

EM-алгоритм и метод релевантных векторов для классификации

Дата: 2 ноября 2011

Метод оптимизации Ньютона

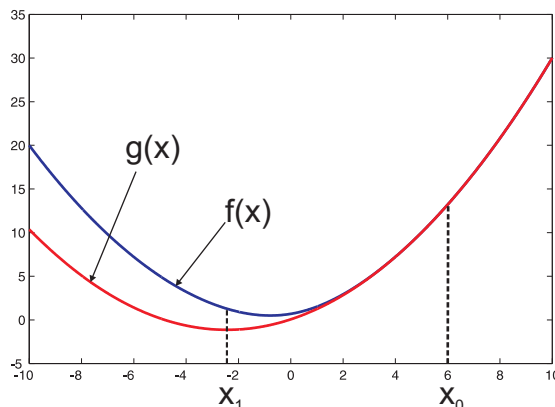


Рис. 1: Иллюстрация одной итерации метода Ньютона для минимизации функции $f(x) = \log(1 + \exp(x)) + x^2/5$ путем аппроксимации с помощью квадратичной функции $g(x)$.

Рассмотрим задачу безусловной оптимизации гладкой функции

$$f(\mathbf{x}) \rightarrow \min_{\mathbf{x}}, \quad \mathbf{x} \in \mathbb{R}^d.$$

Идея метода оптимизации Ньютона заключается в итерационном приближении функции $f(\mathbf{x})$ с помощью квадратичной функции $g(\mathbf{x})$, минимум которой может быть найден аналитически. Приближение функции с помощью квадратичной функции путем приравнивания нулевой, первой и второй производных в точке \mathbf{x}_0 эквивалентно разложению в ряд Тейлора функции f с отбрасыванием всех членов третьего и более высокого порядков:

$$f(\mathbf{x}) \simeq g(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0))^T(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla \nabla f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0).$$

Здесь через ∇f и $\nabla \nabla f$ обозначены градиент и гессиан функции f соответственно. Минимум функции g может быть легко найден путем приравнивания градиента g к нулю:

$$\nabla g(\mathbf{x}_*) = \nabla f(\mathbf{x}_0) + \nabla \nabla f(\mathbf{x}_0)(\mathbf{x}_* - \mathbf{x}_0) = \mathbf{0} \Rightarrow \mathbf{x}_* = \mathbf{x}_0 - (\nabla \nabla f(\mathbf{x}_0))^{-1} \nabla f(\mathbf{x}_0).$$

Таким образом, метод оптимизации Ньютона представляет собой следующий итерационный процесс (см. рис. 1):

$$\mathbf{x}_{t+1} = \mathbf{x}_t - (\nabla \nabla f(\mathbf{x}_t))^{-1} \nabla f(\mathbf{x}_t). \quad (1)$$

Метод Ньютона имеет на порядок более высокую скорость сходимости по сравнению с методом градиентного спуска, но при этом гораздо меньшую область применимости. Сходимость метода Ньютона доказана для выпуклых функций с ограниченным модулем производной (в более общем случае для Липшиц-непрерывных функций с обратимым гессианом). К сожалению, на практике такие функции встречаются нечасто. Тем не менее, во многих случаях метод Ньютона можно применять при условии выбора хорошего начального приближения \mathbf{x}_0 , в окрестности которого функция удовлетворяет условиям выпуклости и Липшиц-непрерывности.

На каждом шаге метода Ньютона (1) требуется вычислять и обращать гессиан оптимизируемой функции. В пространствах большой размерности поиск гессиана может быть связан со значительными вычислительными сложностями. Кроме того, обращение гессиана, близкого к вырожденному,

сопряжено с численными неустойчивостями. В этих случаях используются т.н. квази-ньютоновские методы (например, алгоритм Левенберга-Марквардта, L-BFGS и другие), которые используют различные приближения гессиана и вводят поправки на вырожденный случай.

EM-алгоритм в общем виде

Пусть имеется вероятностная модель, задаваемая совместным распределением $p(X, T|\Theta)$. Здесь X – набор наблюдаемых переменных, T – набор ненаблюдаемых переменных и Θ – набор параметров модели. Рассмотрим задачу обучения модели (поиск параметров Θ по выборке X) с помощью метода максимального правдоподобия:

$$\log p(X|\Theta) = \log \int p(X, T|\Theta) dT \rightarrow \max_{\Theta}. \quad (2)$$

Рассмотрим произвольное вероятностное распределение $q(T)$. Тогда справедлива следующая цепочка равенств:

$$\begin{aligned} \log p(X|\Theta) &= \int q(T) \log p(X|\Theta) dT = \int q(T) \log \frac{p(X, T|\Theta)}{p(T|X, \Theta)} dT = \int q(T) \log \left[\frac{p(X, T|\Theta)}{q(T)} \frac{q(T)}{p(T|X, \Theta)} \right] dT = \\ &= \underbrace{\int q(T) \log p(X, T|\Theta) dT}_{\mathcal{L}(q)} - \underbrace{\int q(T) \log \frac{p(T|X, \Theta)}{q(T)} dT}_{\text{KL}(q||p(T|X, \Theta))}. \end{aligned}$$

Дивергенция $\text{KL}(q||p(T|X, \Theta)) \geq 0$, следовательно

$$\log p(X|\Theta) \geq \mathcal{L}(q). \quad (3)$$

При этом нижняя граница (3) становится точной, если $q(T) = p(T|X, \Theta)$, т.к. дивергенция $\text{KL}(q||p(T|X, \Theta)) = 0$ для тождественных распределений.

