# Машинное обучение.

## Лекция 1.

Тема: Введение в машинное обучение.

① Основные определения и постановки задач.

Машинное обучение — это наука, изучающая способы извлечения закономерностей из ограниченного кол-ва примеров. Есть ещё много близких направлений, к анализу данных, но подо относить любую работу, связанную с извлечением инф. из данных.

Объектом $(x)$ мы будем называть то, для чего хотим сделать предсказание.

Множество всех возможных точек размещения называется пространством объектов и обозначается через $X$.

Величина которую мы хотим определять называется ответом или целевой переменной, а множество её значений — пространством ответов $Y$.

Каждый пример называется обучающим, а все их совокупность — обучающей выборкой

$$X = \{(x_1, y_1), \dots, (x_\ell, y_\ell)\}$$

где $x_1, \dots, x_\ell$ — обучающиеся объекты
а $\ell$ — их кол-во

Объекты — это некие абстрактные сущности которыми компьютеры не умеют оперировать напрямую.

---

Для дальнейшего анализа нам понадобится описать объекты с помощью некоторого харак- котором называют признаками.

Вектор всех признаков объекта $x$ называется признаковым описанием этого объекта.

На работе со слотами данными специализируется активно развивающееся сейчас глубинное обучение (deep learning).

Описанная задача яв-ся примером задачи обучения с учителем (supervised learning) а более конкретно задачей регрессии.

Виды задач с учителем:

1. $Y = \{0, 1\}$ — бинарная классификация
2. $Y = \{1, \dots, K\}$ — многоклассовая классиф-я
3. $Y = \{0, 1\}^K$ — многоклассовая классиф-я с пересекающимися классами
4. Частичное обучение — задача, в которой для одной части объектов обучающей выборки известны и признаки, и ответы, а для другой только признаки.

Существует также обучение без учителя — класс задач, где ответы неизвестны или вообще не существ, и требуется найти некоторые закономерности в данных лишь на основе признаковых описаний:

1. Кластеризация — задача разделения объектов на группы, обладающие некоторыми св-ми.
2. Оценивание плотности — задача приближения распределения объектов.
3. Визуализация — задача изображения много- мерных объектов в двумерном или трёхмерном

---

## пространстве.

4. Понижение размерности — задача генерации таких новых признаков, чтобы их меньше чем исходных, но при этом с их помощью задача решается лучше.

Сложные постановки — обучение с подкреплением (reinforcement learning) где алгоритм на каждом шаге наблюдает какую-то ситуацию, выбирает одно из доступных вид действий, получает некоторую награду и корректирует свою стратегию.

Объекты-признаки
$$X \in R^{\ell \times d}$$
$\ell$ - число объектов
$d$ - число признаков

Построение функции $a : X \to Y$, которая для любого объекта будет предсказывать ответ. Такая функция называется алгоритмом или моделью.

Если функционал устроен так, что его следует минимизировать, то логично назвать его функционал ошибки.

Крайне популярен функционал в задаче регрессии яв-ся среднеквадратичная ошибка:
$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} (a(x_i) - y_i)^2$$

Линейные модели
$$A = \{a(x) = w_0 + w_1 x_1 + \dots + w_d x_d \mid w_0, w_1, \dots, w_d \in R\}$$

---

$$MSE: \frac{1}{\ell} \sum_{i=1}^{\ell} \left( w_0 + \sum_{j=1}^{d} w_j x_{ij} - y_i \right)^2 \to \min_{w_0, w_1 \dots w_d}$$

Процесс поиска оптимального алгоритма называется обучением.

Контроль ёмкости семейства алгоритмов— чем меньше у нас данных для обучения, тем более простое семейство следует выбирать.

Этапы решения задачи машинного обучения:

1. Постановка задачи.
2. Выделение признаков.
3. Формирование выборки.
4. Выбор функционала ошибки.
5. Подготовка данных.
6. Построение модели.
7. Вычисление качества модели.

## Лекция 2.

Тема: Линейная регрессия.

① Линейные модели.

Модели сводятся к суммированию значений признаков с некоторыми весами:
$$a(x) = w_0 + \sum_{j=1}^{d} w_j x_j$$

Параметрами модели яв-ся веса или коэффициенты $w_j$.

