происходит в точке  $(0; 0,5)$ .

$$\varphi(v_j(n)) = \frac{1}{1 + \exp(-av_j(n))}, \quad a > 0, \quad -\infty < v_j(n) < +\infty \quad (4.1.1)$$

$a$ —параметр скорости возрастания  $S$ -образной кривой. При  $a \rightarrow \infty$   $S$ -образная кривая приближается к функции Хевисайда.

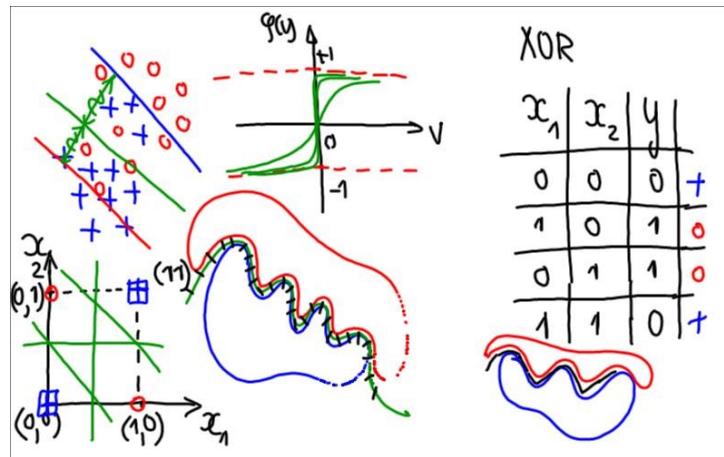


Рис. 11: Выбор функции активации и т.д.

Основное преимущество логистической функции — простота вычисления производной: она равна произведению фактической реакции персептронного нейрона  $y_j(n)$  на дополнение ее до 1:

$$\varphi'(v_j(n)) = a \cdot y_j(n) \cdot [1 - y_j(n)] \quad (4.1.2)$$

Поскольку вычисление фактической реакции производится в любом случае, то для вычисления  $\varphi'(v_j(n))$  дополнительных вычислений почти не происходит.

**Функция гиперболического тангенса** Нижняя и верхняя горизонтальная её асимптота отсекает на оси ординат значения  $\pm 1$ :

$$\varphi(v_j(n)) = a \cdot \tanh(b \cdot v_j(n)), \quad (a, b) > 0 \quad (4.1.3)$$

где  $a, b$ —положительные константы, имеющие геометрический смысл. Рекомендуется выбирать  $a = 1,7159, b = 2/3$ .

Затраты на вычисление также минимальны, поскольку производная функции активации есть разность квадратов:

$$\varphi'(v_j(n)) = \frac{b}{a} [a - y_j(n)] [a + y_j(n)] \quad (4.1.4)$$

Для „оптимального“  $a = 1,7159$  и  $b = 2/3$  функция гиперболического тангенса антисимметрична (обладает радиальной симметрией относительно начала координат) и представлена на рисунке 12. Также выполняется:

$$\begin{aligned}\varphi(-1) &= -1, & \varphi(1) &= 1 \\ \varphi'(0) &= ab = 1,1424 \approx 1 \text{ — угол наклона линейной части } \approx 45^\circ \\ \varphi''(0) &\rightarrow \max \text{ при } v = 1\end{aligned}$$

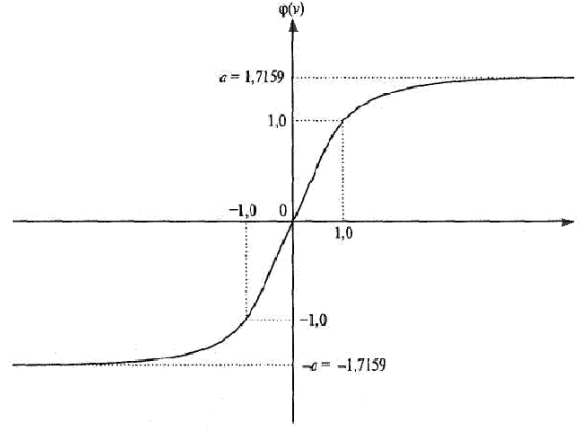


Рис. 12: Антисимметричная функция активации

Поскольку при  $|v_j(n)| \gg 1$  наблюдается насыщение функции активации, что приводит к уменьшению значения производной функции активации и, соответственно, корректирующего значения до 0 (коррекции не происходит), необходимо стараться, чтобы обучающий пример выводил значение локального индуцированного поля в диапазон  $[-1, 1]$ . В этом случае производная функции активации имеет наибольшее значение и это обеспечивает эффективный шаг поиска.

#### 4.1.2 Выбор компромиссного режима обучения персептронных нейронов и сети в целом

Формирование финального значения локального градиента (и направления поиска) производится по текущему обучающему примеру и нескольким предыдущим — комбинированный режим обучения с конечной памятью.

