

Машинное обучение

Семинар 3. Применение вероятностного аппарата в машинном обучении.



23 января 2018

Вероятностная задача

Возьмём какую-то болезнь. Она достаточно редкая — ей болеет примерно один человек на миллион. Есть аппарат, диагностирующий эту болезнь, который для здоровых людей с вероятностью 99.9% даёт отрицательный результат, а для больных с вероятностью 99.9% положительный. Во время ежегодного медицинского обследования человек проходит анализ на этом аппарате и получает положительный результат (диагностирована болезнь). Какова вероятность того, что он болен?

Формула Байеса

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}, \text{ где}$$

$P(A)$ — априорная вероятность события A ;

$P(A | B)$ — вероятность события A при условии события B ;

$P(B | A)$ — вероятность события B при условии события A ;

$P(B)$ — априорная вероятность наступления события B .

Формула Байеса

Если есть полный набор событий $(A_i)_{i=1}^n$, то:

$$P(B) = \sum_{i=1}^n P(B | A_i)P(A_i)$$

$$P(A_k | B) = \frac{P(B | A_k)P(A_k)}{\sum_{i=1}^n P(B | A_i)P(A_i)}$$

$$P(A | B) = \frac{P(B | A)P(A)}{P(B | A)P(A) + P(B | \bar{A})P(\bar{A})}$$

Формула Байеса

В нашей задаче A — человек болен, B — аппарат диагностировал болезнь, нужно найти $P(A | B)$.

$$P(A) = ?$$

$$P(B | A) = ?$$

Формула Байеса

В нашей задаче A — человек болен, B — аппарат диагностировал болезнь, нужно найти $P(A | B)$.

$$P(A) = 10^{-6}$$

$$P(B | A) = 1 - 10^{-3}$$

Формула Байеса

В нашей задаче A — человек болен, B — аппарат диагностировал болезнь, нужно найти $P(A | B)$.

$$P(A) = 10^{-6}$$

$$P(B | A) = 1 - 10^{-3}$$

$$P(A | B) = \frac{(1 - 10^{-3}) 10^{-6}}{(1 - 10^{-3}) 10^{-6} + 10^{-3} (1 - 10^{-6})} \approx 10^{-3}$$

Формула Байеса в повседневности

Если есть полный набор гипотез $(A_i)_{i=1}^n$, и вы оцениваете их шансы как

$$P(A_1): P(A_2): \dots : P(A_n).$$

Далее вы наблюдаете событие B , из-за чего соотношение шансов гипотез меняется, по формуле Байеса оно будет таким:

$$P(A_1)P(B|A_1): P(A_2)P(B|A_2): \dots : P(A_n)P(B|A_n).$$

То есть шаны умножаются на условные вероятности.

Стандартные понятия

Функция распределения случайной величины ξ — это функция $F(x) = P(\xi \leq x)$.

Плотность случайной величины ξ — это функция $f(x) = F'(x)$.

Математическое ожидание случайной величины ξ —
$$E\xi = \int_{-\infty}^{\infty} x f(x) dx.$$

Дисперсия случайной величины ξ — $D\xi = E\xi^2 - (E\xi)^2$

Когда это может пригодиться?

Вы делаете умную автобусную остановку. Она будет говорить сколько ещё осталось ждать автобуса. К сожалению, наладить расписание не удалось, поэтому автобусы приходят на остановку случайно, но с некоторой интенсивностью. Какие проблемы могут возникнуть при предсказании такой величины?

Когда это может пригодиться?

Вы делаете умную автобусную остановку. Она будет говорить сколько ещё осталось ждать автобуса. К сожалению, наладить расписание не удалось, поэтому автобусы приходят на остановку случайно, но с некоторой интенсивностью. Какие проблемы могут возникнуть при предсказании такой величины?

Время между автобусами будет распределено экспоненциально:

$f(x) = \lambda e^{-\lambda x} I\{x \geq 0\}$, где λ — параметр отвечающий за интенсивность

Для экспоненциально распределённой величины ξ значение $P(\xi - s \geq t \mid \xi \geq s)$ не зависит от s . Какие сложности это может вызвать?

Максимизация правдоподобия в задачах ML

Допустим, у нас есть выборка X , целевые переменные y и набор параметров θ . Логично, что нам бы хотелось выбрать наиболее вероятный вектор параметров, при условии того, что мы какие-то данные пронаблюдали:

$$P(\theta | X, y) \rightarrow \max_{\theta}.$$

По формуле Байеса:

$$P(\theta | X, y) \propto P(y | X, \theta)P(\theta | X) = P(y | X, \theta)P(\theta)$$

Последнее равенство выполнено, так как априорные представления о параметрах не зависят от данных, которые мы пронаблюдали.

