

1. Как выглядит бинарный линейный классификатор? (Формула для отображения из множества объектов в множество классов.)

- $y = \text{sign}(\langle wx \rangle - w_0)$, где y предсказываемый класс, x - признаки, w - веса.

1. Что такое отступ алгоритма на объекте? Какие выводы можно сделать из знака отступа?

- $M = y_i \cdot (\langle wx \rangle - w_0)$, y_i - реальный класс объекта. Отступ показывает насколько далеко объект от разделяющей гиперплоскости, если знак отрицательный, то класс предсказывается ошибочно. Функции отступа могут быть различны, но $M > 0$, если $a(x) = y_i$ и $M \leq 0$, если $a(x) \neq y_i$.

1. Как классификаторы вида $a(x) = \text{sign}(\langle w, x \rangle - w_0)$ сводят к классификаторам вида $a(x) = \text{sign}(\langle w, x \rangle)$?

- $x := (x, -1)$ - добавляем -1 к x как новую координату, $w := (w, w_0)$

1. Как выглядит запись функционала эмпирического риска через отступы? Какое значение он должен принимать для «наилучшего» алгоритма классификации?

- $Q(w) = \sum_{i=1}^n I[M_i \leq 0]$.
- $Q(w) \geq 0$, $Q(w) \rightarrow \min \implies$ для "наилучшего" $Q = 0$.

2. Если в функционале эмпирического риска (риск с пороговой функцией потерь) всюду написаны строгие неравенства ($M_i < 0$) можете ли вы сразу придумать параметр w для алгоритма классификации $a(x) = \text{sign}(\langle w, x \rangle)$, минимизирующий такой функционал?
 - $w = (0, 0, 0, \dots, 0)$
3. Запишите функционал аппроксимированного эмпирического риска, если выбрана функция потерь $L(M)$.
 - $Q(w) = \sum_{i=1}^n L(M_i)$
4. Что такое функция потерь, зачем она нужна? Как обычно выглядит ее график?
 - Обычно это "сглаживание" пороговой функции потерь для применения градиентных методов минимизации.
 - При $M < 0 : L(M) \geq 1$, при $M \rightarrow +\infty : L(M) \rightarrow 0$
5. Приведите пример негладкой функции потерь.
 - $L(M) = I[M > 0]$
6. Что такое регуляризация? Какие регуляризаторы вы знаете?
 - Метод добавления некоторой дополнительной информации к условию с целью решить некорректно поставленную задачу или предотвратить переобучение. (например ограничение на норму весов).
 - Новая задача : $Q(w) + \gamma V(w) \rightarrow \min$
 - $L1(\sum w_i), L2(\sum w_i^2)$ - регуляризации
7. Как связаны переобучение и обобщающая способность алгоритма? Как влияет регуляризация на обобщающую способность?

- Обратнопорпорционально...
- Регуляризация увеличивает обобщающую способность, так как препятствует переобучению.

1. Как связаны острые минимумы функционала аппроксимированного эмпирического риска с проблемой переобучения?

- Острые минимумы функционала эмпирического риска говорят о том, что функция риска сильно меняется даже при небольшом изменении весов, а, значит, и признаков, что говорит об отсутствии устойчивости модели, т.е. о переобучении.

2. Что делает регуляризация с аппроксимированным риском как функцией параметров алгоритма?

- препятствует неограниченному увеличению или уменьшению весов.

3. Для какого алгоритма классификации функционал аппроксимированного риска будет принимать большее значение на обучающей выборке: для построенного с регуляризацией или без нее? Почему?

- Для "без регуляризации" меньше. Увеличивая веса можно добиться увеличения положительных отступов и тем самым функционала риска.

4. Для какого алгоритма классификации функционал риска будет принимать большее значение на тестовой выборке: для построенного с оправдывающей себя регуляризацией или вообще без нее? Почему?

- Без "регуляризации" больше, т.к. в отсутствии регуляризации будет переобучение.

5. Что представляют собой метрики качества Accuracy, Precision и Recall?

- Accuracy – доля правильных ответов
- Precision – $\frac{TP}{TP+FP}$
- Recall – $\frac{TP}{P}$

6. Что такое метрика качества AUC и ROC-кривая?

- $FPR = FP/N$, $TPR = TP/P$, тогда ROC – кривая: $TPR(FPR)$. AUC – площадь под графиком.

7. Как построить ROC-кривую (нужен алгоритм), если например, у вас есть правильные ответы к домашнему заданию про фамилии и ваши прогнозы?

- В результате работы алгоритма получим множество $(FPR_i, TPR_i)_{i=0}^l$, по которому строится ROC-кривая:
- $l+$ – количество фамилий в выборке X^l , $l-$ – количество фамилий в выборке.

Сортируем выборку X^l по значениям дискриминантной функции $f(x_i, w)$.