EM-алгоритм для решения задачи (2) представляет собой аналог метода Ньютона для оптимизации произвольной функции, где вместо квадратичного приближения в текущей точке используется нижняя оценка (3). Пусть фиксировано некоторое значение параметров Θ_{old} . Сначала (на E-шаге) находится распределение $q(T)$ как распределение значений скрытых переменных при данных параметрах:

$$q(T) = p(T|X, \Theta_{old}) = \frac{p(X, T|\Theta_{old})}{\int p(X, T|\Theta_{old}) dT}. \quad (4)$$

При таком выборе $q(T)$ нижняя оценка (3) является точной при $\Theta = \Theta_{old}$. Затем (на M-шаге) новые значения параметров Θ находятся путем максимизации нижней границы $\mathcal{L}(q)$, что эквивалентно решению задачи

$$\mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta) \rightarrow \max_{\Theta}, \quad (5)$$

так как энтропия $-\mathbb{E}_q \log q(T)$ распределения $q(T)$ не зависит от Θ . Шаги E и M повторяются в цикле до сходимости. Очевидно, что в процессе EM-итераций нижняя оценка (3), а также значение правдоподобия $p(X|\Theta)$ не убывают. Монотонное неубывание правдоподобия $p(X|\Theta)$ и его ограниченность сверху гарантируют общую сходимость EM-итераций.

Итерационный процесс в EM-алгоритме проиллюстрирован на рис. 2. По аналогии со многими итерационными методами оптимизации, EM-алгоритм позволяет находить только локальный максимум правдоподобия.

Заметим, что во многих практических случаях решение задачи (5) намного проще, чем решение задачи (2). В частности, в рассматриваемой ниже задаче разделения гауссовской смеси задача оптимизации на M шаге может быть решена аналитически.

Вычисление значения функции правдоподобия $p(X|\Theta)$ в фиксированной точке Θ требует интегрирования по пространству T и в ряде случаев может представлять собой вычислительно трудную задачу. Заметим, что эта величина правдоподобия необходима также для вычисления апостериорного распределения $p(T|X, \Theta_{old})$ на E шаге. Однако, распределение $p(T|X, \Theta_{old})$ используется затем только для вычисления математического ожидания логарифма полного правдоподобия на M шаге. Как правило, здесь не требуется знать все апостериорное распределение целиком, а

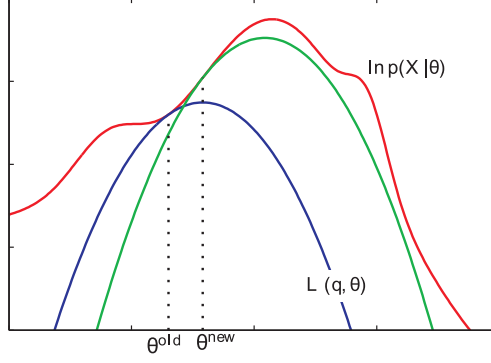


Рис. 2: Иллюстрация итерационного процесса в EM-алгоритме. Нижняя оценка (3) обозначена через $L(q, \theta)$.

достаточно знать лишь несколько статистик этого распределения (например, только мат.ожидания отдельных компонент $\mathbb{E}_{T|X, \Theta_{old}} t_n$ и парные ковариации $\mathbb{E}_{T|X, \Theta_{old}} t_n t_k$). Поэтому EM-алгоритм может быть вычислительно эффективен даже в тех случаях, когда вычисление значения правдоподобия $p(X|\Theta)$ в одной точке затруднено.

EM-алгоритм можно рассматривать как покоординатный подъем для максимизации нижней границы (3). Сначала при фиксированных параметрах Θ_{old} нижняя граница максимизируется по распределению $q(T)$ (эта задача имеет аналитическое решение $q(T) = p(T|X, \Theta_{old})$). Затем при фиксированном $q(T)$ нижняя граница максимизируется по параметрам Θ . Если апостериорное распределение $p(T|X, \Theta_{old})$ не поддается вычислению, то можно ограничить семейство распределений $q(T)$ параметрическим или факторизованным семейством и решать задачу максимизации нижней границы в рамках выбранного семейства распределений. В этом случае будет обеспечиваться монотонный рост нижней границы правдоподобия, но, вообще говоря, не самого значения правдоподобия $p(X|\Theta)$. Такой подход будет обсуждаться в рамках вариационного подхода к решению задачи приближенного байесовского вывода. В случае невозможности точного решения задачи на M шаге можно ограничиться лишь движением в сторону возрастания нижней границы.