Максимизация правдоподобия в задачах ML

Таким образом, задача

$$P(\theta | X, y) \rightarrow \max_{\theta}$$

эквивалентна

$$P(y | X, \theta)P(\theta) \rightarrow \max_{\theta}.$$

Поскольку все объекты выборки независимы, то

$$P(y | X, \theta) = \prod_{i=1}^n P(y_i | x_i, \theta).$$

Подставив в оптимизационную задачу и прологарифмировав получим:

$$\sum_{i=1}^n \log P(y_i | x_i, \theta) + \log P(\theta) \rightarrow \max_{\theta}.$$

Максимизация правдоподобия в задачах ML

Итоговая оптимизационная задача:

$$\sum_{i=1}^n \log P(y_i | x_i, \theta) + \log P(\theta) \rightarrow \max_{\theta}.$$

Сравните с минимизацией регуляризованного эмпирического риска:

$$\sum_{i=1}^n L(y_i, a(x_i)) + R(a) \rightarrow \min_a.$$

Таким образом, различные функции потерь и регуляризаторы приобретают вероятностный смысл условного и априорного распределений соответственно.

Максимизация правдоподобия в задачах ML

В линейных моделях, если положить $P(w) \sim N(0, \tau)$, то

$$\log P(w) = \log \left(\frac{1}{\sqrt{2\pi\tau}} e^{-\frac{\|w\|^2}{2\tau}} \right) = -\text{const}(w) - \frac{1}{2\tau} \|w\|^2,$$

что при максимизации по w соответствует l_2 регуляризации.

- ▶ Какое априорное распределение соответствует l_1 регуляризации?
- ▶ Какое условное распределение соответствует MSE функции потерь?
- ▶ Какое условное распределение соответствует MAE функции потерь?
- ▶ Какая функция потерь будет в задаче про автобусную остановку?

Максимизация правдоподобия в задачах ML

В линейных моделях, если положить $P(w) \sim N(0, \tau)$, то

$$\log P(w) = \log \left(\frac{1}{\sqrt{2\pi\tau}} e^{-\frac{\|w\|^2}{2\tau}} \right) = -\text{const}(w) - \frac{1}{2\tau} \|w\|^2,$$

что при максимизации по w соответствует l_2 регуляризации.

- ▶ Распределение Лапласа
- ▶ $y_i \sim N(a(x_i), 1)$
- ▶ $y_i \sim \text{Lap}(1, a(x_i))$
- ▶ $-\ln a(x_i) + a(x_i)y_i$

Байесовский классификатор

Дан вектор признаков $(x_j)_{j=1}^d$, по формуле Байеса:

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)}$$

То есть выбор наиболее вероятного класса эквивалентен задаче

$$P(x | y)P(y) \rightarrow \max_y,$$

так как знаменатель не зависит от y .

Наивный Байесовский классификатор

Наивный Байесовский классификатор предполагает, что

$$P(x | y) = \prod_{j=1}^d P(x^j | y)$$

так как знаменатель не зависит от y . Таким образом, задача определения класса становится:

$$P(y) \prod_{j=1}^d P(x^j | y) \rightarrow \max_y,$$

Наивный Байесовский классификатор

Если по выборке оценить условные вероятности и априорную вероятность, то мы сможем решать данную оптимизационную задачу для новых объектов x , подставив в формулу вместо вероятностей их оценки:

$$\hat{P}(y) \prod_{j=1}^d \hat{P}(x^j|y) \rightarrow \max_y,$$

однако, как найти эти оценки?

- ▶ $\hat{P}(y)$ — частота класса y в обучающей выборке
- ▶ $\hat{P}(x^j|y)$ — зависит от того, как устроены признаки

Наивный Байес в Sklearn

MultinomialNB

- ▶ Признаки категориальные
- ▶ $\hat{P}(x^j | y) = \frac{N_{yj} + \alpha}{N_y + d\alpha}$, где $N_{yj} = \sum_{i=1}^n x_i^j$, $N_y = \sum_{j=1}^d N_{yj}$, а α —
сглаживающий параметр
- ▶ Может интерпретироваться как линейная модель.

