

Обучение с учителем. Классификация. Дискриминантный анализ. Логистическая регрессия. Метод опорных векторов. Выбор модели с помощью кросс-валидации. Метод стохастического градиента.

Белкова Анна, Редкокош Кирилл, Лобанова Полина, гр. 21.М03-мм

11 декабря 2022 г.

Содержание

1	Классификация	2
1.1	Постановка задачи	2
1.2	Метрики качества классификации	2
1.2.1	Матрица ошибок	2
1.2.2	Accuracy	3
1.2.3	Precision, recall и F-мера	3
1.2.4	AUC-ROC и AUC-PR	4
1.3	Модификации датасета для выравнивания соотношения классов	4
1.3.1	Случайная наивная избыточная выборка	4
1.3.2	SMOTE	5
1.3.3	Tomek Links	7
2	Дискриминантный анализ	8
2.1	Байесовский классификатор	8
2.2	Линейный дискриминантный анализ	9
2.3	Канонические переменные	9
2.4	Значимость канонических переменных	10
2.5	Квадратичный дискриминантный анализ	10
2.6	Оценка параметров	10
2.7	Regularized Discriminant Analysis	11
2.8	Наивный байесовский классификатор	11
3	SVM. Метод опорных векторов.	11
3.1	SVM. Hard-margin SVM	11
3.2	SVM.Slack variables	13
3.3	Kernel trick	13
3.4	Сложность SVM	14
3.5	Мультиклассовый SVM	15
4	Логистическая регрессия	15
4.1	Логистическая регрессия. Подход через минимизацию функции потерь	15
4.2	Логистическая регрессия. Вероятностный подход	16
4.3	Линейная и логистическая регрессия	17
4.4	Логистическая регрессия. Регуляризация	17
4.5	Многоклассовая логистическая регрессия	18
4.6	Логистическая регрессия. Преимущества и недостатки	18

1 Классификация

1.1 Постановка задачи

Дано:

- X – матрица признаков, где $\mathbf{x}_i \in \mathbb{R}^p$ – i вектор этой матрицы (признаки i индивида).
- \mathcal{Y} – конечное множество номеров (имён, меток) классов, где $y_i \in \mathcal{Y}$, а \mathbf{y} – вектор меток классов для матрицы X .

Задача:

- По выборке $(X_{train}, \mathbf{y}_{train})$, построить классификатор $f: \mathbb{R}^p \rightarrow \mathcal{Y}$, который по выборке $(X_{test}, \mathbf{y}_{test})$ предскажет метку класса.
- Хотим, чтобы на f достигается минимальная ошибка классификации в некотором смысле.

На генеральном языке:

- $\xi \in \mathbb{R}^p$ – случайный вектор признаков.
- $\eta \in \mathcal{Y}$ – дискретная случайная величина, метка класса.
- $P(\xi, \eta)$ – их совместное распределение.

Дано:

Выборка $(X_{train}, \mathbf{y}_{train})$ – N реализаций случайного вектора (ξ, η) , по выборке необходимо построить классификатор

$$f: \mathbb{R}^p \rightarrow \mathcal{Y}.$$

Линейная модель классификации в общем случае: $f(\mathbf{x}_i, \mathbf{w}) = \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle$, где \mathbf{w} вектор весов признаков, а w_0 некоторый сдвиг, который позволяет нам не получить 0 значение.

Мы будем рассматривать классификатор в следующей форме, добавив вектор из единиц:

$$f(\mathbf{x}_i, \mathbf{w}) = \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle.$$

Отступ (margin): $M_i = \langle (\mathbf{w}, \mathbf{x}_i), \mathbf{x}_i \rangle y_i$ – "расстояние" между реальным и предсказанным значением. В случае отрицательного значения, считаем, что объект не принадлежит классу.

Функцию потерь $\mathcal{L}(M)$ – неотрицательная функция, характеризующая величину ошибки предсказания. Пороговая функция потерь: $[M(x_i) < 0]$.

Тогда задача классификации можно свести к минимизации функции потерь:

$$\mathcal{L}(M) \rightarrow \min$$

Далее будем рассматривать отступ, со знаком минус, как штраф за неверную классификацию.

1.2 Метрики качества классификации

Часто возникает необходимость в изучении различных аспектов качества уже обученного классификатора. Обсудим подробнее распространённые подходы к измерению качества моделей.

1.2.1 Матрица ошибок

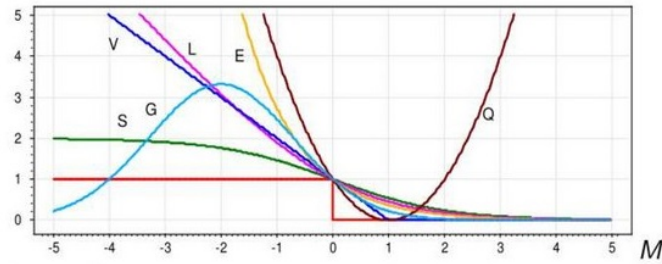
Перед переходом к самим метрикам необходимо ввести важную концепцию для описания этих метрик в терминах ошибок классификации – confusion matrix (матрица ошибок).

Допустим, что у нас есть два класса и алгоритм, предсказывающий принадлежность каждого объекта одному из классов, тогда матрица ошибок классификации будет выглядеть следующим образом:

	$y=1$	$y=0$
$\hat{y} = 1$	True Positive (TP)	False Positive (FP)
$\hat{y} = 0$	False Negative (FN)	True Negative (TN)

Где \hat{y} – это ответ алгоритма на объекте, а y – истинная метка класса на этом объекте. Таким образом, ошибки классификации бывают двух видов: False Negative (FN) и False Positive (FP).

Функции потерь $\mathcal{L}(M)$ в задачах классификации на два класса



$E(M) = e^{-M}$ — экспоненциальная (AdaBoost);
 $L(M) = \log_2(1 + e^{-M})$ — логарифмическая (LogitBoost);
 $G(M) = \exp(-cM(M + s))$ — гауссовская (BrownBoost);
 $Q(M) = (1 - M)^2$ — квадратичная;
 $S(M) = 2(1 + e^M)^{-1}$ — сигмоидная;
 $V(M) = (1 - M)_+$ — кусочно-линейная (SVM);

1.2.2 Accuracy

Интуитивно понятной, очевидной и почти неиспользуемой метрикой является ассурасу — доля правильных ответов алгоритма:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Эта метрика бесполезна в задачах с неравными классами, и это легко показать на примере.

