

Support Vector Machines

Д. Корчемкин, В. Агеев
622 группа

26 ноября 2017 г.

1 SVM

Будем рассматривать задачу классификации в рамках обучения с учителем.

Имеется выборка $\{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^p$, $y_i \in \{-1, 1\}$; задачей является построение классифицирующего правила $f: \mathbb{R}^p \rightarrow \{-1, 1\}$.

Отметим, что никаких дополнительных ограничений на распределение не требуется.

1.1 Hard-margin SVM

Предположим, что присутствует линейная разделимость, т.е. существует гиперплоскость (определяемая уравнением $x^\top \beta - \beta_0 = 0$ ($x, \beta \in \mathbb{R}^p$; $\beta_0 \in \mathbb{R}$), такая, что точки, соответствующие разным классам лежат в различных полупространствах относительно гиперплоскости.

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

$$\begin{cases} x_i^\top \beta - \beta_0 < 0 & y_i = -1 \\ x_i^\top \beta - \beta_0 > 0 & y_i = 1 \end{cases} \Leftrightarrow (x_i^\top \beta - \beta_0) y_i > 0$$

В таком случае, классифицирующим правилом разумно принять

$$g(x) = \text{sign}(x^\top \beta - \beta_0)$$

Ясно, что в случае 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 x_i ; это расстояние будем называть зазором (margin).

Легко видеть, что 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

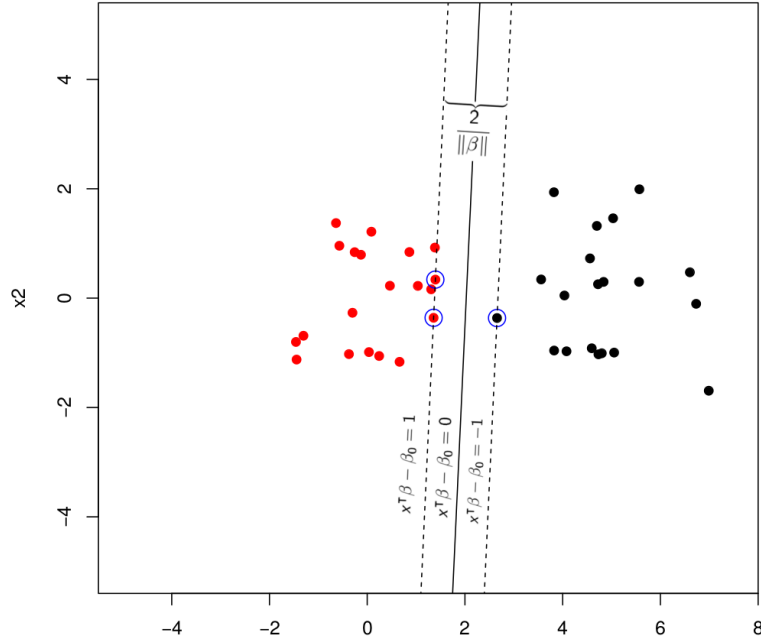


Рис. 1: Hard-margin SVM

увеличить, сместив гиперплоскость, которой не принадлежит ни одной точки); точки, которые принадлежат одной из гиперплоскостей — будем называть опорными векторами.

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

$$\begin{aligned} x^T \beta - \beta_0 &= -1 \\ x^T \beta - \beta_0 &= 1 \end{aligned}$$

а расстояние между ними составит $\frac{2}{\|\beta\|}$ (см. поясняющий рисунок 1).

Принадлежность точек обучающей выборки полупространства описывается уравнениями

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

Таким образом, в случае линейно разделимой выборки, задача выбора оптимальной гиперплоскости сводится к следующей задаче квадратичного программирования с линейными ограничениями:

$$\begin{cases} \frac{1}{2} \|\beta\|_2^2 \rightarrow \min_{\beta, \beta_0} \\ (x_i^T \beta - \beta_0) y_i \geq 1, \forall i \end{cases} \quad (1)$$

Пользуясь принципом Лагранжа, из (1) получаем задачу:

$$\begin{cases} \inf_{\beta, \beta_0} \frac{1}{2} \|\beta\|_2^2 - \sum_{i=1}^n \alpha_i (y_i (x_i^\top \beta - \beta_0) - 1) \rightarrow \max_{\alpha_1, \dots, \alpha_n} \\ \alpha_i \geq 0, \forall i \end{cases} \quad (2)$$