BernoulliNB

- ▶ Признаки бинарные
- ▶ Оценки аналогично MultinomialNB, но при применении
 $P(x^j | y) = P(j | y)x^j + (1 - P(j | y))(1 - x^j)$

Наивный Байес в Sklearn

BernoulliNB

- ▶ Признаки бинарные
- ▶ Оценки аналогично MultinomialNB, но при применении $P(x^j | y) = P(j | y)x^j + (1 - P(j | y))(1 - x^j)$

GaussianNB

- ▶ Признаки вещественные
- ▶ Предполагается, что признаки распределены нормально, но на практике это не критично
- ▶ $P(x^j | y) = \frac{1}{\sqrt{2\pi\sigma_{yj}^2}} e^{-\frac{(x^j - \mu_{yj})^2}{2\sigma_{yj}^2}}$, где μ_{yj} и σ_{yj}^2 — параметры нормального распределения, находятся методом максимального правдоподобия по обучающей выборке

Метод максимального правдоподобия и байесовские методы

Метод максимального правдоподобия для оценки параметра θ по выборке X :

$$P(\theta | X) \rightarrow \max_{\theta}.$$

То есть на самом деле мы ищем не всё распределение $P(\theta | X)$, а только его моду. В байесовских методах ищут непосредственно всё распределение $P(\theta | X)$. А уже потом, в зависимости от задачи, берут его моду, среднее или другую статистику. Это позволяет более детально учесть получаемую информацию, но вычислительно гораздо сложнее.

Байесовские методы

Применим формулу Байеса:

$$P(\theta | X) = \frac{P(X | \theta)P(\theta)}{\int P(X | \theta)P(\theta)d\theta}.$$

Заметьте, что так как нас интересует точное распределение, то мы не можем пренебречь знаменателем, как делали это раньше. С этим интегралом и связана основная проблема, так как числитель вычисляется легко в силу независимости объектов выборки.

Байесовские методы. Примеры

- ▶ Считаем счётчики при обучении MultinomialNB
- ▶ Подпросили монетку n раз, выпало k орлов, нужно оценить вероятность орла
- ▶ Оцениваем конверсию показа в клик рекламного блока

Во всех этих случаях, если данных немного мы бы применяли сглаживание. На самом деле у этой операции есть байесовское обоснование.

Байесовские методы. Сопряжённое распределение

$$P(\theta | X) = \frac{P(X | \theta)P(\theta)}{\int P(X | \theta)P(\theta)d\theta}.$$

В некоторых случаях интеграл в знаменателе можно найти аналитически. Для этого необходимо, чтобы априорное распределение на параметр θ было связано с условным распределением $P(X | \theta)$ определённым способом.

Если апостериорное распределение $P(\theta|X)$ принадлежит тому же семейству вероятностных распределений, что и априорное распределение $P(\theta)$ (т.е. имеет тот же вид, но с другими параметрами), то это семейство распределений называется сопряжённым семейству функций правдоподобия $p(X|\theta)$.

Байесовские методы. Пример сопряжённых распределений

Для случайной величины, распределённой по закону Бернулли (приведённые примеры) с неизвестным параметром q в качестве сопряжённого априорного распределения обычно выступает бета-распределение с плотностью вероятности:

$$f(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1-q)^{\beta-1},$$

где $B(\alpha, \beta)$ — бета функция т.е. $\int_0^1 q^{\alpha-1} (1-q)^{\beta-1} dq$

Байесовские методы. Пример сопряжённых распределений

Итак, $P(p)$ — бета распределение, а $P(x|q) = q^x(1 - q)^{(1-x)}$.
Тогда

$$\begin{aligned} P(X|q)P(q) &= P(q) \prod_{i=1}^n P(x_i|q) = \\ &= \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1 - q)^{\beta-1} \prod_{i=1}^n q^{x_i} (1 - q)^{(1-x_i)} = \\ &= \frac{1}{B(\alpha, \beta)} q^{\alpha-1+\sum_{i=1}^n x_i} (1 - q)^{\beta-1+n-\sum_{i=1}^n x_i} \end{aligned}$$

$$P(X|q)P(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1+\sum_{i=1}^n x_i} (1-q)^{\beta-1+n-\sum_{i=1}^n x_i}$$

Это означает, что

$$\begin{aligned} \int_0^1 P(X|q)P(q) dq &= \frac{1}{B(\alpha, \beta)} \int_0^1 q^{\alpha-1+\sum_{i=1}^n x_i} (1-q)^{\beta-1+n-\sum_{i=1}^n x_i} dq = \\ &= \frac{B(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i)}{B(\alpha, \beta)} \end{aligned}$$

Это означает, что $P(q | X)$ будет иметь Бета распределение с параметрами $\alpha + \sum_{i=1}^n x_i$ и $\beta + n - \sum_{i=1}^n x_i$

Байесовские методы. Пример сопряжённых распределений

По формуле математического ожидания Бета распределения имеем, что

$$E(q | X) = \frac{\alpha + \sum_{i=1}^n x_i}{\alpha + \beta + n},$$

то есть фактически то сглаживание, которое мы делали.