Допустим, мы хотим оценить работу спам-фильтра почты. У нас есть 100 не-спам писем, 90 из которых наш классификатор определил верно (True Negative = 90, False Positive = 10), и 10 спам-писем, 5 из которых классификатор также определил верно (True Positive = 5, False Negative = 5). Тогда ассурасу:

$$accuracy = \frac{5 + 90}{5 + 90 + 10 + 5} = 86,4\%$$

Однако если мы просто будем предсказывать все письма как не-спам, то получим более высокую ассурасу:

$$accuracy = \frac{0 + 100}{0 + 100 + 0 + 10} = 90,9\%$$

При этом, наша модель совершенно не обладает никакой предсказательной силой, так как изначально мы хотели определять письма со спамом. Преодолеть это нам поможет переход с общей для всех классов метрики к отдельным показателям качества классов.

1.2.3 Precision, recall и F-мера

Для оценки качества работы алгоритма на каждом из классов по отдельности введем метрики precision (точность) и recall (полнота).

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

Precision можно интерпретировать как долю объектов, названных классификатором положительными и при этом действительно являющимися положительными, а recall показывает, какую долю объектов положительного класса из всех объектов положительного класса нашел алгоритм.

Именно введение precision не позволяет нам записывать все объекты в один класс, так как в этом случае мы получаем рост уровня False Positive. Recall демонстрирует способность алгоритма обнаруживать данный класс вообще, а precision — способность отличать этот класс от других классов.

Precision и recall не зависят, в отличие от accuracy, от соотношения классов и потому применимы в условиях несбалансированных выборок.

Существует несколько различных способов объединить precision и recall в агрегированный критерий качества. F-мера (в общем случае F_β) — среднее гармоническое precision и recall :

$$F_\beta = (1 + \beta^2) \cdot \frac{precision \cdot recall}{(\beta^2 \cdot precision) + recall}$$

β в данном случае определяет вес точности в метрике, и при $\beta = 1$ это среднее гармоническое (с множителем 2, чтобы в случае precision = 1 и recall = 1 иметь $F_1 = 1$) F-мера достигает максимума при полноте и точности, равными единице, и близка к нулю, если один из аргументов близок к нулю.

1.2.4 AUC-ROC и AUC-PR

Одним из способов оценить модель, является AUC-ROC (или ROC AUC) — площадь (Area Under Curve) под кривой ошибок (Receiver Operating Characteristic curve). Данная кривая представляет из себя линию от (0,0) до (1,1) в координатах True Positive Rate (TPR) и False Positive Rate (FPR):

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

TPR— это полнота, а FPR показывает, какую долю из объектов negative класса алгоритм предсказал неверно. В идеальном случае, когда классификатор не делает ошибок ($FPR = 0$, $TPR = 1$) мы получим площадь под кривой, равную единице; в противном случае, когда классификатор случайно выдает вероятности классов, AUC-ROC будет стремиться к 0.5, так как классификатор будет выдавать одинаковое количество TP и FP.

Каждая точка на графике соответствует выбору некоторого порога. Площадь под кривой в данном случае показывает качество алгоритма.

Precision и recall также используют для построения кривой и, аналогично AUC-ROC, находят площадь под ней:

1.3 Модификации датасета для выравнивания соотношения классов

Одним из распространенных способов решения проблемы несбалансированных данных является избыточная выборка. Чрезмерная выборка относится к различным методам, которые направлены на увеличение количества экземпляров из недопредставленного класса в наборе данных.

1.3.1 Случайная наивная избыточная выборка

Самый простой способ сделать это - случайным образом выбрать наблюдения из класса меньшинства и добавить их в набор данных, пока мы не достигнем баланса между большинством и классом меньшинства.

Одна проблема со случайной наивной избыточной выборкой заключается в том, что она просто дублирует уже существующие данные. Поэтому, хотя алгоритмы классификации подвергаются большему количеству наблюдений из класса меньшинства, они не узнают больше о том, как отличить наблюдения одного класса от другого. Новые данные не содержат больше информации о характеристиках класса, чем старые данные.

