

Познавач на Височина
Решение с Линејна Регресија

Николај Пашов

30 мај 2022 г.

Съдържание

1	Описание на решението	3
1.1	Псевдокод	4
2	Инструкции за компилиране	4
3	Примерни резултати	4
3.1	Експериментиране	4
4	Анализ	6

1 Описание на решението

Същината на задачата е чрез линейна регресия да се намери функция $f : \mathbb{R}^3 \rightarrow \mathbb{R}$, която изчислява височината на човек, разполагайки с възрастта и тежестта. Знаем, че:

$$f(age, weight) = a_0 + a_1 * age + a_2 * weight$$

Дефинираме функция на грешката:

$$L(y, gt) = (y - gt)^2$$

, където $y = f(age, weight)$ за случайни $age, weight$, а gt (ground truth) е истинската стойност от базата данни, която съответства на дадените $age, weight$.

Задачата се свежда до оптимизационна задача на функцията L чрез параметрите a_0, a_1, a_2 . Методът за оптимизация е градиентно спускане.

Алгоритъмът се състои от множество итерации през данните за трениране, в които се случва оптимизация на параметрите посредством техните частни производни спрямо L . Добавя се множител lr (learning rate), с който да се контролира оптимизационната стъпка:

$$a_x \leftarrow \frac{\partial L}{\partial a_x} * lr$$

С прости сметки се вижда, че:

$$\frac{\partial L}{\partial a_x} = \frac{\partial L}{\partial f} * \frac{\partial f}{\partial a_x} = 2 * (y - gt) * a_x, x \neq 0$$

$$\frac{\partial L}{\partial a_x} = \frac{\partial L}{\partial f} * \frac{\partial f}{\partial a_x} = 2 * (y - gt) * 1, x = 0$$

Съответната сметка се прави по време на итериране през данните за трениране, откъдето се взимат стойностите на y и gt .

Крайната оценка на тренирането се извършва с друга функция на грешката:

$$Acc = \frac{1}{n} * \sum_{i=0}^n abs(prediction(test[i][age], test[i][weight]), test[i][height])$$

1.1 Псевдокод

```
(train, test) ← input_data.split(0.8, 0.2)
[a0, a1, a2] ← [random(a0, a1, a2)]
f(age, weight) = a0 + age * a1 + weight * a2
iteration ← 0
lr ← 2e - 6
while it ≠ length(train) do
    prediction ← f(train[it][age], train[it][weight])
    a0 ← a0 - 2*(prediction - data[it][height]) * lr
    a1 ← a1 - 2 * (prediction - data[it][height]) * a1 * lr
    a2 ← a2 - 2 * (prediction - data[it][height]) * a2 * lr
```

While цикълът се изпълнява произволен брой пъти. Броят изпълнение е хиперпараметър, който подлежи на експериментиране

2 Инструкции за компилиране

Имплементацията е реализирана на Python 3.10. Очаква се да бъде изпълнима и на по-стари версии, заради backward compatibility, но не е тествана и не се гарантира от автора. Изисква се инсталиран Python и python командата в терминала да извиква интерпретатор на Python 3.10. За изпълнение на програмата:

```
cd <code_directory>
python <filename>
```

Изисква се .csv файлът с данните за трениране да е в същата директория и да се казва "age_weight.csv"

3 Примерни резултати

3.1 Експериментиране

В програмата е реализирано търсене на оптимална скорост на трениране. То работи по следния начин:

- Задава фиксирани скорости (lr) на трениране в масив.
- Разцепва train множеството на train и validation (Пропорции - 4:1).

```
Iteration: 5
Loss: 11.071495360682619
Iteration: 10
Loss: 2.0210486032708066
Iteration: 15
Loss: 0.567680734963195
Iteration: 20
Loss: 0.18482616343072647
Iteration: 25
Loss: 0.06609095636239821
Iteration: 30
Loss: 0.026488231595684617
Iteration: 35
Loss: 0.012792851433954238
Iteration: 40
Loss: 0.007989335978525691
Iteration: 45
Loss: 0.006280523841797584
Iteration: 50
Loss: 0.005672013276050359
Iteration: 55
Loss: 0.00545695258561716
Iteration: 60
Loss: 0.005382976238350585
Iteration: 65
Loss: 0.005355711791346578
Iteration: 70
Loss: 0.0053493785655751806
Iteration: 75
Loss: 0.005345048548668659
Iteration: 80
Loss: 0.005346277603118689
Iteration: 85
Loss: 0.005346434412933582
Iteration: 90
Loss: 0.005345498738160399
Iteration: 95
Loss: 0.00534650502383333
Iteration: 100
Loss: 0.005346689071300991
```

Фигура 1: Трениране върху тестово множество със 100 итерации. $lr = 6e - 04$. Loss е функцията L която се пресмята на всеки 5 обучителни цикъла (обхождане с трениране на параметрите върху множеството за трениране). Стойността е осреднена върху елементите от множеството за трениране.

```
lr: 2e-05 Acc: 0.05489679759376924
lr: 5e-05 Acc: 0.0131378065704029
lr: 7e-05 Acc: 0.06421912394091132
lr: 0.0001 Acc: 0.06136979687598651
lr: 0.0004 Acc: 0.023950766288360213
lr: 0.0007 Acc: 0.024044680138715156
lr: 0.001 Acc: 0.02397123441187671
lr: 0.003 Acc: 0.02407774423056882
lr: 0.005 Acc: 0.27407520285230946
lr: 0.01 Acc: 0.6296062331128828
Best learning rate according to hyperparameter search: 5e-05
Actual best learning rate: 7e-05
```

Фигура 2: Експеримент с търсене на скорост на учене.

- Извършва трениране за всяко lr и пресмята Acc за съответното lr .
- Предикторите от всичките тренировки се тестват спрямо множеството за тестване и се пресмята кое lr е в действителност най-доброто.

4 Анализ

В настоящия проект беше реализиран оптимизационен алгоритъм за линейна функция чрез линейна регресия. Намаляващата функция на грешката L при трениране индицира, че обучение наистина се случва и алгоритъмът се подобрява. Крайните стойности на функцията на грешката Acc варираха между 0.1 и 0.003, което са добри стойности в сравнение с диапазона на числата, които се познават (височината на хората).

Функцията за търсене на най-добра скорост на трениране невинаги открива най-оптималната скорост. Възможна причина за това е малкото количество данни. Въпреки това, може да се заключи, че при големи стойности на lr ($lr > 3e - 3$) функцията на грешката е осезателно по-голяма. Очакваната причина за това е, че с толкова голяма стъпка локалните минимума на функцията лесно се "прескачат".