Вес $w_0$ также называется свободным коэф-м или сдвигом (bias).

Линейная модель в более компактном виде:
$$a(x) = w_0 + \langle w, x \rangle$$
где $w = (w_1, ..., w_d)$ — вектор весов.

Достаточно часто используется след приём, позволяющий упростить запись ещё сильнее. Добавим к признаковому описанию каждого объекта $(d+1)$-й признак равный единице. Вес при этом признаке как раз будет иметь смысл свободного коэф-та и необходимость в слагаемом $w_0$ отпадает:
$$a(x) = \langle w, x \rangle$$

② Области применимости линейных моделей.

Категориальные признаки
Допустим категориальный признак $f_j(x)$ принимает значение из множества $C = \{c_1, ..., c_m\}$. Заменим его на $m$ бинарных признаков $b_1(x), ..., b_m(x)$:
$$b_i(x) = [f_j(x) = c_j] \quad \text{one-hot}$$

Отметим, что признаки $b_1(x), ..., b_m(x)$ яв-ся линейно зависимыми: для любого объекта выполнено
$$b_1(x) + ... + b_m(x) = 1$$

Если мы применим линейную модель к данным после one-hot кодирования признака
$$a(x) = w_1[f(x)=c_1] + ... + w_m[f(x)=c_m] + [\text{взаимодействие с другими признаками}]$$

Работа с текстом.
Найдём все слова, которые есть в нашей выборке текстов и пронумеруем их: $\{c_1, ..., c_m\}$. Будем кодировать текст по признакам $b_1(x), ..., b_m(x)$ где $b_i(x)$ равен как-бы вхождений слова $c_i$ в текст.

Линейная модель над токен признаками:
$$a(x) = w_1 b_1(x) + ... + w_m b_m(x) + ...,$$

Бинаризация числовых признаков
Чтобы линейная модель смогла сделать подходящую для модели бинаризуем признак. Для этого выберем некоторую сетку точек $\{t_1, ..., t_m\}$. Это потом быть равномерная сетка между минимальным и max значением признака. Чем сетка из эмпирических квантилей.

Добавим края точки $t_0 = -\infty$ и $t_{m+1} = +\infty$
$$b_i(x) = [t_{i-1} < x_j \le t_i], \quad i \in 1, ..., m+1$$

Линейная модель над этими признаками:
$$a(x) = w_1[t_0 < x_j \le t_1] + ... + w_m[t_m < x_j \le t_{m+1}] + ...,$$

③ Измерение ошибки в задачах регрессии
MSE. Основной способ измерить отклонение — посчитать квадрат разности
$$L(y, a) = (a - y)^2$$

Основанный на ней функционал называется среднеквадратичным отклонением (CO)
$$MSE(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} (a(x_i) - y_i)^2$$

Величина CO плохо интерпретируется поскольку не сохраняет ед. измерение — так, если мы учу в рублях, то MSE будет измеряется в квадратах рублей. Чтобы избежать этого используют корень из CO.
$$RMSE(a, X) = \sqrt{\frac{1}{\ell} \sum_{i=1}^{\ell} (a(x_i) - y_i)^2}$$

Среднеквадратичная ошибка подходит для сравнение всех моделей или для контроля качества во время обучения, но не позволяет сделать выводы о том, насколько хорошо данная модель решает задачу.

В таких ситуациях вместо CO полезно использовать коэф детерминации
$$R^2(a, X) = 1 - \frac{\sum_{i=1}^{\ell}(a(x_i)-y_i)^2}{\sum_{i=1}^{\ell}(y_i - \bar{y})^2}$$
где $\bar{y} = \frac{1}{\ell}\sum_{i=1}^{\ell} y_i$ — среднее значение целевой переменной.