Такой режим реализуется на основе сглаживания текущих оценок параметров персептронного нейрона для нескольких обучающих примеров (текущего и нескольких предыдущих). точки зрения подавления шумовой компоненты необходимо выбирать конечный временной интервал сглаживания (размер памяти).

Один из эффективных алгоритмов комбинированного режима обучения нейронов, предусматривающий учёт текущего обучающего примера и нескольких предыдущих (сглаживание зашумлённых текущих оценок параметра нейрона) — (*обобщённое дельта-правило*) — опирается на модель экспоненциальной скользящей средней (ЕМА).

Специфика адаптации модели состоит в формулировке модели ЕМА в величинах коррекции параметров персептронного нейрона

$$\Delta w_{ji}(n) = w_{ji}(n) - w_{ji}(n-1) \quad (4.1.5)$$

В соответствии с формулой ЕМА:

$$\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) + (1 - \eta) \Delta w_{ji}(n-1) \quad (4.1.6)$$

Для получения очередной коррекции необходима фактическая реакция нейронов предыдущего слоя, взвешенная локальным градиентом  $\delta_j(n)$  и параметром скорости обучения  $\eta$ , к которой добавляется предыдущая коррекция  $\Delta w_{ji}(n-1)$ , что соответствует структуре ЕМА (из предыдущего семестра):

$$y(n) = \alpha x(n) + (1 - \alpha) y(n-1), \quad \alpha = \frac{2}{K+1} \quad (4.1.7)$$



В этой модели размер эффективной памяти  $K$  управляется параметром  $\alpha$ , который является весом штрафа за точность описания закономерности, скрытой в данных.

Параметр скорости обучения  $\eta$  может быть интерпретирован как параметр сглаживания, определяющий величину памяти, т.е.:

$$\eta = \frac{2}{N_{smo} + 1} \quad (4.1.8)$$

где  $N_{smo}$  — количество итераций обучения.

Однако, в теории ИНС принята другая формулировка процедуры экспоненциального сглаживания за счёт добавления *момента инерции*  $\alpha$  получим *обобщённое дельта-правило*:

$$\Delta w_{ji}(n) = \eta \delta_j(n) y_i(n) + \alpha \Delta w_{ji}(n-1) \quad (4.1.9)$$

При  $\eta \rightarrow 0$  память увеличивается, значит, чем ближе постоянная момента  $\alpha$  к 1, тем больше память — большее количество обучающих примеров алгоритм „помнит“ и тем „сильнее“ сглаживание.

При  $\alpha = 0$  получим обычное дельта-правило Уидроу-Хоффа (режим последовательного обучения), коррекция осуществляется после каждого обучающего примера.

При  $\alpha \rightarrow 1$  получим режим пакетного обучения.

#### 4.1.3 Максимизация информативности

Для сокращения количества итераций обучения (улучшения его сходимости) рекомендуется выбирать обучающие примеры, обеспечивающие наибольшие отклонения фактических реакций от указаний учителя. В этом случае локальный градиент будет иметь наибольшее значение, что приведёт к эффективному значительному шагу коррекции.

Также необходимо стараться выбирать обучающие примеры так, чтобы не „попадать“ в ветви насыщения функции активации — работать на линейном или квазилинейном участке функции активации.

Для улучшения сходимости алгоритма обратного распространения ошибок необходимо как минимум реализовать линейное масштабирование входных данных — *нормализацию входов*. Рекомендуется в 0-м слое реализовать *декоррелирующее* преобразование входных данных.

На практике получила популярность комбинация 2 нейросетевых парадигм:

1. Сеть главных компонент — реализация декоррелирующего преобразования (преобразование Карунена-Лоэва).
2. Многослойный персептрон.

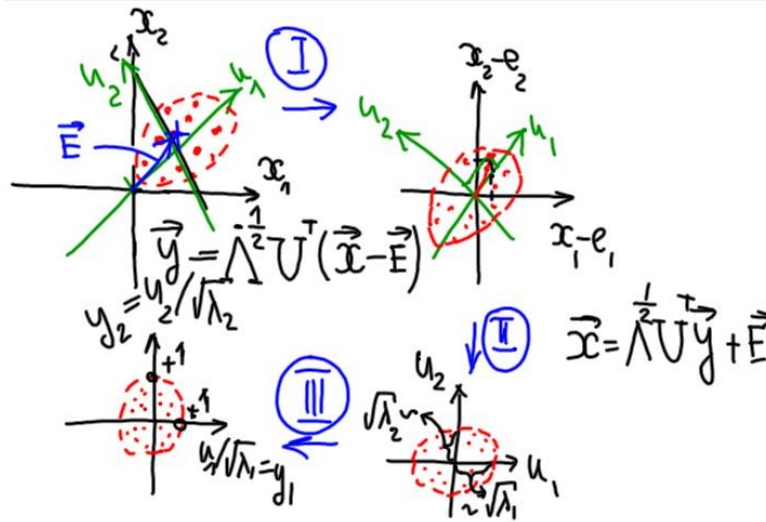


Рис. 13: Схема декоррелирующего преобразования

Рассмотрим схему декоррелирующего преобразования 13. В общем случае обучающие примеры обладают ненулевым центром группировки (отлично от нуля математическое ожидание в пространстве исходных признаков  $\vec{E}$ ) и имеется значимая корреляция.

