

MO2

March 2021

1 Ядровые методы

Данные: $x = (x_1, \dots, x_m)$

Базисные функции: $\phi(x_1, \dots)$

Модель принимает вид: $a(x) = \sum_{j=1}^m w_j \phi_j(x)$

Для хорошего качества нужно много базисных функций \rightarrow Ядровые методы позволяют не перебирать большое количество базисных функций

- Быстрое обучение

Ядровые методы

1. Двойственное представление для линейной регрессии

$$Q(w) = \frac{1}{2} \sum_{i=1}^l (\sum_{j=1}^m (w_j \phi_j(x_i) - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2 = \frac{1}{2} \|\Phi w - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$\Phi = \begin{pmatrix} \phi_1(x_1) & \dots & \phi_m(x_1) \\ \dots & \dots & \dots \\ \phi_1(x_l) & \dots & \phi_m(x_l) \end{pmatrix}$$

$$\nabla_w Q = \Phi^T (\Phi w - y) + \lambda w \rightarrow w = -\frac{1}{\lambda} \Phi^T (\Phi w - y) \rightarrow w = \Phi^T a$$

w является линейной комбинацией строк $\Phi \rightarrow$ Решение можно искать из $w = \Phi^T a$

$$Q(a) = \frac{1}{2} \|\Phi \Phi^T a - y\|^2 + \frac{\lambda}{2} a^T \Phi \Phi^T a \rightarrow \min_a$$

$\Phi \Phi^T$ - матрица Грама (попарных скалярных произведений объектов)

Можно записать $Q(w)$ так, что он зависит только от скалярных произведений объектов

2. SVM

$$\begin{cases} \sum_{i=1}^l \lambda_i - \frac{1}{2} \sum_{i,j=1}^l \lambda_i \lambda_j y_i y_j < x_i, x_j > \rightarrow \max_{\lambda} \\ 0 \geq \lambda_i \leq C \\ \sum_{i=1}^l \lambda_i y_i = 0 \end{cases}$$

Такая формулировка задачи зависит от скалярных произведений объектов

3. Алгоритм

- (a) Добавляем новые признаки
- (b) $x, z \in X$
- (c) Делаем это так, что $\langle \phi(x), \phi(z) \rangle$ выражается через $\langle x, z \rangle$
- (d) Используем метод, который использует скалярные произведения объектов
- (e) В этом методе $\langle x, z \rangle \rightarrow \langle \phi(x), \phi(z) \rangle$ (*Kernel trick*)

4. Ядро - функция $K(x, z) = \langle \phi(x), \phi(z) \rangle$, где $\phi : X \rightarrow H$

- (a) H - спрямляющее пространство
- (b) ϕ - спрямляющее отображение

5. Теорема Мерсера

- (a) $K(x, z)$ - ядро $\leftrightarrow \begin{cases} K(x, z) = K(z, x) \\ K \text{ неотрицательно определенная} \end{cases}$
- (b) $\text{НО} \rightarrow \forall l, \forall (x_1, \dots, x_l) \in R^d \rightarrow (K(x_i, x_j))_{i,j=1}^l \text{ НО}$
- (c) На практике теорема Мерсера слишком сложна для применения

6. Теорема 1

- (a) Если
 - i. $K_1(x, z), K_2(x, z)$ - ядра, $x, z \in X$
 - ii. $f^{(x)}$ - вещественная функция на X
 - iii. $\phi : X \rightarrow R^n$
 - iv. K_3 - ядро заданное на R^n
- (b) То следующие функции являются ядрами:
 - i. $K(x, z) = K_1(x, z) + K_2(x, z)$
 - ii. $K(x, z) = \alpha K_1(x, z)$

- iii. $K(x, z) = K_1 K_2$
- iv. $K(x, z) = f^{(x)} f^{(z)}$
- v. $K(x, z) = K(\phi(x), \phi(z))$

7. Теорема 2

- (a) Если:
 - i. $K_1(x, z), K_2(x, z), \dots$ - последовательность ядер
 - ii. $\exists K(x, z) = \lim_{n \rightarrow \infty} K_n(x, z), \forall x, z$
- (b) То:
 - i. K - ядро