$(FPR_0, TPR_0) = (0, 0)$

- for i in $\{1, \dots, l\}$:
- if $y_i == -1$:
- $(FPR_i, TPR_i) = (FPR_{i-1} + 1/l-, TPR_{i-1})$
- else:
- $(FPR_i, TPR_i) = (FPR_{i-1}, TPR_{i-1} + 1/l+)$

2 Вероятностный смысл регуляризаторов

Покажите, что регуляризатор в задаче линейной классификации имеет вероятностный смысл априорного распределения параметров моделей. Какие распределения задают ℓ_1 -регуляризатор и ℓ_2 -регуляризатор?

Решение

Допустим, множество $X \times Y$ является вероятностным пространством, и задана параметрическая модель совместной плотности распределения объектов и классов $p(x, y|w)$. Введем параметрическое семейство априорных распределений $p(w; \gamma)$, где γ — неизвестная и не случайная величина (гиперпараметр).

Тогда будем считать, что выборка может быть порождена каждой из плотностей $p(x, y|w)$ с параметризованной γ вероятностью $p(w; \gamma)$.

Приходим к принципу максимума совместного правдоподобия данных и модели:

$$L_\gamma(w, X^l) = \ln p(X^l, w; \gamma) = \sum_{i=1}^l \ln p(x_i, y_i|w) + \underbrace{\ln p(w; \gamma)}_{\text{регуляризатор}} \rightarrow \max_w.$$

Вспомним, что этот принцип эквивалентен принципу минимизации аппроксимированного эмпирического риска

$$Q(w; X^l) = \sum_{i=1}^l \mathcal{L}(y_i f(x_i, w)) + \underbrace{\gamma V(w)}_{\text{регуляризатор}} \rightarrow \min_w,$$

если положить

$$\begin{aligned} -\ln p(x_i, y_i|w) &= \mathcal{L}(y_i f(x_i, w)), \\ \ln p(w; \gamma) &= \gamma V(w). \end{aligned}$$

Таким образом, получаем, что регуляризатор $V(w)$ соответствует параметрическому семейству априорных распределений плотностей $p(w; \gamma)$ — параметров моделей.

ℓ_1 -регуляризатор

Пусть $w \in \mathbb{R}^n$ имеет n -мерное распределение Лапласа:

$$p(w; C) = \frac{1}{(2C)^n} \exp\left(-\frac{\|w\|_1}{C}\right), \quad \|w\|_1 = \sum_{j=1}^n |w_j|,$$

т.е. все веса независимы, имеют нулевое матожидание и равные дисперсии; C — гиперпараметр.

Логарифмируя, получаем регуляризатор по ℓ_1 -норме:

$$-\ln p(w; C) = \frac{1}{C} \sum_{j=1}^n |w_j| + \text{const}(w).$$

ℓ_2 -регуляризатор

Пусть $w \in \mathbb{R}^n$ имеет n -мерное гауссовское распределение:

$$p(w; \sigma) = \frac{1}{(2\pi\sigma)^{n/2}} \exp\left(-\frac{\|w\|^2}{2\sigma}\right), \quad \|w\|^2 = \sum_{j=1}^n w_j^2,$$

т.е. все веса независимы, имеют нулевое матожидание и равные дисперсии σ ; σ — гиперпараметр.

Логарифмируя, получаем регуляризатор по ℓ_2 -норме:

$$-\ln p(w; \sigma) = \frac{1}{2\sigma} \|w\|^2 + \text{const}(w).$$

3 SVM и максимизация разделяющей полосы

Покажите, как получается условная оптимизационная задача, решаемая в SVM из соображений максимизации разделяющей полосы между классами. Можно отталкиваться от линейно разделимого случая, но итоговое выражение должно быть для общего. Как эта задача сводится к безусловной задаче оптимизации?

Решение

Рассмотрим задачу классификации на два непересекающихся класса, в которой объекты опи-

сываются n -мерными вещественными векторами: $X = \mathbb{R}^n, Y = \{-1, +1\}$.

Будем строить линейный пороговый классификатор:

$$a(x) = \text{sign} \left(\sum_{j=1}^n w_j x^j - w_0 \right) = \text{sign} (\langle w, x \rangle - w_0),$$

где $x = (x^1, \dots, x^n)$ — признаковое описание объекта x , вектор $w = (w^1, \dots, w^n) \in \mathbb{R}^n$ и скалярный порог $w_0 \in \mathbb{R}$ — параметры алгоритма.