Наличие значимой корреляции в геометрическом смысле означает, что в простейшем случае морфология диаграммы рассеяния представляет собой гиперэллипсоид. Чем больше он вытянут вдоль главных осей (главных информативных признаков), тем более эти главные оси значимы.

Декоррелирующее преобразование реализуется в 3 этапа.

### **I Параллельный перенос СК исходных признаков в центр группировки**

Получаем новые оси  $\vec{x} - \vec{E}$  — признак минус его математическое ожидание. При этом морфология диаграммы рассеяния всё ещё демонстрирует наличие явной корреляции, поскольку главные оси повернуты относительно новых осей  $\vec{x} - \vec{E}$ .

### **II Поворот системы координат до совпадения с главными информативными осями**

Выполняется умножение на транспонированную матрицу собственных векторов  $\Lambda^T$ . Новые оси представляют собой главные и шумовые признаки.

В этой системе координат эллипсоид ещё имеет различный масштаб вдоль осей, но корреляции уже нет. В этом случае корреляционная матрица имеет диагональный вид — на главной диагонали стоят собственные значения (аналоги дисперсий).

### **III Масштабирование вдоль главных осей**

Каждый признак масштабируется — делится на сингулярное число  $\sqrt{\lambda}$  (аналог СКО). Получаем новую морфологию диаграммы рассеяния в виде гиперболы, где, по сути, реализуется метрика Евклида.

#### **4.1.4 Инициализация. Выбор начального приближения для биаса и синаптических весов**

Рациональным является следующий подход для начального выбора биаса и синаптических весов персептронных нейронов.

Начальное приближение рекомендуют выбирать на основе выборки псевдослучайных чисел, полученных из равномерного распределения с нулевым математическим ожиданием и СКО, обратно пропорциональным квадратному корню из количества входов выбранного нейрона. Т.е., необходимо реализовать процедуру статистического моделирования.

*Ермаков, Михайлов. Курс статистического моделирования*

Представлены алгоритмы и формулы формирования псевдослучайных формул и стандартных распределений.

## Лекция 07.05

### 5.1 Сеть главных компонент

#### 5.1.1 Метод главных компонент (МГК)

Lec\_Neuro\_04.docx

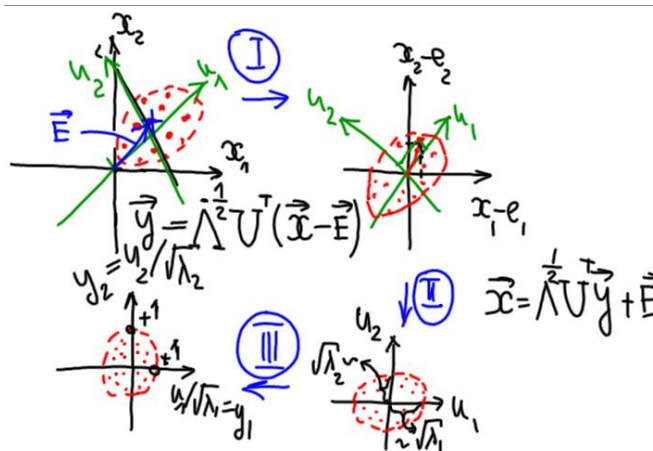


Рис. 14: Метод главных компонент

При решении задач распознавания образов на первом этапе формируется пространство исходных признаков, которое обычно обладает следующими недостатками:

1. Избыточность исходного списка признаков, поскольку он формулируется экспертным (эмпирическим) методом.
2. Значимая корреляция исходных признаков (содержание неравных долей информации об объекте анализа)

Одним из инструментов анализа морфологии такого пространства является анализ диаграммы рассеяния, которая может демонстрировать явные направления, вдоль которых степень свободы (рассеяния) значима, т.е., точки в пространстве образуют эллипсоид. Э является признаком значимой корреляции.

Для анализа пространства применяется ковариационная матрица исходных признаков, в которой прежде всего интересны собственные вектора и собственные значения.

Существует 2 группы собственных значений:

1. **Значимые** ( $\gg 0$ ). Образуют *сигнальное подпространство*, размерность (количество степеней свободы) которого определяется рангом  $M$  ковариационной матрицы обучающих примеров. Также определяют энергию сигнального подпространства, поскольку собственные значения являются аналогом дисперсии, и количество степеней свободы.
2. **Малозначимые** ( $\approx 0$ ). Образуют *шумовое подпространство*, порождённое ошибками измерений и вычислений.

Согласно МГК собственные вектора матрицы ковариации исходных признаков определяют новый базис (как было показано Пирсоном), а ошибка описания исходных признаков с помощью МГК определяется суммой малозначимых собственных значений.