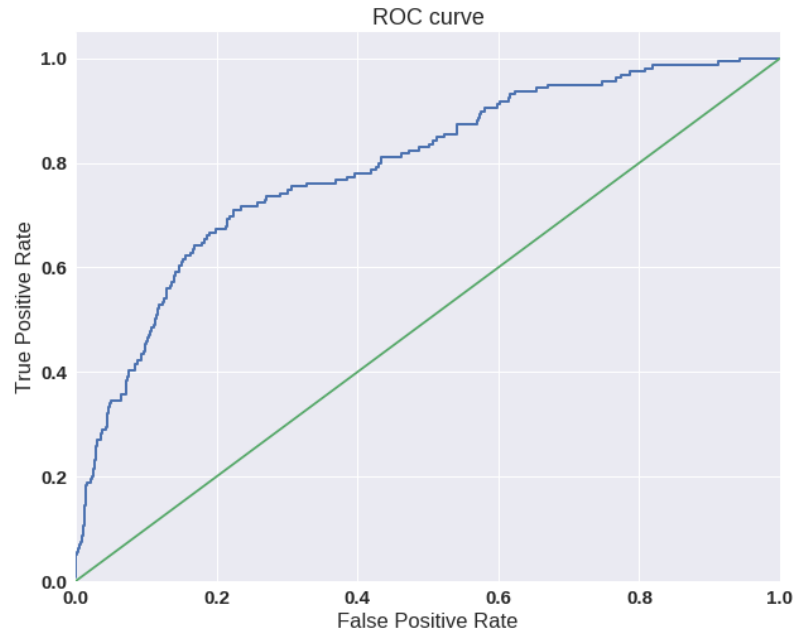


Рис. 1: ROC-кривая

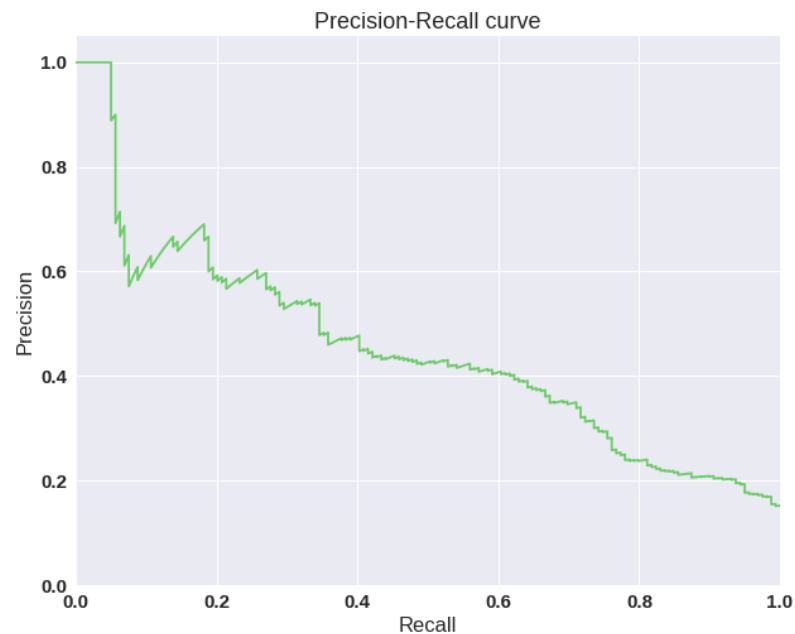


Рис. 2: PR-кривая

1.3.2 SMOTE

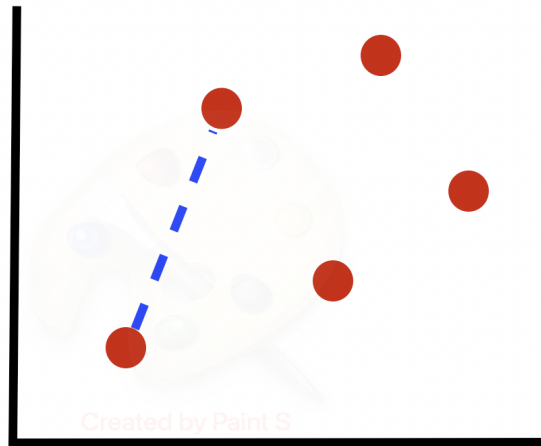
Метод увеличения числа примеров миноритарного класса (Synthetic Minority Over-sampling Technique, SMOTE) — это алгоритм предварительной обработки данных, используемый для устранения дисбаланса классов в наборе данных.

В общих чертах этот алгоритм можно описать следующим образом. Он находит разность между данным образцом и его ближайшим соседом. Эта разность умножается на случайное число в интер-

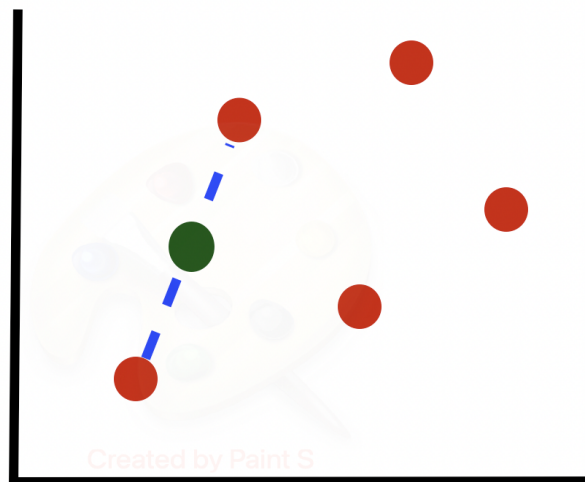
вале от 0 до 1. Полученное значение добавляется к данному образцу для формирования нового синтезированного образца в пространстве признаков. Подобные действия продолжаются со следующим ближайшим соседом, до заданного пользователем количества образцов.

Проиллюстрируем алгоритм более подробно. Предположим, у нас есть несбалансированный набор данных (индивидов одного класса гораздо больше, чем другого).

Берем индивида и вычисляем k -ближайших соседей. Затем выбираем случайного ближайшего соседа из k -ближайших соседей.



Вычисляем разность между двумя точками и умножаем ее на случайное число от 0 до 1. Получаем синтезированный образец вдоль линии между двумя точками.

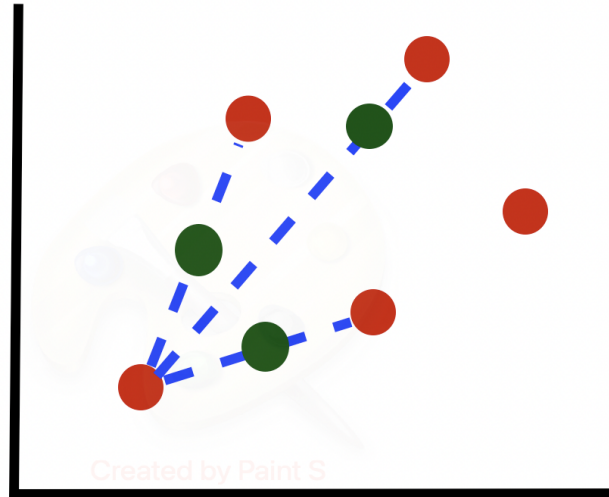


Выбираем m соседей для этого индивида и повторяем процедуру дублирования. m подбирается из соотношения классов.

Для каждого из исходных индивидов повторяем весь алгоритм для достижения равного количества индивидов в классах.

SMOTE Расширения

Как и в большинстве алгоритмов, есть несколько расширений SMOTE. Они нацелены на улучшение SMOTE путем добавления его функциональности или уменьшения его слабых сторон. Примеры расширений SMOTE, которые можно найти в `imblearn`, включают:



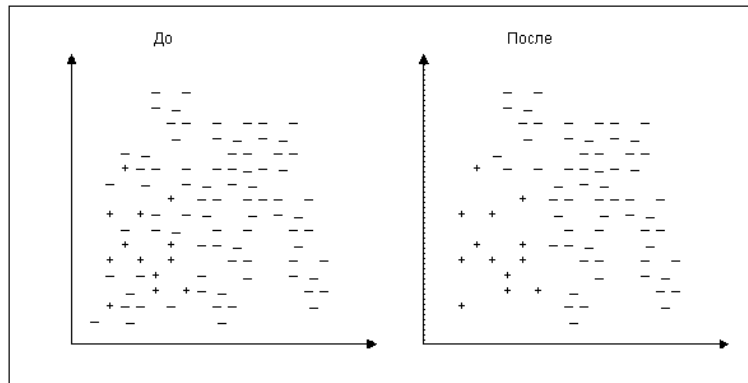
- BorderlineSMOTE: Вместо избыточной выборки между всеми наблюдениями меньшинств, BorderlineSMOTE стремится увеличить количество наблюдений меньшинств, которые граничат с наблюдениями большинства. Цель здесь - дать классификатору возможность более четко различать эти пограничные наблюдения.
- SVMSMOTE: SVMSMOTE, как следует из его названия, использует алгоритм машины опорных векторов для генерации новых наблюдений меньшинства вблизи границы между классами большинства и меньшинства.

1.3.3 Tomek Links

Пусть индивиды E_i и E_j принадлежат к различным классам, $d(E_i, E_j)$ – расстояние между указанными примерами. Пара (E_i, E_j) называется связью Томека, если не найдется ни одного примера E_l такого, что будет справедлива совокупность неравенств:

$$\begin{cases} d(E_i, E_l) < d(E_i, E_j), \\ d(E_j, E_l) < d(E_i, E_j) \end{cases}$$

Согласно данному подходу, все индивиды из большей группы, входящие в связи Томека, должны быть удалены из набора данных. Этот способ хорошо удаляет записи, которые можно рассматривать в качестве «зашумляющих». Далее визуально показан набор данных в двумерном пространстве признаков до и после применения поиска связей Томека.