EM-алгоритм можно применять также для решения задачи обучения вероятностной модели со скрытыми переменными $p(X, T|\Theta)$ с помощью максимизации апостериорного распределения:

$$p(\Theta|X) \rightarrow \max_{\Theta} \Leftrightarrow \log p(X|\Theta) + \log p(\Theta) \rightarrow \max_{\Theta}.$$

Здесь $p(\Theta)$ — априорное распределение на параметры модели. В этом случае нижней оценкой для оптимизируемого функционала является следующая величина:

$$\log p(X|\Theta) + \log p(\Theta) \geq \mathcal{L}(q) + \log p(\Theta).$$

Итерационная оптимизация данной нижней границы по распределению $q(T)$ и по параметрам Θ приводит к E шагу (4) и M шагу

$$\mathbb{E}_{T|X, \Theta_{old}} \log p(X, T|\Theta) + \log p(\Theta) \rightarrow \max_{\Theta}.$$

EM-алгоритм для разделения гауссовской смеси

Рассмотрим вероятностную модель смеси нормальных распределений:

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k), \quad \sum_{k=1}^K w_k = 1, \quad w_k \geq 0. \quad (6)$$

Модель смеси распределений (не обязательно нормальных) можно рассматривать как модель со скрытой переменной t , которая обозначает номер компоненты смеси:

$$p(t = k) = w_k, \quad (7)$$

$$p(\mathbf{x}|t = k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \Sigma_k). \quad (8)$$

Легко показать, что маргинальное распределение $p(\mathbf{x}) = \sum_k p(\mathbf{x}|t = k)p(t = k)$ в этой модели совпадает с распределением (6). В этом смысле модели (7)-(8) и (6) эквивалентны.

Интерпретация вероятностной модели смеси распределений как модели со скрытой переменной позволяет генерировать выборку из модели смеси следующим образом. Сначала с вероятностями, равными \mathbf{w} , генерируется номер компоненты смеси, из которой затем генерируется точка \mathbf{x} .

Для аппроксимации выборки $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ с помощью модели смеси из K гауссиан можно воспользоваться методом максимального правдоподобия:

$$p(X|\mathbf{w}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}) = \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{w}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}) = \prod_{n=1}^N \left(\sum_k w_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k) \right) \rightarrow \max_{\mathbf{w}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}},$$

$$\sum_k w_k = 1, w_k \geq 0,$$

$$\Sigma_k = \Sigma_k^T, \Sigma_k \succ 0. \quad (9)$$

Данная задача условной оптимизации может быть эффективно решена с помощью EM-алгоритма, описанного выше. Заметим, что количество компонент смеси K не может быть найдено аналогичным образом с помощью максимизации правдоподобия, т.к. значение правдоподобия данных тем выше, чем больше компонент K используется. Для поиска оптимального значения K можно воспользоваться скользящим контролем, где критерием качества аппроксимации тестовых данных является значение правдоподобия.

При использовании EM-алгоритма для решения задачи (9) вероятностная модель смеси распределений интерпретируется как вероятностная модель со скрытыми переменными. Вычислим значение мат.ожидания логарифма полного правдоподобия, необходимого для решения задачи оптимизации на M шаге:

$$\mathbb{E}_q \log p(X, T|\mathbf{w}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}) = \sum_{n=1}^N \sum_{k=1}^K q(t_n = k) (\log w_k + \log \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \Sigma_k)). \quad (10)$$

Заметим, что это выражение зависит только от вероятностей отдельных скрытых переменных $q(t_n = k)$. Нетрудно показать, что эти величины можно найти следующим образом:

$$\gamma_{nk} \triangleq q(t_n = k) = p(t_n = k|\mathbf{x}_n, \mathbf{w}^{old}, \{\boldsymbol{\mu}_k^{old}\}, \{\Sigma_k^{old}\}) = \frac{w_k^{old} \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k^{old}, \Sigma_k^{old})}{\sum_{j=1}^K w_j^{old} \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j^{old}, \Sigma_j^{old})}. \quad (11)$$

Также нетрудно показать, что задача максимизации критерия (10) при ограничениях $\sum_{k=1}^K w_k = 1$, $w_k \geq 0$ может быть решена аналитически:

$$w_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}, \quad (12)$$

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}}, \quad (13)$$

$$\Sigma_k = \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}}. \quad (14)$$

Заметим, что решение для Σ_k (14) удовлетворяет условию симметричности и положительной определенности. Кроме того, формулы (13),(14) соответствуют оценкам максимального правдоподобия для многомерного нормального распределения, в которых каждый объект \mathbf{x}_n берется с весом γ_{nk} .

Таким образом, EM-алгоритм для смеси нормальных распределений заключается в итерационном применении формул (11) и (12)-(14). Этот процесс имеет простую интерпретацию. Величина γ_{nk} показывает степень соответствия между объектом \mathbf{x}_n и компонентой k (определяет вес объекта \mathbf{x}_n для компоненты k). Эти веса затем используются на M шаге для вычисления новых значений параметров компонент. Иллюстрация применения EM-алгоритма для разделения нормальной смеси с двумя компонентами показана на рис. 3.