Относительная ошибка определяется как:

$$\frac{\lambda_{M+1} + \dots + \lambda_{M+1}}{\lambda_1 + \dots + \lambda_M} \quad (5.1.1)$$

То есть, задавшись допустимым значением ошибок описания, можно получить критерий отбора собственных значений.

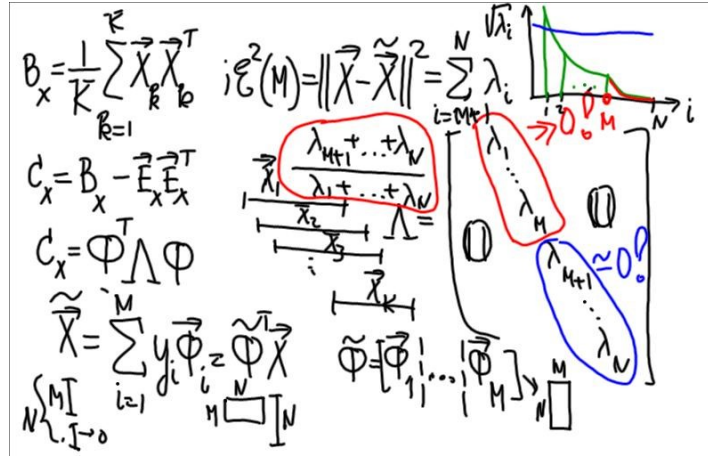


Рис. 15: Метод главных компонент (продолжение)

### 5.1.2 Сеть главных компонент (введение)

Сеть главных компонент в процессе обучения по набору немеченных обучающих примеров должна выдать направляющие косинусы осей сигнального подпространства (компоненты собственных векторов), а совокупность фактических реакций нейронов сети главных компонент дать оценки дисперсий (оценки собственных значений), что позволит ортогонально спроецировать исходные данные на оси сигнального подпространства. Задача упрощается в случае, если корреляция значима, поскольку тогда эффективный ранг будет на порядок меньше ранга пространства исходных признаков.

Аналогично МГК функционирует и в задачах прогнозирования.

В основе сети главных компонент лежат выборочные оценки направляющих косинусов собственных векторов сигнального подпространства и оценки собственных значений, определяющихся ковариационной матрицей (второй центральный момент). В литературе ковариация и корреляция в русской и иностранной литературе альтернативны (ковариация — correlation). Основой выборочного оценивания является корреляционное произведение.

### 5.1.3 Гипотеза Хебба (принцип самоусиления (корреляции))

Правило коррекции синаптических весов, устанавливающее связь между входным сигналом и его реакцией в виде корреляционного произведения в простейшей форме:

$$\Delta w_{kj}(n) = \eta \cdot y_k(n) \cdot x_j(n) \quad (5.1.2)$$

иначе называемое *правилом обучения синапса Хебба*. Здесь  $x_j$  — сигнал на входе линейного нейрона от  $j$ -го исходного признака,  $y_k$  — сигнал на выходе  $k$ -го нейрона,  $\eta \in [0, 1]$  — параметр скорости обучения.

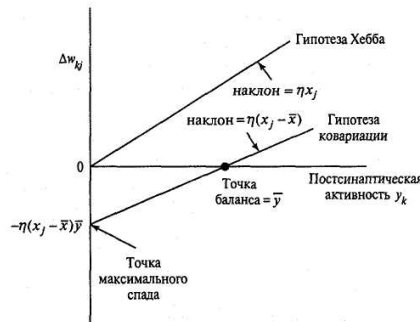


Рис. 16: Иллюстрация принципов Хебба и ковариации

Зависимость корректирующего значения  $\Delta w_{kj}(n)$  от выходного сигнала  $y_k$  представлена на 16. Угловой коэффициент для гипотезы Хебба равен  $\eta \cdot x_j$ . Видим, что правило Хебба не вводит ограничение на значение синаптических весов, которые могут возрастать бесконечно, что приводит к состоянию *насыщения сети* (состояние, при которой выходные значения ИНС нечувствительны к изменению исходных признаков).

#### 5.1.4 Принцип ковариации

Для стабилизации весов необходимо ввести обратную связь и вместо *принципа самоусиления (корреляции)* ввести принцип конкурентности — *принцип ковариации*.

Для этого правило Хебба модифицируют центрированием обучающих примеров — происходит параллельный перенос системы координат исходных признаков в центр группировки обучающих примеров (см. прошлое занятие). Тогда получаем:

$$\Delta w_{kj} = \eta (x_j - \bar{x}) (y_k - \bar{y}) \quad (5.1.3)$$

где  $\bar{x}$  и  $\bar{y}$  — усредненные по времени значения предсинаптического  $x_j$  и постсинаптического  $y_k$  сигналов.

Теперь возможны и отрицательные корректирующие значения и существует *точка баланса*, что видно для соответствующей прямой на 16.

## 5.2 Метод настройки весов для сети главных компонент

Линейный нейрон, являющийся структурным элементом сети главных компонент, называют *фильтром Хебба*. Несмотря на простоту, в синаптических весах фильтра Хебба сохраняются направляющие косинусы осей в пространстве главных компонент.

