

Постановка задачи классификации

Задана $D_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ — выборка объёма m , где $y \sim P(y|\mathbf{x})$. Найдя $P(y|\mathbf{x})$, можно для каждого \mathbf{x} давать предсказания значения y следующим способом

$$\hat{y} = \operatorname{argmax}_{y \in \{-1, 1\}} P(y|\mathbf{x})$$

Чтобы восстановить $P(y|\mathbf{x})$ введём семейство функций $\mathcal{F} = \{f(\mathbf{x}, \mathbf{w}) | \mathbf{w} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{R}^n\}$ таких, что $f(\mathbf{x}, \mathbf{w}) = p(y|\mathbf{x})$. Задана логистическая регрессия, поэтому

$$f(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{x}^T \mathbf{w})}$$

а вектор оптимальных параметров \mathbf{w} вычисляется как доставляющий максимум правдоподобия данных. Совместное правдоподобие:

$$P(D_m, \mathbf{w}) = p(D_m | \mathbf{w}) p(\mathbf{w})$$

где $p(\mathbf{w})$ — априорное распределение векторов параметров, которое является нормальным $\mathcal{N}(\boldsymbol{\mu}, \mathbf{A})$. Можно вычислить вектор оптимальных параметров $\hat{\mathbf{w}}$

$$\mathcal{L}(D_m, \mathbf{w}) = \sum_{i=1}^m y_i \ln P(y_i | \mathbf{x}_i, \mathbf{w}) + \ln p(\mathbf{w}) = - \sum_{i=1}^m y_i \ln(1 + \exp(-\mathbf{x}_i^T \mathbf{w})) + \frac{1}{2C} \|\mathbf{w}\|_2^2 \rightarrow \max_{\mathbf{w} \in \mathbb{R}^n}$$

Для неизвестного объекта \mathbf{x} предсказываем целевую переменную следующим образом

$$y = \operatorname{argmax}_{y \in \{0, 1\}} \frac{1}{1 + \exp(-\mathbf{x}^T \hat{\mathbf{w}})}$$

Постановка задачи определения оптимального объёма выборки

Предположим задано m объектов выборки D , порождённой распределением $P(\mathbf{x}, y)$. Объём m заранее недостаточен для обучения устойчивой логистической регрессии. Необходимо, имея m объектов, спрогнозировать такое m^* , что каждый новый объект после m^* не будет давать новой информации о распределении параметров логистической регрессии, обученной на m^* объектах. Обозначим $\mathbb{P}(\mathbf{w} | m)$ распределение весов классификатора, обученного на выборке D_m порождённой $P(\mathbf{x}, y)$. Задача прогнозирования оптимального объёма выборки имеет вид

$$m^* = \min m \in \mathbb{N} \forall k \in \mathbb{N} \rightarrow \rho(P(\mathbf{w} | m), P(\mathbf{w} | m + k)) < \varepsilon$$

где ρ — некоторая функция расстояния между распределениями, а ε — заданный порог.

Базовый метод

В качестве расстояния между распределениями мы используем расстояние Кульбака-Лейблера. Предполагая, что выборка простая и поскольку целевая переменная порождается распределением Бернулли, распределение векторов параметров является нормальным [ссылка Стрижов]. Удаляя из выборки по l элементов, получим набор подвыборок размера $m - l$. Обучив на этих подвыборках логистические регрессии, получим набор оптимальных векторов параметров. Этот набор порождён нормальным распределением. Выбрав два

различных объёма m_1 и m_2 , мы генерируем два нормальных распределения векторов параметров с матожиданиями $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2$ и матрицами ковариации $\mathbf{A}_1, \mathbf{A}_2$. Таким образом

$$\rho(\mathbf{P}(\mathbf{w}|m), \mathbf{P}(\mathbf{w}|m+k)) = KL(\mathcal{N}(\hat{\mathbf{w}}_1, \mathbf{A}_1) || \mathcal{N}(\hat{\mathbf{w}}_2, \mathbf{A}_2)) =$$

$$\frac{1}{2} \left(\text{tr}(\mathbf{A}_2^{-1} \mathbf{A}_1) + (\hat{\mathbf{w}}_1 - \hat{\mathbf{w}}_2)^T \mathbf{A}_2^{-1} (\hat{\mathbf{w}}_1 - \hat{\mathbf{w}}_2) - n + \ln \left(\frac{\det(\mathbf{A}_2)}{\det(\mathbf{A}_1)} \right) \right)$$

Чтобы оценить оптимальный объём выборки поочерёдно добавляем к выборке по 1 элементу пока объём выборки не достигнет m . из выборки меняя l от $m-k$ до 1. Получим $m-k$ распределений $\mathbf{P}(\mathbf{w}|m-l)$. Чем больше $m-l$, тем точнее логистическая регрессия восстанавливает распределение $\mathbf{P}(y|\mathbf{x})$, поскольку обучена на большем количестве объектов. Обозначи m^* оптимальный объём выборки, т. е. такой, что по такой выборке распределение $\mathbf{P}(y|\mathbf{x})$ восстанавливается точно. Чем меньше расстояние $KL(\mathbf{P}(\mathbf{w}|m-l) || \mathbf{P}(\mathbf{w}|m^*))$, тем ближе $m-l$ к оптимальному размеру выборки. Поскольку m^* — оптимальный объём, $KL(\mathbf{P}(\mathbf{w}|m^*+p) || \mathbf{P}(\mathbf{w}|m^*))$ не меняется существенно при $p \rightarrow \infty$. Таким образом график зависимости $KL(\mathbf{P}(\mathbf{w}|m) || \mathbf{P}(\mathbf{w}|m^*))$ от m выходит на плато на m^* . По полученным $m-k$ распределениям аппроксимируем эту зависимость. Имея аппроксимирующую функцию, можно найти m^* как точку, на которой график функции выходит на плато.

*ЗАМЕТКА чтобы получить прогноз необходимого объёма, так или иначе, нужно разбираться с тем, как ведет себя матожидание и матрица ковариации как случайная величина, зависящая от размера выборки. Идея в том, что мы эмперически попробуем аппроксимировать функцию w от m , матрицу ковариации.