

Дано:

множество объектов X , множество классов Y ;

$X^\ell = \left\{ \begin{array}{l} x_1, \dots, x_\ell \\ y_1, \dots, y_\ell \end{array} \right\}$ — размеченная выборка (labeled data);

$X^k = \{x_{\ell+1}, \dots, x_{\ell+k}\}$ — неразмеченная выборка (unlabeled data).

Два варианта постановки задачи:

- *Частичное обучение* (semi-supervised learning):
построить алгоритм классификации $a: X \rightarrow Y$.
- *Трансдуктивное обучение* (transductive learning):
зная **все** $\{x_{\ell+1}, \dots, x_{\ell+k}\}$, получить метки $\{y_{\ell+1}, \dots, y_{\ell+k}\}$.

Типичные приложения:

классификация и каталогизация текстов, изображений, и т. п.

Линейный классификатор на два класса $Y = \{-1, 1\}$:

$$a(x) = \text{sign}(\langle w, x \rangle - w_0), \quad w, x \in \mathbb{R}^n, \quad w_0 \in \mathbb{R}.$$

Отступ объекта x_i :

$$M_i(w, w_0) = (\langle w, x_i \rangle - w_0) y_i.$$

Задача обучения весов w, w_0 по размеченной выборке:

$$Q(w, w_0) = \sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}.$$

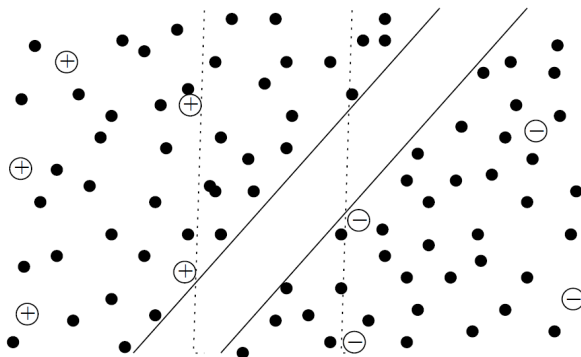
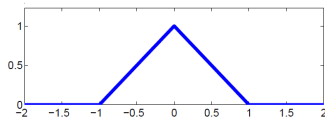
Функция $\mathcal{L}(M) = (1 - M)_+$ штрафует за уменьшение отступа.

Идея!

Функция $\mathcal{L}(M) = (1 - |M|)_+$ штрафует за попадание в зазор, $|M_i|$ не зависит от y_i и определён на неразмеченных объектах.

Функция потерь для трансдуктивного SVM

Функция потерь $\mathcal{L}(M) = (1 - |M|)_+$ штрафует за попадание объекта внутрь разделяющей полосы.



Обучение весов w, w_0 по частично размеченной выборке:

$$Q(w, w_0) = \sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2C} \|w\|^2 + \\ + \gamma \sum_{i=\ell+1}^{\ell+k} (1 - |M_i(w, w_0)|)_+ \rightarrow \min_{w, w_0}.$$

Достоинства и недостатки TSVM:

- ⊕ как и в обычном SVM, можно использовать ядра;
- ⊕ имеются эффективные реализации для больших данных;
- ⊖ решение неустойчиво, если нет области разреженности;
- ⊖ требуется настройка двух параметров C, γ ;

Sindhwani, Keerthi. Large scale semisupervised linear SVMs. SIGIR 2006.

Линейный классификатор на конечное множество классов $|Y|$:

$$a(x) = \arg \max_{y \in Y} \langle w_y, x \rangle, \quad x, w_y \in \mathbb{R}^n.$$

Вероятность того, что объект x_i относится к классу y :

$$P(y|x_i, w) = \frac{\exp \langle w_y, x_i \rangle}{\sum_{c \in Y} \exp \langle w_c, x_i \rangle}.$$

Задача максимизации регуляризованного правдоподобия:

$$Q(w) = \sum_{i=1}^{\ell} \log P(y_i|x_i, w) - \frac{1}{2C} \sum_{y \in Y} \|w_y\|^2 \rightarrow \max_w,$$

Оптимизация $Q(w)$ — методом стохастического градиента.

Теперь учтём неразмеченные данные $X^k = \{x_{\ell+1}, \dots, x_{\ell+k}\}$. Пусть $b_j(x)$ — бинарные признаки, $j = 1, \dots, m$.

Оценим вероятности $P(y|b_j(x) = 1)$ двумя способами:

1) эмпирическая оценка по размеченным данным X^ℓ :

$$\hat{p}_j(y) = \frac{\sum_{i=1}^{\ell} b_j(x_i) [y_i = y]}{\sum_{i=1}^{\ell} b_j(x_i)};$$

2) оценка по неразмеченным данным X^k и линейной модели:

$$p_j(y, w) = \frac{\sum_{i=\ell+1}^{\ell+k} b_j(x_i) P(y|x_i, w)}{\sum_{i=\ell+1}^{\ell+k} b_j(x_i)}.$$

Будем минимизировать расстояние между $\hat{p}_j(y)$ и $p_j(y, w)$, в качестве расстояния между распределениями возьмём дивергенцию Кульбака–Лейблера.

Минимизация KL-дивергенции между $\hat{p}_j(y)$ и $p_j(y, w)$:

$$\text{KL}(\hat{p}_j(y) \parallel p_j(y, w)) = \sum_y \hat{p}_j(y) \log \frac{\hat{p}_j(y)}{p_j(y, w)} \rightarrow \min_w.$$

Вычтем сумму KL-дивергенций по всем признакам $j = 1, \dots, m$ из функционала регуляризованного правдоподобия $Q(w)$:

$$\begin{aligned} \tilde{Q}(w) = & \sum_{i=1}^{\ell} \log P(y_i | x_i, w) - \frac{1}{2C} \sum_{y \in Y} \|w_y\|^2 + \\ & + \gamma \sum_{j=1}^m \sum_{y \in Y} \hat{p}_j(y) \log \left(\frac{\sum_{i=\ell+1}^{\ell+k} b_j(x_i) P(y | x_i, w)}{\sum_{i=\ell+1}^{\ell+k} b_j(x_i)} \right) \rightarrow \max_w, \end{aligned}$$

где γ — коэффициент регуляризации.

Mann, McCallum. Simple, robust, scalable semi-supervised learning via expectation regularization. ICML 2007.

- ❶ Оптимизация $\tilde{Q}(w)$ — методом стохастического градиента.
- ❷ Возможные варианты задания переменных b_j :
 - $b_j(x) \equiv 1$, тогда $P(y|b_j(x) = 1)$ — априорная вероятность класса y (label regularization) — хорошо подходит для задач с несбалансированными классами;
 - $b_j(x) = [\text{термин } j \text{ содержится в тексте } x]$ — для задач классификации и каталогизации текстов.
- ❸ метод слабо чувствителен к выбору C и γ ,
- ❹ устойчив к погрешностям оценивания $\hat{p}_j(y)$,
- ❺ не требует большого числа размеченных объектов ℓ ,
- ❻ хорошо подходит для категоризации текстов,
- ❼ в экспериментах показывает высокую точность.

Mann, McCallum. Simple, robust, scalable semi-supervised learning via expectation regularization. ICML 2007.

- Задача SSL занимает промежуточное положение между классификацией и кластеризацией, но не сводится к ним.
- Простые методы-обёртки требуют многократного обучения, что вычислительно неэффективно.
- *Методы кластеризации* легко адаптируются к SSL путём введения ограничений (constrained clustering), но, как правило, вычислительно трудоёмки.
- *Методы классификации* адаптируются сложнее, но приводят к более эффективному частичному обучению.