*Lec\_Neuro\_14.docx*

#### 5.2.1 Правила обучения линейного нейрона

Как было отмечено выше, синаптические веса линейного нейрона по сути своей являются направляющими косинусами оси в пространстве исходных признаков.

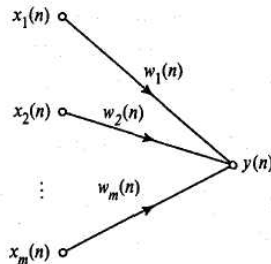


Рис. 17: Граф линейного нейрона

Рассмотрим модель линейного нейрона (представлена на 17):

$$y = \sum_{i=1}^m w_i x_i \quad (5.2.1)$$

где  $w_i$  — синаптические веса (всего  $m$  штук),  $n$  — итерация обучения.

**Важно! Здесь отсутствует сигнал стандартного смещения +1 в отличие от персептронного нейрона.**

В рамках гипотезы корреляции (согласно правилу Хебба):

$$w_i(n) = w_i(n-1) + \eta \cdot y(n-1) \cdot x_i(n-1), \quad i = 1, 2, \dots, m \quad (5.2.2)$$

где  $n$  — дискретное время,  $\eta$  — параметр скорости обучения.

Здесь реализован только *принцип самоусиления*, что нас не устраивает. Необходимо ввести *эффект конкуренции* между синаптическими весами: между собой они могут соревноваться, но при этом длина вектора синаптических весов всегда должна быть равна 1, т.е. геометрически синаптические веса приобретают смысл направляющих косинусов.

Для этого после каждой итерации будем проводить нормировку, если совокупность весов не образует вектор единичной длины:

$$w_i(n) = \frac{w_i(n-1) + \eta \cdot y(n-1) \cdot x_i(n-1)}{\sqrt{\sum_{i=1}^m [w_i(n-1) + \eta \cdot y(n-1) \cdot x_i(n-1)]^2}} \quad (5.2.3)$$

В результате нормировки правило обучения становится нелинейным. Поэтому применим методы линеаризации: аппроксимируем нелинейную зависимость рядом (Маклорена или Тейлора), степень аргумента которого будет убывать достаточно быстро. В качестве аргумента используем параметр скорости обучения  $\eta \in [0, 1]$ .

Необходимо вычислить набор частных производных по  $\eta$  и их значения при  $\eta = 0$ , чтобы получить полином, у которого отбросим степени, старше первой (для ряда Маклорена). Итого получим:

$$w_i(n) = w_i(n-1) + \eta \cdot y(n-1) [x_i(n-1) - y(n) \cdot w_i(n-1)] \quad (5.2.4)$$

Это полученное линеаризованное правило обучения синаптических весов линейных нейронов, реализующее оба принципа (самоусиления и конкуренции), иногда называют *правилом Ойя*.

Данное правило обучения в итоге выдаёт наибольшее собственное значение ковариационной матрицы обучающих примеров, а веса являются направляющими косинусами. Для понимания этого реализуем матричную формулировку алгоритма.

### 5.2.2 Матричная формулировка алгоритма Ойя

Для входного вектора исходных признаков  $\vec{x}(n)$  и вектора синаптических весов  $\vec{w}(n)$  введём следующие обозначения:

$$\vec{x}(n) = [x_1(n), x_2(n), \dots, x_m(n)]^T \quad (5.2.5)$$

$$\vec{w}(n) = [w_1(n), w_2(n), \dots, w_m(n)]^T \quad (5.2.6)$$

Тогда выходная реакция линейного нейрона и правило его обучения:

$$y(n) = \vec{x}^T(n) \cdot \vec{w}(n) = \vec{w}^T(n) \cdot \vec{x}(n) \quad (5.2.7)$$

$$\vec{w}(n) = \vec{w}(n-1) + \eta \cdot y(n-1) [\vec{x}(n-1) - y(n) \cdot \vec{w}(n-1)] \quad (5.2.8)$$

Подставим в правило Ойя 5.2.8 фактическую реакцию нейрона (различные представления для двух вхождений в формулу 5.2.7) и получим:

$$\vec{w}(n) = \vec{w}(n-1) + \eta \cdot y(n-1) [\vec{x}(n-1) \cdot \vec{x}^T(n-1) \cdot \vec{w}(n-1) - \vec{w}^T(n-1) \cdot \vec{x}(n-1) \cdot \vec{x}^T \cdot \vec{w}(n-1) \cdot \vec{w}(n-1)] \quad (5.2.9)$$

Здесь  $\vec{x}(n) \cdot \vec{x}^T(n)$  — *элементарная корреляционная матрица*, квадратная матрица размерности  $m \times m$ . Введём замены:

$$\vec{h}(\vec{w}(n), \vec{x}(n)) = [\vec{x}(n) \cdot \vec{x}^T(n) \cdot \vec{w}(n) - \vec{w}^T(n) \cdot \vec{x}(n) \cdot \vec{x}^T \cdot \vec{w}(n) \cdot \vec{w}(n)] \quad (5.2.10)$$

$$\vec{\eta}(n) = \eta \cdot y(n) \quad (5.2.11)$$

Получим:

$$\vec{w}(n) = \vec{w}(n-1) + \vec{\eta}(n-1) \cdot \vec{h}(\vec{w}(n-1), \vec{x}(n-1)), \quad n = 0, 1, 2, \dots \quad (5.2.12)$$

Это нелинейное стохастическое разностное уравнение.