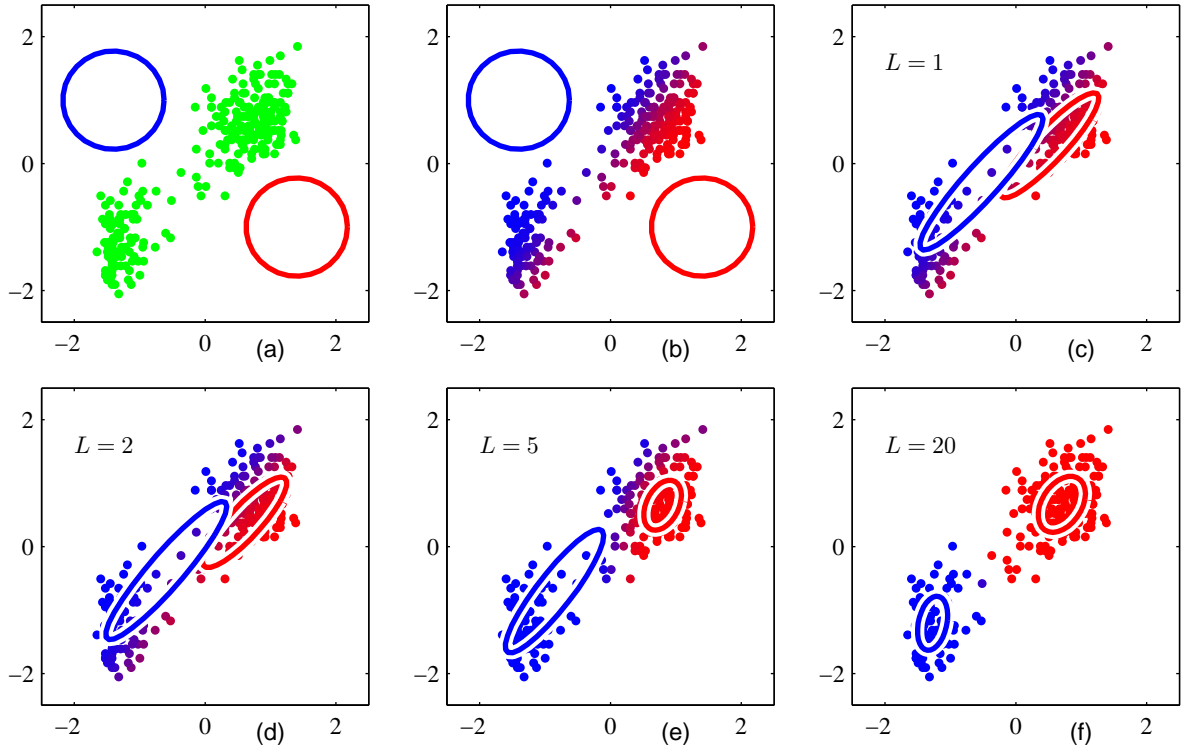


Рис. 3: Иллюстрация применения EM-алгоритма для разделения смеси нормальных распределений с двумя компонентами. На рис. а показана исходная выборка и начальное приближение для двух компонент. На рис. б показан результат E шага. При этом цвета объектов соответствуют значениям γ_{nk} . На рис. с-f показаны результаты вычислений после 1, 2, 5 и 20 итераций.

Одним из применений смеси нормальных распределений является решение задачи кластеризации на K кластеров. В этом случае номер кластера для объекта \mathbf{x}_n определяется величиной

$$k_n = \arg \max_k \gamma_{nk}.$$

Такая схема кластеризации является вероятностным обобщением известного метода кластеризации K -средних.

В заключение этого раздела заметим, что восстановление плотности по данным (в частности, смеси нормальных распределений) является простейшим способом решения задачи идентификации. Для этого сначала для всех объектов обучающей выборки, обладающих заданным свойством, восстанавливается плотность распределения. Затем, для нового объекта \mathbf{x} решение о наличии у него заданного свойства принимается, если значение восстановленной плотности $p(\mathbf{x})$ выше определенного порога.

Логистическая и мультиномиальная регрессия

Рассмотрим задачу классификации на два класса. Имеется обучающая выборка из N объектов $(T, X) = \{t_n, \mathbf{x}_n\}_{n=1}^N$, где $\mathbf{x}_n \in \mathbb{R}^d$ – вектор вещественных признаков для объекта n , а $t_n \in \{-1, 1\}$ – его метка класса. Задача заключается в том, чтобы определить метку класса t для нового объекта, представленного только своим вектором признаков \mathbf{x} .

Рассмотрим решение данной задачи в классе линейных решающих правил. Для этого введем функцию

$$f(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x}).$$

Здесь $w_j \in \mathbb{R}$ – некоторые веса, а $\phi_j : \mathbb{R}^d \rightarrow \mathbb{R}$ – фиксированные базисные функции, задаваемые пользователем. Примеры базисных функций:

- Исходные признаки: $\phi_j(\mathbf{x}) = x_j$;
- Полиномиальные базисные функции: $\phi_j(\mathbf{x}) = x_{j_1} x_{j_2} \dots x_{j_l}$, $j_1, \dots, j_l \in \{1, \dots, d\}$, $1 \leq l \leq L$. L – степень полинома.
- Радиальные базисные функции: $\phi_j(\mathbf{x}) = \exp(-\gamma \|\mathbf{x} - \hat{\mathbf{x}}_j\|^2)$, $\gamma > 0$. Здесь $\hat{\mathbf{x}}$ – некоторые фиксированные вектора в пространстве \mathbb{R}^d , например, объекты обучающей выборки или центры кластеров.