Так как все функции гладкие, то \inf достигается в точке, в которой выполнены необходимые условия экстремума:

$$\begin{aligned} \frac{\partial}{\partial \beta} : \quad \beta &= \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial}{\partial \beta_0} : \quad 0 &= \sum_{i=1}^n \alpha_i y_i \end{aligned}$$

подставляя эти равенства в (2) получаем:

$$\begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k x_i^\top x_k \rightarrow \max_{\alpha_1, \dots, \alpha_n} \\ \alpha_i \geq 0, \forall i \\ \sum_{i=1}^n \alpha_i y_i = 0 \\ \alpha_i (y_i (x_i^\top \beta - \beta_0) - 1) = 0 \end{cases} \quad (3)$$

Из условия регулярности ККТ:

$$\alpha_i (y_i (x_i^\top \beta - \beta_0) - 1) = 0 \quad (4)$$

в оптимальной точке, т.е. либо $\alpha_i = 0$, либо x_i является опорным вектором (принадлежит одной из пары плоскостей, описанных выше).

Таким образом, решающее правило строится на основе «сложных» для классификации наблюдений, а остальные наблюдения — не влияют (явным образом) на расположение разделяющей гиперплоскости.

1.2 Soft-margin SVM

Понятно, что требование линейной разделимости классов слишком сильное для реальной применимости SVM как метода классификации.

Позволим для этого каждому из ограничений в задаче (1) быть несколько ослабленным:

$$\begin{cases} \frac{1}{2} \|\beta\|_2^2 \rightarrow \min_{\beta, \beta_0} \\ (x_i^\top \beta - \beta_0) y_i \geq 1 - \xi_i, \forall i \\ \sum_{i=1}^n \xi_i \leq t \\ \xi_i \geq 0 \forall i \end{cases} \quad (5)$$

(t — параметр алгоритма; линейно-разделимый случай соответствует $t = 0$).

Далее, повторяя все рассуждения для линейно-разделимого случая, используем принцип Лагранжа:

$$\left\{ \begin{array}{l} \inf_{\beta, \beta_0, \xi_1, \dots, \xi_n} \frac{1}{2} \|\beta\|_2^2 + \lambda \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i \left(y_i (x_i^\top \beta - \beta_0) - (1 - \xi_i) \right) - \sum_{i=1}^n \gamma_i \xi_i \rightarrow \max_{\alpha_1, \dots, \alpha_n, \gamma_1, \dots, \gamma_n, \lambda} \\ 0 \leq \alpha_i \leq t, \forall i \\ \gamma_i \geq 0, \forall i \\ \lambda \geq 0 \end{array} \right. \quad (6)$$

Опять же, ввиду гладкости, \inf достигается в точке, в которой выполнены необходимые условия экстремума:

$$\begin{aligned} \frac{\partial}{\partial \beta} : \quad & \beta = \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial}{\partial \beta_0} : \quad & 0 = \sum_{i=1}^n \alpha_i y_i \\ \frac{\partial}{\partial \xi_i} : \quad & \alpha_i = \lambda - \gamma_i \end{aligned}$$

используя эти равенства в (6), получаем:

$$\left\{ \begin{array}{l} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k x_i^\top x_k \rightarrow \max_{\alpha_1, \dots, \alpha_n} \\ 0 \leq \alpha_i \leq t, \forall i \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{array} \right. \quad (7)$$

Из условия регулярности ККТ опять же следует классификация векторов на опорные, аутлаеры и не участвующие в построении правила предсказания точки.

Параметр t , определяющий допустимое нарушение ограничений, позволяет варьировать количество опорных векторов (которые, вследствие ККТ-условий, соответствуют векторам, участвующим в построении классифицирующего правила; пример изменения разделяющей гиперплоскости и набора опорных векторов показан на рисунке 2).

1.3 Методы оптимизации

В зависимости от соотношения размерностей n и p , может быть выгодно (с точки зрения вычислительных затрат) решать прямую (5), или двойственную 7 задачу.

При этом, несмотря на то, что обе задачи являются задачами квадратичного программирования с линейными ограничениями неравенства, использование методов общего вида осложняется экспоненциальной (в худшем случае) сложностью.