По мере обучения объёма обучающих примеров переходим к непрерывному времени и разностное уравнение превращается в систему линейных ДУ:

$$\frac{d}{dt} \vec{w}(t) = \langle \vec{h}(\vec{w}(t)) \rangle \quad (5.2.13)$$

При этом используем обычную выборочную среднюю по матрице набора обучающих примеров  $\vec{X}$  (по всем реализациям  $\vec{x}$ ):

$$\langle \vec{h}(\vec{w}) \rangle = \lim_{n \rightarrow \infty} E \left[ h(\vec{w}, \vec{X}) \right] \quad (5.2.14)$$

Функция коррекции  $\vec{h}(\vec{w}, \vec{x})$  для фильтра Хебба имеет вид:

$$\vec{h}(\vec{w}, \vec{x}) = \vec{x}(n) \cdot y(n) - y^2(n) \cdot \vec{w}(n) = \vec{x}(n) \cdot \vec{x}^T(n) \cdot \vec{w}(n) - [\vec{w}^T(n) \cdot \vec{x}(n) \cdot \vec{x}^T(n) \cdot \vec{w}(n)] \vec{w}(n) \quad (5.2.15)$$

В квадратных скобках 5.2.15 — квадратичная форма (по определению положительна).

Усредняя функцию коррекции  $\vec{h}(\vec{w}, \vec{x})$  5.2.15 по столбцам матрицы обучающих примеров:

$$\langle \vec{h} \rangle = \lim_{n \rightarrow \infty} E \left[ \vec{X}(n) \cdot \vec{X}^T(n) \cdot \vec{w}(n) - \left( \vec{w}^T(n) \cdot \vec{X}(n) \cdot \vec{X}^T(n) \cdot \vec{w}(n) \right) \vec{w}(n) \right] = \mathbf{R} \cdot \vec{w}(\infty) - [\vec{w}^T(\infty) \cdot \mathbf{R} \cdot \vec{w}(\infty)] \vec{w}(\infty) \quad (5.2.16)$$

где  $\mathbf{R}$  — матрица корреляции исходных данных,  $\vec{w}(\infty)$  — предельное значение вектора синаптических весов

Итого получим систему ЛДУ:

$$\frac{d}{dt} \vec{w}(t) = \langle \vec{h}(\vec{w}(t)) \rangle = \mathbf{R} \cdot \vec{w}(t) - [\vec{w}^T(t) \cdot \mathbf{R} \cdot \vec{w}(t)] \vec{w}(t) \quad (5.2.17)$$

решение которой известно:

$$\lim_{n \rightarrow \infty} \vec{w}(n) = \vec{q}_1 \quad (5.2.18)$$

где  $\vec{q}_1$  — оценка первого собственного вектора. Считаем свойство фильтра Хебба доказанным.

Для существования решения накладываются ограничения на выбор параметра скорости обучения  $\eta$ , который от итерации к итерации должен изменяться, чтобы выполнялись условия (*стохастическая аппроксимация Роббенса-Монро*):

$$\sum_{n=1}^{\infty} \eta(n) = \infty \quad (5.2.19)$$

$$\sum_{n=1}^{\infty} \eta^p(n) < \infty, \quad p > 1 \quad (5.2.20)$$

$$\eta(n) \rightarrow 0 \text{ при } n \rightarrow \infty \quad (5.2.21)$$

Обучение линейного нейрона по правилу Ойя гарантирует в конце обучения следующие результаты:

### Свойства фильтра Хебба

1. Вектор синаптических весов  $\vec{w}(n)$  является оценкой первого собственного вектора  $\vec{q}_1$  корреляционной матрицы  $\mathbf{R}$  входных данных. Синаптические веса также являются направляющими косинусами.
2. Дисперсия совокупности фактических реакций линейного нейрона, обученного по правилу Ойя, даёт выборочную оценку дисперсии ковариационной матрицы обучающих примеров — её наибольшее собственное значение  $\lambda_1$ .

## Лекция 14.05 (не проводилась)



## Лекция 21.05

### 7.1 Архитектура сети главных компонент — обобщённый алгоритм Хебба

Фильтр Хебба на основе одного линейного нейрона оценивает первую главную компоненту входного сигнала. Сеть прямого распространения с одним слоем линейных нейронов способна оценивать все главные компоненты входного сигнала произвольной размерности.

