

Отчет Лабораторная работа №3

Введение

В этой лабораторной работе мы реализовали и исследовали методы стохастической оптимизации поиска минимума многомерных функций, а также сравнили их с методами из прошлых лабораторных работы. В работе были использованы Python 3 и библиотека `scipy-optimize`.

1 Описание методов

SGD

Реализация SGD

Листинг 1: Реализация метода SGD

```
1 def sgd(  
2     X, y,  
3     lr_fun,  
4     lr0=0.01,  
5     n_epochs=800,  
6     batch=1,  
7     decay=0.003,  
8     momentum=0,  
9     reg_type=None,  
10    tol=1e-4,  
11    patience=20,  
12    log_every=1  
13 ):  
14     m, n = X.shape  
15     w = np.zeros((n, 1))  
16     v = np.zeros_like(w)  
17  
18     loss_hist, iter_hist = [], []  
19     best_loss, wait, it = np.inf, 0, 0  
20  
21     for epoch in range(1, n_epochs + 1):  
22         lr = lr_fun(epoch - 1, lr0, decay)  
23         perm = np.random.permutation(m)  
24  
25         for i in range(0, m, batch):  
26             xi = X[perm[i:i + batch]]  
27             yi = y[perm[i:i + batch]]  
28  
29             g = calc_gradient(xi, yi, w, reg_type)  
30             v = momentum * v - lr * g  
31             w += v  
32             it += 1  
33  
34             if it % log_every == 0:  
35                 cur_loss = calc_loss(X, y, w, reg_type)  
36                 loss_hist.append(cur_loss)  
37                 iter_hist.append(it)  
38  
39                 if best_loss - cur_loss > tol:  
40                     best_loss = cur_loss  
41                     wait = 0  
42                 else:  
43                     wait += 1  
44                     if wait >= patience:  
45                         return w, iter_hist, loss_hist  
46     return w, iter_hist, loss_hist
```

2 Графики

Реализуем отображение графиков на Python, который:

- отображает визуализацию 2D
- отрисовывает траекторию градиентного спуска

Используемые библиотеки

- `numpy` — работа с массивами данных
- `matplotlib.pyplot` — создание 2D-графиков

3 Описание результатов

Результаты для разных размеров батчей и $\text{learning rate} = 10^5$

3.1 Реализация SGD с разными batch size

| batch size | Iterations | Mem (MiB) | MSE | MAE | R^2 |
|------------|------------|-----------|-------|-------|-------|
| 1 | 297 | 1.3 | 0.425 | 0.526 | 0.350 |
| 10 | 402 | 0.9 | 0.394 | 0.508 | 0.397 |
| 100 | 449 | 1.0 | 0.391 | 0.503 | 0.402 |
| 500 | 466 | 1.5 | 0.390 | 0.503 | 0.403 |
| 1000 | 465 | 1.0 | 0.390 | 0.503 | 0.403 |

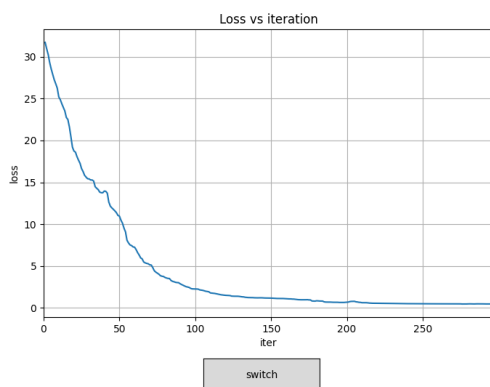


Рис. 1: batch size 1

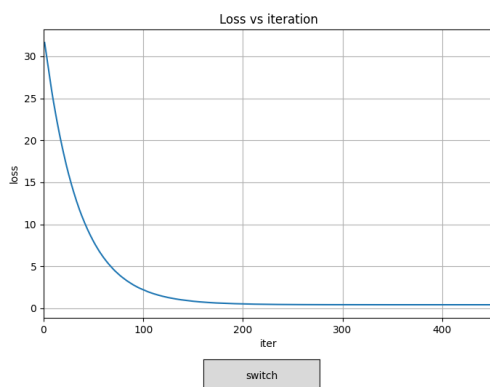
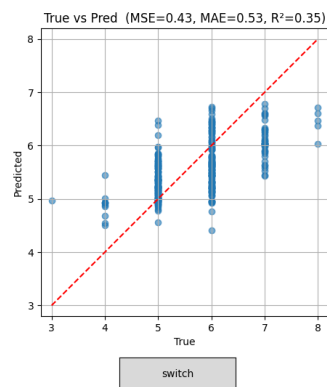