MAE. Заменим квадрат отклонения на модуль:
$$L(y, a) = |a - y|$$
Соответствующий функционал называется средним абсолютным отклонением MAE
$$MAE(a, X) = \frac{1}{\ell}\sum_{i=1}^{\ell}|a(x_i) - y_i|$$
Если мы выбрали MSE в качестве функционала ошибки то получаем след задачу.

$$\frac{1}{\ell}\sum_{i=1}^{\ell}(a - y_i)^2 \to \min_a$$

Минимум достигается на среднем значении всех ответов.
$$a^*_{MSE} = \frac{1}{\ell}\sum_{i=1}^{\ell} y_i$$

Функционал MAE:
$$\frac{1}{\ell}\sum_{i=1}^{\ell}|a - y_i| \to \min_a$$

Медиана ответов:
$$a^*_{MAE} = median\{y_i\}_{i=1}^{\ell}$$

Huber loss. Функция потерь Хубера:
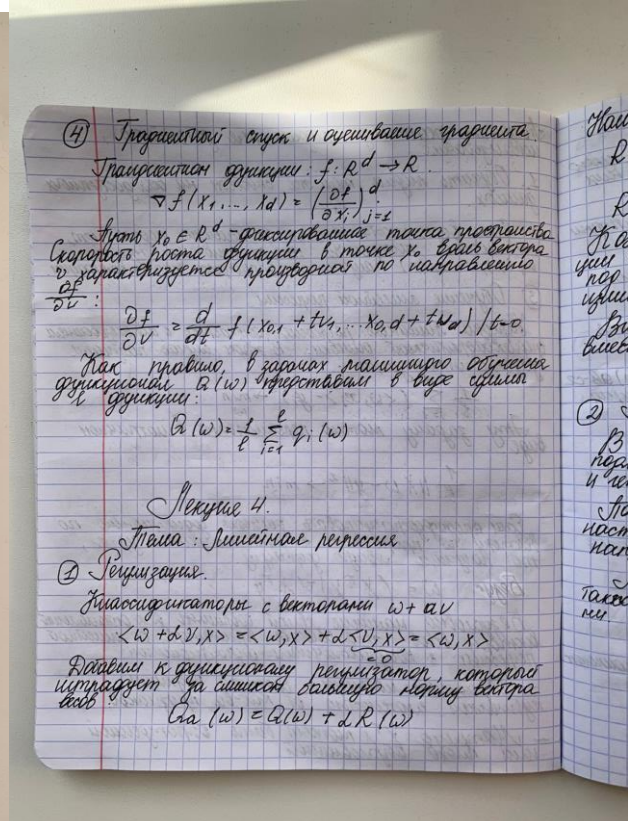$$L_\delta(y,a) = \begin{cases} \frac{1}{2}(y-a)^2, & |y-a| < \delta \\ \delta(|y-a| - \frac{1}{2}\delta), & |y-a| \ge \delta \end{cases}$$

Log-Cosh. Функция потерь log-cosh:
$$L(y,a) = \log\cosh(a - y)$$

MSLE. Среднеквадратичная логарифмическая ошибка
$$L(y,a) = (\log(a+1) - \log(y+1))^2$$

MAPE и SMAPE. Средний абсолютный процентной ошибки
$$L(y,a) = \left|\frac{y - a}{y}\right|$$

симметричная модификация
$$L(y,a) = \frac{|y-a|}{(|y|+|a|)/2}$$

Квантильная функция потерь
$$Q(a, X^\ell) = \sum_{i=1}^{\ell} \rho_\tau (y_i - a(x_i)).$$

где:
$$\rho_\tau(z) = (\tau-1)[z<0]z + \tau[z\geq 0]z = (\tau - \tfrac{1}{2})z + \tfrac{1}{2}|z|$$

Кросс валидация – это метод, предназначенный для оценки качества работы модели, широко применяемый в машинном обучении.
Он помогает сравнить между собой различные модели и выбрать наилучшую для конкретной задачи.

### Лекция 3.
Тема: Линейная регрессия.

① Переобучение

Нередко в машинном обучении модель оказывается переобученной – ее качество на новых данных существенно хуже качества на обучающей выборке.

Одномерная выборка, значения единственного признака $x$ в которой генерируются равномерно на отрезке $[0,1]$, а значения целевой переменной вычисляются по формуле
$$y = \cos(1{,}5\pi x) + N(0{,}001)$$

где $N(\mu, \sigma^2)$ – нормальное распределение со средним $\mu$ и дисперсией $\sigma^2$.

② Оценивание качества моделей.