Для формулировки прямой задачи в виде

$$\sum_{i=1}^n \max \left\{ 0, 1 - y_i (x_i^\top \beta - \beta_0) \right\} + \eta \|\beta\|_2^2 \rightarrow \min_{\beta, \beta_0}$$

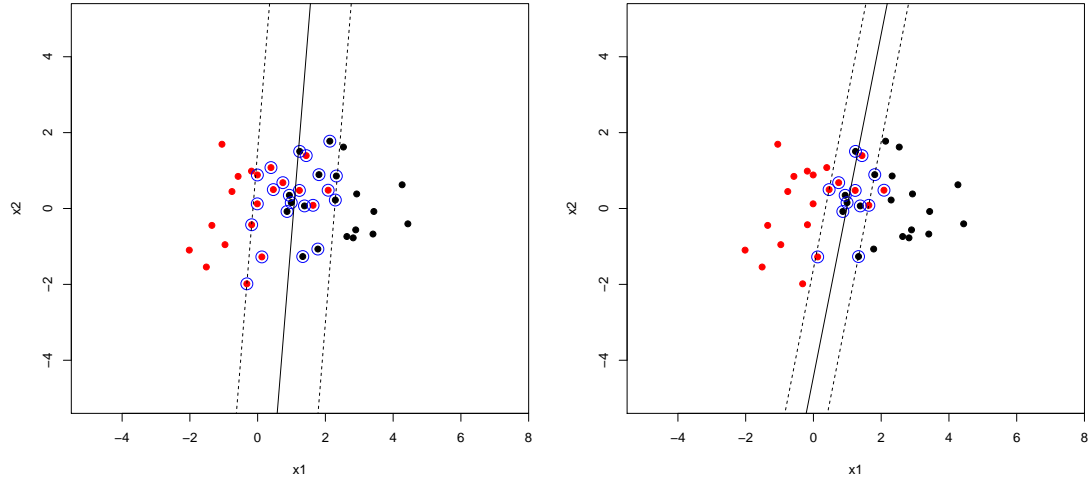


Рис. 2: Влияние максимально допустимой ошибки на soft-margin SVM

предлагается использовать (стохастический) градиентный спуск со специфическим (для задачи SVM) выбором величины шага (Pegasos: primal estimated sub-gradient solver for SVM); в этом случае можно показать, что решение с точностью ε достигается за $O\left(\frac{1}{\varepsilon}\right)$ итераций; при этом количество наблюдений не влияет на асимптотику количества итераций.

Для двойственной задачи

$$\mathcal{F}(\alpha_1, \dots, \alpha_n) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^\top x_j \rightarrow \max_{\alpha_i}$$

предлагается (A Dual Coordinate Descent Method for Large-scale Linear SVM) использовать покоординатный градиентный спуск:

- Каждый α_i изменяется в направлении $\frac{\partial \mathcal{F}}{\partial \alpha_i}$
- Очередное решение проецируется на множество допустимых
- Поддерживается необходимая для вычисления $x^\top \beta$ информация

Последовательность решений сходится как минимум линейно

1.4 SVM как частный случай Empirical Risk Minimization

Можно доказать, что задача SVM эквивалентна задаче

$$C \sum_{i=1}^n \max \{1 - (x^\top \beta + \beta_0) y_i, 0\} + \frac{1}{2} \|\beta\|_2^2 \rightarrow \min_{\beta_0, \beta}$$

где C – параметр алгоритма; домножив на $\frac{1}{nC}$ получается эквивалентная задача

$$\frac{1}{n} \sum_{i=1}^n \max \{1 - (x^\top \beta + \beta_0) y_i, 0\} + \frac{t}{2} \|\beta\|_2^2 \rightarrow \min_{\beta_0, \beta}$$

В такой формулировке можно сравнить SVM, логистическую регрессию и LDA как непрерывные (SVM) и гладкие (LDA, логистическая регрессия) аппроксимации ошибки классификации (в смысле методов минимизации [регуляризованного] эмпирического риска) следующими функциями потерь (графики функций потерь приведены на рисунке 3):

- SVM (hinge loss): $\max \{1 - (x^\top \beta + \beta_0) y_i, 0\}$
- LDA¹: $(1 - y_i (\beta_0 + \beta^\top x_i))^2$
- Logistic regression: $\log (1 + e^{-y_i (\beta^\top x_i + \beta_0)})$

Можно заметить, что LDA и логистическая регрессия имеют штраф и за правильно классифицированные наблюдения; для логистической регрессии этот штраф убывает при удалении от разделяющей гиперплоскости, а для LDA — начинает увеличиваться при удалении от центра класса.