Рис. 2: batch size 100

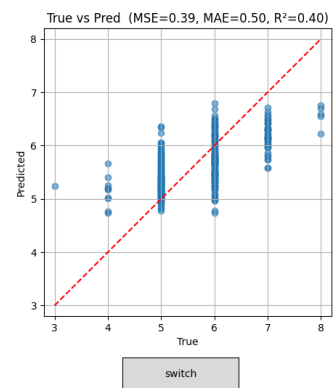
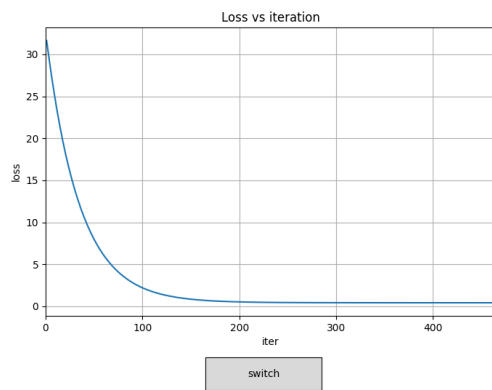


Рис. 3: batch size 1000

3.2 Выбор шага и регуляция

| Method | Iterations | Mem (MiB) | MSE | MAE | R ² |
|--------------------|------------|-----------|-------|-------|----------------|
| exp | 427 | 1.6 | 0.400 | 0.512 | 0.388 |
| step | 938 | 0.4 | 0.429 | 0.513 | 0.344 |
| const | 413 | 0.7 | 0.401 | 0.513 | 0.386 |
| exp + l1(0.1) | 437 | 2.0 | 0.401 | 0.513 | 0.387 |
| step + l1(0.1) | 933 | 0.5 | 0.429 | 0.513 | 0.344 |
| const + l1(0.1) | 426 | 1.2 | 0.401 | 0.513 | 0.387 |
| exp + l2(0.1) | 430 | 2.3 | 0.400 | 0.513 | 0.387 |
| step + l2(0.1) | 929 | 0.4 | 0.429 | 0.513 | 0.344 |
| const + l2(0.1) | 386 | 0.7 | 0.402 | 0.512 | 0.385 |
| exp + elasticnet | 430 | 2.6 | 0.400 | 0.512 | 0.388 |
| step + elasticnet | 932 | 1.0 | 0.429 | 0.513 | 0.344 |
| const + elasticnet | 374 | 0.1 | 0.403 | 0.513 | 0.384 |

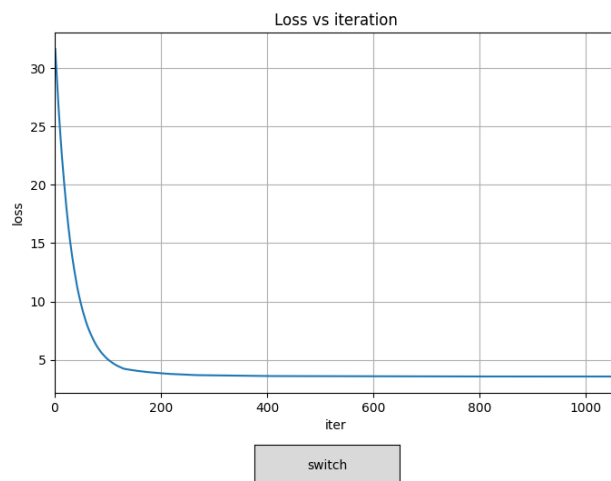


Рис. 4: step + l1(0.1)