8. Полиномиальные ядра

- (a) $p(v)$ - многочлен с неотриц. коэфф
- (b) $K(x, z) = w_0 + w_1 \langle x, z \rangle + w_2 \langle x, z \rangle^2 + \dots$
- (c) Является ядром по теореме 1
- (d) $K(x, z) = (\langle x, z \rangle + R)^m = \sum_{i=0}^m C_m^i R^{m-i} \langle x, z \rangle^i$
 - i. Если расписать все $\langle x, z \rangle^i$, то получим все мономы степени i от исходных признаков
 - ii. Зачем R ? \rightarrow коэффициент при мономе $= \sqrt{C_m^i R^{m-i}}$
 - iii. Сравним веса при мономах 1 и $(m-1)$ $\sqrt{\frac{C_m^{m-1} R}{C_m^1 R^{m-1}}} = \sqrt{\frac{1}{R^{m-2}}}$
 - iv. R больше - мономы высоких степеней имеют низкий вклад
 - v. Конечномерное спрямляющее пространство, но можно сделать линейно разделимое пространство

9. Гауссовы ядра

- (a) Позволяет перевести в бесконечномерное спрямляющее пространство
- (b) $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$
 - i. $\exp(\langle x, z \rangle) = \sum_{k=0}^{\infty} \frac{\langle x, z \rangle^k}{k!}, \forall x, z = \lim_{n \rightarrow \infty} \sum_{k=0}^{\infty} \frac{\langle x, z \rangle^k}{k!}$
 - А. Разложение через ряд Тейлора
 - В. Ядро, как последовательность ядер
 - ii. $\frac{\exp(\langle x, z \rangle)}{2\sigma^2}$ - ядро, аналогично
 - iii. $\exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\langle x-z, x-z \rangle}{2\sigma^2}\right) = \exp\left(-\frac{\langle x, x \rangle - \langle x, z \rangle - \langle z, z \rangle + \langle x, z \rangle}{2\sigma^2}\right) =$
 $\frac{\exp(\langle x, z \rangle / \sigma^2)}{\exp(\|x\|^2 / \sigma^2) \exp(\|z\|^2 / \sigma^2)}$

$$\text{iv. } \exp(< x, z > / \sigma^2) = K(x, z) = < \phi(x), \phi(z) >$$

$$\text{v. } \phi(\tilde{x}) = \frac{\phi(x)}{\|\phi(x)\|} = \frac{\phi(x)}{\sqrt{K(x, x)}}$$

$$\text{vi. } < \phi(\tilde{x}), \phi(\tilde{z}) > = \frac{< \phi(x), \phi(z) >}{\sqrt{K(x, x)K(z, z)}}$$

(с) Какое спрямляющее пространство? - бесконечная сумма всех мономов

(d) Утверждение: x_1, \dots, x_l - различные векторы из \mathbb{R}^d

Тогда:

$$G = (\exp(-\frac{\|x-z\|^2}{2\sigma^2}))_{i,j=1}^l - \text{ невырожденная при } \sigma^2 > 0$$

(e) $x_1, \dots, x_l \in \mathbb{R}^d$ - их матрица Грамма невырождена $\rightarrow \phi(x_1, \dots, x_l)$
ЛНЗ \rightarrow бесконечное количество ЛНЗ векторов \rightarrow бесконечномерное пространство

10. Ядровой SVM

$$(a) \begin{cases} \frac{1}{2}\|w\|^2 + C \sum_{i=1}^l \xi_i \rightarrow \min_{w,b,\xi} \\ y_i(< w, x_i > + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

$$L(w, b, \xi, \lambda, \mu) = \frac{1}{2}\|w\|^2 + C \sum_{i=1}^d \xi_i - \sum_{i=1}^l \lambda_i (y_i(< w, x_i > + b) - 1 + \xi_i) - \sum_{i=1}^l \mu_i \xi_i$$

В точке оптимума $\nabla_w L = 0$

$$\nabla_w L = w - \sum_{i=1}^l \lambda_i y_i x_i = 0 \rightarrow w = \sum_{i=1}^l \lambda_i y_i x_i$$

$$\nabla_b L = \sum_{i=1}^l \lambda_i y_i = 0$$

$$\nabla_{\xi_i} L = C - \lambda_i - \mu_i = 0 \rightarrow \lambda_i + \mu_i = C$$