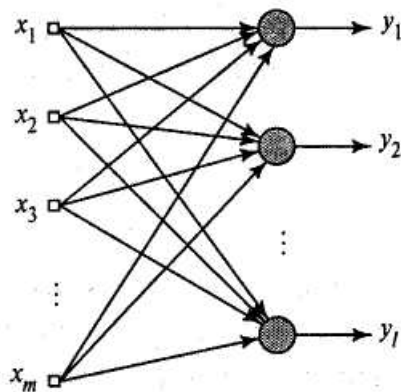


Рис. 18: Архитектура сети главных компонент

Lec\_Neuro\_14.docx

Сеть главных компонент состоит из  $m$  входных нейронов, обозначенных  $x_i$ . Имеется единственный слой из  $l$  нейронов  $y_j$ , каждый из которых обучается правилом Ойа. В синаптических весах каждого нейрона сохраняются направляющие косинусы соответствующего собственного вектора ковариационной матрицы. Матрица этих весов формирует базис размерностью  $l < m$ , поворну́тый относительно базиса входных признаков.

Размерность сигнального подпространства (пространства информативных признаков) определяется рангом ковариационной матрицы данных (*в электронной версии это  $l$* ). МГК позволяет обнаружить наличие сильной корреляции в исходных признаках и, в силу этого, размерность сигнального подпространства значительно меньше пространства входных признаков ( $l \ll m$ ). Каждый нейрон позволяет получить оценку дисперсии (собственные значения ковариационной матрицы данных) — масштаба соответствующей оси.

В основе алгоритма обучения сети главных компонент лежит *обобщенный алгоритм обучения Хебба* (Generalized Hebbian Algorithm, GHA). Алгоритм устроен следующим образом:

$$\Delta w_{ji}(n) = \eta \left[ y_j(n)x_i(n) - y_j(n) \sum_{k=1}^j w_{ki}(n)y_k(n) \right], \quad i = 1, \dots, m, \quad j = 1, \dots, l \quad (7.1.1)$$

где  $\Delta w_{ji}(n)$  устанавливает связь между  $x_i$  и  $y_j$  и является *величиной коррекции* пред- и постсинаптических сигналов  $x_i$  и  $y_j$  соответственно.

#### 7.1.1 Последовательное обучение

Для упрощения анализа на первом этапе примем предположим, что реализуется последовательное во времени обучение (настройка параметров) линейных нейронов: сначала обучается 1-й линейный нейрон, за ним 2-й линейный нейрон и т.д.

Рассмотрим обучение 1-го линейного нейрона. В этом случае  $j = 1$ , а поскольку остальные нейроны не используются, сумма в 7.1.1 вырождается в  $\sum_{k=1}^1$  — алгоритм вырождается в правило Ойа, реализующий 1-й и 2-й принципы самоорганизации (самоусиления [гипотеза корреляции, 1-я компонента] и конкуренции).

Результатом является формирование в пространстве признаков оси, единичный орт которой указывает направление первого главного собственного вектора ковариационной матрицы обучающего примера. При проекции этой оси на ось сигнального подпространства получим порцию полезной информации, которую даёт  $j$ -я ось сигнального подпространства. Вдоль оси с  $j = 1$  степень рассеяния обучающих примеров наибольшая.

Далее из вектора обучающего примера  $\vec{x}$  вычитаем фактическую реакцию, масштабированную вектором синаптических весов — проекцию его на ось сигнального подпространства:

$$\vec{x}'(n) = \vec{x}(n) - \vec{w}(n)y_1(n) \quad (7.1.2)$$

По сути, происходит "обеление" обучающего примера, из которого "удаляется" первая главная (информативная) компонента. Таким образом, второй нейрон ( $j = 2$ ) обрабатывает уже модифицированный вектор входных компонентов  $\vec{x}''(n)$ . Разобьём сумму в 7.1.1 на 2 компоненты:

$$\Delta w_{2i}(n) = \eta \cdot y_2(n) [x_i(n) - y_1(n)w_{1i}(n) - w_{2i}(n)] \quad (7.1.3)$$

Введём новое обозначение:

$$x'_{2i} = x_i(n) - w_{1i}(n)y_1(n) \quad (7.1.4)$$

Тогда 7.1.1 для  $j = 2$  с учётом 7.1.4 становится аналогичным правилу Ойа:

$$\Delta w_{2i}(n) = \eta \cdot y_2(n) [x'_{2i}(n) - w_{2i}(n)y_2(n)] \quad (7.1.5)$$

Удаление первой информативной компоненты 7.1.2 приводит 7.1.1 к виду, аналогичному правилу обучения Ойа.