Регрессионные кривые для признаков набора различной сложности.

Ridge
MSE = 4.08e-01 +/- 4.25e-01



На данной идее основан подход с отложенной выборкой.
На обучающей выборке, как это следует из названия, модель обучается, а на контрольной выборке проверяется ее качество. Если значение функционала на контрольной выборке оказалось удовлетворительным, то можно считать, что модель смогла извлечь закономерности при обучении.

С помощью кросс-валидации, размеченные данные разбиваются на $k$ блоков $X_1,\ldots,X_k$ примерно одинакового размера.
$$CV = \frac{1}{k}\sum_{i=1}^{k} a(a_i(x), X_i).$$

Получить финальную модель для дальнейшего использования.

1. Обучить модель тем же способом на всех доступных данных.

2. Если возможности обучить финальную модель нет, то можно построить композицию из моделей $a_1(x),\ldots,a_k(x)$ полученных в процессе кросс-валидации.

③ Обучение линейной регрессии.

Нередко линейная регрессия обучается с использованием среднеквадратичной ошибки. В этом случае получаем задачу оптимизации:
$$\frac{1}{\ell}\sum_{i=1}^{\ell} (\langle w, x_i\rangle - y_i)^2 \to \min_w$$

Эту задачу можно переписать в матричном виде
$$\frac{1}{\ell}\|Xw - y\|^2 \to \min_w$$

Если продифференцировать данный функционал по вектору $w$, приравнять к нулю и решить ур-х, то получим явную формулу.
$$w = (X^T X)^{-1} X^T y$$

Безусловно, наличие явной формулы для оптимального вектора весов – это большое преимущество линейной регрессии с квадратичными функциями.

Ограничение матрицы – слишком большая емкость от кол-ва признаков.
- Матрица $X^T X$ может быть вырожденной или плохо обусловленной.

④ Градиентный спуск и оценивание градиента.

Градиент функции: $f: R^d \to R$.
$$\nabla f(x_1,\ldots, x_d) = \left(\frac{\partial f}{\partial x_j}\right)_{j=1}^d$$

Пусть $x_0 \in R^d$ – фиксированная точка пространства. Скорость роста функции в точке $x_0$ вдоль вектора $v$ характеризуется производной по направлению $\frac{\partial f}{\partial v}$:
$$\frac{\partial f}{\partial v} = \frac{d}{dt} f(x_{01} + tv_1,\ldots, x_{0,d} + tv_d)\big|_{t=0}$$

Как правило, в задачах машинного обучения функционал $Q(w)$ представим в виде суммы функций:
$$Q(w) = \frac{1}{\ell}\sum_{i=1}^{\ell} q_i(w)$$

### Лекция 4.
Тема: Линейная регрессия

① Регуляризация.

Классификаторы с векторами $w + \alpha v$
$$\langle w + \alpha v, x\rangle = \langle w, x\rangle + \alpha\langle v, x\rangle = \langle w, x\rangle$$

Добавим к функционалу регуляризатор, который штрафует за слишком большую норму вектора весов
$$Q_\alpha(w) = Q(w) + \alpha R(w)$$

Наиболее распространёнными яв-ся $L_2$ и $L_1$ регуляризаторы

$$R(\omega) = \|\omega\|_2 = \sum_{i=1}^{d} \omega_i^2$$

$$R(\omega) = \|\omega\|_1 = \sum_{i=1}^{d} |\omega_i|$$

Коэффициент $\lambda$ называется параметр регуляризации и контролирует баланс между подгонкой под обучающую выборку и штрафом за излишнюю сложность.

Для решения при использовании $L_2$-регуляризации вместе со среднеквадратичной ошибкой:

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

② Гиперпараметры

В машинном обучении принято разделять подлежащие настройке величины на параметры и гиперпараметры.

Параметрами называют величины, которые настраиваются по обучающей выборке — например, веса линейной регрессии.

Линии уровня функционала качества, а также ограничение, задаваемое $L_2$ и $L_1$ регуляризаторами



---

③ Разреженные модели.

Модели, в которых некоторые веса равны нулю, называют разреженными, по скольку прогноз в них зависит лишь от части признаков.

1. Может быть заранее известно, что различающими яв-ся не все признаки.