2 Расширения SVM

2.1 Kernel trick

Описанный выше алгоритм может быть использован лишь в случае, когда данные относительно похожи на линейно-разделимые.

Предположим, что существует некоторое отображение $\Phi : \mathbb{R}^p \rightarrow V$ где V – некоторое гильбертово пространство.

Тогда, несложно заметить, что применяя SVM к образам исходных векторов, мы получаем задачу квадратичного программирования:

$$\begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \langle \Phi(x_i), \Phi(x_k) \rangle_V \rightarrow \max_{\alpha_1, \dots, \alpha_n} \\ 0 \leq \alpha_i \leq t, \forall i \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

и классифицирующая функция

$$h(x) = \text{sign} \left[\sum_{i=1}^n \alpha_i y_i \langle \Phi(x_i), \Phi(x) \rangle_V - \beta_0 \right]$$

¹В случае $\sum y_i = 0$ можно показать следующую цепочку эквивалентных переходов: LDA \Leftrightarrow FDA \Leftrightarrow CCA (с 1 переменной в одном из наборов признаков) \Leftrightarrow OLS $\sum (y_i - (\beta_0 + \beta^\top x_i))^2 \rightarrow \min_{\beta, \beta_0} \Leftrightarrow$ ERM

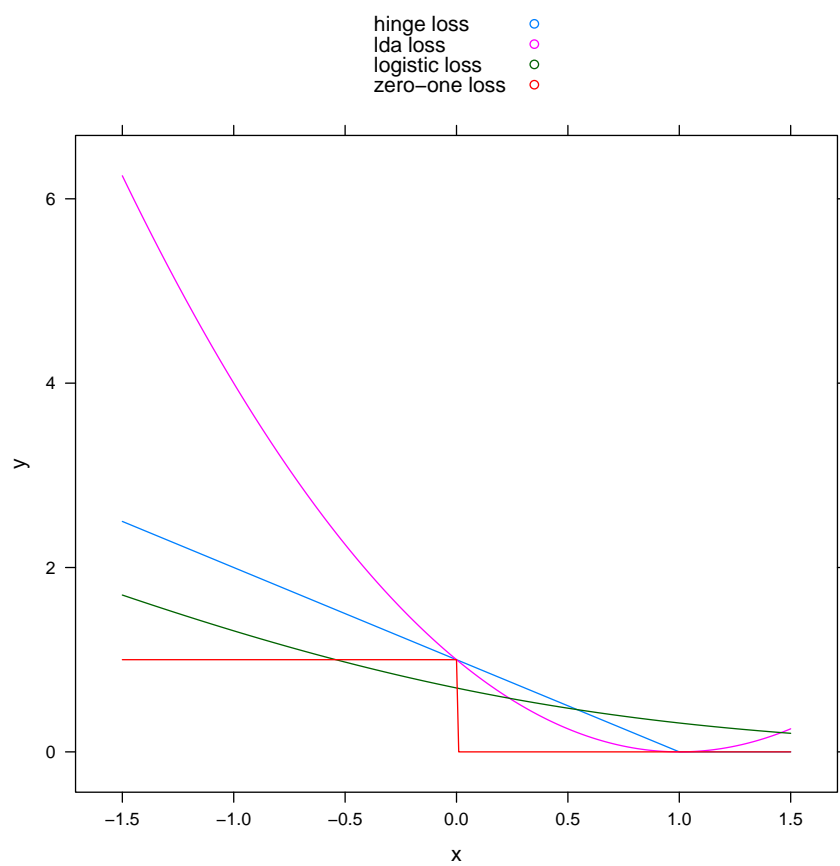


Рис. 3: Различные функции потерь как непрерывные аппроксимации ошибки классификации

(правило классификации остаётся прежним — $\text{sign}(h(x))$)