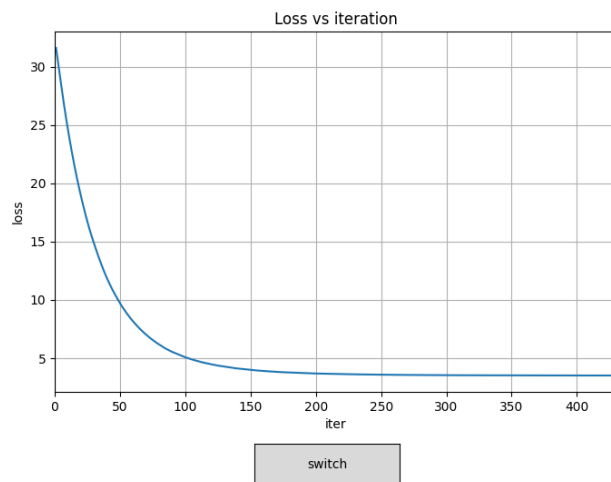


Рис. 5: exp + elasticnet

Регуляция меняет не так много, выбор шага реально влияет на результат, что наглядно видно. Так, exp хорошо себя показал в малом количестве итераций, но достаточно много тратил памяти. Противоположно же работал метод step. Чем то средним между двумя этими методами был const.

3.3 Keras и torch

| Method | MSE | MAE | R^2 |
|--------|-------|-------|--------|
| KERAS | 0.691 | 0.665 | -0.057 |
| TORCH | 0.482 | 0.539 | 0.263 |

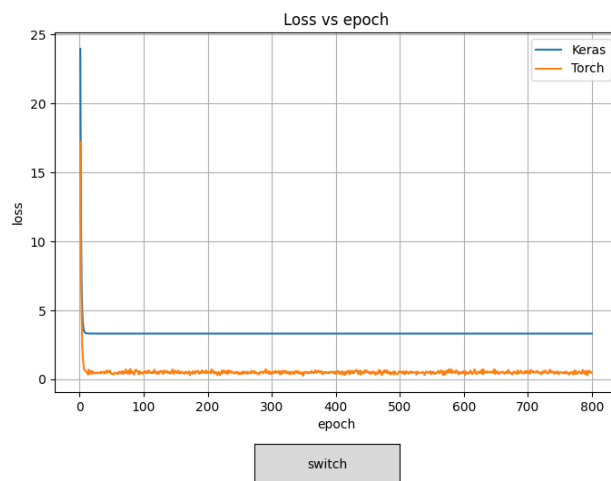


Рис. 6: Keras и Torch

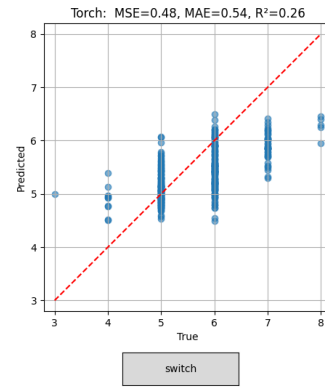
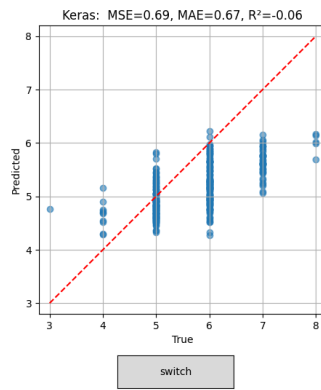


Рис. 7: batch size 100

Ключевое различие состоит в том, что в Keras вы просто описываете архитектуру модели и вызываете метод `fit`, передавая нужные параметры (число эпох, размер батча, разделение на тренировочную и валидационную выборки)

В PyTorch же вам нужно вручную реализовать этот цикл: в каждом проходе по эпохам формировать батчи, вычислять выходы модели, считать лосс, выполнять `backward()` и `step()`, то есть напрямую управлять всеми этапами оптимизации.

3.4 Momentum

Реализуем моментум

| Momentum | Iterations | Mem (MiB) | MSE | MAE | R ² |
|----------|------------|-----------|-------|-------|----------------|
| 0 | 417 | 1.9 | 0.401 | 0.512 | 0.387 |
| 0.3 | 460 | 2.0 | 0.401 | 0.514 | 0.386 |
| 0.6 | 434 | 2.2 | 0.400 | 0.512 | 0.388 |
| 0.9 | 453 | 1.8 | 0.399 | 0.512 | 0.389 |

Как видим моментум незначительно улучшил наш результат

4 Поставим задачу регрессии с ε -insensitive-потерями

Пусть $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$ — обучающая выборка. Будем искать линейный предсказатель в (возможно) «скрытом» пространстве

$$f(x) = w^\top \varphi(x) + b, \quad w \in \mathcal{H}, b \in \mathbb{R},$$

где $\varphi : \mathbb{R}^d \rightarrow \mathcal{H}$ — фиксированное отображение (ядровой трюк позволит не выписывать φ явно).