2. К модели могут выдвигаться ограничения по скорости построения предсказаний.

3. В обучающей выборке объектов может быть существенно меньше, чем признаков

Можно показать, что если функционал $Q(\omega)$ яв-ся выпуклым, то задача безусловной минимизации функции $Q(\omega) + \lambda\|\omega\|$ эквивалента задаче условной оптимизации.

$$\begin{cases} Q(\omega) \to \min \\ \|\omega\|_1 \leqslant C \end{cases}$$

Найдём изменение $L_2$ и $L_1$ норм вектора при увеличении первой компоненты на некоторое положительное число $0 < \varepsilon$:

$$\|\omega - (\delta, 0)\|_2^2 = 1 - 2\delta + \delta^2 + \varepsilon^2$$

$$\|\omega - (\delta, 0)\|_1 = 1 - \delta + \varepsilon$$

Вычислим то же самое для изменения 2 компонент

$$\|\omega - (0, \delta)\|^2 = 1 - 2\varepsilon\delta + \delta^2 + \varepsilon^2$$

$$\|\omega - (0, \delta)\|_1 = 1 - \delta + \varepsilon$$

---

Проксимальные методы — это класс методов оптимизации, которые хорошо подходят для функционалов с негладкими слагаемыми.

$$\omega^{(k)} = S_{\lambda\alpha}\left(\omega^{(k-1)} - \alpha^2 \omega \nabla F(\omega^{(k-1)})\right)$$

$$F(\omega) = \|X\omega - y\|^2 - \text{функционал ошибки без регуляризатора}$$

$$S_{\lambda\alpha}(\omega_i) = \begin{cases} \omega_i - \lambda\alpha, & \omega_i > \lambda\alpha \\ 0, & |\omega_i| < \lambda\alpha \\ \omega_i + \lambda\alpha, & \omega_i < -\lambda\alpha \end{cases}$$

④ Преобразование признаков

С помощью линейной регрессии можно восстанавливать нелинейные зависимости, если провести преобразование признакового пространства.

$$x = (x_1, \ldots x_d) \to \psi(x) = (\varphi_1(x) \ldots \varphi_m(x)).$$

Переход к квадратичным признакам:

$$\psi(x) = (x_1, \ldots x_d, x_1^2, \ldots x_d^2, x_1, \ldots x_{d-1} x_d)$$

Другие преобразования

$\log x_i$ — для признаков с тяжёлыми хвостами

$\exp(\|x - \mu\|^2 / \sigma)$ — для измерения близости до некоторой точки

$\sin(x_j / T)$ — для задач с периодическими зависимостями.

---

## Лекция 5.

### Тема: Линейная классификация

① Линейные модели классификации.

Пусть $X = \mathbb{R}^d$ — пространство объектов, $Y = \{-1, +1\}$ множество допустимых ответов, $X = \{(x_i, y_i)\}_{i=1}^{\ell}$ обучающая выборка.

Линейная модель классификации определяется:

$$a(x) = \text{sign}(\langle \omega, x \rangle + \omega_0) = \text{sign}\left(\sum_{j=1}^{d} \omega_j x_j + \omega_0\right)$$

где $\omega \in \mathbb{R}^d$ — вектор весов

$\omega_0 \in \mathbb{R}$ — сдвиг

Соответствующий функционал называется долей правильных ответов

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i]$$

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) \neq y_i] = \frac{1}{\ell} \sum_{i=1}^{\ell} [\text{sign}\langle \omega, x_i \rangle \neq y_i]$$

Функционал несколько видоизменим

$$Q(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [y_i \langle \omega, x_i \rangle < 0] \to \min_{\omega}$$

$$M_i = y_i \langle \omega, x_i \rangle - \text{отступ}$$

Функционал оценивает ошибку алгоритма на объекте $x$ с помощью пороговой функции потерь $L(M) = [M < 0]$, где аргумент функции яв-ся отступ $M = y\langle \omega, x \rangle$.

$$Q(a,X) \leq \frac{1}{\ell} \sum_{i=1}^{\ell} \tilde{L}(y_i \langle \omega, x_i \rangle) \to \min_{\omega}$$