Можно заметить, что во всех выражениях результат применения Φ используется только для использования в скалярном произведении с результатом применения Φ к другому вектору из \mathbb{R}^p , что позволяет (используя теорему Мерсера) использовать произвольную симметричную положительно определённую функцию $k(u, v) : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ (ядро) в качестве скалярного произведения в некотором векторном пространстве; в этом случае преобразование Φ может соответствовать отображению, состоящему из собственных функций k :

$$\begin{aligned} k(u, v) &= \sum_{i=1}^{\infty} \theta_i \varphi_i(u) \varphi_i(v) \\ \Phi(x) &= [\varphi_1(x), \dots, \varphi_k(x), \dots,] \end{aligned} \quad (8)$$

Данное соображение позволяет применять SVM к данным в достаточной степени линейно разделимым в некотором гильбертовом пространстве (в том числе — бесконечномерном); в том числе — без предъявления в явном виде отображения из исходного пространства в данное.

Так как непосредственная проверка положительной определённости ядра представляет собой проблему, можно воспользоваться следующими операциями, приводящими к получению новых ядер:

- Скалярное произведение в векторном пространстве
- Положительная константа
- Произведение ядер: $K(u, v) = K_1(u, v) K_2(u, v)$
- Произведение отображений: $K(u, v) = \varphi(u) \varphi(v), \varphi : x \rightarrow \mathbb{R}$
- Линейная комбинация с положительными коэффициентами: $K(u, v) = \alpha_1 K_1(u, v) + \alpha_2 K_2(u, v), \alpha_{1,2} > 0$
- Композиция ядра и отображения: $K(u, v) = K_1(\varphi(u), \varphi(v))$
- Степенной ряд:

$$f : \mathbb{R} \rightarrow \mathbb{R}$$

сходящийся степенной ряд с положительными коэффициентами, тогда

$$K(u, v) = f(K_1(u, v))$$

является ядром

Часто используемые ядра:

- RBF (radial basis functions): $k(u, v) = e^{-\gamma \|u-v\|_2^2}$
- Полиномиальное (степеней $\leq d$): $k(u, v) = (\langle u, v \rangle + 1)^d$

Используя (8) можно показать, что использование полиномиального ядра степени d соответствует отображению в C_{p+d}^d -мерное пространство, гиперплоскостям в котором будут соответствовать поверхности порядка d в исходном пространстве.

Рассматривая RBF-ядро, можно убедиться, что $k(u, v) \xrightarrow{\gamma \rightarrow \infty} \mathbf{1}_{u=v}$; что позволяет добиться линейной разделимости для произвольного набора данных (предполагая отсутствие наблюдений с одинаковыми значениями признаков, принадлежащим разным классам).

2.2 Изменение регуляризации

С целью отбора признаков можно изменить регуляризацию в одной из эквивалентных формулировок SVM:

$$C \sum_{i=1}^n \max \left\{ 0, 1 - y_i (x_i^T \beta - \beta_0) \right\} + \Phi(\beta) \rightarrow \min_{\beta, \beta_0}$$

- LASSO SVM: $\Phi(\beta) = \|\beta\|_1$
 - Чем меньше C , тем больше влияние ℓ_1 -регуляризации
 - Может попеременно отбрасывать и шумовые и значимые признаки при варьировании C
 - Зависимые признаки не группируются
- Doubly-regularized SVM: $\Phi(\beta) = \alpha \|\beta\|_1 + \frac{1}{2} \|\beta\|_2^2$
 - Чем больше α , тем больше влияние ℓ_1 -регуляризации
 - Присутствует эффект группировки
 - Может попеременно отбрасывать и шумовые и значимые признаки при варьировании C, α
- Support Features Machine: $\Phi(\beta) = \sum_{i=1}^p \max \{ 2\mu\beta_i, \mu^2 + \beta_i^2 \}$
 - μ – параметр «селективности»
 - Присутствует эффект группировки
 - Значимые ($\|\beta_j\| > \mu$) признаки группируются и входят в решение совместно
 - Шумовые признаки ($\|\beta_j\| < \mu$) подавляются
- Relevance Features Machine: $\Phi(\beta) = \sum_{i=1}^p \ln \left(\beta_i^2 + \frac{1}{\mu} \right)$
 - μ – параметр «селективности»
 - Присутствует эффект группировки

– Лучше выбирает набор значимых признаков

(α, μ — дополнительные параметры; коэффициент soft-margin SVM t , эквивалентным преобразованием перемещён к \sum)

2.3 Support Vector Regression

В секциях 1.1-1.2 было показано, что классифицирующее правило в SVM зависит только от «сложно»-классифицируемых наблюдений (опорных векторов).