-insensitive-loss игнорирует ошибки радиусом до $\varepsilon > 0$:

$$|y - f(x)|_\varepsilon := \max(0, |y - f(x)| - \varepsilon).$$

Её график — «плоская» центральная полка шириной 2ε .

Прямая задача SVR

Чтобы свести абсолютную величину к линейным ограничениям, вводим неотрицательные slack-переменные ξ_i^+ , ξ_i^- для превышений сверху/снизу трубы.

$$\min_{w, b, \xi^+, \xi^-} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i^+ + \xi_i^-)$$

$$\begin{cases} y_i - (w^\top \varphi(x_i) + b) \leq \varepsilon + \xi_i^+, \\ (w^\top \varphi(x_i) + b) - y_i \leq \varepsilon + \xi_i^-, \\ \xi_i^+, \xi_i^- \geq 0, \quad i = 1, \dots, m. \end{cases}$$

$\frac{1}{2}\|w\|^2$ — регуляризация, эквивалентная максимизации ширины зазора $2/\|w\|$; $C > 0$ — жёсткость модели: чем больше C , тем дороже ошибаться маленькое $C \rightarrow$ модель более плоская и допускает больше выбросов.

Лагранжиан и условия ККТ

Введём множители $\alpha_i, \alpha'_i \geq 0$ для первых двух ограничений и $\eta_i^+, \eta_i^- \geq 0$ для $\xi_i^\pm \geq 0$. Лагранжиан:

$$\begin{aligned} \mathcal{L} = & \frac{1}{2}\|w\|^2 + C \sum_i (\xi_i^+ + \xi_i^-) - \sum_i \alpha_i (\varepsilon + \xi_i^+ - y_i + w^\top \varphi_i + b) \\ & - \sum_i \alpha'_i (\varepsilon + \xi_i^- + y_i - w^\top \varphi_i - b) - \sum_i \eta_i^+ \xi_i^+ - \sum_i \eta_i^- \xi_i^-, \end{aligned}$$

где $\varphi_i = \varphi(x_i)$.

Условия стационарности $\partial \mathcal{L} / \partial (\cdot) = 0$:

$$\begin{cases} w = \sum_i (\alpha_i - \alpha'_i) \varphi_i, \\ \sum_i (\alpha_i - \alpha'_i) = 0, \\ 0 = C - \alpha_i - \eta_i^+, \quad 0 = C - \alpha'_i - \eta_i^-. \end{cases}$$

$$\eta_i^\pm \geq 0 \Rightarrow 0 \leq \alpha_i, \alpha'_i \leq C.$$

Двойственная задача

Подставляя выражение для w в \mathcal{L} и используя $K(x_i, x_j) = \varphi_i^\top \varphi_j$, получаем:

$$\begin{aligned} \max_{\alpha, \alpha'} & -\frac{1}{2} \sum_{i,j} (\alpha_i - \alpha'_i)(\alpha_j - \alpha'_j) K(x_i, x_j) - \varepsilon \sum_i (\alpha_i + \alpha'_i) + \sum_i (\alpha_i - \alpha'_i) y_i \\ \text{s.t.} & \sum_i (\alpha_i - \alpha'_i) = 0, \quad 0 \leq \alpha_i, \alpha'_i \leq C. \end{aligned}$$

Размерность задачи — $2m$ вместо d ; ядро $K(\cdot, \cdot)$ позволяет работать в бесконечномерном \mathcal{H} .

Разреженность. В оптимуме большинство α_i, α'_i обнулены. Точки с ненулевыми коэффициентами называются опорными (SV) и полностью определяют модель.

Восстановление прямого решения

$$w^* = \sum_{i \in SV} (\alpha_i - \alpha'_i) \varphi(x_i), \quad b^* = y_k - w^\top \varphi(x_k) - \text{sgn}(\alpha_k - \alpha'_k) \varepsilon,$$

где x_k — любой SV, находящийся внутри трубы ($0 < \alpha_k < C$ или $0 < \alpha'_k < C$). Прогноз для новой точки x :

$$f(x) = \sum_{i \in SV} (\alpha_i - \alpha'_i) K(x_i, x) + b^*.$$