Решение о принадлежности объекта \mathbf{x} к классу t принимается по знаку линейной функции f :

$$t(\mathbf{x}) = \begin{cases} +1, & \text{если } f(\mathbf{x}, \mathbf{w}) > 0, \\ -1, & \text{иначе.} \end{cases}$$

Метод классификации «Логистическая регрессия» представляет собой следующую вероятностную модель:

$$\begin{aligned} p(T, \mathbf{w}|X, \alpha) &= p(T|X, \mathbf{w})p(\mathbf{w}|\alpha) = p(\mathbf{w}|\alpha) \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}), \\ p(t|\mathbf{x}, \mathbf{w}) &= \frac{1}{1 + \exp(-tf(\mathbf{x}, \mathbf{w}))}, \\ p(\mathbf{w}|\alpha) &= \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}I). \end{aligned} \quad (15)$$

Здесь соответствие между меткой класса t и значением линейной функции f на объекте \mathbf{x} вводится с помощью логистической функции. Гауссовское априорное распределение вводится для предотвращения переобучения, а параметр $\alpha \geq 0$ играет роль параметра регуляризации. Обучение в модели (15) (поиск весов \mathbf{w}) с помощью максимизации апостериорного распределения приводит к следующей задаче оптимизации:

$$\begin{aligned} p(\mathbf{w}|T, X, \alpha) \rightarrow \max_{\mathbf{w}} &\Leftrightarrow \log p(T|X, \mathbf{w}) + \log p(\mathbf{w}|\alpha) \rightarrow \max_{\mathbf{w}} \Leftrightarrow \\ & - \sum_{n=1}^N \log(1 + \exp(-t_n \sum_{j=1}^M w_j \phi_j(\mathbf{x}_n))) - \frac{\alpha}{2} \|\mathbf{w}\|^2 \rightarrow \max_{\mathbf{w}}. \end{aligned} \quad (16)$$

Можно показать, что оптимизируемая функция в этой задаче является вогнутой по \mathbf{w} (ее гессиан является неположительно определенным). Поэтому для решения задачи (16) можно воспользоваться методом оптимизации Ньютона (1). Параметр регуляризации α является структурным параметром, который может быть определен, например, с помощью скользящего контроля. Скалярный параметр в семействе радиальных или полиномиальных базисных функций также может быть определен с помощью скользящего контроля.

Пример применения логистической регрессии для задачи классификации с двумя признаками приведен на рис. 4а–д.

Рассмотрим задачу классификации на K классов. Имеется обучающая выборка $(T, X) = \{t_n, \mathbf{x}_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^d$ – вектор признаков для объекта n , $t_n \in \{1, \dots, K\}$ – его метка класса. Задача состоит в определении метки класса t для нового объекта, представленного только своим вектором признаков \mathbf{x} .

Для решения данной задачи введем K линейных функций

$$f_k(\mathbf{x}, \mathbf{w}^k) = \sum_{j=1}^M w_j^k \phi_j(\mathbf{x}), \quad k = 1, \dots, K.$$

Здесь w^k – веса k -ой линейной функции, а ϕ – общие базисные функции. Решение о принадлежности объекта \mathbf{x} к классу t принимается по максимуму значений линейных функций f_1, \dots, f_K :