2 Дискриминантный анализ

Суть дискриминантного анализа заключается в том, чтобы смоделировать распределение X в каждом из классов отдельно, а затем использовать теорему Байеса, чтобы получить $P(Y = i | X = x)$

2.1 Байесовский классификатор

В качестве меры ошибки предсказания введем функцию потерь. Рассмотрим матрицу \mathbf{L} размера $K \times K$, где $K = \text{card}(\mathcal{Y})$. На диагонали \mathbf{L} стоят нули, а $\mathbf{L}(i, j) = \lambda_{ij}$ – цена ошибки отнесения элемента класса Y_i к классу Y_j . Часто используется 0-1 функция потерь, где каждая ошибка оценивается единицей.

Математическое ожидание функции потерь (средний риск):

$$R(a) = \mathbb{E}(\mathbf{L}(\eta, a(\xi))) = \mathbb{E}_\xi \sum_{k=1}^K L(Y_k, a(\xi)) P(Y_k | \xi).$$

Отсюда получаем функцию классификации:

$$f(x) = \arg \min_{Y \in \mathcal{Y}} \sum_{k=1}^K L(Y_k, Y) P(Y_k | \xi = x).$$

Если подставим сюда 0-1 функцию потерь, получим

$$f(x) = \arg \min_{Y \in \mathcal{Y}} 1 - P(Y | \xi = x).$$

Или, что то же самое

$$f(x) = \arg \max_{Y \in \mathcal{Y}} P(Y | \xi = x) = \arg \max_{Y \in \mathcal{Y}} P(Y) P(\xi | \eta = Y).$$

Это решение называется байесовским классификатором, а такой подход – принципом максимума апостериорной вероятности.

Для построения байесовского классификатора, нам необходимо знать апостериорные вероятности $P(Y | \xi = x)$.

Обозначим $p_i(x) = P(\xi = x | \eta = Y_i)$ условные плотности классов, $\pi_i = P(\eta = Y_i)$ – априорные вероятности, $\sum_{i=1}^K \pi_i = 1$. По теореме Байеса получим:

$$P(Y = i | X = x) = \frac{p_i(x) \pi_i}{\sum_{i=1}^K p_i(x) \pi_i}.$$

Поэтому в качестве классифицирующих функций берут

$$f_i(x) = \frac{p_i(x) \pi_i}{\sum_{j=1}^k p_j(x) \pi_j}.$$

Так как знаменатель у всех f_i одинаковый, его можно отбросить, и итоговые классифицирующие функции будут выглядеть как $f_i(x) = P(x | C_i) \pi_i = p_i(x) \pi_i$.

Возникает вопрос: откуда брать априорные вероятности?

1. Равномерно, $\forall i \in 1 : k \pi_i = 1 / k$.
2. По соотношениям в обучающей выборке: $\pi_i = n_i / \sum_{j=1}^k n_j$.
3. На основе другой дополнительной информации о данных (результаты предыдущих исследований, etc.)

Построенный метод классификации $\text{predict}(x) = \arg \max_i \pi_i p_i(x)$ минимизирует среднюю апостериорную ошибку:

$$\sum_{i=1}^k \pi_i P(\text{predict}(x) \neq i | Y_i).$$

2.2 Линейный дискриминантный анализ

Модель: ξ — дискретная с.в., принимающая значения $\{Y_i\}_{i=1}^k$, $\mathcal{P}(\eta | \xi = Y_i) = \mathcal{N}(\mu_i, \Sigma)$. (Предполагаем, что классы имеют нормальное распределение с одинаковой ковариационной матрицей)

Тогда плотность в точке x :

$$p_i(x) = p(x | \xi = Y_i) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x^T - \mu_i)\Sigma^{-1}(x - \mu_i)\right),$$

и классифицирующая функция $f_i(x) = \pi_i p(x | \xi = Y_i)$, где π_i — априорная вероятность наблюдения попасть в i -ю группу. Для упрощения вычислений можно переписать классифицирующую функцию через возрастающее монотонное преобразование как

$$g_i(x) = \log f_i(x) = \log \pi_i - \frac{1}{2} \log |\Sigma| - \frac{1}{2}(x^T - \mu_i)\Sigma^{-1}(x - \mu_i).$$

Сократив часть, не зависящую от номера класса, получаем линейные классифицирующие функции:

$$h_i(x) = -\frac{1}{2}\mu_i^T \Sigma^{-1} \mu_i + \mu_i^T \Sigma^{-1} x + \log \pi_i.$$

2.3 Канонические переменные

Задача: найти линейное преобразование $\mathbf{Z} = A^T \mathbf{X}$, в результате которого получаются признаки наилучшим образом разделяющие группы. Хотелось бы, чтобы эти признаки оказались ортогональны. Далее опишем эту задачу более формально.

Вычислим внутриклассовую ковариационную матрицу:

$$\mathbf{E} = \frac{1}{n - K} \sum_{i=1}^K \sum_{j: y_j = Y_i} (x_j - \hat{\mu}_i)^T (x_j - \hat{\mu}_i)$$

Вычисляем межклассовую ковариационную матрицу (с точностью до коэффициента):

$$\mathbf{H} = \sum_{i=1}^K n_i (\hat{\mu}_i - \hat{\mu})^T (\hat{\mu}_i - \hat{\mu}).$$

Пусть $\zeta = A\xi$ — новый признак, тогда распределение $P(\zeta | \eta = Y_k) = \mathcal{N}_p(A^T \mu_k, A^T \Sigma_k A)$.

На выборочном языке новые признаки $\mathbf{Z} = A^T \mathbf{X}$. Выборочная ковариационная матрица (с точностью до коэффициента) новых признаков имеет вид:

$$A^T \mathbf{T} A = A^T (\mathbf{E} + \mathbf{H}) A = A^T \mathbf{E} A + A^T \mathbf{H} A,$$

где \mathbf{T} — total covariance matrix, первое слагаемое — оценка внутригрупповых отклонений, а второе — оценка межгрупповых отклонений. Воспользовавшись критерием Фишера перейдем к обобщенной задаче на собственные числа и собственные вектора:

$$\frac{A^T \mathbf{H} A}{A^T \mathbf{E} A} \rightarrow \max_A.$$

Пусть $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ — собственные числа матрицы $\mathbf{E}^{-1} \mathbf{H}$, а A_1, \dots, A_d — соответствующие им собственные вектора. Тогда максимум выше равен λ_1 и достигается на A_1 . При этом $A_i^T \mathbf{E} A_j = 0$. Далее

$$\max_{A, A \perp A_1} \frac{A^T \mathbf{H} A}{A^T \mathbf{E} A} = \lambda_2,$$

достигается на A_2 и так далее.