Повторив аналогично действия, удалим вторую главную компоненту из набора обучающих примеров для перехода к обучению 3-го нейрона, т.е. в гиперплоскости, образованной 1-м и 2-м собственными векторами, имея масштабированные вектора информативных признаков, их сумму удаляем из вектора обучающих примеров.

$$\vec{x}'(n) = \vec{x}(n) - \vec{w}(n)y_1(n) - \vec{w}(n)y_1(n) \quad (7.1.6)$$

Аналогичные действия повторяем и для третьего нейрона, что позволяет перейти к обучению 4-го линейного нейрона. Т.е., в общем случае:

$$\Delta w_{ji}(n) = \eta \cdot y_j(n) [x'_{ji}(n) - w_{ji}(n)y_j(n)], \quad i = 1, \dots, m, \quad j = 1, \dots, l \quad (7.1.7)$$

$$x'_{ji}(n) = x_i(n) - \sum_{k=1}^{j-1} w_{ki}(n)y_k(n) \quad (7.1.8)$$

Итого получим, что в исходном пространстве сформирован базис сигнального подпространства

$$\vec{w}_1, \dots, \vec{w}_l,$$

на который ортогонально спроецирован вектор обучающего примера, т.е. Реализована *пространственная декомпозиция*. Остаток (значение вектора обучающего примера после вычета всех главных компонент) — *ошибка реконструкции данных*, — вектор, который целиком содержится в шумовом подпространстве.

Анализ показывает, что такого рода процедура представляет собой декомпозицию исходных признаков на информативные признаки в виде произведений собственных чисел ковариационной матрицы (совпадают с  $y_i$ ) на собственные вектора ковариационной матрицы  $\vec{q}_i$ .

Для получения оценки дисперсии ошибки реконструкции из каждого обучающего примера необходимо вычесть собственные вектора и просуммировать длины полученных векторов-остатков, либо получить относительное значение, пользуясь значениями собственных чисел ковариационной матрицы:

$$\frac{\sum_{k=l+1}^m \lambda_k}{\sum_{k=1}^m \lambda_k} \quad (7.1.9)$$

### 7.1.2 Параллельное обучение

Рассмотрим новые объекты анализа.

Сформируем матрицу единичных векторов синаптических весов сети главных компонент размерности  $l \times m$ :

$$\mathbf{W}(n) = [\vec{w}_1(n), \vec{w}_2(n), \dots, \vec{w}_l(n)]^T \quad (7.1.10)$$

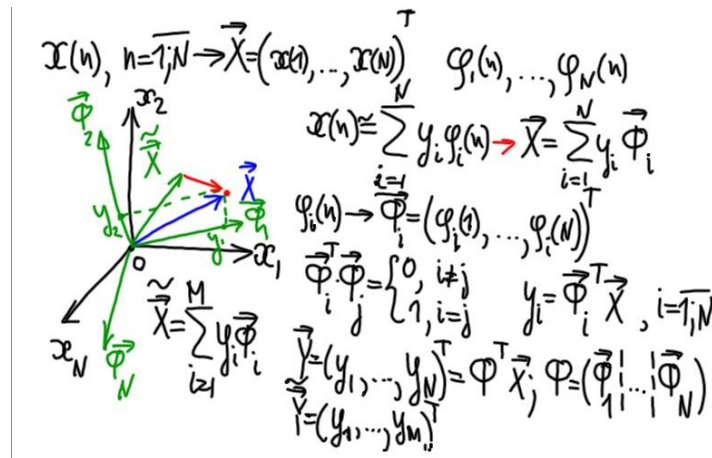


Рис. 19: hyperflat 12 1

Заметим, что:

$$\mathbf{W}^T(n) = \{w_{ji}(n)\}, \quad i = 1, \dots, m, \quad j = 1, \dots, l$$

Вводится величина коррекции (обобщенный алгоритм Хебба в матричном виде):

$$(8.91) \Delta \mathbf{W}(n) = \eta(n) \{ \vec{y}(n) \cdot \vec{x}^T(n) - LT [\vec{y}(n) \cdot \vec{y}^T(n)] \cdot \mathbf{W}(n) \} \quad (7.1.11)$$

$\vec{y}(n)$  – вектор-столбец с  $l$  компонент,  $\vec{x}(n)$  – вектор-столбец с  $m$  компонент,  $LT \{ \vec{y}(n) \cdot \vec{y}^T(n) \}$  – оператор, устанавливающий все элементы матрицы  $\vec{y}(n) \cdot \vec{y}^T(n)$ , расположенные выше ее главной диагонали, в нуль (формирует нижнюю треугольную матрицу). Нижняя треугольная матрица необходима для реализации ГНА. По мере обучения матрица и вектора обновляются.

Не уверен, но  $\eta$  – матрица

Если параметр скорости обучения  $\eta(n)$  удовлетворяет условиям:

$$\lim_{n \rightarrow \infty} \eta(n) = 0, \quad \sum_{n=0}^{\infty} \eta(n) = \infty$$

то обобщенный алгоритм Хебба будет сходиться к фиксированной точке, а  $\mathbf{W}(n)$  сходится к матрице, строки которой являются первыми  $l$  собственными векторами матрицы ковариации размерности  $m \times m$  входных векторов размерности  $m \times 1$ , упорядоченных по убыванию собственных значений.

## Лекция 28.05 (показаны новые наработки в 514м)