Для начала предположим, что выборка линейна разделима: найдутся w, w_0 , задающие разделяющую гиперплоскость $\langle w, x \rangle = w_0$, при которых функционал числа ошибок

$$Q(w, w_0) = \sum_{i=1}^l [y_i (\langle w, x_i \rangle - w_0) \leq 0] = 0.$$

Найдем, как оптимальнее расположить разделяющую гиперплоскость. Для простоты выполним нормировку параметров алгоритма: домножим w и w_0 на такую константу, что

$$\min_{i=1, l} y_i (\langle w, x_i \rangle - w_0) = 1.$$

Хочется максимизировать ширину разделяющей полосы. Тогда на границе разделяющей полосы будут лежать точки из обучающей выборки: x_- и x_+ , принадлежащие соответственно -1 и $+1$ классам. Ширина полосы

$$\begin{aligned} \left\langle (x_+ - x_-), \frac{w}{\|w\|} \right\rangle &= \frac{\langle w, x_+ \rangle - \langle w, x_- \rangle}{\|w\|} = \left| y_i (\langle w, x_i \rangle - w_0) = 1, i \in \{+, -\} \right| = \\ &= \frac{(w_0 + 1) - (w_0 - 1)}{\|w\|} = \frac{2}{\|w\|}. \end{aligned}$$

Получаем, что ширина полосы максимальна, когда норма вектора w минимальна. Значит, можно сформулировать следующую задачу оптимизации:

$$\left\{ \begin{array}{l} \frac{1}{2} \|w\|^2 \rightarrow \min; \\ w, w_0 \end{array} \right.$$

$$\bigwedge y_i (\langle w, x_i \rangle - w_0) \geq 1, i = 1, \dots, l.$$

Если работать с линейно неразделимой выборкой, $y_i (\langle w, x_i \rangle - w_0)$ не обязательно будет не меньше 1. Ослабим эти ограничения и введем в минимизируемый функционал штраф за суммарную ошибку:

$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i \rightarrow \min_{w, w_0, \xi}; \\ M_i(w, w_0) = y_i (\langle w, x_i \rangle - w_0) \geq 1 - \xi_i, i = 1, \dots, l; \\ \xi_i \geq 0, i = 1, \dots, l. \end{cases}$$

ξ_i в этих условиях будет показывать величину ошибки на x_i объекте.

Преобразуем условия на ξ_i :

$$\begin{cases} \xi_i \geq 1 - M_i(w, w_0) \\ \xi_i \geq 0 \end{cases}$$

Значит, минимум ξ_i будет при $\xi_i = (1 - M_i(w, w_0))_+$.

Получаем эквивалентную задачу безусловной оптимизации:

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

4 Kernel trick

Придумайте ядро, которое позволит линейному классификатору с помощью Kernel Trick построить в исходном пространстве признаков разделяющую поверхность $x_1 + 2x_2 = 3$. Какой будет размерность спрямляющего пространства?

Возьмем квадратичное ядро: $K(x, y) = \langle x, y \rangle^2 = (x_1 y_1 + x_2 y_2)^2 = x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2 = \langle (x_1^2, x_2^2, \sqrt{2}x_1 x_2), (y_1^2, y_2^2, \sqrt{2}y_1 y_2) \rangle$.

Получим отображение в спрямляющее пространство $H = \mathbb{R}^3$:

$$\psi : (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1 x_2).$$

Тогда линейная поверхность в H будет иметь вид: $\langle (x_1^2, x_2^2, \sqrt{2}x_1 x_2), (w_1, w_2, w_3) \rangle + w_0 = w_1 x_1^2 + w_2 x_2^2 + w_3 \sqrt{2}x_1 x_2 + w_0 = 0$.
 $\left| w = (1, 2, 0), w_0 = -3 \right| = x_1^2 + 2x_2^2 - 3 = 0$.

(Вообще-то говоря, w ищется из условия $\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^l \lambda_i y_i \psi(x_i) = 0$ после того, как найдены λ_i в задаче (4.1) со скалярным произведением $K(x_i, x_j)$ вместо $\langle x_i, x_j \rangle$.)

w_0 можно найти, подставив в выражение $w_0 = \langle w, \psi(x_i) \rangle - y_i$ произвольный опорный граничный вектор x_i .)

5 l1-регуляризация

Покажите с помощью теоремы Куна-Таккера, что ограничение l1-нормы вектора весов числом и добавление штрафа с его l1-нормой приводят к построению одного и того же алгоритма. Можно считать, что регуляризатор добавляется по существу, т.е. меняет итоговый ответ по сравнению с оптимизационной задачей без регуляризатора.

Лассо Тибширани:

$$\begin{cases} Q(\alpha) = \|F\alpha - y\|^2 \rightarrow \min_{\alpha}; \\ \sum_{j=1}^n |\alpha_j| \leq \kappa \Leftrightarrow \sum_{j=1}^n |\alpha_j| - \kappa \leq 0. \end{cases}$$

По теореме Куна-Таккера:

$$\begin{cases} Q(\alpha) = \|F\alpha - y\|^2 + \lambda \left(\sum_{j=1}^n |\alpha_j| - \kappa \right) = \|F\alpha - y\|^2 + \lambda \left(\sum_{j=1}^n |\alpha_j| \right) + \text{const} \rightarrow \min_{\alpha}; \\ \lambda \geq 0; \\ \lambda \left(\sum_{j=1}^n |\alpha_j| - \kappa \right) = 0 \Leftrightarrow \lambda = 0 \text{ или } \sum_{j=1}^n |\alpha_j| = \kappa. \end{cases}$$

In []: