

Лекция 3

Линейные модели в задачах регрессии и классификации

Макаренко В.А., Габдуллин Р.А.

МГУ им. М.В. Ломоносова

24 января 2023

X – множество объектов,

Y – множество ответов,

$y : X \rightarrow Y$ – неизвестная зависимость.

Дано:

$\{x_1, x_2, \dots, x_\ell\} \subset X$ – обучающая выборка,

$y_i = y(x_i)$, $i = 1, \dots, \ell$ – известные ответы.

Найти:

$a : X \rightarrow Y$ – решающая функция, приближающая y на всём X .

Линейная модель регрессии

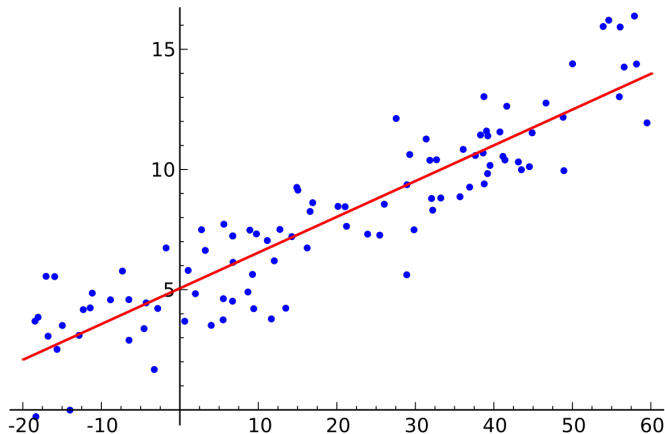


Рис.: Источник: [Википедия](#)

Линейная модель регрессии

- Семейство алгоритмов:

$$A = \{a(x, \theta) | \theta \in \mathbb{R}^{n+1}\},$$

$$a(x, \theta) = \theta_0 + \sum_{j=1}^n \theta_j f_j(x) = \sum_{j=0}^n \theta_j f_j(x),$$

если положить $f_0(x) \equiv 1$.

- Эмпирический риск:

$$Q(\theta, \mathbb{X}) = \sum_{i=1}^{\ell} w_i \cdot \left(y_i - a(x_i, \theta) \right)^2,$$

где w_i – вес, степень важности объекта i -го объекта.

- Метод наименьших квадратов (МНК):

$$\theta^* = \underset{\theta}{\operatorname{argmin}} Q(\theta, \mathbb{X}).$$

Метод максимального правдоподобия

- Вероятностная модель:

$$y_i = a(x_i, \theta) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2),$$

где $\{\varepsilon_i\}$ – независимые нормальные случайные величины.

- Функция правдоподобия ответов:

$$L(y_1, \dots, y_\ell | \theta) = \prod_{i=1}^{\ell} \frac{1}{\sqrt{2\pi\sigma_i^2}} \cdot \exp\left(-\frac{(y_i - a(x_i, \theta))^2}{2\sigma_i^2}\right).$$

- Метод максимального правдоподобия:

$$L(y_1, \dots, y_\ell | \theta) \rightarrow \max_{\theta} \iff -\ln L(y_1, \dots, y_\ell | \theta) \rightarrow \min_{\theta},$$

$$\ln L(y_1, \dots, y_\ell | \theta) = \text{const} + \frac{1}{2} \cdot \sum_{i=1}^{\ell} \frac{(y_i - a(x_i, \theta))^2}{\sigma_i^2},$$

$$\sum_{i=1}^{\ell} w_i \cdot (y_i - a(x_i, \theta))^2 \rightarrow \min_{\theta}, \quad w_i = \frac{1}{\sigma_i^2}.$$

- Многомерная линейная регрессия:

$$a(x, \theta) = \sum_{j=0}^n \theta_j f_j(x).$$

- Матричная запись:

$$F = \begin{pmatrix} f_0(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_0(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad y = \begin{pmatrix} y_1, \\ \dots, \\ y_\ell \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_0, \\ \dots, \\ \theta_n \end{pmatrix},$$

$$y = F\theta,$$

$$Q(\theta, \mathbb{X}) = \|y - F\theta\|^2 \rightarrow \min_{\theta}.$$

Какой геометрический смысл у МНК?

Аналитическое решение

- Многомерная линейная регрессия:

$$a(x_i, \theta) = \sum_{j=0}^n \theta_j F_{ij}.$$

- Необходимое условие минимума:

$$\frac{\partial Q}{\partial \theta_j} = -2 \sum_{i=1}^{\ell} (y_i - a(x_i, \theta)) \cdot \frac{\partial a(x_i, \theta)}{\partial \theta_j} = -2 \sum_{i=1}^{\ell} (y_i - a(x_i, \theta)) \cdot F_{ij} = 0,$$

то есть

$$\frac{\partial Q}{\partial \theta} = 2F^T(F\theta - y) = 0.$$

- Нормальная система уравнений:

$$F^T F \theta = F^T y.$$

- Решение нормальной системы уравнений:

$$\theta = (F^T F)^{-1} F^T y.$$

Градиентный спуск:

- Выбрать начальное приближение $\theta(0)$.
- Шаг в сторону антиградиента:

$$\theta(i+1) = \theta(i) - \alpha(i) \cdot \frac{\partial Q}{\partial \theta} \Big|_{\theta=\theta(i)} = \theta(i) - \alpha(i) \cdot 2F^T(F\theta(i) - y).$$

- Повторять до сходимости.

Варианты:

- Классический градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по всей выборке.
- Стохастический градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по одному наблюдению.
- Mini-batch градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по части выборки.

Численное решение

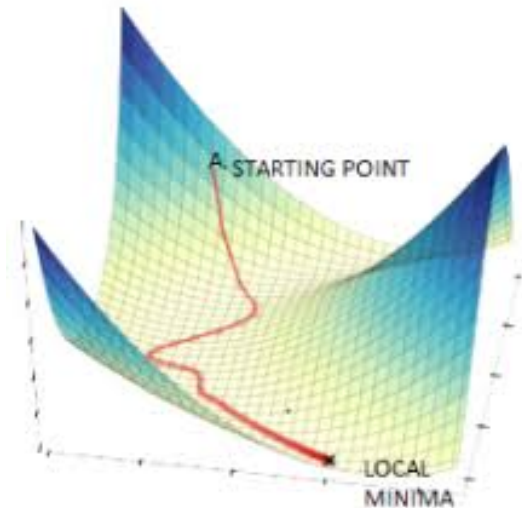


Рис.: Источник: datasciencecentral.com

Теорема (Гаусса-Маркова)

Пусть выполнены следующие условия:

- $y_i = a(x_i, \theta) + \varepsilon_i$.
- $\text{rank}(F) = n + 1$.
- $\mathbb{E}\varepsilon_i = 0$.
- $\text{cov}(\varepsilon) = \sigma^2 I$.

Тогда

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^{\ell} \left(y_i - a(x_i, \theta) \right)^2$$

является оптимальной оценкой в классе линейных оценок.

Проблема мультиколлинеарности

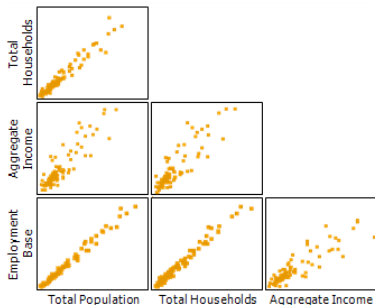


Рис.: Источник: medium.com

- Два или более признаков почти линейно зависимы.
- Решение получается неустойчивым.
- Неустойчивое решение ведет к переобучению.

Гребневая (Ridge) регрессия (L_2 -регуляризация)

- Эмпирический риск:

$$Q(\theta, \mathbb{X}) = \|y - F\theta\|^2 + \lambda\|\theta\|^2, \quad \lambda > 0.$$

- Необходимое условие минимума:

$$\frac{\partial Q}{\partial \theta} = 2F^T(F\theta - y) + 2\lambda\theta = 2((F^T F + \lambda I)\theta - F^T y) = 0,$$

где I – единичная матрица.

- Решение в явном виде:

$$\theta = (F^T F + \lambda I)^{-1} F^T y.$$

- Собственные числа матрицы $F^T F$ становятся больше на λ , при этом собственные векторы – те же самые.

Вероятностная интерпретация гребневой регрессии

- Вероятностная модель:

$$\theta \sim \mathcal{N}(0, \tau^2 I),$$

$$y_i = a(x_i, \theta) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2),$$

где $\{\varepsilon_i\}$ – независимые нормальные случайные величины.

- Апостериорное распределение весов:

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}.$$

- Максимум апостериорного распределения:

$$p(\theta|y) \rightarrow \max_{\theta} \iff -\ln p(y|\theta)p(\theta) \rightarrow \min_{\theta},$$

$$-\ln p(y|\theta)p(\theta) = \text{const} + \frac{1}{2} \cdot \left(\sum_{i=1}^{\ell} \frac{(y_i - a(x_i, \theta))^2}{\sigma^2} + \sum_{j=0}^n \frac{\theta_j^2}{\tau^2} \right),$$

то есть $\lambda = \frac{\sigma^2}{\tau^2}$.

Lasso-регрессия (L_1 -регуляризация)

- Эмпирический риск:

$$Q(\theta, \mathbb{X}) = \|y - F\theta\|^2 + \lambda \sum_{j=0}^n |\theta_j|, \quad \lambda > 0.$$

- Вероятностная интерпретация: веса независимы и имеют одно и то же распределение Лапласа.

Эквивалентные задачи поиска условного минимума.

- Для *Ridge*-регрессии (L_2):

$$\|y - F\theta\|^2 \rightarrow \min_{\theta}, \quad \sum_{j=0}^n \theta_j^2 \leq \kappa_1.$$

- Для *Lasso*-регрессии (L_1):

$$\|y - F\theta\|^2 \rightarrow \min_{\theta}, \quad \sum_{j=0}^n |\theta_j| \leq \kappa_2.$$

Ridge vs. Lasso

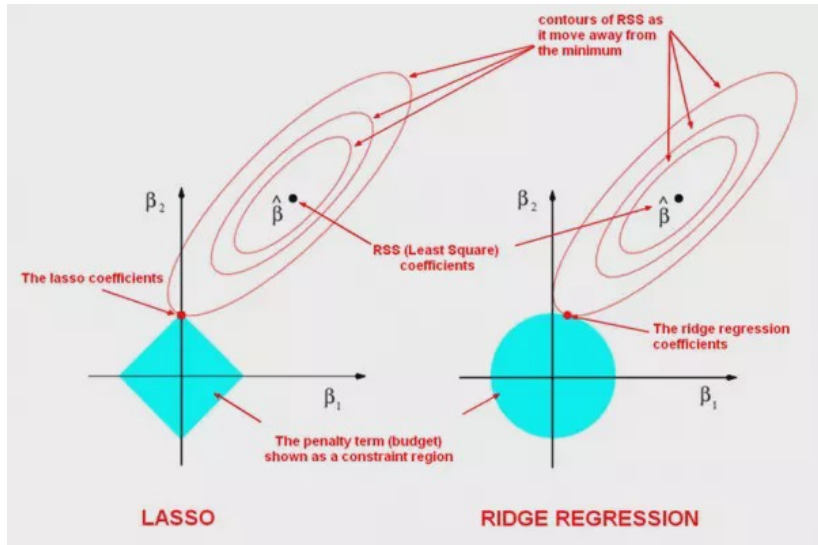


Рис.: Источник: medium.com

Ridge vs. Lasso

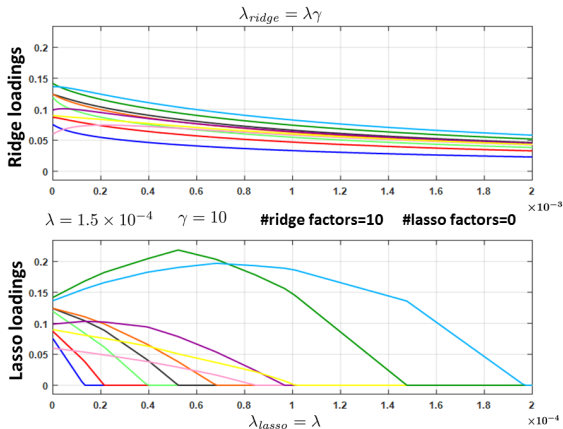


Рис.: Источник: arpm.co

Обязательно ли штрафовать за отклонение коэффициентов именно от нуля?

Масштабирование вещественных признаков

Веса модели чувствительны к сдвиг-масштабным преобразованиям признаков:

$$\tilde{f}_j(x) = \alpha_j + \beta_j f_j(x),$$

$$\begin{aligned}\theta_0 + \sum_{j=1}^n \theta_j \cdot \frac{\tilde{f}_j(x) - \alpha_j}{\beta_j} &= \left(\theta_0 - \sum_{j=1}^n \frac{\alpha_j \theta_j}{\beta_j} \right) + \sum_{j=1}^n \frac{\theta_j}{\beta_j} \cdot \tilde{f}_j(x) = \\ &= \tilde{\theta}_0 + \sum_{j=1}^n \tilde{\theta}_j \cdot \tilde{f}_j(x).\end{aligned}$$

Часто признаки преобразовывают, приводя их к единой шкале:

$$\tilde{f}_j(x) = \frac{f_j(x) - \mathbb{E}f_j(x)}{\sqrt{\mathbb{D}f_j(x)}} \quad \text{или} \quad \tilde{f}_j(x) = \frac{f_j(x) - \min_i f_j(x_i)}{\max_i f_j(x_i) - \min_i f_j(x_i)}.$$

Преобразование категориальных признаков

Пусть признак f_j может принимать одно из K возможных значений:

$$f_j(x) \in \{1, \dots, K\}.$$

One-hot кодирование. Признак f_j «разбивается» на $K - 1$ признаков:

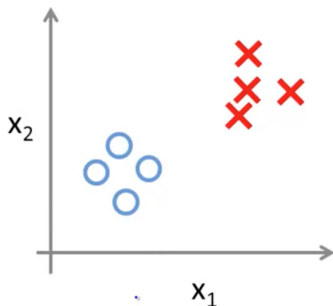
$$\tilde{f}_{j,k}(x) = [f_j(x) = k], \quad k = 1, \dots, K - 1.$$

Конструирование новых признаков на основе имеющихся:

- Применение функций к признакам (степени, логарифм, экспонента, ...).
- Добавление взаимодействий между признаками (перемножение, деление, ...).

Задача классификации

Binary classification:



Multi-class classification:

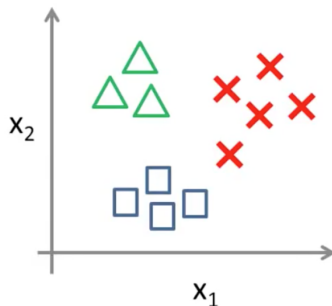


Рис.: Источник: medium.com

Модель бинарной классификации

- Множество ответов:

$$Y = \{-1, 1\}.$$

- Семейство вещественных дискриминантных функций:

$$S = \{s(x, \theta) | \theta \in \Theta\}.$$

- Семейство алгоритмов:

$$a(x, \theta) = \text{sign } s(x, \theta).$$

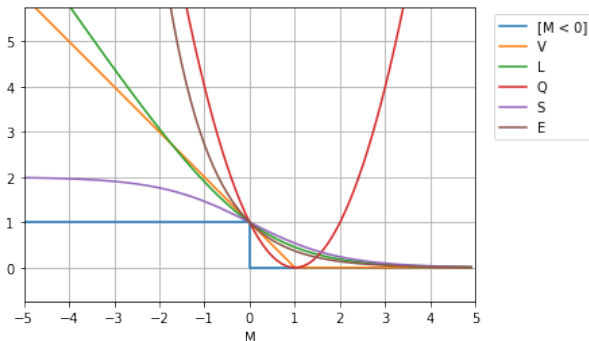
- Эмпирический риск:

$$Q(\theta, \mathbb{X}) = \sum_{i=1}^{\ell} [M(x_i, \theta) < 0] \equiv \sum_{i=1}^{\ell} [y_i \cdot s(x_i, \theta) < 0].$$

- Минимизация мажоранты эмпирического риска:

$$Q(\theta, \mathbb{X}) = \sum_{i=1}^{\ell} [M(x_i, \theta) < 0] \leq \sum_{i=1}^{\ell} L(M(x_i, \theta)) \rightarrow \min_{\theta}.$$

Мажоранты эмпирического риска



Часто используемые функции потерь L :

- $V(M) = (1 - M)_+$
- $L(M) = \log_2(1 + e^{-M})$
- $Q(M) = (1 - M)^2$
- $S(M) = 2(1 + e^M)^{-1}$
- $E(M) = e^{-M}$

Вероятностная модель бинарной классификации

Постановка задачи.

- Объекты: $\{x_i\}_{i=1}^{\ell}$.
- Ответы:

$$\begin{aligned} y_1, y_2, \dots, y_{\ell} &- \text{н.с.в.}, \\ y_i &\sim \text{Be}(p(x_i, \theta)), \quad i = 1, \dots, \ell, \quad \theta \in \Theta, \\ \mathbb{P}(y_i = y | \theta) &= p(x_i, \theta)^y \cdot (1 - p(x_i, \theta))^{1-y} \end{aligned}$$


- Оценить параметр $\theta \in \Theta$.

Логарифм функции правдоподобия ответов:

$$\ln L(y_1, \dots, y_{\ell} | \mathbb{X}, \theta) = \sum_{i=1}^{\ell} \left(y_i \ln p(x_i, \theta) + (1 - y_i) \ln(1 - p(x_i, \theta)) \right).$$

Задача оптимизации:

$$\sum_{i=1}^{\ell} L(y_i, p(x_i, \theta)) \rightarrow \min_{\theta},$$

где $L(y, a) = -(y \ln a + (1 - y) \ln(1 - a))$ – функция потерь Log Loss. 

Log Loss и оценка вероятностей

Функция потерь Log Loss:

$$L(y, a) = -\left(y \ln a + (1 - y) \ln(1 - a)\right), \quad y \in \{0, 1\}, \quad a \in [0, 1].$$

Пусть $Y \sim \text{Be}(p)$, тогда

$$a^* = \underset{a \in [0, 1]}{\operatorname{argmin}} \mathbb{E} L(Y, a) = \underset{a \in [0, 1]}{\operatorname{argmax}} (p \ln a + (1 - p) \ln(1 - a)).$$

Имеем:

$$\frac{\partial(p \ln a + (1 - p) \ln(1 - a))}{\partial a} = \frac{p}{a} - \frac{1 - p}{1 - a},$$

$$\frac{p}{a^*} - \frac{1 - p}{1 - a^*} = 0 \iff a^* = p.$$

Таким образом, $p(x, \theta^*)$ – оценка $\mathbb{P}(y(x) = 1)$.

Пороговая модель бинарной классификации

- Модель ответов:

$$y_i = [s(x_i, \theta) + \varepsilon_i > 0],$$

где $\{\varepsilon_i\}$ - н.о.р.с.в. с абсолютно непрерывной симметричной ф.р. F_ε ,

$$p(x_i, \theta) = \mathbb{P}(s(x_i, \theta) + \varepsilon_i > 0) = \mathbb{P}(\varepsilon_i < s(x_i, \theta)) = F_\varepsilon(s(x_i, \theta)).$$

- Модель инвариантна относительно масштабирования:

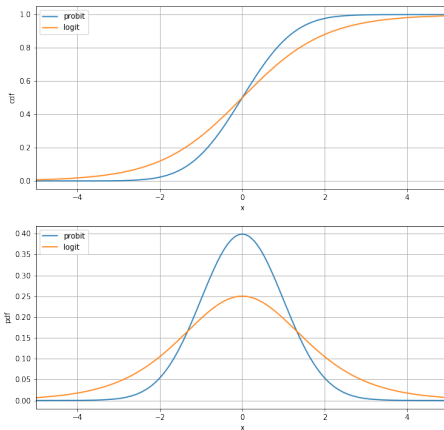
$$y_i = [\alpha \cdot (s(x_i, \theta) + \varepsilon_i) > 0] = [s(x_i, \theta) + \varepsilon_i > 0], \quad \alpha > 0,$$

поэтому можно зафиксировать любой удобный масштаб для $\{\varepsilon_i\}$.

- Задача оптимизации (минимизация Log Loss):

$$-\sum_{i=1}^{\ell} \left(y_i \ln F_\varepsilon(s(x_i, \theta)) + (1 - y_i) \ln(1 - F_\varepsilon(s(x_i, \theta))) \right) \rightarrow \min_{\theta}.$$

Модели остатков



Примеры моделей остатков:

- Probit-модель: $F_{\varepsilon}(x) = \Phi(x)$.
- Logit-модель: $F_{\varepsilon}(x) = \sigma(x) = (1 + e^{-x})^{-1}$.

Logit-модель и эмпирический риск

Минимизация Log Loss:

$$-\sum_{i=1}^{\ell} \left(y_i \log \left(\frac{1}{1 + \exp\{-s(x_i, \theta)\}} \right) + (1 - y_i) \log \left(\frac{1}{1 + \exp\{s(x_i, \theta)\}} \right) \right) \rightarrow \min_{\theta}.$$

Потери на одном наблюдении:

$$L(x_i, \theta) = \begin{cases} \log(1 + \exp\{-s(x_i, \theta)\}), & y = 1 \\ \log(1 + \exp\{s(x_i, \theta)\}), & y = 0 \end{cases} = \log(1 + \exp\{-M(x_i, \theta)\}).$$

Эмпирический риск:

$$Q(\mathbb{X}, \theta) = \sum_{i=1}^{\ell} \log(1 + \exp\{-M(x_i, \theta)\}).$$

Логистическая регрессия

- Семейство дискриминантных функций:

$$s(x, \theta) = \theta_0 + \sum_{j=1}^n \theta_j \cdot f_j(x) = \sum_{j=0}^n \theta_j \cdot f_j(x),$$

если положить $f_0(x) \equiv 1$.

- Семейство алгоритмов:

$$a(x, \theta) = [s(x, \theta) > 0].$$

- Оценка вероятности принадлежности позитивному классу:

$$\mathbb{P}(y(x) = 1) \approx \sigma(s(x, \theta)).$$

- Задача оптимизации (минимизация Log Loss):

$$-\sum_{i=1}^{\ell} \left(y_i \ln \sigma(s(x_i, \theta)) + (1 - y_i) \ln(1 - \sigma(s(x_i, \theta))) \right) \rightarrow \min_{\theta}.$$

Почему лог. регрессию нужно обязательно учить с регуляризатором?

Обучение модели логистической регрессии

Производная сигмоиды:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)).$$

Производные дискриминантной функции по параметрам:

$$\frac{\partial s(x, \theta)}{\partial \theta_j} = f_j(x).$$

Производная Log Loss по предсказанной вероятности:

$$\frac{\partial L(y, a)}{\partial a} = \frac{1 - y}{1 - a} - \frac{y}{a}.$$

Обучение модели логистической регрессии

$$\sigma'(x) = \sigma(x)(1-\sigma(x)), \quad \frac{\partial s(x, \theta)}{\partial \theta_j} = f_j(x), \quad \frac{\partial L(y, a)}{\partial a} = \frac{1-y}{1-a} - \frac{y}{a}.$$

Объединяем:

$$\begin{aligned} \frac{\partial L(y, \sigma(s(x, \theta)))}{\partial \theta_j} &= \frac{\partial L}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial s} \cdot \frac{\partial s}{\partial \theta_j} = \\ &= \left(\frac{1-y}{1-\sigma(s(x, \theta))} - \frac{y}{\sigma(s(x, \theta))} \right) \cdot \sigma(s(x, \theta)) \cdot (1-\sigma(s(x, \theta))) \cdot f_j(x) = \\ &= \left(\sigma(s(x, \theta))(1-y) - (1-\sigma(s(x, \theta)))y \right) \cdot f_j(x) = \left(\sigma(s(x, \theta)) - y \right) \cdot f_j(x). \end{aligned}$$

Обучение модели логистической регрессии

Обозначим:

$$F = \begin{pmatrix} f_0(x_1) & f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots & \dots \\ f_0(x_\ell) & f_1(x_\ell) & \dots & f_n(x_\ell) \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \dots \\ \theta_n \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}$$

Для $u \in \mathbb{R}^m$ обозначим:

$$\sigma(u) = \begin{pmatrix} \sigma(u_1) \\ \sigma(u_2) \\ \dots \\ \sigma(u_m) \end{pmatrix}.$$

Обучение модели логистической регрессии

Эмпирический риск:

$$Q(\mathbb{X}, \theta) = \sum_{i=1}^{\ell} L(y_i, s(x_i, \theta))$$

Производные по параметрам:

$$\frac{\partial Q}{\partial \theta_j} = \sum_{i=1}^{\ell} \left(\sigma(s(x_i, \theta)) - y_i \right) \cdot f_j(x_i)$$

В матричных обозначениях:

$$\frac{\partial Q}{\partial \theta} = F^T (\sigma(F\theta) - y).$$

Ср. с линейной регрессией:

$$\frac{\partial Q}{\partial \theta} = 2F^T (F\theta - y).$$

Обучение модели логистической регрессии

Градиентный спуск:

- Выбрать начальное приближение $\theta^{(0)}$.
- Шаг в сторону антиградиента:

$$\theta^{(i+1)} = \theta^{(i)} - \alpha^{(i)} \cdot \left. \frac{\partial Q}{\partial \theta} \right|_{\theta=\theta^{(i)}} = \theta^{(i)} - \alpha^{(i)} \cdot F^T (\sigma(F\theta^{(i)}) - y).$$

- Повторять до сходимости.

Варианты:

- Классический градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по всей выборке.
- Стохастический градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по одному наблюдению.
- Mini-batch градиентный спуск: на каждой итерации делаем шаг в сторону антиградиента эмпирического риска по части выборки.

Множественная классификация. One-vs-all

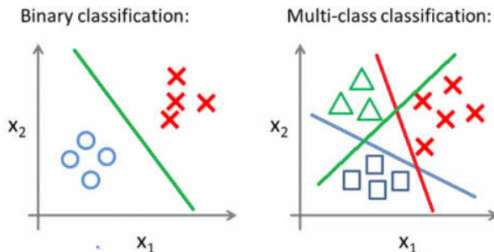


Рис.: Источник: towardsdatascience.com

- Для k -го класса обучаем свою дискриминантную функцию $s(x, \theta_k)$, отделяющую объекты этого класса от всех остальных.
- Итоговый алгоритм:
$$a(x) = \underset{k}{\operatorname{argmax}} s(x, \theta_k).$$

Множественная логистическая регрессия

- Совместно учим K дискриминантных функций. У каждого свой набор весов:

$$s(x, \theta_k) = \sum_{j=0}^n \theta_{k,j} \cdot f_j(x).$$

- Итоговый алгоритм:

$$a(x) = \operatorname{argmax}_k s(x, \theta_k).$$

- Распределение моделируется с помощью Softmax:

$$\operatorname{Softmax}(s_1, \dots, s_K) = \left(\frac{\exp(s_1)}{\sum_k \exp(s_k)}, \dots, \frac{\exp(s_K)}{\sum_k \exp(s_k)} \right).$$

Как изменится Softmax, если ко всем компонентам добавить одно и то же число? Как связаны между собой бинарный случай и множественный с двумя классами?

Множественная логистическая регрессия

Обучение методом максимального правдоподобия ответов:

$$\begin{aligned} \ln L(y_1, \dots, y_\ell | \mathbb{X}, \theta_1, \dots, \theta_K) = \\ = \sum_{i=1}^{\ell} \left(s(x_i, \theta_{y_i}) - \ln \left(\sum_k \exp\{s(x_i, \theta_k)\} \right) \right) \rightarrow \max_{\theta_1, \dots, \theta_K}. \end{aligned}$$

Что делать, если у нас классы не являются
взаимоисключающими (например, предсказываем хэштеги для
для статей) ?

- Линейная модель регрессии
 - МНК и ММП, их связь
 - Аналитическое решение
 - Численное решение. Градиентный спуск
 - Проблема мультиколлинеарности
 - L_1 и L_2 -регуляризации
 - Вероятностный смысл регуляризации
 - Преобразование и конструирование признаков
- Задача классификации
 - Типы задач классификации
 - Постановка задачи. Эмпирический риск
 - Мажоранты эмпирического риска
 - Вероятностная модель бинарной классификации
 - Log Loss
 - Пороговая модель бинарной классификации
 - Probit и Logit
 - Логистическая регрессия
 - Обучение модели логистической регрессии
 - One-vs-all
 - Множественная логистическая регрессия