

Линейные модели регрессии

Викулин Всеволод

v.vikulin@corp.mail.ru

18 марта 2019

Повторение

Обучение — приобретение нужной функциональности посредством опыта.
(<https://postnauka.ru/video/53575>)

Множество объектов X , множество допустимых ответов Y . Существует $y : X \rightarrow Y$
Для $\{x_1, \dots, x_N\} = X_{train}$ известны $\{y_1, \dots, y_N\} = y_{train}$
Построить $a : X \rightarrow Y$, которая приближала бы y .

Пример

Строим рекомендательную систему музыки. Что является объектом?

Объекты описываются набором признаков. Признак объекта x — это результат измерения некоторой его характеристики.

В курсе разберем 2 постановки задачи обучения с учителем:

- Классификация – $Y = \{1, \dots, M\}$
- Регрессия – $Y = \mathbb{R}$

Пример

Строим рекомендательную систему музыки. Какой берем Y ? Какие берем признаки?

Минимизация эмпирического риска

Функция потерь (loss function) $L(a, x, y)$ – неотрицательная функция, показывающая величину ошибки алгоритма a на объекте x с ответом y .

Функционал качества $Q(a, X_{train}, y_{train}) = \frac{1}{N} \sum_{i=1}^N L(a, x_i, y_i), x_i \in X_{train}, y_i \in y_{train}$

Принцип минимизации эмпирического риска:

$$a^* = \underset{a}{\operatorname{argmin}} Q(a, X_{train}, y_{train})$$

Пример

Строим рекомендательную систему музыки. Какую берем функцию потерь?

Не обязательно, что $\underset{a}{\operatorname{argmin}} Q(a, X_{train}, y_{train})$ – полезный алгоритм. Можно просто запомнить обучающую выборку.

Проблема **переобучения** – значения $Q(a, X_{train}, y_{train})$ значительно меньше, чем значение $Q(a, X_{test}, y_{test})$ на **контрольной** выборке.

Если $Q(a, X_{test}, y_{test})$ примерно равна $Q(a, X_{train}, y_{train})$, то говорят, что алгоритм обладает **обобщающей способностью**

Функций $a(x)$, которые идеально описывают обучающую выборку бесконечно много, и обобщающей способности у таких алгоритмов не будет. Нужно сузить функциональное пространство перебора. Будем искомую функцию **параметризовать** – модель $a(x, w)$ описывается вектором весов w .

Тогда задача превращается в поиск весов:

$$w^* = \underset{w}{\operatorname{argmin}} Q(w, X_{train}, y_{train})$$

Параметрическая модель, предсказание которой **линейно** зависит от признаков, называется **линейной моделью**.

Не забываем нашу любимую формулу:

$$\textit{Learning} = \textit{Representation} + \textit{Evaluation} + \textit{Optimization}$$

Читаем статью <https://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf>
(1500 цитирований)

Знакомьтесь, линейная регрессия!

Представление линейной регрессии

Пусть объект описывается D признаками f_1, f_2, \dots, f_D . Тогда модель:

$a(x, w) = w_0 + \sum_{j=1}^D f_j w_j$ называется **линейной моделью**, где w – D -мерный вектор признаков w_1, w_2, \dots, w_D .

Далее будем считать, что в векторе признаков есть тождественно равный единице признак f_0 , тогда формула упростится до:

Representation: $a(x, w) = \sum_{j=0}^D f_j w_j = \mathbf{x} \cdot \mathbf{w}$

Весы модели интерпретируемы. w_i – значение, на которое изменится предсказание, если признак f_i увеличить на единицу.

Evaluation.

- квадратичная функция потерь (Q называется MSE) $L(a, x, y) = (a(x) - y)^2$
- абсолютная функция потерь (Q называется MAE) $L(a, x, y) = |a(x) - y|$
- логарифмическая функция потерь (Q называется $MSLE$)
 $L(a, x, y) = (\log(a(x) + 1) - \log(y + 1))^2$
- абсолютная-процентная функция потерь (Q называется $MAPE$) $L(a, x, y) = \frac{|a(x) - y|}{y}$
- все, что сами придумаете

Обычно используют MSE :

$$Q(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \cdot \mathbf{w} - y_i)^2 = \frac{1}{N} \|\mathbf{X} \cdot \mathbf{w} - \mathbf{y}\|^2$$

Feature Matrix (\mathbf{X})

n_features →

← n_samples

Target Vector (\mathbf{y})

← n_samples

Точное решение

$\nabla_{\mathbf{w}} Q(\mathbf{w})$ – градиент, вектор частных производных. Необходимое условие минимума – градиент равен нулю.

$$Q(\mathbf{w}) = \frac{1}{N} \|\mathbf{X} \cdot \mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{X} \cdot \mathbf{w} - \mathbf{y})^T (\mathbf{X} \cdot \mathbf{w} - \mathbf{y})$$

Два правила векторного дифференцирования:

$$1) \nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{c} = \nabla_{\mathbf{w}} \mathbf{c}^T \mathbf{w} = \mathbf{c}; 2) \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{C} \mathbf{w}) = (\mathbf{C} + \mathbf{C}^T) \mathbf{w}$$

$$\nabla_{\mathbf{w}} Q(\mathbf{w}) = \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Домашнее задание

Покажите, что это действительно минимум.