Условие дополняющей нежесткости:

$$\lambda_i (y_i(< w, x_i > + b) - 1 + \xi_i) = 0 \rightarrow \lambda_i = 0 \text{ или } (y_i(< w, x_i > + b) - 1 + \xi_i) = 0$$

$$\mu_i \xi_i = 0 \rightarrow \mu_i = 0 \text{ или } \xi_i = 0$$

Свойства лагранжиана:

$$\lambda \geq 0, \mu \geq 0$$

(b) Типы объектов

- i. $\lambda_i = 0 \rightarrow \mu_i = C \rightarrow \xi_i = 0 \rightarrow x_i$ лежит с правильной стороны от разделяющей гиперплоскости и на достаточном расстоянии от нее. $w = \sum_{i=1}^l \lambda y_i x_i \rightarrow$ объект не влияет на веса. Называется **периферийный**.
- ii. $0 < \lambda_i < 1 \rightarrow \mu \neq 0 \rightarrow \xi_0 = 0$. x_i не залезает на разделяющую полосу, но $y_i(< w, x_i > +b) = 1 \rightarrow x_i$ лежит прямо на границе. Дает вклад в w . x_i - **опорный граничный**.
- iii. $\lambda_i = C \rightarrow \xi_i > 0$. x_i дает вклад в w . $\xi_i > 0 \rightarrow x_i$ нарушает границу - **Опорные нарушители**.

- (c) Подставляем $w = \sum_{i=1}^l \lambda y_i x_i$ в лагранжиан, учтем ограничения $\sum_{i=1}^l \lambda_i y_i = 0$ и $C - \lambda_i - \mu_i = 0$
Двойственная задача SVM

$$\begin{cases} L = \sum_{i=1}^l \lambda_i - \frac{1}{2} \sum_{i,j=1}^l \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \max_{\lambda} \\ \sum_{i=1}^l \lambda_i y_i = 0 \\ 0 \leq \lambda_i \leq C \end{cases}$$

- (d) Если λ - решение, то $w = \sum_{i=1}^l \lambda_i y_i x_i$ - решение исходной задачи
- (e) Задача зависит от объектов только через скалярное произведение \rightarrow можно заменить его на ядро
- (f) Находим b Берем $x_i : 0 < \lambda_i < C \rightarrow \xi_i = 0 \rightarrow y_i(< w, x_i > +b) = 1 \rightarrow b = y_i - \langle w, x_i \rangle$
- (g) Минусы ядрового SVM
- i. Сложно контролировать переобучение
 - ii. Необходимо хранить в памяти матрицу Грамма
 - iii. Нельзя менять функцию потерь

11. Применение ядерной модели

- (a) $a(x) = \text{sign}(\langle w, x \rangle + b) = \text{sign}(\langle \sum_{i=1}^l \lambda y_i x_i, x \rangle + b) = \text{sign}(\sum_{i=1}^l \lambda_i y_i \langle x_i, x \rangle + b)$

1.1 Семинар: Задачи условной оптимизации

Учебник: *Boyd, Convex Optimization*

$$\begin{cases} f_0(x) \rightarrow \min_{x \in R^d} \\ f_i(x) \leq 0, i = 1, \dots, m \\ h_i(x) = 0, i = 1, \dots, p \end{cases}$$

$$G(x) = f_0(x) + \sum_{i=1}^m I_-(f_i(x)) + \sum_{i=1}^p I_0(h_i(x)) \rightarrow \min$$

Штрафы за нарушение ограничений:

$$I_-(z) = \begin{cases} 0, z \leq 0 \\ +\infty, z > 0 \end{cases}$$

$$I_0 = \begin{cases} 0, z = 0 \\ +\infty, z \neq 0 \end{cases}$$

$G(x) \rightarrow \infty$ в точках где не выполняется условие

Проблема: Недифференцируема

Заменяем функции на их аппроксимации ($\hat{I}_- = ax$)

Лагранжиан:

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

$$\lambda_i \geq 0$$

x - прямые (primal) переменные

λ, ν - двойственные переменные

Двойственная функция

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu)$$

- Двойственная функция всегда вогнутая
- Дает нижнюю оценку на минимум функции в прямой задаче
 x' - допустимая точка (все условия выполнены)

$$L(x', \lambda, \nu) = f_0(x') + \sum_{i=1}^m \lambda_i f_i(x') + \sum_{i=1}^p \nu_i h_i(x')$$

$$f_i(x) \leq 0, h_i(x) = 0 \rightarrow L(x', \lambda, \nu) \leq f_0(x')$$

$$\inf_x L(x, \lambda, \nu) \leq \inf_{x'} L(x', \lambda, \nu) \leq \inf_{x'} f_0(x')$$

\uparrow - это и есть решение исходной задачи

$$g(\lambda, \nu) \leq f_0(x_*)$$

$$g(\lambda, \nu) \rightarrow \max_{\lambda, \nu}, \lambda_i \geq 0$$

λ^*, ν^* - решение двойственной задачи

$g(\lambda^*, \nu^*) \leq f_0(x_*)$ - слабая двойственность

$g(\lambda^*, \nu^*) = f_0(x_*)$ - сильная двойственность

Достаточное условие сильной двойственности (Условие Слейтера)

– Задача выпуклая:

f_0, f_1, \dots, f_m - выпуклые

h_1, \dots, h_p - линейные

– $\exists x'$, что все ограничения выполнены строго

Пусть имеет место сильная двойственность:

$$g(\lambda^*, \nu^*) = f_0(x_*)$$

$$g(\lambda^*, \nu^*) = \inf_x (f_0(x) + \sum \lambda^* f_i(x) + \sum \nu^* h_i(x)) \leq f_0(x_*) + \sum \lambda^* f_i(x_*) + \sum \nu^* h_i(x_*) \leq f_0(x_*)$$

Все неравенства являются равенствами:

- Если решить безусловную задачу при подставлении λ^*, ν^* , то получим решение прямой задачи
- $\lambda_i^* f_i(x^*) = 0$ - условие дополняющей нежесткости

Теорема Куна-Такера

Необходимые условия для

$$\begin{cases} \nabla_x L(x_*, \lambda^*, \nu^*) = 0 \\ f_i(x) \leq 0 \\ h_i(x) = 0 \\ \lambda_i \geq 0 \\ \lambda_i f_i(x_*) = 0 \\ \text{Сильная двойственность} \end{cases} \Leftrightarrow x_*, \lambda^*, \nu^* \text{ решения}$$

2 Аппроксимации ядер, ЕМ алгоритм

Скалярные произведения тяжело хранить из-за размера матрицы.

Есть ли возможность построить $\tilde{\phi}(x) \rightarrow \langle \tilde{\phi}(x_i), \tilde{\phi}(x_j) \rangle \approx K(x_i, x_j)$

2.1 Метод случайных признаков Фурье

$$K(x, z) = K(x - z)$$

K - непрерывная функция

Теорема Бохнера

$$K(x - z) \rightarrow \exists p(w) \rightarrow K(x - z) = \int_{R^d} p(w) e^{iw^T(x-z)} dw$$

Используем:

$$K(x-z) = \int_{R^d} p(w) e^{iw^T(x-z)} dw \xrightarrow{\text{Формула Эйлера}^1} \int_{R^d} p(w) \cos(w^T(x-z)) + i \int_{R^d} p(w) \sin(w^T(x-z)) dw$$

$$\xrightarrow{K(x-z) - \text{веществ.}} \text{Комплексная часть} = 0 \rightarrow K(x-z) = \int_{R^d} p(w) \cos(w^T(x-z)) dw$$

$$\xrightarrow{\text{Монте-Карло}^2} K(x-z) \approx \{w_j \sim p(w)\} : \frac{1}{n} \sum_{i=1}^n \cos w_j^T(x-z)$$

$$= \frac{1}{n} \sum_{i=1}^n \cos w_j^T x \cos w_j^T z + \sin w_j^T x \sin w_j^T z$$

$$\tilde{\phi}(x) = \frac{1}{n} (\cos w_1^T x, \dots, \cos w_n^T x, \sin w_1^T x, \dots, \sin w_n^T x)$$

$$K(x-z) = \langle \tilde{\phi}(x), \tilde{\phi}(z) \rangle$$

Для гауссова ядра:

$$p(w) = \mathcal{N}(0, 1)$$

¹ $e^{ix} = \cos x + i \sin x$

² $\int_a^b f(x) dx = \frac{b-a}{n} \sum_{i=1}^N f(u_i)$

2.2 ЕМ алгоритм

Смесь распределений:

$$\begin{cases} p(x) = \sum_{k=0}^K \pi_k p_k(x) \\ \sum \pi_k = 1 \end{cases}$$

Вероятностный эксперимент:

Выбираем K из $[\pi_1, \dots, \pi_K]$, выбираем x из $p_{i_k}(x)$

Z - скрытые переменные

$$Z = \{0, 1\}^K, \sum Z_k = 1$$

$$p(Z_k = 1) = \pi_k$$

$$p(z) = \prod_{k=1}^K \pi_k^{Z_k}$$

$$p(x \mid Z_k = 1) = p_k(x)$$

$$p(x \mid z) = \prod_{k=1}^K (p_k(x)^{Z_k})$$

$$p(x, z) = p(x \mid z)p(z) = \prod_{k=1}^K (\pi_k p_k(x))^{Z_k}$$

$$p(x) = \sum_{k=1}^K p(x, z = k) = \sum_{k=1}^K \pi_k p_k(x)$$

Вероятностная кластеризация:

$p_k(x)$ - распределение k -го кластера

$$x \rightarrow (p_1(x), \dots, p_K(x))$$

Хотим описать X смесью распределений

$$p(x) = \sum_{k=1}^K \pi_k \phi(x \mid \theta_k), \phi(x \mid \theta_k) \sim \mathcal{N}(\mu, \Sigma), \theta = (\mu, \Sigma)$$

Неполное правдоподобие:

$$\ln(P(X \mid \Theta)) = \sum_{i=1}^l \log \sum_{k=1}^K \pi_k \phi(x_i \mid \theta_k) \rightarrow \max_{\theta}$$

Логарифм многооптимальная функция - просто оптимизировать ее сложно

Используем функцию полного правдоподобия

$$\log(P, X \mid \Theta) = \sum_{i=1}^l \log \sum_{k=1}^K (\pi_k \phi(x_i \mid \theta_k))^{Z_k}$$

$$\sum_{i=1}^l \sum_{k=1}^K Z_{ik} (\log \pi_k + \log \phi(x_i \mid \theta_k)) \rightarrow \max_{\Theta}$$

Известно аналитическое решение для нормального распределения.
Не знаем Z_{ik}

$$\Theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$$

Используем метод ALS для поиска Z, Θ

1. Оптимизация по скрытым переменным

Апостериорное распределение: $p(Z \mid X, \Theta) = \frac{P(X, Z \mid \Theta)}{p(X \mid \Theta)}$

$$Z^* = \arg \max_Z p(Z \mid X, \Theta)$$

2. Оптимизировать по Θ

$$\log p(X, Z^* \mid \Theta) \rightarrow \max_{\Theta}$$

3. Повторять до сходимости

Можно лучше. Не гарантирует сходимости

ЕМ-алгоритм - метод обучения моделей со скрытыми переменными

ЕМ-алгоритм

1. Е-шаг - вычисляем $p(Z \mid X, \Theta)$ и запоминаем
2. М-шаг

$$E_{Z \sim p(Z \mid X, \Theta)} \log p(X, Z \mid \Theta) = \sum_Z p(Z \mid X, \Theta) \log p(X, Z \mid \Theta) \rightarrow \max_{\Theta}$$

Вывод EM-алгоритма

$$\log p(X | \Theta) = Z(q, \Theta) + KL(q || p)$$

$$L(q, \Theta) = \sum_Z q(Z) \log \frac{p(X, Z | \Theta)}{q(Z)}$$

$$KL(q || p) = - \sum_Z q(Z) \log \frac{p(Z | X, \Theta)}{q(Z)}$$
$$\forall q(Z)$$

$L(q, \Theta)$ - нижняя оценка

Берем $q(Z) = p(Z | X, \Theta)$ - получаем E-шаг

$L(q, \Theta) = \sum_{Z \sim q(Z)} p(Z) \log(\dots)$ - M-шаг

EM-алгоритм дает гарантии на рост правдоподобия