$$t(\mathbf{x}) = \arg \max_k f_k(\mathbf{x}, \mathbf{w}^k). \quad (17)$$

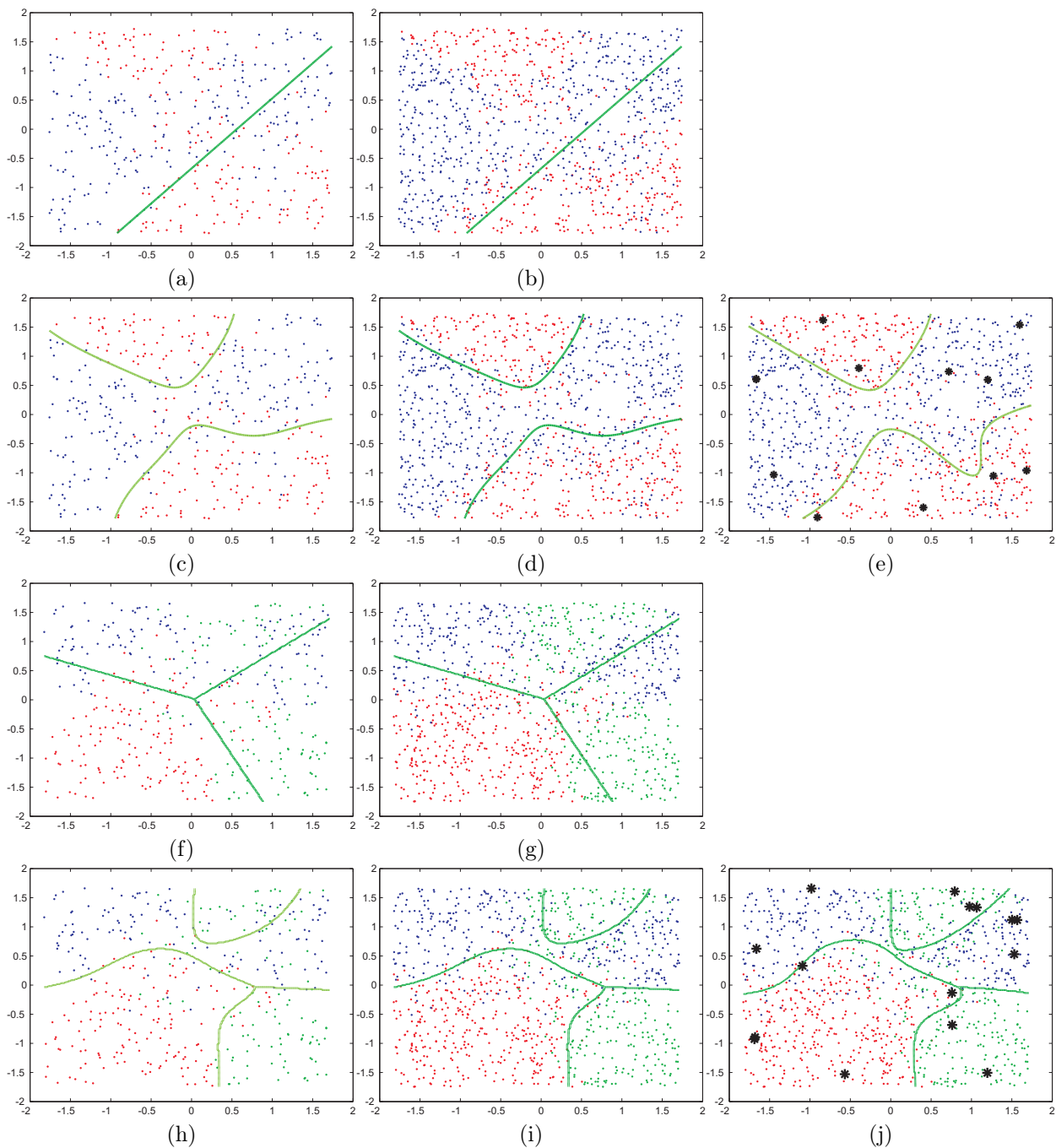


Рис. 4: Иллюстрация работы логистической и мультиномиальной регрессии, а также метода релевантных векторов. (a) – обучающая выборка с двумя классами, логистическая регрессия, исходные признаки, (b) – тестовая выборка с двумя классами, логистическая регрессия, исходные признаки, (c) – обучающая выборка с двумя классами, логистическая регрессия, радиальные базисные функции, (d) – тестовая выборка с двумя классами, логистическая регрессия, радиальные базисные функции, (e) – тестовая выборка с двумя классами, метод релевантных векторов, радиальные базисные функции, черными звездочками показаны релевантные центры радиальных базисных функций, (f) – обучающая выборка с тремя классами, мультиномиальная регрессия, исходные признаки, (g) – тестовая выборка с тремя классами, мультиномиальная регрессия, исходные признаки, (h) – обучающая выборка с тремя классами, мультиномиальная регрессия, радиальные базисные функции, (i) – тестовая выборка с тремя классами, мультиномиальная регрессия, радиальные базисные функции, (j) – тестовая выборка с тремя классами, многоклассовый метод релевантных векторов, радиальные базисные функции, черными звездочками показаны релевантные центры радиальных базисных функций.

По аналогии с логистической регрессией будем строить вероятностную модель, в которой функция соответствия между меткой класса t и набором значений линейных функций f_1, \dots, f_K выглядит как

$$p(t = k | \mathbf{x}, W) = p(t = k | f_1(\mathbf{x}, \mathbf{w}^1), \dots, f_K(\mathbf{x}, \mathbf{w}^K)) = \frac{\exp(f_k(\mathbf{x}, \mathbf{w}^k))}{\sum_{j=1}^K \exp(f_j(\mathbf{x}, \mathbf{w}^j))}.$$

Здесь через $W \in \mathbb{R}^{K \times M}$ обозначена матрица весов всех линейных функций. Введенная функция называется мультиномиальной функцией, а основанная на ней вероятностная модель — «мультиномиальной регрессией»:

$$p(T, W | X, \alpha) = p(T | X, W) p(W | \alpha) = p(W | \alpha) \prod_{n=1}^N p(t_n | \mathbf{x}_n, W) = p(W | \alpha) \prod_{n=1}^N \frac{\exp(f_{t_n}(\mathbf{x}_n, \mathbf{w}^{t_n}))}{\sum_{j=1}^K \exp(f_j(\mathbf{x}_n, \mathbf{w}^j))},$$

$$p(W | \alpha) = \prod_{k=1}^K \mathcal{N}(\mathbf{w}^k | \mathbf{0}, \alpha^{-1} I).$$

Здесь, как и в случае логистической регрессии, вводится гауссовское априорное распределение со скалярным параметром регуляризации α для предотвращения переобучения. Заметим, что модель мультиномиальной регрессии переходит в модель логистической регрессии для случая $K = 2$ и $f_2 = -f_1$.

Обучение в мультиномиальной регрессии (поиск весов W) по методу максимизации апостериорного распределения может быть осуществлено с помощью метода Ньютона решения задачи

$$\sum_{n=1}^N \left[f_{t_n}(\mathbf{x}_n, \mathbf{w}^{t_n}) - \log \left(\sum_{j=1}^K \exp(f_j(\mathbf{x}_n, \mathbf{w}^j)) \right) \right] - \frac{\alpha}{2} \sum_{j=1}^K \|\mathbf{w}^j\|^2 \rightarrow \max_W.$$

Параметр регуляризации α может быть настроен с помощью скользящего контроля. Пример применения мультиномиальной регрессии для задачи классификации с двумя признаками показан на рис. 4f-i.

Метод релевантных векторов для задачи классификации

Метод релевантных векторов (RVM) является надстройкой над вероятностной моделью логистической (для двух классов) и мультиномиальной (для многих классов) регрессии, позволяющей в процессе обучения отбирать информативные (релевантные) признаки (базисные функции).

Рассмотрим сначала метод RVM для случая двух классов. Вероятностная модель RVM выглядит следующим образом:

$$p(T, \mathbf{w} | X, \alpha) = p(T | X, \mathbf{w}) p(\mathbf{w} | \alpha) = p(\mathbf{w} | \alpha) \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}),$$

$$p(t | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-t f(\mathbf{x}, \mathbf{w}))}, \quad (18)$$

$$p(\mathbf{w} | \alpha) = \prod_{j=1}^M \mathcal{N}(w_j | 0, \alpha_j^{-1}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \text{diag}(\alpha)^{-1}).$$

Как легко заметить, модель (18) полностью повторяет модель логистической регрессии (15), в которой матрица ковариации в гауссовском априорном распределении заменена на диагональную с [различными элементами на диагонали](#). Таким образом, здесь вводится свой параметр регуляризации $\alpha_j \geq 0$ для каждого веса w_j . Обучение в модели (18) (поиск параметров α) производится с помощью метода максимального правдоподобия (обоснованности):

$$p(T | X, \alpha) = \int p(T | X, \mathbf{w}) p(\mathbf{w} | \alpha) d\mathbf{w} = \int Q(\mathbf{w}, \alpha) d\mathbf{w} \rightarrow \max_{\alpha}. \quad (19)$$

Здесь через $Q(\mathbf{w}, \alpha)$ обозначена подинтегральная функция. Интеграл в выражении для обоснованности (19) не вычисляется аналитически. Для приближенного оценивания данного интеграла в

RVM используется приближение Лапласа, при котором подынтегральная функция Q заменяется на свое ненормированное гауссовское приближение, интеграл от которого вычисляется аналитически:

$$p(T|X, \boldsymbol{\alpha}) \simeq \frac{Q(\mathbf{w}_{MP}, \boldsymbol{\alpha}) \sqrt{2\pi}^{-M}}{\sqrt{\det \nabla_{\mathbf{w}}^2 \log Q(\mathbf{w}_{MP}, \boldsymbol{\alpha})}}, \quad \mathbf{w}_{MP} = \arg \max_{\mathbf{w}} Q(\mathbf{w}, \boldsymbol{\alpha}). \quad (20)$$

Здесь через $\nabla_{\mathbf{w}}^2$ обозначен гессиан функции по \mathbf{w} . Считая \mathbf{w}_{MP} фиксированным и приравняв градиент полученного приближения обоснованности по $\boldsymbol{\alpha}$ к нулю, получаем следующую формулу пересчета для $\boldsymbol{\alpha}$:

$$\alpha_j^{new} = \frac{1 - \alpha_j^{old} \Sigma_{jj}}{w_{MP,j}^2}, \quad \Sigma = (-\nabla_{\mathbf{w}}^2 \log Q(\mathbf{w}_{MP}, \boldsymbol{\alpha}))^{-1}. \quad (21)$$

Таким образом, мы получаем итерационный алгоритм максимизации обоснованности (19). На каждой итерации сначала при текущем значении $\boldsymbol{\alpha}$ находится \mathbf{w}_{MP} и строится приближение для обоснованности (20), а затем при текущем \mathbf{w}_{MP} новое значение $\boldsymbol{\alpha}$ находится по формуле пересчета (21).

В процессе максимизации обоснованности (19) часть компонент вектора $\boldsymbol{\alpha}$ стремится к плюс бесконечности. В результате априорное распределение для соответствующих весов \mathbf{w} переходит в дельта-функцию с центром в нуле. Как следствие, часть базисных функций (признаков) исключается (признается нерелевантными) из решающего правила, т.к. их веса обнуляются.

Согласно последовательному байесовскому подходу, прогноз метки класса t_{new} для тестового объекта \mathbf{x}_{new} в RVM соответствует поиску распределения

$$p(t_{new}|\mathbf{x}_{new}, T, X, \boldsymbol{\alpha}) = \int p(t_{new}|\mathbf{x}_{new}, \mathbf{w})p(\mathbf{w}|T, X, \boldsymbol{\alpha})d\mathbf{w}.$$

По аналогии с обоснованностью (19) данный интеграл не вычисляется аналитически. Простейшей схемой аппроксимации здесь является замена апостериорного распределения $p(\mathbf{w}|T, X, \boldsymbol{\alpha})$ на дельта-функцию $\delta(\mathbf{w} - \mathbf{w}_{MP}(\boldsymbol{\alpha}))$, где $\mathbf{w}_{MP}(\boldsymbol{\alpha}) = \arg \max_{\mathbf{w}} Q(\mathbf{w}, \boldsymbol{\alpha})$. Тогда

$$p(t_{new}|\mathbf{x}_{new}, T, X, \boldsymbol{\alpha}) \simeq p(t_{new}|\mathbf{x}_{new}, \mathbf{w}_{MP}(\boldsymbol{\alpha})).$$

Более точной схемой аппроксимации для прогнозного распределения является использование приближения Лапласа. Пример применения метода релевантных векторов показан на рис. 4е.