Недостатки точного решения:

- Можно написать только для очень редких функций потерь
- Обращение матрицы – кубическая сложность
- Матрица $X^T X$ может быть плохо обусловленной, если признаки линейно зависимы

Итеративные методы оптимизации:

- Нулевого порядка: золотое сечение, метод парабол, имитация отжига, генетические алгоритмы и т.д.
- Первого порядка: градиентный спуск, метод сопряженных градиентов, квазиньютоновские методы и т.д.
- Второго порядка: ньютоновские методы.

В курсе разбираем градиентный спуск.

Антиградиент функции показывает направления наискорейшего убывания функции.

❶ выбрать начальную длину шага α_0 , начальное приближение \mathbf{w}_0

❷ $\mathbf{w}_{new} = \mathbf{w}_{old} - \alpha \nabla_{\mathbf{w}} Q(\mathbf{w}_{old})$

❸ $\alpha = f(k), k = k + 1$

❹ Повторять (2), (3) до сходимости

$f(k) = \alpha_0, f(k) = \frac{\alpha_0}{k}, f(k) = \frac{\alpha_0}{k^p}, \dots$

В задачах машинного обучения оптимизируемая функция специального вида:

$$Q(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{w}, \mathbf{x}_i, y_i)$$

$$\nabla_{\mathbf{w}} Q(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}_i, y_i)$$

Пример

Для $L(\mathbf{w}, \mathbf{x}_i, y_i) = (\mathbf{x}_i \cdot \mathbf{w} - y_i)^2$ сможете посчитать $\nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}_i, y_i)$?

- 1 выбрать начальную длину шага α_0 , начальное приближение \mathbf{w}_0 , размер батча n
- 2 выбрать случайно $\{j_1, j_2, \dots, j_n\}$
- 3 оценить градиент $\nabla_{\mathbf{w}} Q^*(\mathbf{w}_{old}) = \frac{1}{n} \sum_{j=1}^n \nabla_{\mathbf{w}} L(\mathbf{w}_{old}, \mathbf{x}_j, y_j)$
- 4 $\mathbf{w}_{new} = \mathbf{w}_{old} - \alpha \nabla_{\mathbf{w}} Q^*(\mathbf{w}_{old})$
- 5 $\alpha = f(k), k = k + 1$
- 6 Повторять (2 - 5) до сходимости

В зависимости от n выделяют:

- $n = N$ – градиентный спуск (Gradient Descent, GD)
- $n = 1$ – стохастический градиентный спуск (Stochastic Gradient Descent, SGD)
- $n > 1$ и $n < N$ – мини-батч градиентный спуск (Mini-Batch Gradient Descent, MBGD)

Почему это вообще работает? Потому что такая оценка не является смещенной:

$$\mathbb{E}_{j_k} \nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}_{j_k}, y_{j_k}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{x}_i, y_i)$$

Градиентного спуск более точный и быстрее сходится (если мерить числом итераций до сходимости), но у стохастического спуска каждая итерация намного быстрее, то есть по времени он может сойтись быстрее!

Линейная регрессия невероятно популярный алгоритм машинного обучения. Плюсы алгоритма:

- Быстро учится
- Быстро предсказывает
- Легко интерпретируется
- Легко хранить в памяти
- Легко применять с дифференцируемой функцией потерь

Весомый минус - не способен учитывать нелинейные зависимости в данных.

Пример

Строим рекомендательную систему музыки. Где могут появиться нелинейные зависимости?

Уменьшаем минусы, добавляем плюсы!

Правильные признаки для линейной регрессии

Категориальные признаки кодируем:

- One-hot кодирование - категориальный признак с k значениями превращаем в k бинарных признаков!
- Кодирование через целевую переменную (нельзя включать переменную самого объекта)
- Кодирование через вещественные признаки

Для вещественных признаков применяем нелинейные функции - возводим в степень, берем синус и т.д.

Учитываем взаимодействия:

- Пару вещественных перемножаем, делим и т.д.
- Для пары бинарных используем логические операции

Невозможно сделать правильное признаковое пространство без понимания самой задачи!