1. $\tilde{L}(M) = \log(1 + e^{-M})$ — логистическая функция потерь

2. $\tilde{L}(M) = (1-M)_+ = \max(0, 1-M)$ — кусочно-линейная функция потерь

3. $\tilde{L}(M) = (-M)_+ = \max(0, -M)$ — кусочно-линейная функция потерь

4. $\tilde{L}(M) = e^{-M}$ — экспоненциальная функция потерь

5. $\tilde{L}(M) = 2/(1 + e^M)$ — сигмоидная функция потерь

## ② Метрики качества классификации

Наиболее очевидной мерой качества в задаче класс-ии яв-ся доля правильных ответов, которую мы уже упоминали:

$$\text{accuracy}(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i]$$

Доля правильных ответов:

$$\text{accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

---

|  | $y = 1$ | $y = -1$ |
|---|---|---|
| $a(x) = 1$ | Tru Positive (TP) | False Positive (FP) |
| $a(x) = -1$ | False negative (FN) | True negative (TN) |

Гараздо более инф критериями яв-ся точность и полнота

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

Существует несколько способов получить один из критериев качества на основе точности и полноты. Один из них — F мера, гармоническое среднее точности и полноты

$$F = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Для измерения качества ранжирования нередко используют среднюю точность:

$$AP = \frac{1}{\ell_+} \sum_{k=1}^{\ell} [y_{(k)} = 1] \, \text{precision} @ k$$

Простота концентрации

$$\text{lift} = \frac{\text{precision}}{(TP + FN)/\ell}$$

Широко используется также интегральная метрика качества семейства, как площадь под ROC — кривой.

---

$$FPR = \frac{FP}{FP + TN} \quad ; \quad TPR = \frac{TP}{TP + FN}$$

## Лекция 6.

### Тема: Линейные классификации

Линейный классификатор основанный на минимизации верхней оценки:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i \langle \omega, x_i \rangle) \to \min_{\omega}$$

1. $L(y, z) = \log(1 + \exp(-y z))$ — логическая функция потерь

2. $L(y, z) = (1 - y z)_+$ — кусочно-линейная функция потерь

### ① Логическая функция потерь.

Пусть в каждой точке пространства объектов $x \in \mathbb{R}$ задана вероятность $p(y = +1 | x)$ того что объект $x$ будет принадлежать классу $+1$.

Если в выборке объект $x$ встречается $n$ раз с ответами $\{y_1 \ldots y_n\}$

$$\arg\min_{b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^{n} L(y_i, b) \approx p(y = +1 | x)$$

при стремлении $n$ к бесконечности

$$\arg\min_{b \in \mathbb{R}} E[L(y, b) | x] = p(y = +1 | x).$$

С точки зрения алгоритма вероятность того что $x$ в выборке встретится с классом $y$ равна $b(x_i)^{[y_i = +1]}(1 - b(x_i))^{[y_i = -1]}$

---

$$Q(a, X) = \prod_{i=1}^{\ell} b(x_i)^{[y_i = +1]}(1 - b(x_i))^{[y_i = -1]}$$

$$- \sum_{i=1}^{\ell} ([y_i = 1] \log b(x_i) + [y_i = -1] \log(1 - b(x_i)))$$

Матожидание функции потерь в точке $x$

$$E[L(y, b) | x] = E[-[y = +1] \log b - [y = -1] \log(1 - b) | x] =$$
$$= -p(y = +1 | x) \log b - (1 - p(y = +1 | x)) \log(1 - b)$$

Продифференцируем по $b$:

$$\frac{\partial}{\partial b} E[L(y, b) | x] = -\frac{p(y = +1 | x)}{b} + \frac{1 - p(y = +1 | x)}{1 - b} = 0$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$ Мы будем использовать сигмоидную функцию

$$p(y = 1 | x) = \frac{1}{1 + \exp(-\langle \omega \times \rangle)}$$

Выразим из нее скалярное произведение

$$\langle \omega, y \rangle = \log \frac{p(y = +1 | x)}{p(y = -1 | x)}$$

### ② Метод опорных векторов.

Пусть задан некоторый классификатор $a(x) = \text{sign}(\langle \omega \times \rangle + b)$, если одновременно умножить параметры $\omega$ и $b$ на одну и ту же положительную константу то классификатор не изменится