Вектора A_i называют каноническими коэффициентами, а новые признаки Z_i — каноническими переменными, Z_i ортогональны.

2.4 Значимость канонических переменных

Возникает вопрос: сколько канонических переменных нам окажется достаточно взять? Другими словами, нужно проверить гипотезу:

$$H_0 : A_i, i = \ell, \dots, d \text{ не описывают отличия.}$$

Введем статистику Λ — *prime* (Wilks' Lambda):

$$\Lambda_\ell^p = \prod_{i=\ell}^d \frac{1}{1 + \lambda_i}.$$

Тогда гипотезу выше можно переформулировать так

$$H_0 : \Lambda_\ell^p = 1 \Leftrightarrow \lambda_\ell = \dots = \lambda_d = 0 \Leftrightarrow \text{rank} \mathbf{B} = \ell - 1.$$

Критерий:

$$t = \Lambda_\ell^p \sim \Lambda_{\nu_{\mathbf{B}} + (\ell-1), \nu_{\mathbf{W}} - (\ell-1)}.$$

Другими вариантами статистиками для проверки гипотезы могут являться:

- Roy's greatest root

$$r_1^2 = \frac{\lambda_1}{1 + \lambda_1};$$

- Pillai's trace

$$V = \text{trace}(\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1});$$

- Hotelling-Lawley trace

$$V = \text{trace}(\mathbf{H}\mathbf{E}^{-1}).$$

2.5 Квадратичный дискриминантный анализ

Модель: ξ — дискретная с.в., принимающая значения $\{Y_i\}_{i=1}^k$, $\mathcal{P}(\eta \mid \xi = Y_i) = \mathcal{N}(\mu_i, \Sigma_i)$. (Предполагаем, что каждый класс имеет многомерное нормальное распределение с различными ковариационными матрицами)

Тогда плотность в точке x :

$$p(x \mid \xi = Y_i) = \frac{1}{(2\pi)^{p/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (x^T - \mu_i) \Sigma_i^{-1} (x - \mu_i) \right),$$

и классифицирующая функция $f_i(x) = \pi_i p(x \mid \xi = Y_i)$. Применяем возрастающее монотонное преобразование и оставляем в классифицирующей функции только члены, отличающиеся в разных группах:

$$g_i(x) = \log f_i(x) = \log \pi_i - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (x^T - \mu_i) \Sigma_i^{-1} (x - \mu_i),$$

получаем квадратично зависящую от x классифицирующую функцию.

2.6 Оценка параметров

На практике параметры распределений классов нам не известны, поэтому предлагается использовать следующие оценки максимального правдоподобия параметров нормальных плотностей классов.

- Среднее $\bar{\mu}_i = \frac{1}{n_i} \sum_{j: y_j = Y_i} x_j$,
- Ковариационная матрица класса $\hat{\Sigma}_i = \frac{1}{n_i - 1} \sum_{j: y_j = Y_i} (x_j - \bar{\mu}_i)^T (x_j - \bar{\mu}_i)$,
- Pooled ковариационная матрица $\hat{\Sigma} = \sum_{j=1}^K \frac{n_i - 1}{n - K} \hat{\Sigma}_i$.

2.7 Regularized Discriminant Analysis

Оценка ковариационной матрицы $\hat{\Sigma}_i$ может оказаться выраженной или плохо обусловленной. Опишем компромис между LDA и QDA, а так же борьбу с мультиколлинеарностью.

- Regularized Discriminant Analysis. Рассматривается матрица $\hat{\Sigma}_i(\alpha) = \alpha \hat{\Sigma}_i + (1 - \alpha) \hat{\Sigma}$, где $\hat{\Sigma}$ – pooled ковариационная матрица. Здесь $\alpha \in [0, 1]$ порождает континуум моделей между LDA и QDA, выбирается скользящим контролем.
- Дополнительно к предыдущему методу можно похожим образом модифицировать pooled ковариационную матрицу и рассматривать $\hat{\Sigma}(\gamma) = \gamma \hat{\Sigma} + (1 - \gamma) \sigma^2 \mathbf{I}_p$, где γ определяет вид ковариационной матрицы и выбирается скользящим контролем.

2.8 Наивный байесовский классификатор

Предположим, что признаки независимы внутри групп и имеют нормальное распределение:

$$p_i(x) = \prod_{j=1}^p p_{ij}(x_j), \quad p_{ij}(x_j) = \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-\frac{(x_j - \mu_{ij})^2}{2\sigma_{ij}^2}}.$$

Отсюда классифицирующую функцию можно представить в виде:

$$\delta_i(x) = -\frac{1}{2} \sum_{j=1}^p \frac{(x_j - \mu_{ij})^2}{2\sigma_{ij}^2} + \log(\pi_i).$$

Аналогично подходам выше, можно подбирать ковариационную матрицу скользящим контролем в виде:

$$\hat{\Sigma}_i(\alpha) = \alpha \hat{\Sigma}_i + (1 - \alpha) \text{diag}(\sigma_{i1}^2, \dots, \sigma_{ip}^2), \alpha \in [0, 1].$$

Такой подход может быть полезен, когда признаков очень много и оценивать плотности классов оказывается сложно. Плотности p_{ki} можно оценивать по отдельности, а если признак дискретный, для этого можно использовать гистограмму.

Не смотря на такое оптимистичное предположение, наивный байесовский классификатор часто превосходит более сложные методы.

3 SVM. Метод опорных векторов.

Входные данные: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, $\mathbf{x}_i \in \mathbb{R}^p$, $y_i \in \{-1, 1\}$; Задача: построение классифицирующего правила $f: \mathbb{R}^p \rightarrow \{-1, 1\}$.

Предположим, что присутствует линейная разделимость, т.е. существует гиперплоскость (определяемая уравнением $\mathbf{x}^T \mathbf{w} - w_0 = 0$ ($\mathbf{x}, \mathbf{w} \in \mathbb{R}^p$; $w_0 \in \mathbb{R}$).