Проведем модельный эксперимент. Обучающая выборка состоит из 100 объектов (по 50 объектов из каждого класса) и пятидесяти признаков. При этом только два признака (пятый и семнадцатый) являются информативными с точки зрения разделения двух классов (см. рис. 5а). Остальные признаки являются шумовыми. Тестовая выборка сгенерирована аналогичным образом и состоит из 1000 объектов (по 500 объектов из каждого класса), см. рис. 5б. При обучении логистической регрессии все веса в линейном решающем правиле оказываются отличными от нуля (см. рис. 5с). Ошибка на тестовой выборке при этом составляет 16%. При обучении метода релевантных векторов только три веса оказываются отличными от нуля (см. рис. 5д), в том числе веса для пятого и семнадцатого признака. В результате ошибка на тестовой выборке сокращается до 8%, что соответствует байесовскому уровню ошибки для данной задачи.

Рассмотрим теперь метод RVM для случая K классов. Для этого модифицируем вероятностную модель мультиномиальной регрессии следующим образом:

$$p(T, W, \boldsymbol{\alpha}|X) = p(T|X, W)p(W|\boldsymbol{\alpha}) = p(W|\boldsymbol{\alpha}) \prod_{n=1}^N p(t_n|\mathbf{x}_n, W),$$

$$p(t_n|\mathbf{x}_n, W) = \frac{\exp(f_{t_n}(\mathbf{x}_n, \mathbf{w}^{t_n}))}{\sum_{j=1}^K \exp(f(\mathbf{x}_n, \mathbf{w}^j))},$$

$$p(W|\boldsymbol{\alpha}) = \prod_{j,m=1}^{K,M} \mathcal{N}(w_m^j|0, \alpha_m^{-1}).$$

Заметим, что в этой модели скалярный параметр регуляризации α_j вводится для всех весов из j -го столбца матрицы W . В результате устремление α_j к плюс бесконечности соответствует обнулению

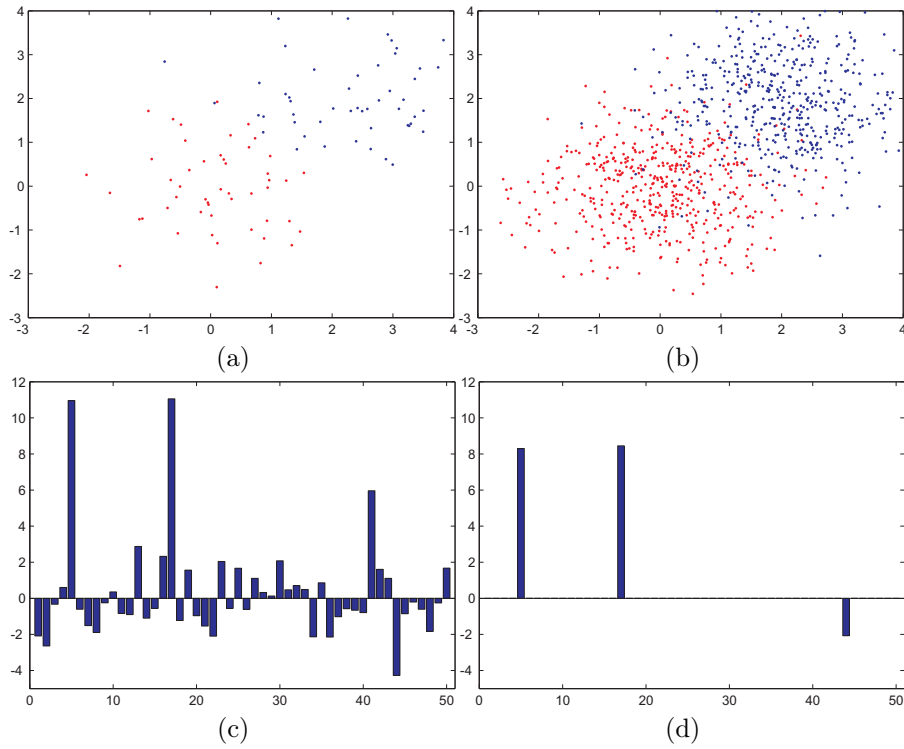


Рис. 5: Модельный эксперимент. (a) – два информативных признака в обучающей выборке, (b) – два информативных признака в тестовой выборке, (c) – веса обученной логистической регрессии, (d) – веса обученного метода релевантных векторов.

весов j -ой базисной функции сразу для всех линейных функций f_1, \dots, f_K , т.е. исключению данного признака из рассмотрения.

Схема обучения и принятия решения в модели RVM для многих классов полностью повторяет схему для двух классов. На этапе обучения вектор α настраивается путем максимизации обоснованности $p(T|X, \alpha)$, в которой используется приближение Лапласа. На этапе принятия решения для нового объекта в простейшем случае можно ограничиться вычислением весов $W_{MP}(\alpha)$ и применением решающего правила (17) для этих весов. Пример применения многоклассового метода релевантных векторов показан на рис. 4j.