Масштаб признаков очень важен для скорости сходимости градиентного спуска.

Пример

$f(x, y, z) = (x - 1)^2 + (y - 1)^2 + (z - 1)^2$. Стартуем из $(0, 0, 0)$ с шагом 0.5
А если так $f(x, y, z) = 1000(x - 1)^2 + 200(y - 1)^2 + (z - 1)^2$?

Два основных вида нормализации:

- Стандартизация. $f_j = \frac{f_j - \text{mean}(f_j)}{\text{std}(f_j)}$
- Min-max нормализация. $f_j = \frac{f_j - \min(f_j)}{\max(f_j) - \min(f_j)}$

Хотим еще сузить функциональное пространство $a(x)$, чтобы улучшить обобщающую способность. Наложим доп. штраф, если решение удаляется от нашего представления о правильном решении.

$$Q_r(\mathbf{w}) = Q(\mathbf{w}) + \alpha R(\mathbf{w}),$$

$R(\mathbf{w})$ - регуляризатор, α - параметр регуляризации.

Пример

Строим рекомендательную систему музыки. Какие регуляризаторы стоит сделать?

Базовые регуляризаторы, штрафующие большие значения весов:

- L_1 регуляризация (Ridge регрессия). $R(\mathbf{w}) = \sum_{j=1}^D |w_j|$; С ней Q_r не будет гладким!
- L_2 регуляризация (Lasso регрессия). $R(\mathbf{w}) = \sum_{j=1}^D w_j^2$.

Домашнее задание

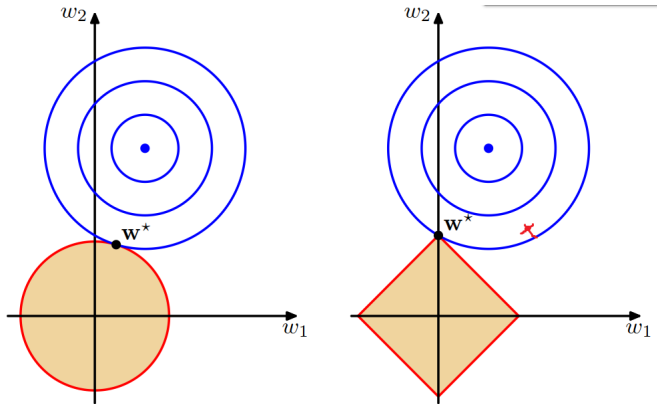
Показать, что с L_2 точное решение линейной регрессии станет $\mathbf{w} = (X^T X + \alpha I)^{-1} X^T \mathbf{y}$.
Что это дает для нас?

Домашнее задание

Выписать правило обновления весов для SGD линейной регрессии с L_2 .

Регуляризация

Можно показать, что L_1 приводит к занулению весов, то есть автоматически отбирает важные признаки!



Источник: Bishop.

Стохастическая оптимизация позволяет не хранить данные в оперативной памяти! Взяли батч, обновили веса и сразу же его забыли. Позволяет легко обучаться на терабайтах данных. Рекомендуем ознакомиться с библиотекой **Vowpal Wabbit**.

Online learning – обучение, когда прецеденты поступают потоком.

Пример

Строим рекомендательную систему музыки. Как в ней устроим online learning?

Учим линейную модель предсказывать спам/не спам почтовых писем. Делаем бинарные признаки f_i факт наличия слова $word_i$. Сколько таких будет признаков?

Чему равен градиент по j компоненте $\nabla_{\mathbf{w}}(\mathbf{x}_i \cdot \mathbf{w} - y_i)^2$ в точках $x_{i,j} = 0$?

Разреженный формат данных - для каждого объекта краним словарь вида
{номер ненулевого признака : значение}

Пример

Строим рекомендательную систему музыки. Как сделаем признаковое пространство?

Регрессия и причинно-следственные связи

Предсказываем сколько прослужит машина по ее цвету. Получили, что вес ненулевой: желтые машины дольше работают?

Предсказываем, результат физических упражнений по среднему баллу. Зависимостей нет, коэффициент около 0. А если добавить признак, поступил ли студент на военную кафедру МГУ?

Линейная регрессия позволяет находить причинно-следственные связи если:

- Линейная регрессия должна содержать все признаки, являющиеся причинами признакам
- Динейная регрессия не должна содержать признаки, которые являются следствиями одновременно самого признака и целевой переменной.

Рекомендую курс: www.coursera.org/learn/stats-for-data-analysis

Неправильные интерпретации:

<https://alexanderdyakonov.files.wordpress.com/2015/07/dyakonovfunnydm.pdf>

Спасибо за внимание!