Оказывается, такую же идею можно использовать и для задач регрессии:

$$\begin{cases} \frac{1}{2}\beta^T\beta \rightarrow \min_{\beta, \beta_0} \\ |y_i - (x_i^T\beta + \beta_0)| \leq \varepsilon \end{cases}$$

(где ε — параметр алгоритма); добавляя возможность нарушения ограничений, переходим к задаче

$$\begin{cases} \frac{1}{2}\beta^T\beta + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \rightarrow \min_{\beta, \beta_0} \\ y_i - (x_i^T\beta + \beta_0) \leq \varepsilon + \xi_i^+ \\ -y_i + (x_i^T\beta + \beta_0) \leq \varepsilon + \xi_i^- \\ \xi_i^+ \geq 0 \\ \xi_i^- \geq 0 \end{cases}$$

Можно показать, что двойственной к ней является задача

$$\begin{cases} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i^+ - \alpha_i^-) (\alpha_j^+ - \alpha_j^-) x_i^T x_j + \varepsilon \sum_{i=1}^n (\alpha_i^+ + \alpha_i^-) + \sum_{i=1}^n y_i (\alpha_i^- - \alpha_i^+) \rightarrow \min \\ \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) = 0 \\ 0 \leq \alpha_i^+ \leq C \\ 0 \leq \alpha_i^- \leq C \\ \beta = \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) x_i \end{cases}$$

При этом, в силу ККТ-условий, многие $\alpha_i = 0$, т.е. в построении функции регрессии как и в SVM участвует лишь ограниченное количество опорных наблюдений.

Аналогично SVM, двойственная задача регрессии допускает нелинейное обобщение с использованием kernel trick.

2.4 Multi-class SVM

В предъявленном построении SVM рассматривался лишь случай классификации с двумя классами.

Для распространения идеи SVM на классификацию с большим количеством классов существует несколько подходов, в частности:

- Классификация с использованием сравнений вида “один со многими”
- Классификация с использованием сравнений вида “каждый с каждым”

2.4.1 Сравнения “один со многими”

Для классификации с N классами строится N классифицирующих правил $h_i(x)$; кодирующих принадлежность i -му классу за 1, а принадлежность любому другому классу за -1 .

В качестве результирующего решающего правила используется

$$h(x) = \operatorname{argmax}_i h_i(x)$$

2.4.2 Сравнения “каждый с каждым”

Для классификации с N классами строится $\frac{N(N-1)}{2}$ классифицирующих правил, производящих классификацию для каждой возможной пары классов.

Обозначив за N_i количество сравнений, в которых элемент x был классифицирован как принадлежащий i -ому классу; в качестве классифицирующего правила предлагается использовать

$$h(x) = \operatorname{argmax}_i N_i$$

3 Выбор параметров

Ввиду наличия свободы выбора значения параметра регуляризации, ядра (или семейства ядер), необходимо предъявить процедуру сравнения построенных классификаторов. Так как в предлагаемой процедуре не используются никакие предположения о распределении $P(x, y)$, использование информационных критериев (BIC, AIC, ...) для этих целей невозможно.

Предлагается рассмотреть процедуры выбора параметров, основанные на (общей) идее кросс-валидации и специфичную для SVM оценку эмпирического риска на основе комбинаторной размерности.

3.1 Кросс-валидация

Предполагая, что выборка является выборкой из распределения $\langle x_i, y_i \rangle \sim \mathcal{P}(x, y)$, с точки зрения выбора параметров или типа алгоритма классификации (или регрессии), хотелось бы получать не только выборочную оценку качества классификации (регрессии), но и оценку пригодности выбранного алгоритма в смысле качества классификации (регрессии) на всём распределении \mathcal{P} , а не только конкретной выборки из него.

Для этого предлагается построить несколько моделей по различным подмножествам датасета, оценить для каждой из моделей качество классификации (регрессии) по части датасета, не участвовавшей в оценке параметров, после чего получить оценку математического ожидания качества классификации (регрессии).

Обозначив за $N = \{1, \dots, n\}$ множество индексов элементов выборки, выберем K подмножеств (примерно одинакового размера) $N_k \subset N$.

Построим K классификаторов (функций регрессии) $f_k(x) = f(x, \hat{\theta}(N \setminus N_k))$.

3.1.1 Классификация

Для каждой из них вычислим эмпирическую ошибку классификации

$$\hat{\varepsilon}_k = \frac{1}{\#N_k} \sum_{i \in N_k} \mathbf{1}_{f_k(x_i) \neq y_i}$$

Используя данную оценку, построим оценку ошибки классификации

$$\hat{\varepsilon} = \frac{1}{n} \sum_{i=1}^K \#N_k \hat{\varepsilon}_k$$

Также можно оценить «разброс» ошибки классификации как

$$\hat{\sigma}_\varepsilon = \sqrt{\frac{1}{K-1} \sum_{i=1}^K (\hat{\varepsilon}_k - \hat{\varepsilon})^2}$$

3.1.2 Регрессия

Для каждой из построенных моделей вычислим ошибку регрессии

$$\hat{\varepsilon}_k = \frac{1}{\#N_k} \sum_{i \in N_k} (f_k(x_i) - y_i)^2$$

И используя данные оценки, построим оценку дисперсии остатка регрессии для семейства f :

$$\hat{\varepsilon} = \frac{1}{n} \sum_{i=1}^K \#N_k \hat{\varepsilon}_k$$

3.1.3 Общие соображения

Получив оценки $\hat{\varepsilon}$ для всех интересующих параметрических семейств классификаторов (функций регрессии), следует выбрать семейство, для которого $\hat{\varepsilon}$ минимальна. После этого имеет смысл повторить оценку параметров, используя все доступные данные.

Стоит отметить, что оценки $\hat{\varepsilon}$ обычно консервативны (хуже, чем производительность наилучшего представителя семейства с параметрами, оцененными по всем данным).

Существует несколько подходов к разделению выборки для cross-validation:

- K-fold cross-validation: набор индексов N разбивается на K примерно-равных дизъюнктивных подмножества N_k
- Leave-one-out cross-validation: N -fold cross-validation; рассматривается набор N_k , соответствующий всем одноэлементным подмножествам N

3.2 Оценка через комбинаторную размерность

Рассмотрим некоторое параметрическое семейство классификаторов $f_\alpha : \mathbb{R}^p \rightarrow \{-1, 1\}$; α — параметры классификатора.

Предполагая, что выборка является выборкой из распределения $\langle x_i, y_i \rangle \sim \mathcal{P}(x, y)$ можно определить риск как

$$R(\alpha) = \int \mathbf{1}_{y \neq f_\alpha(x)} d\mathcal{P}(x, y)$$

Определив эмпирический риск (являющийся случайной величиной) как

$$R_{emp}(\alpha) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{y_i \neq f_\alpha(x_i)}$$

можно показать, что (независимо от \mathcal{P}) с вероятностью $1 - \eta$ выполнено неравенство:

$$R(\alpha) \leq R_e(\alpha) + \sqrt{\frac{h \left(1 + \log \frac{2n}{h} \right) - \log \frac{\eta}{4}}{n}}$$

где h — комбинаторная размерность семейства классификаторов, определяемая как максимальное количество точек, которые при любом их расположении и разделении на классы при некотором α будут безошибочно классифицированы.

Например, для линейного классификатора в p -мерном пространстве VC-размерность составляет $p + 1$, для полиномиального ядра степени d — C_{d+p}^d (ввиду числа мономов, являющихся собственными функциями ядра), для RBF VC-размерность бесконечна (см. описанную в 2.1 конструкцию, разделяющую произвольное множество точек на основе RBF).

Можно отметить, что при конечной VC-размерности и $n \rightarrow \infty$ дополнительное слагаемое стремится к 0, т.е. $R_{emp}(\alpha) \rightarrow R(\alpha)$ по вероятности (и эта сходимости наблюдается вне зависимости от конкретного распределения \mathcal{P}).