$$\min_{x \in X} |\langle \omega, x \rangle + b| = 1$$

Расстояние от произвольной точки $x_0 \in \mathbb{R}^d$ до гиперплоскости, определяемой данным классификатором, равно

$$p(x_0, a) = \frac{|\langle w, x \rangle + b|}{\|w\|}$$

$$\min_{x \in X} \frac{|\langle w, x \rangle + b|}{\|w\|} = \frac{1}{\|w\|} \min_{x \in X} |\langle w, x \rangle + b| = \frac{1}{\|w\|}$$

Приходим к оптимизационной задаче, соответствующей поиску оптимальных векторов для линейно неразделимой выборки

$$\begin{cases} \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \to \min_{w, b, \xi} \\ y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad i = 1 \dots \ell \\ \xi_i \geq 0, \quad i = 1, \dots \ell \end{cases}$$

Задача безусловной оптимизации

$$\begin{cases} \xi_i \geq 1 - y_i(\langle w, x_i \rangle + b) \\ \xi_i \geq 0 \end{cases}$$

Поскольку при этом в функционале требуется чтобы штрафа $\xi_i$ были как можно меньше то менее то можно получить след. явную формулу

$$\xi_i = \max(0, 1 - y_i(\langle w, x_i \rangle + b))$$

Данной выражение для $\xi_i$ уже учитывает в себе ограничения задачи

$$\frac{1}{2}\|w\|^2 + C \sum_{i=1}^{\ell} \max(0, 1 - y_i(\langle w, x_i \rangle + b)) \to \min_{w, b}$$

---

Лекция 7

Тема: Многоклассовая классификация и категориальные признаки

① Многоклассовая классификация

В данном разделе будем считать, что каждый объект относится к одному из $k$ классов.

$$Y = \{1 \dots K\}$$

Обучим $K$ линейных классификаторов $b_1(x)$, ... $b_k(x)$, выдающих оценки принадлежности классам $1, \dots K$ соответственно.

$$b_k(x) = \langle w_k, x \rangle + w_{0k}$$

Классификатор с номером $k$ будем обучать по выборке $(x_i, 2[y_i = k] - 1)_{i=1}^{\ell}$

Итоговый классификатор будет выдавать класс, соответствующий самому уверенному из бинарных алгоритмов:

$$a(x) = \arg\max_{k \in \{1 \dots K\}} b_k(x)$$

Обучим $C_K^2$ классификатор $a_{ij}(x)$, $i, j = 1, K$, $i \neq j$

$$b_k(x) = sgn(\langle w_k, x \rangle + w_{0k})$$

Классификатор $a_{ij}(x)$ будем настраивать по подвыборке $X_{ij} \subset X$, содержащей только объектов классов $i$ и $j$.

---

$$X_{ij} = \{(x_n, y_n) \in X \mid [y_n = i] = 1 \text{ или } [y_n = j] = 1\}$$

Оператор Soft Max $(z_1 \dots z_k)$

$$SoftMax(z_1 \dots z_k) = \left( \frac{\exp(z_1)}{\sum_{k=1}^{K} \exp(z_k)} \dots \frac{\exp(z_k)}{\sum_{k=1}^{K} \exp(z_k)} \right)$$

В этом случае вероятность $k$-го класса будет выражается как

$$P(y = k \mid x, w) = \frac{\exp(\langle w_k, x \rangle + w_{0k})}{\sum_{j=1}^{K} \exp(\langle w_j, x \rangle + w_{0j})}$$

Настроить $K$ набор параметров $w_1 \dots w_k$ и итоговый алгоритм определим как

$$a(x) = \arg\max_{k \in \{1 \dots k\}} \langle w_k, x \rangle$$

Рассмотрим след функцию потерь:

$$\max_k \{\langle w_k, x \rangle + 1 - [k = y(x)]\} - \langle w_{y(x)}, x \rangle$$

Анализ отступа для многоклассового случая об. норма Фробениуса матрица $w$, $k$-я строка которой совпадает с $w_k$.

$$\mathcal{P} = \frac{1}{\|w\|} = \frac{1}{\sum_{i=1}^{K} \sum_{j=1}^{\ell} w_{ki}^2}$$