Как и прежде у нас есть линейный классификатор:

$$f(\mathbf{x}_i) = \text{sign}(\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0),$$

где $\mathbf{x}, \mathbf{w} \in \mathbb{R}^p$; $w_0 \in \mathbb{R}$. И кусочно линейна функция потерь, где $M_i = -\langle (\mathbf{w}, \mathbf{x}_i), \mathbf{x}_i \rangle y_i$:

$$L(M_i) = \max\{0, 1 + M_i\}$$

3.1 SVM. Hard-margin SVM

Предположим, что присутствует линейная разделимость, т.е. существует гиперплоскость (определяемая уравнением $\mathbf{x}^T \mathbf{w} - w_0 = 0$), такая, что точки, соответствующие разным классам лежат в различных полу-пространствах относительно гиперплоскости.

Факт принадлежности наблюдений из разных классов разным полупространствам можно (возможно, изменив знаки β, β_0) описать уравнениями:

$$\begin{cases} \mathbf{x}^T \mathbf{w} - w_0 < 0 & y_i = -1 \\ \mathbf{x}^T \mathbf{w} - w_0 > 0 & y_i = 1 \end{cases} \Leftrightarrow (\mathbf{x}^T \mathbf{w} - w_0) y_i > 0$$

В таком случае, классифицирующим правилом разумно принять $f(x) = \text{sign}(x^T w - w_0)$

В линейно разделимых данных может существовать более одной гиперплоскости, разделяющей данные. Введём критерий оптимальности: максимальное расстояние между двумя гиперплоскостями, параллельных данной и симметрично расположенных относительно неё, при котором между ними не находится ни одна из точек; это расстояние будем называть зазором (margin).

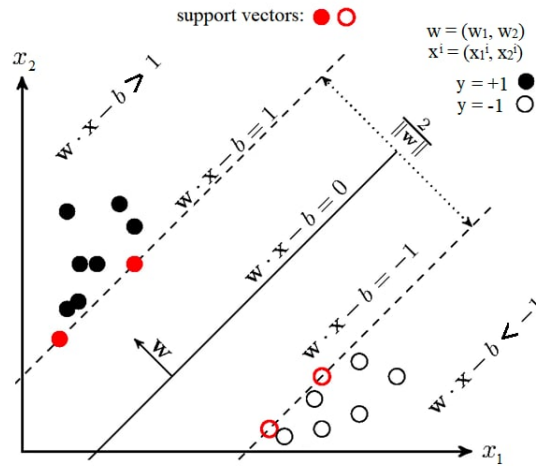
Для каждой из двух параллельных гиперплоскостей будет принадлежать некоторое количество точек из соответствующего класса (иначе, так как количество точек в выборке конечно, то расстояние между гиперплоскостями можно увеличить, сместив гиперплоскость, которой не принадлежит ни одной точки); точки, которые принадлежат одной из гиперплоскостей — будем называть опорными векторами.

С точностью до нормировки вектора w эта пара гиперплоскостей может быть описана парой уравнений:

$$\begin{cases} x^T w - w_0 = -1 \\ x^T w - w_0 = 1 \end{cases}$$

а расстояние между ними составит $\frac{2}{\|w\|}$ (см. рисунок) Принадлежность точек обучающей выборки полу-пространства описывается уравнениями

$$\begin{cases} x_i^T \beta - \beta_0 \leq 1 & y_i = -1 \\ x_i^T \beta - \beta_0 \geq 1 & y_i = 1 \end{cases} \Leftrightarrow (x_i^T \beta - \beta_0) y_i \geq 1$$



Тогда задачу можно свести к задаче квадратичного программирования с линейными ограничениями:

$$\begin{cases} \frac{1}{2} \langle w, w \rangle \rightarrow \min_w \\ \langle x_i, w \rangle - w_0 \geq 1 \end{cases}$$

Воспользуемся методом множителей Лагранжа:

$$\begin{cases} \inf_{w, w_0} \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^n \alpha_i [y_i (x_i^T w + w_0) - 1] \rightarrow \max_{\alpha_i} \\ \alpha_i \geq 0, \forall i \\ y_i (x_i, w) - w_0 \geq 1 \end{cases}$$

Так как оптимизируемая функция гладкая, можно воспользоваться необходимыми условиями экстремума

$$\frac{\partial}{\partial w} : w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial}{\partial w_0} : 0 = \sum_{i=1}^n \alpha_i y_i$$

Двойственная задача Вольфа:

$$\begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \rightarrow \max_{\alpha_i} \\ \alpha_i \geq 0, \forall i \\ y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0) \geq 1 \\ \mathbf{w}_0 = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \end{cases}$$

В точке оптимума выполнены условия Каруша-Куна-Такера, в частности:

$$\alpha_i [1 - y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0)] = 0 \forall i$$

Т.е. либо

- $\alpha_i = 0$ т.е. наблюдение не влияет на \mathbf{w} , w_0
- $\alpha_i > 0 \Rightarrow y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0) = 1$ – такое наблюдение будем называть опорным вектором

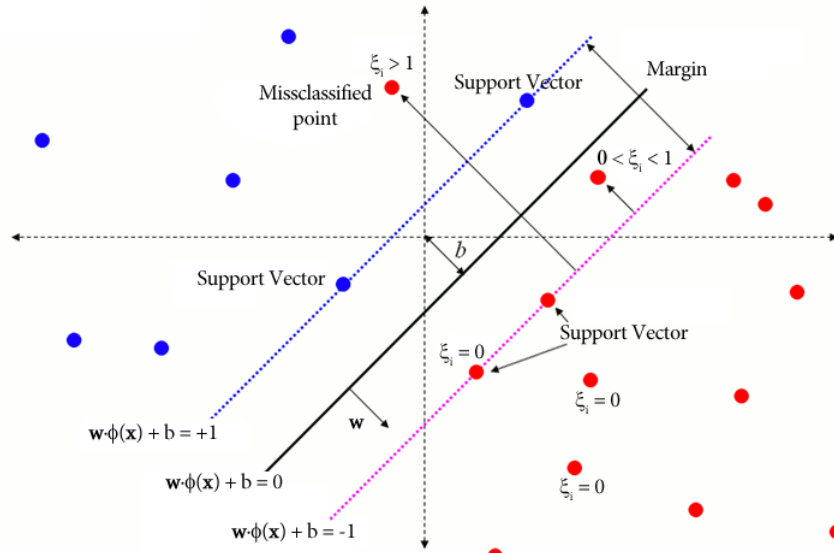
3.2 SVM.Slack variables

Поскольку в случае линейно неразделимой выборки по определению любой линейный классификатор будет ошибаться, условие $(\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0) y_i \geq 1$ не может быть выполнено для всех i . Введём ошибки $\xi \geq 0$ алгоритма и штрафы за эти ошибки в минимизируемую функцию следующим образом:

$$\begin{cases} \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C \sum_{i=1}^n \xi_i \rightarrow \min_{\mathbf{w}, w_0} \\ (\langle \mathbf{x}_i, \mathbf{w} \rangle - w_0) y_i \geq 1 - \xi_i \end{cases},$$

где C задает размер штрафа за ошибки.

Мы опять получили задачу линейного программирования.



3.3 Kernel trick

Чтобы применять SVN в нелинейном случае, строилось спрямляющее пространство. В основе этого лежит очень простая и очень красивая идея: если в каком-то исходном пространстве признаков классы не являются линейно разделимыми, то может быть можно отобразить это пространство признаков в какое-то новое, в котором классы уже будут линейно разделимы.

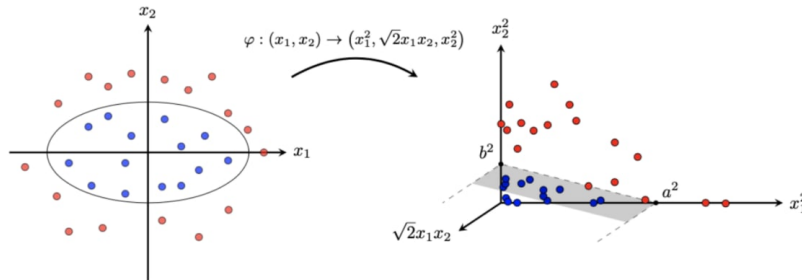
Пусть $\phi(\mathbf{x}_i)$ - спрямляющее отображение. Тогда, можем записать SVM в спрямляющем отображении, используя следующее скалярное произведение:

$$\mathbf{x}_i \rightarrow \phi(\mathbf{x}), \mathbf{w} \rightarrow \phi(\mathbf{w}), \langle \mathbf{w}, \mathbf{x}_i \rangle \rightarrow \langle \phi(\mathbf{w}), \phi(\mathbf{x}_i) \rangle$$

Чтобы получить нелинейную разделимость в исходном пространстве - зададим скалярное произведение следующего вида:

$$K(\mathbf{w}, \mathbf{x}_i) = \langle \phi(\mathbf{w}), \phi(\mathbf{x}_i) \rangle$$

K - симметричная нелинейная функция



Наиболее часто используемые ядра

- Линейное ядро:

$$K(\mathbf{w}, \mathbf{x}_i) = \langle \mathbf{w}, \mathbf{x}_i \rangle$$

- Полиномиальное ядро:

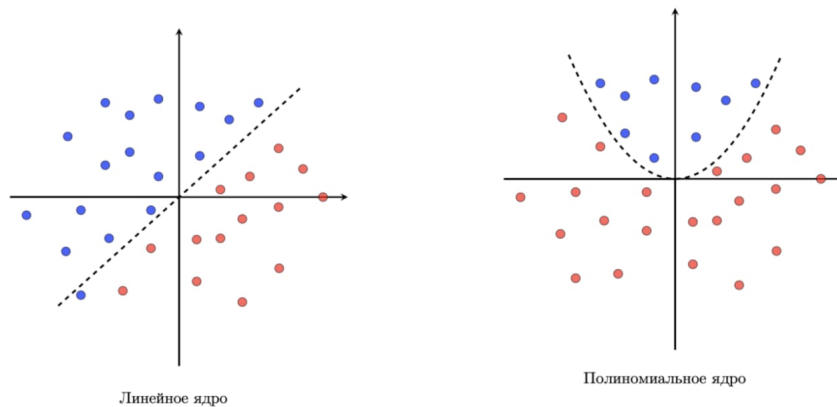
$$K(\mathbf{w}, \mathbf{x}_i) = (\langle \mathbf{w}, \mathbf{x}_i \rangle + r)^d$$

- Радиальное ядро:

$$K(\mathbf{w}, \mathbf{x}_i) = \exp(-\gamma ||(\mathbf{w} - \mathbf{x}_i)||^2)$$

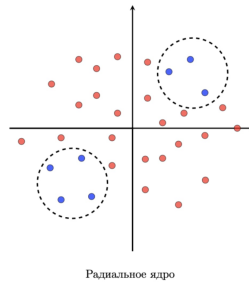
- Сигмовидное ядро.

Примеры:



3.4 Сложность SVM

Так как для SVM нужно решать задачу квадратичного программирования, то сложность варьируется между $O(p \cdot n^2)$ и $O(p \cdot n^3)$, где p - количество признаков, n - количество индивидов, зависимости от набора данных.



Радiallyе ядро

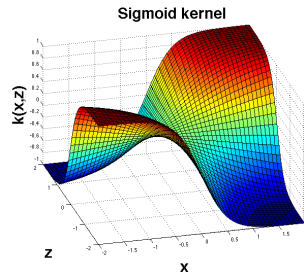


Рис. 3: Сигмовидное ядро

3.5 Мультиклассовый SVM

- Сравнение "один со многими". Строим N классифицирующих правил $f_i(\mathbf{x})$, кодирующих принадлежность к i -му классу за 1, -1 иначе. В качестве решающего правила используется

$$f(\mathbf{x}) = \underset{i}{\operatorname{argmax}} f_i(\mathbf{x})$$

- Сравнение "каждый с каждым". Строим $\frac{N(N-1)}{2}$ классифицирующих правил, производящих классификацию для каждой возможной пары классов. Обозначим за N_i количество сравнений, в которых элемент x был классифицирован как принадлежащий к i -му классу. В качестве решающего правила используется

$$f(\mathbf{x}) = \underset{i}{\operatorname{argmax}} N_i$$

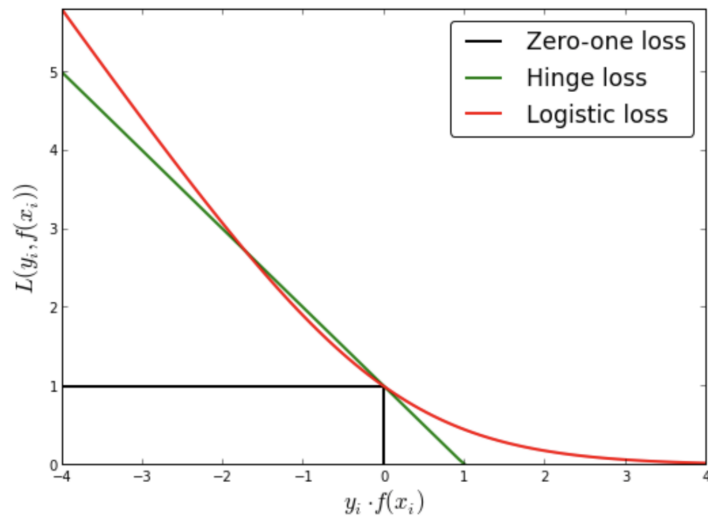
4 Логистическая регрессия

4.1 Логистическая регрессия. Подход через минимизацию функции потерь

Линейная модель классификации:

- $f(x, \theta) = \operatorname{sign}\langle \theta, x \rangle$, $x, \theta \in \mathbb{R}^p$
- $M = \langle \theta, x \rangle y$ — отступ.