Получаем:

$$\begin{cases} \frac{1}{2}\|w\|^2 \to \min \\ \langle w_{y_i}, x_i \rangle + [y_i = k] \langle w_k, x_i \rangle \geq 1, \quad i = 1 \dots \ell \end{cases}$$

---

Получим

$$\begin{cases} \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \to \min_{w, \xi} \\ \langle w_{y_i}, x_i \rangle + [y_i = k] - \langle w_k, x_i \rangle \geq 1 - \xi_i, \quad i = 1, \dots \ell \\ \xi_i \geq 0, \quad i = 1, \dots \ell \end{cases}$$

② Классификация с пересекающимися классами.

Для учета корреляций между классами можно воспользоваться следующим несложным подходом. Разобьем обучающую выборку $X$ на две части $X_1, X_2$. На первой части обучим $k$ независимых классификаторов $b_1(x), \dots b_k(x)$.

$$X_{ik}' = b_k(x_i), \quad x_i \in X_2$$

Существуют подходы, которые пытаются в рамках одной модели учитывать взаимосвязи между классами. Один из них предлагает преобразовать пространство ответов так, что классы оказались как можно менее зависимыми.

$$Y = U \Sigma V^T$$

Обозначим через $V_M$ матрицу состоящую из тех $M$ столбцов матрицу $V$, которые соответствуют наибольшим сингулярным числам. Спроецируем с ее помощью матрицу $Y$:

$$Y V_M = Y' \in R^{\ell \times M}$$

В таком ложной мере ошибки будут сильнее расставлены между этими множеством, то есть классов, факт принадлежности которым угадан неверно.

$$hamming(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \frac{|Y_i \setminus Z_i| + |Z_i \setminus Y_i|}{K}$$

Стандартные метрики качества классификации можно обобщить на multilabel-задачу так же как и на случай с непересекающимися классами — через микро- или макро-усреднение. Есть и несколько другой подход к обобщению основных метрик качества:

$$accuracy(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \left| \frac{Y_i \cap Z_i}{Y_i \cup Z_i} \right|$$

$$precision(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \left| \frac{Y_i \cap Z_i}{Z_i} \right|$$

$$recall(a, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \left| \frac{Y_i \cap Z_i}{Y_i} \right|$$

③ Категориальные признаки

Простейший способ — создать $n$ индикаторов, каждый из которых будет отвечать за одно из возможных значений признака. Таким образом, мы формируем $n$ бинарных признаков с условиями, но формируем $n$ бинарных признаков $g_1(x) \ldots g_n(x)$, которые определяются как

$$g_i(x) = [f(x) = u_i]$$

Определим наш способ кодирования. Вычислим для каждого значения $u$ категориального признака $(K+1)$ величин:

$$counts(u, X) = \sum_{(x,y) \in X} [f(x) = u]$$

$$successes_k(u, X) = \sum_{(x,y) \in X} [f(x) = u][y = k]$$

После того, как данные величины посчитаны, заменим наш категориальный признак $f(x)$ на $K$ вещественных $g_1(x) \ldots g_k(x)$:

$$g_k(x, X) = \frac{successes(f(x), X) + C_k}{counts(f(x), X) + \sum_{m=1}^{} C_m}, \quad k = 1 \ldots k$$

При использовании счётчиков нередко используют след. приёмы:

1. К признакам можно добавить не только дроби $g_k(x)$, но и значение $counts(f(x), X)$ и $successes_k(f(x), X)$

2. Можно генерировать парные категориальные признаки, т.е. для каждой пары категориальных признаков $f_i(x)$ и $f_j(x)$ создать новый признак $f_{ij}(x) = (f_i(x), f_j(x))$

3. Если у категориальных признаков много возможных значений, то хранение статистик $counts(u, X)$ и $successes(u, X)$ может потребовать сущ. кол-ва памяти.

4. Можно вычислять несколько счётчиков для разных значений параметров $C_1, \ldots C_k$.

5. Можно все редкие значение категориального признака объединить в одно, поскольку скорее всего, для редких значений не получится качественно оценить статистику $successes_k$ и $counts$.