В качестве аппроксимации пороговой функции потерь берется логарифмическая функция потерь $L(M) = \log(1 + e^{-M})$.



Задача 1. $Q(X_n, \theta) = \sum_{i=1}^n \log(1 + \exp(-y_i \langle \theta, x_i \rangle)) \rightarrow \min_{\theta}$

Методы решения задачи минимизации:

- метод стохастического градиента
- метод Ньютона-Рафсона

4.2 Логистическая регрессия. Вероятностный подход

$P(y = 1|x, \theta) = \sigma_{\theta}(M) = \frac{1}{1 + e^{-\langle x, \theta \rangle y}}$ — сигмоидная функция.

Свойства $\sigma(z)$:

- $\sigma(z) \in [0, 1]$, задана на $(-\infty, +\infty)$
- $\sigma(z) \rightarrow 1, z \rightarrow +\infty; \sigma(z) \rightarrow 0, z \rightarrow -\infty$
- $\sigma(z) + \sigma(-z) = 1$
- $\sigma'(z) = \sigma(z)\sigma(-z)$

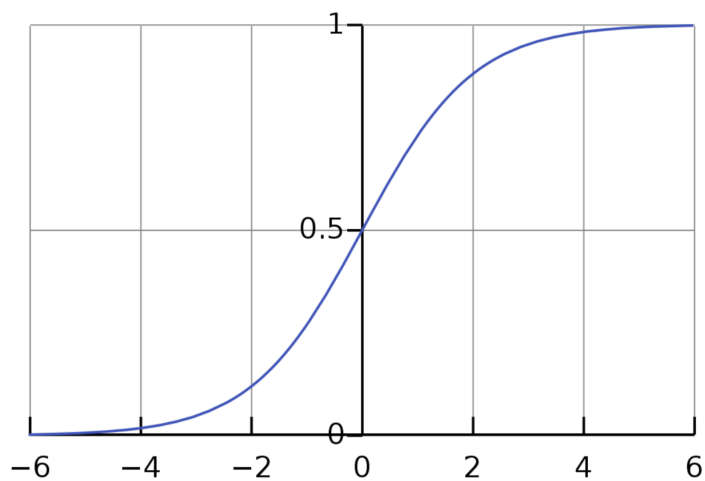


Рис. 4: Сигмоидная функция

Пусть $Y = \{0, 1\}$.

- $P(y_i = 1|x; \theta) = \sigma_\theta(x)$
- $P(y_i = 0|x; \theta) = 1 - \sigma_\theta(x)$

Тогда $P(y|x; \theta) = (\sigma_\theta(x))^y (1 - \sigma_\theta(x))^{1-y}$.

Функция правдоподобия:

$$\begin{aligned} Q(X_n, \theta) &= -\log L(\theta) = -\log \prod_{i=1}^n (\sigma_\theta(x_i))^{y_i} (1 - \sigma_\theta(x_i))^{1-y_i} = \\ &= -\sum_{i=1}^n [y_i \log(\sigma_\theta(x_i)) + (1 - y_i) \log(1 - \sigma_\theta(x_i))] \rightarrow \min_{\theta} \end{aligned}$$

4.3 Линейная и логистическая регрессия

Существуют примеры данных, для которых логистическая регрессия показывает лучшие результаты, чем линейная.

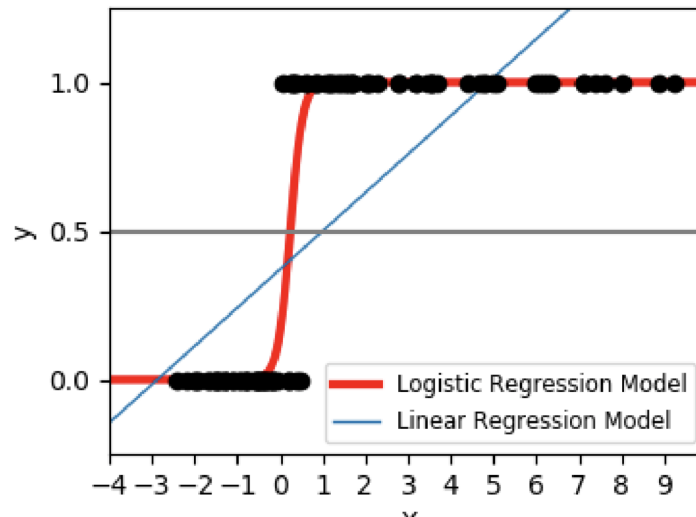


Рис. 5: Линейная и логистическая регрессия

4.4 Логистическая регрессия. Регуляризация

$$Q(\theta) = -\sum_{i=1}^n [y_i \log(\sigma_\theta(x_i)) + (1 - y_i) \log(1 - \sigma_\theta(x_i))]$$

Регуляризация в логистической регрессии:

- **L2:** $Q_\tau(\theta) = Q(\theta) + \frac{\tau}{2} \sum_{j=1}^p \theta_j^2 \rightarrow \min_{\theta}$
- **L1:** $Q_\tau(\theta) = Q(\theta) + \tau \sum_{j=1}^p |\theta_j| \rightarrow \min_{\theta}$

Параметр τ можно подбирать с помощью кросс-валидации.

Методы решения задачи минимизации:

- метод стохастического градиента
- метод Ньютона-Рафсона.

4.5 Многоклассовая логистическая регрессия

Линейный классификатор при произвольном числе классов $Y = \{1, \dots, K\}$:

$$\hat{f}(x, \theta) = \arg \max_{y \in Y} \langle \theta_y, x \rangle, \quad x, \theta_y \in \mathbb{R}^p$$

Вероятность того, что объект x относится к классу i :

$$P(y = i|x; \theta) = \frac{\exp \langle \theta_i, x \rangle}{\sum_{z \in Y} \exp \langle \theta_z, x \rangle} = \frac{e^{\theta_i^T x}}{\sum_{k=1}^K e^{\theta_k^T x}}$$

Задача:

$$Q(X_n, \theta) = - \sum_{i=1}^n \log P(y_i|x_i; \theta) \rightarrow \min_{\theta}$$

4.6 Логистическая регрессия. Преимущества и недостатки

Плюсы:

1. Позволяет оценить вероятности принадлежности объектов к классу
2. Достаточно быстро работает при больших объемах выборки
3. Применима в случае отсутствия линейной разделимости, если на вход подать полиномиальные признаки

Минусы:

1. Плохо работает в задачах, в которых зависимость сложная, нелинейная