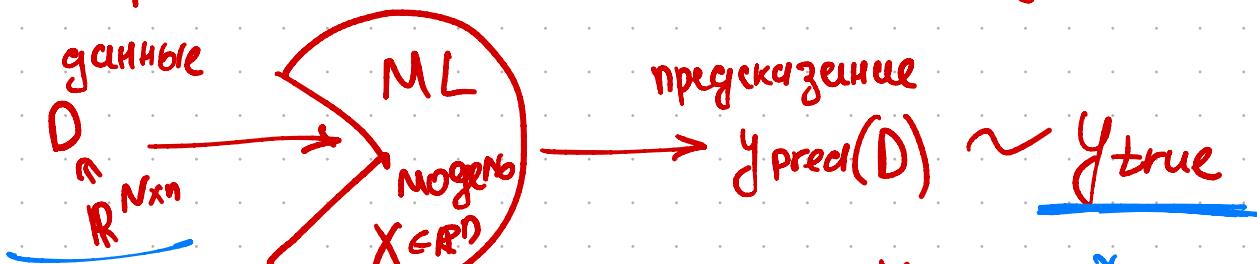


Градиентные методы в реальной жизни

$$f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x) \rightarrow \min_{x \in \mathbb{R}^p}$$

N - # данных в батче

p - # обучаемых параметров модели.



$$\text{Loss} = \frac{1}{N} \sum_{i=1}^N (Y_{\text{pred}}(d_i) - Y_{\text{true}}^i)^2 \rightarrow \min_{X \in \mathbb{R}^p}$$

$$Y_{\text{pred}}^i = X^T \cdot d_i \rightarrow f(x) = \frac{1}{N} \sum_{i=1}^N (X^T d_i - Y_{\text{true}}^i)^2.$$

GD

$$x_{k+1} = x_k - \alpha_k \cdot \nabla f(x_k)$$

2. p параметров
Memory

$$p \sim 10^0 \sim 10^1$$

$$\nabla f(x_k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_k)$$

$$2 \cdot 10^0 \cdot 16 \text{ байт} \sim$$

$$\sim 10^9 \cdot 320 \frac{\text{байт}}{\text{байт}}$$

$$\sim 320 \text{ Гбайт} \cdot 40 \text{ Гб.}$$

$$N = 10^7$$

Классическое ускорение:

MOMENTUM

- Как ускорить обучение модели с помощью SGD?

▼ Beyond gradient descent, vol. 2

Концепция оптимальных методов, метод тяжёлого шарика ускоренный метод Нестерова

Схема получения оценок снизу на сложность методов и задач

- Фиксируем класс функций \mathcal{F}
- Фиксируем класс методов оптимизации \mathcal{M}
- Ищем настолько плохую функцию из класса \mathcal{F} , что любой метод из класса \mathcal{M} сходится не лучше, чем некоторая оценка
- Такая оценка называется оценкой снизу

▼ Пример

Фиксируем класс методов

Рассмотрим такие методы, что

$$x_{k+1} = x_0 + \text{span}\{f'(x_0), \dots, f'(x_k)\}$$

- Далее в рамках этого семинара для краткости только такие методы будем называть методами первого порядка
- Весь последующий анализ **НЕ** применим, если

$$x_{k+1} = x_0 + G(f'(x_0), \dots, f'(x_k)),$$

где G - некоторая нелинейная функция

- С такими методами мы познакомимся на одном из ближайших занятий

▼ Фиксируем класс функций

Выпуклые функции с липшицевым градиентом

Теорема. Существует выпуклая функция с Липшицевым градиентом, такая что

$$f(x_t) - f^* \geq \frac{3L\|x_0 - x^*\|_2^2}{32(t+1)^2},$$

нижняя оценка для методов 1-го порядка convex

где $x_k = x_0 + \text{span}\{f'(x_0), \dots, f'(x_{k-1})\}$, $1 \leq k \leq t$

Привести пример такой функции и доказать эту теорему Вам надо в домашнем задании.

Сильно выпуклые функции с липшицевым градиентом

Теорема. Существует сильно выпуклая функция с Липшицевым градиентом, такая что

$$f(x_t) - f^* \geq \frac{\mu}{2} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2t} \|x_0 - x^*\|_2^2,$$

где $\kappa = \frac{L}{\mu}$ и $x_k = x_0 + \text{span}\{f'(x_0), \dots, f'(x_{k-1})\}$, $1 \leq k \leq t$

нижняя оц.

μ сильно вып.

$$\kappa = \frac{L}{\mu}$$

Привести пример такой функции и доказать эту теорему Вам надо в домашнем задании.

▼ Оценки сходимости известных методов

Оценки сходимости для градиентного спуска: напоминание

- Пусть
 - $f(x)$ дифференцируема на \mathbb{R}^n
 - $f(x)$ выпукла
 - $f'(x)$ удовлетворяет условию Липшица с константой L
 - $\alpha = \frac{1}{L}$

Тогда

$$f(x_k) - f^* \leq \frac{2L\|x_0 - x^*\|_2^2}{k+4}$$

$$f(x_k) - f^* \leq \frac{2L\|x_0 - x^*\|_2^2}{k+4}$$

convex

$$\alpha = \frac{1}{L}$$

- Пусть

- $f(x)$ дифференцируема на \mathbb{R}^n ,
- градиент $f(x)$ удовлетворяет условию Липшица с константой L
- $f(x)$ является сильно выпуклой с константой μ
- $\alpha = \frac{2}{\mu + L}$

Тогда для градиентного спуска выполнено:

$$\|x_k - x^*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k \|x_0 - x^*\|_2$$

$$f(x_k) - f^* \leq \frac{L}{2} \left(\frac{\kappa - 1}{\kappa + 1} \right)^{2k} R^2$$

где $\kappa = \frac{L}{\mu}$

$$\alpha = \frac{1}{\mu + L}$$

$$\|x_k - x^*\|^2 \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k R^2$$

Оценки сходимости для метода сопряжённых градиентов: напоминание

Для сильно выпуклой квадратичной функции

$$f(x) = \frac{1}{2} x^\top A x - b^\top x$$

и метода сопряжённых градиентов справедлива следующая оценка сходимости

$$\sqrt{2(f_k - f^*)} = \|x_k - x^*\|_A \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|x_0 - x^*\|_A,$$

где $\kappa(A) = \frac{\lambda_1(A)}{\lambda_n(A)} = \frac{L}{\mu}$ - число обусловленности матрицы A , $\lambda_1(A) \geq \dots \geq \lambda_n(A) > 0$ - собственные значения матрицы A

Can we do better?

Существует ли метод, который сходится в соответствии с нижними оценками для

- произвольной сильно выпуклой функции с липшицевым градиентом (не только квадратичной)?
- произвольной выпуклой функции с липшицевым градиентом?

▼ Метод тяжёлого шарика (heavy-ball method)



- Предложен в 1964 г. Б.Т. Поляком

- Для квадратичной целевой функции зигзагообразное поведение градиентного спуска обусловлено неоднородностью направлений
- Давайте учитывать предыдущие направления для поиска новой точки
- Метод тяжёлого шарика

$$x_{k+1} = x_k - \alpha_k f'(x_k) + \beta_k (x_k - x_{k-1})$$

- Помимо параметра шага вдоль антиградиента α_k появился ещё один параметр β_k

Геометрическая интерпретация метода тяжёлого шарика

$$\beta_k = \beta = \text{const}$$

$$V_k = \beta_k V_{k-1} + g_k$$

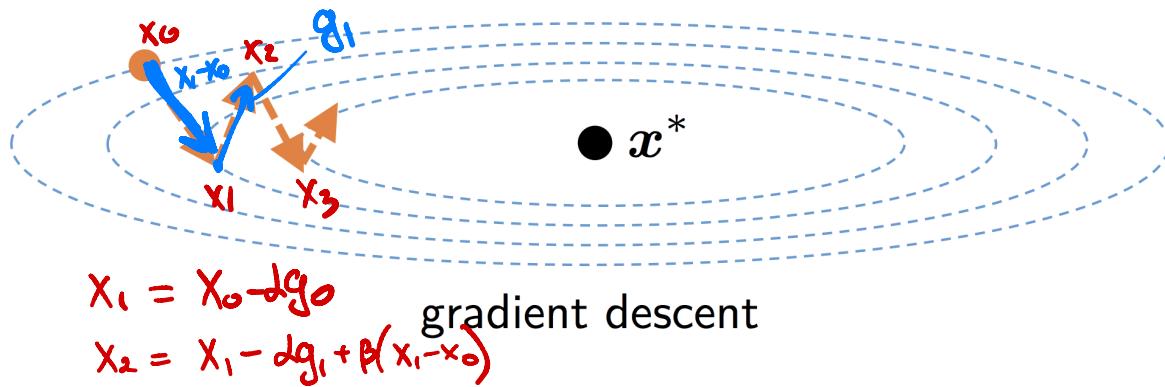
$$0 < \beta < 1$$

$$X_{k+1} = X_k - \alpha_k V_k$$

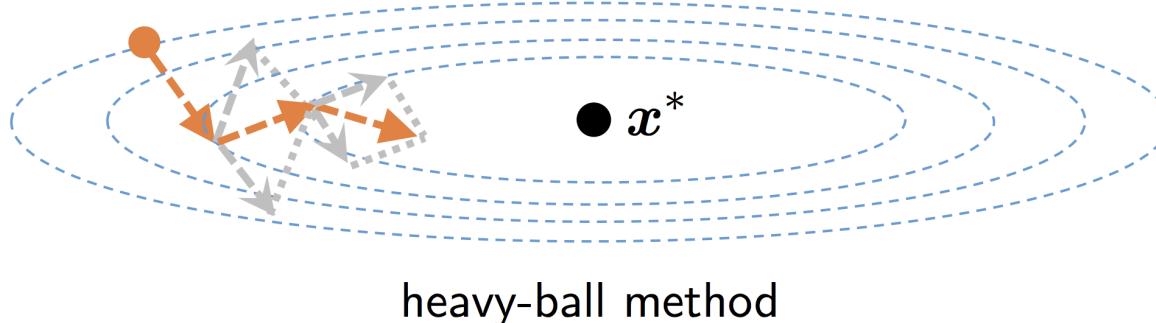
и где:

$$X_{k+1} = X_k - \alpha_k (\beta_k V_{k-1} + g_k) =$$

$$= X_k - \alpha_k (\beta (\beta V_{k-2} + g_{k-1}) + g_k) = X_k - \alpha_k (\beta^2 g_{k-2} + \beta g_{k-1} + g_k)$$



Картинка [отсюда](#)



Теорема сходимости

Пусть f сильно выпукла с Липшицевым градиентом. Тогда для

$$\alpha_k = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}$$

и

$$\beta_k = \max(|1 - \sqrt{\alpha_k L}|^2, |1 - \sqrt{\alpha_k \mu}|^2)$$

справедлива следующая оценка сходимости

$$\left\| \begin{bmatrix} x_{k+1} - x^* \\ x_k - x^* \end{bmatrix} \right\|_2 \leq \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \left\| \begin{bmatrix} x_1 - x^* \\ x_0 - x^* \end{bmatrix} \right\|_2$$

- Совпадает с оценкой снизу для методов первого порядка!
- Оптимальные параметры α_k и β_k определяются через **неизвестные** константы L и μ

▼ Схема доказательства

- Перепишем метод как

$$\begin{bmatrix} x_{k+1} \\ x_k \end{bmatrix} = \begin{bmatrix} (1 + \beta_k)I & -\beta_k I \\ I & 0 \end{bmatrix} \begin{bmatrix} x_k \\ x_{k-1} \end{bmatrix} + \begin{bmatrix} -\alpha_k f'(x_k) \\ 0 \end{bmatrix}$$

- Используем теорему из анализа

$$\begin{bmatrix} x_{k+1} - x^* \\ x_k - x^* \end{bmatrix} = \underbrace{\begin{bmatrix} (1 + \beta_k)I - \alpha_k \int_0^1 f''(x(\tau))d\tau & -\beta_k I \\ I & 0 \end{bmatrix}}_{=A_k} \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix},$$

где $x(\tau) = x_k + \tau(x^* - x_k)$

- В силу интегральной теоремы о среднем $A_k(x) = \int_0^1 f''(x(\tau))d\tau = f''(z)$, поэтому L и μ ограничивают спектр $A_k(x)$
- Сходимость зависит от спектрального радиуса матрицы итераций A_k
- Получим оценку на спектр A_k

$$A_k = \begin{bmatrix} (1 + \beta_k)I - \alpha_k A(x_k) & -\beta_k I \\ I & 0 \end{bmatrix}$$

- Пусть $A(x_k) = U\Lambda(x_k)U^\top$, поскольку гессиан - симметричная матрица, тогда

$$\begin{bmatrix} U^\top & 0 \\ 0 & U^\top \end{bmatrix} \begin{bmatrix} (1 + \beta_k)I - \alpha_k A(x_k) & -\beta_k I \\ I & 0 \end{bmatrix} \begin{bmatrix} U & 0 \\ 0 & U \end{bmatrix} = \begin{bmatrix} (1 + \beta_k)I - \alpha_k \Lambda(x_k) & -\beta_k I \\ I & 0 \end{bmatrix} = \hat{A}_k$$

- Ортогональное преобразование не меняет спектральный радиус матрицы
- Далее сделаем перестановку строк и столбцов так, чтобы

$$\hat{A}_k \simeq \text{diag}(T_1, \dots, T_n),$$

где $T_i = \begin{bmatrix} 1 + \beta_k - \alpha_k \lambda_i & -\beta_k \\ 1 & 0 \end{bmatrix}$ и \simeq обозначает равенство спектральных радиусов поскольку матрица перестановки является ортогональной

- Покажем как сделать такую перестановку на примере матрицы 4×4

$$\left[\begin{array}{cccc} a & 0 & c & 0 \\ 0 & b & 0 & c \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{array} \right] \rightarrow \left[\begin{array}{cccc} a & 0 & c & 0 \\ 1 & 0 & 0 & 0 \\ 0 & b & 0 & c \\ 0 & 1 & 0 & 0 \end{array} \right] \rightarrow \left[\begin{array}{ccccc} a & c & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & b & c & 0 \\ 0 & 0 & 1 & 0 & 0 \end{array} \right]$$

- Свели задачу к оценке спектрального радиуса блочно-диагональной матрицы \hat{A}_k

- $$\bullet \quad \rho(\hat{A}_k) = \max_{i=1,\dots,n} \{ |\lambda_1(T_i)|, |\lambda_2(T_i)| \}$$

- Характеристическое уравнение для T_i

$$\beta_k - u(1 + \beta_k - \alpha_k \lambda_i - u) = 0 \quad u^2 - u(1 + \beta_k - \alpha_k \lambda_i) + \beta_k = 0$$

- Дальнейшее изучение распределения корней и их границ даёт оценку из условия теоремы

▼ Эксперименты

Тестовая задача 1

$$\nabla f = \dots$$

$$f(x) = \frac{1}{2}x^\top Ax - b^\top x \rightarrow \min_x,$$

где матрица A плохо обусловлена, но положительно определена!

```

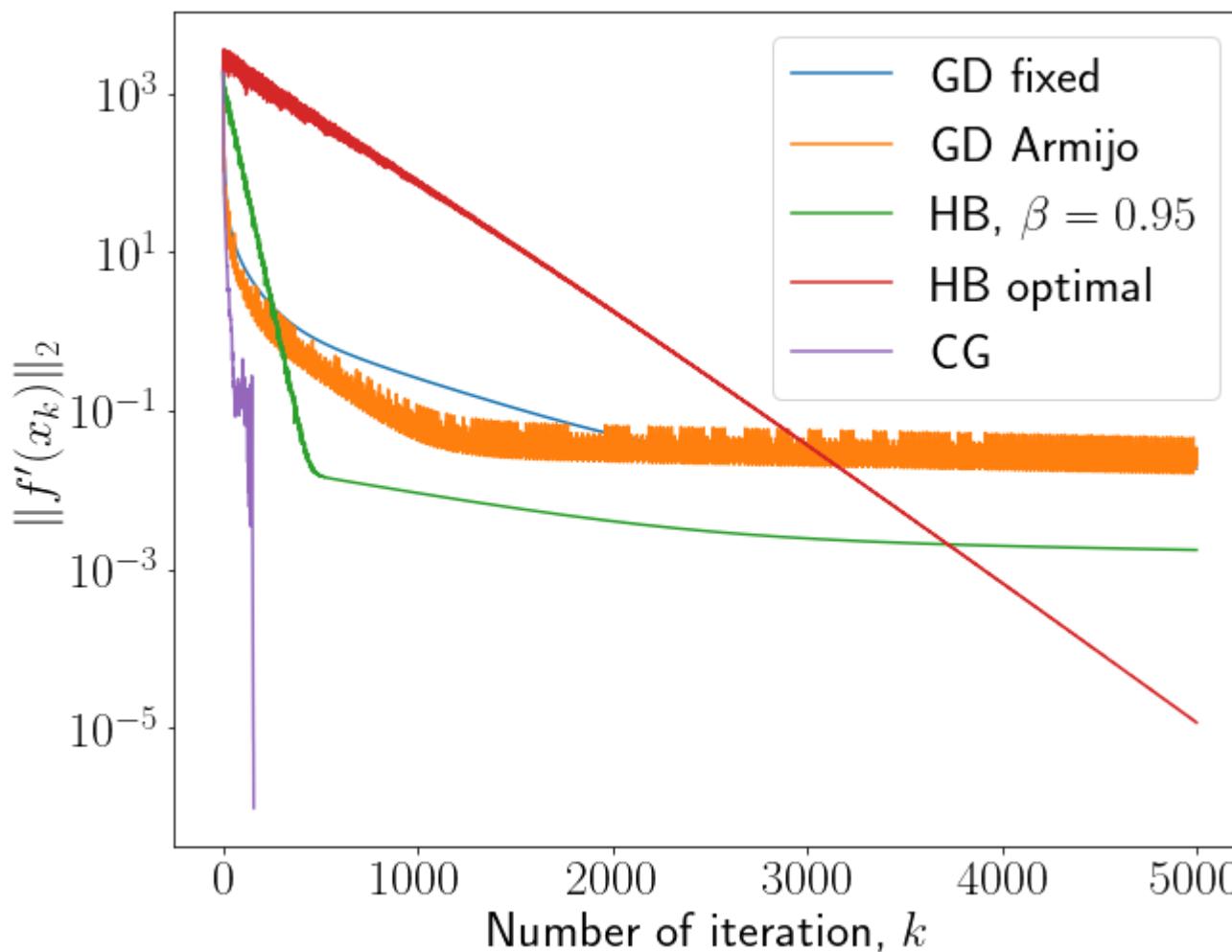
17 print(alpha_opt, beta_opt)
18
388.0536004941595 0.0017707060804411218
Condition number = 219151.89922287207
0.010263957266938104 0.9914918736025333

1 methods = {
2     "GD fixed": fo.GradientDescent(f, grad, ss.ConstantStepSize(1 / L)),
3     "GD Armijo": fo.GradientDescent(f, grad,
4         ss.Backtracking("Armijo", rho=0.5, beta=0.1, init_alpha=1.)),
5     r"HB, $\beta = {}$".format(beta_test): HeavyBall(f, grad, ss.ConstantStepSize(1 / L), beta=beta_test),
6     "HB optimal": HeavyBall(f, grad, ss.ConstantStepSize(alpha_opt), beta = beta_opt),
7     "CG": fo.ConjugateGradientQuad(A, b)
8 }
9 x0 = np.random.randn(n)
10 max_iter = 5000
11 tol = 1e-6

1 for m in methods:
2     _ = methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

1 figsize = (10, 8)
2 fontsize = 26
3 plt.figure(figsize=figsize)
4 for m in methods:
5     plt.semilogy([np.linalg.norm(grad(x)) for x in methods[m].get_convergence()], label=m)
6 plt.legend(fontsize=fontsize, loc="best")
7 plt.xlabel("Number of iteration, $k$", fontsize=fontsize)
8 plt.ylabel(r"$\|f'(x_k)\|_2$", fontsize=fontsize)
9 plt.xticks(fontsize=fontsize)
10 _ = plt.yticks(fontsize=fontsize)

```



```

1 for m in methods:
2     print(m)
3     %timeit methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

GD fixed
518 ms ± 87.2 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
GD Armijo
5.42 s ± 499 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
HB, $\beta = 0.95$ 
322 ms ± 44.7 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
HB optimal
299 ms ± 24 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
CG
22.6 ms ± 1.66 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)

```

$$f(w) = \frac{1}{2} \|w\|_2^2 + C \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle x_i, w \rangle)) \rightarrow \min_w$$

```

1 n = 300
2 m = 1000
3 import sklearn.datasets as skldata
4 import jax
5 import jax.numpy as jnp
6 import scipy.optimize as scopt
7 from jax.config import config
8 config.update("jax_enable_x64", True)
9
10 X, y = skldata.make_classification(n_classes=2, n_features=n, n_samples=m, n_informative=n//3)
11 C = 1
12
13 @jax.jit
14 def f(w):
15     return jnp.linalg.norm(w)**2 / 2 + C * jnp.mean(jnp.logaddexp(jnp.zeros(X.shape[0]), -y * (X @ w)))
16
17 # def grad(w):
18 #     denom = scspec.expit(-y * X.dot(w))
19 #     return w - C * X.T.dot(y * denom) / X.shape[0]
20
21 autograd_f = jax.jit(jax.grad(f))
22 x0 = jnp.ones(n)
23 print("Initial function value = {}".format(f(x0)))
24 print("Initial gradient norm = {}".format(jnp.linalg.norm(autograd_f(x0))))
25
WARNING:absl:No GPU/TPU found, falling back to CPU. (Set TF_CPP_MIN_LOG_LEVEL=0 and rerun for more info.)
Initial function value = 165.04873928789436
Initial gradient norm = 18.646187895005486

1 alpha_test = 5e-3
2 beta_test = 0.9
3
4 methods = {
5     "GD, $\alpha_k = {}$".format(alpha_test): fo.GradientDescent(f, autograd_f, ss.ConstantStepSize(alpha_test)),
6     "GD Armijo": fo.GradientDescent(f, autograd_f,
7         ss.Backtracking("Armijo", rho=0.7, beta=0.1, init_alpha=1.)),
8     "HB, $\beta = {}$".format(beta_test): HeavyBall(f, autograd_f, ss.ConstantStepSize(alpha_test), beta=beta_test),
9
10 }
11 # x0 = np.random.rand(n)
12 # x0 = jnp.zeros(n)
13 x0 = jax.random.normal(jax.random.PRNGKey(0), (X.shape[1],))
14 max_iter = 400
15 tol = 1e-3

1 for m in methods:
2     print(m)
3     _ = methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)
4
5 scopt_cg_array = []
6 def callback(x, arr):
7     arr.append(x)
8
9 scopt_cg_callback = lambda x: callback(x, scopt_cg_array)
10 x = scopt.minimize(f, x0, tol=tol, method="CG", jac=autograd_f, callback=scopt_cg_callback, options={"maxiter": max_iter})
11 x = x.x

GD, $\alpha_k = 0.005$  

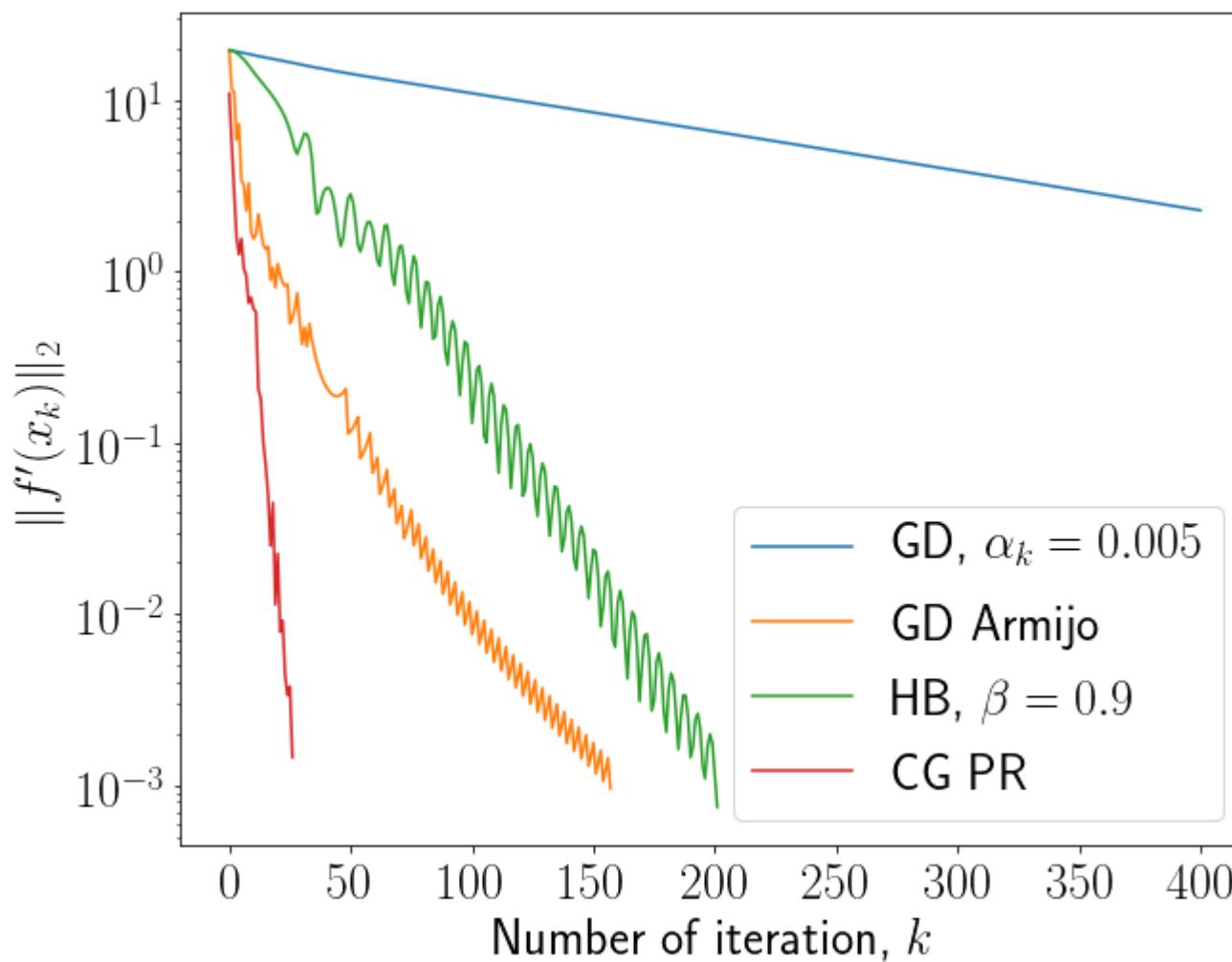
GD Armijo  

HB, $\beta = 0.9$  

1 figsize = (10, 8)
2 fontsize = 26
3 plt.figure(figsize=figsize)
4 for m in methods:
5     plt.semilogy([np.linalg.norm(autograd_f(x)) for x in methods[m].get_convergence()], label=m)
6 plt.semilogy([np.linalg.norm(autograd_f(x)) for x in scopt_cg_array], label="CG PR")
7 plt.legend(fontsize=fontsize, loc="best")
8 plt.xlabel("Number of iteration, $k$", fontsize=fontsize)
9 plt.ylabel(r"$\| f'(x_k) \|_2$", fontsize=fontsize)
10 plt.xticks(fontsize=fontsize)
11 _ = plt.yticks(fontsize=fontsize)

```



```

1 for m in methods:
2     print(m)
3     %timeit methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

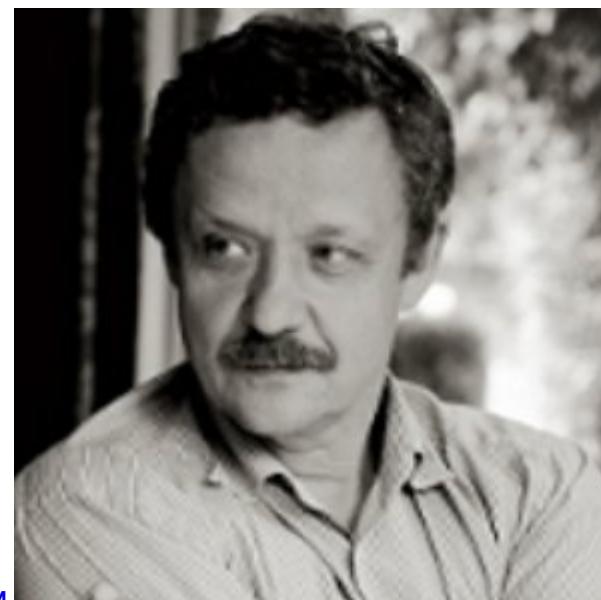
GD, $alpha_k = 0.001$
53.1 ms ± 11.1 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)
GD Armijo
613 ms ± 32.3 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
HB, $beta = 0.9$
57.2 ms ± 11.9 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)

```

Главное про метод тяжёлого шарика

- Двухшаговый метод
- Не обязательно монотонный
- Параметры зависят от неизвестных констант
- Решает проблему осцилляций для плохо обусловленных задач
- Сходимость для сильно выпуклых функций совпадает с оптимальной оценкой

▼ Ускоренный метод Нестерова



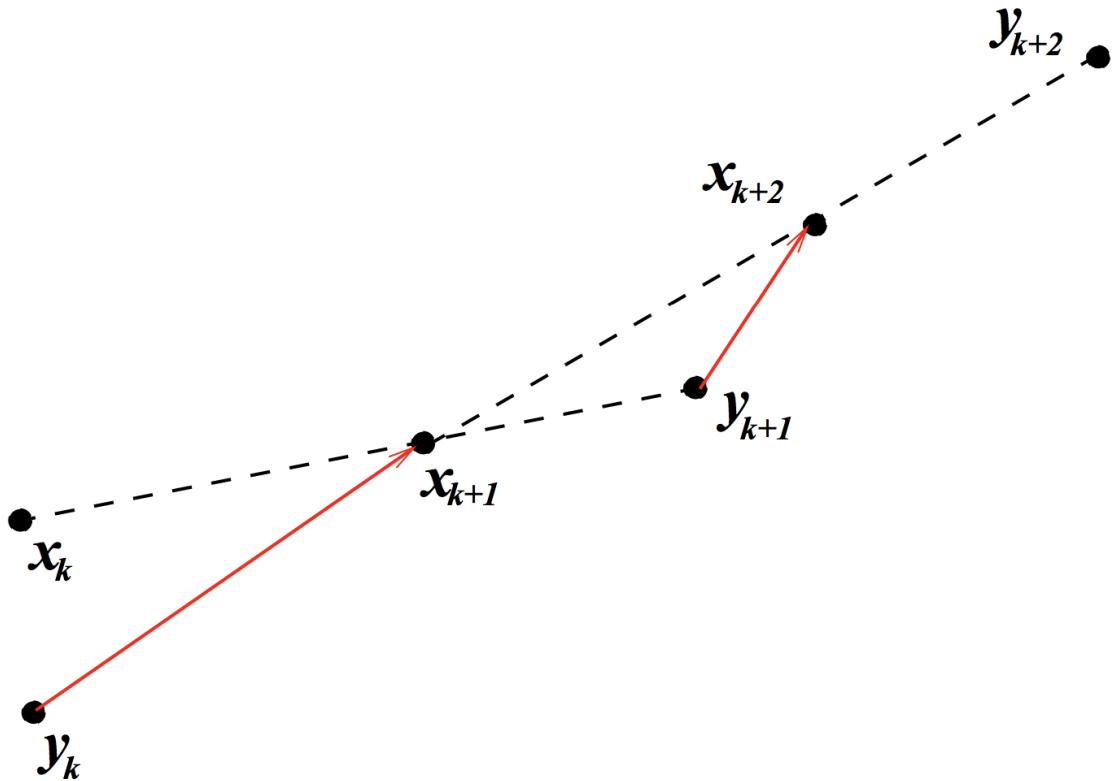
- Предложен в 1983 г. Ю.Е. Нестеровым
- Одна из возможных форм записи

$$\begin{aligned}
 y_0 &= x_0 \\
 x_{k+1} &= y_k - \alpha_k f'(y_k) \\
 y_{k+1} &= x_{k+1} + \frac{k}{k+3}(x_{k+1} - x_k)
 \end{aligned}$$

- Сравните с методом тяжёлого шарика
- Также не обязательно монотонен

- Для любителей геометрии есть альтернативный метод под названием [geometric descent](#) с такой же скоростью сходимости

Геометрическая интерпретация ускоренного метода Нестерова



▼ Теорема сходимости

- Пусть f выпукла с Липшицевым градиентом, а шаг $\alpha_k = \frac{1}{L}$. Тогда ускоренный метод Нестерова сходится как

$$f(x_k) - f^* \leq \frac{2L\|x_0 - x^*\|_2^2}{(k+1)^2}$$

- Пусть f сильно выпукла с липшицевым градиентом. Тогда ускоренный метод Нестерова при шаге $\alpha_k = \frac{1}{L}$ сходится как

$$f(x_k) - f^* \leq L\|x_k - x_0\|_2^2 \left(1 - \frac{1}{\sqrt{k}}\right)^k$$

▼ Тестовая задача

$$f(x) = \frac{1}{2}x^\top Ax - b^\top x \rightarrow \min_x$$

где матрица A положительно полуопределённая, то есть функция **НЕ** является сильно выпуклой

```

1 import numpy as np
2 import liboptpy.unconstr_solvers.fo as fo
3 import liboptpy.step_size as ss
4 import liboptpy.base_optimizer as base
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7
8 np.random.seed(42)
9 n = 100
10 A = np.random.randn(n, n)
11 A = A.T.dot(A)
12 A_eigvals = np.linalg.eigvalsh(A)
13 mu = np.min(A_eigvals)
14 A = A - (mu - 1e-6) * np.eye(n)
15 x_true = np.random.randn(n)
16 b = A.dot(x_true)
17 f = lambda x: 0.5 * x.dot(A.dot(x)) - b.dot(x)
18 grad = lambda x: A.dot(x) - b
19 A_eigvals = np.linalg.eigvalsh(A)
20 L = np.max(A_eigvals)
21 mu = np.min(A_eigvals)
22 print(L, mu)

```

388.0518307880791 9.99999702771782e-07

```

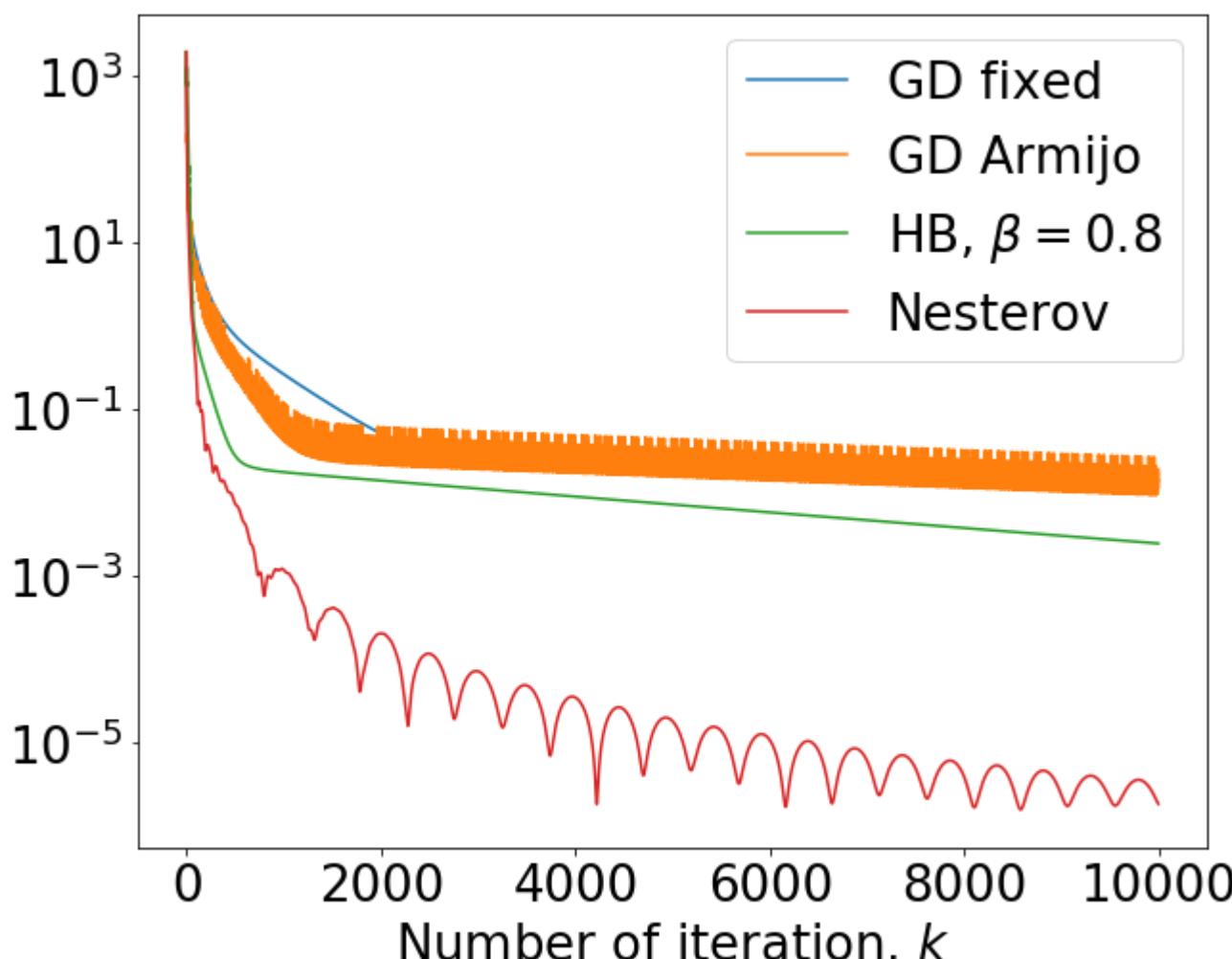
1 class HeavyBall(base.LineSearchOptimizer):
2     def __init__(self, f, grad, step_size, beta, **kwargs):

```

```

3     super().__init__(f, grad, step_size, **kwargs)
4     self._beta = beta
5
6     def get_direction(self, x):
7         self._current_grad = self._grad(x)
8         return -self._current_grad
9
10    def _f_update_x_next(self, x, alpha, h):
11        if len(self.convergence) < 2:
12            return x + alpha * h
13        else:
14            return x + alpha * h + self._beta * (x - self.convergence[-2])
15
16    def get_stepsize(self):
17        return self._step_size.get_stepsize(self._grad_mem[-1], self.convergence[-1], len(self.convergence))
18
19 beta_test = 0.8
20 methods = {
21     "GD fixed": fo.GradientDescent(f, grad, ss.ConstantStepSize(1 / L)),
22     "GD Armijo": fo.GradientDescent(f, grad,
23         ss.Backtracking("Armijo", rho=0.5, beta=0.1, init_alpha=1.)),
24     r"HB, $\beta = {}$".format(beta_test): HeavyBall(f, grad, ss.ConstantStepSize(1 / L), beta=beta_test),
25     "Nesterov": fo.AcceleratedGD(f, grad, ss.ConstantStepSize(1 / L)),
26 }
27
28 x0 = np.random.randn(n)
29 max_iter = 10000
30 tol = 1e-6
31
32
33 for m in methods:
34     _ = methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)
35
36
37 figsize = (10, 8)
38 fontsize = 26
39 plt.figure(figsize=figsize)
40 for m in methods:
41     plt.semilogy([np.linalg.norm(grad(x)) for x in methods[m].get_convergence()], label=m)
42 #     plt.semilogy([np.linalg.norm(x - x_true) for x in methods[m].get_convergence()], label=m)
43 #     plt.semilogy([f(x) - f(x_true) for x in methods[m].get_convergence()], label=m)
44 plt.legend(fontsize=fontsize, loc="best")
45 plt.xlabel("Number of iteration, $k$", fontsize=fontsize)
46 # plt.ylabel(r"$\| f'(x_k) \|_2^2$", fontsize=fontsize)
47 plt.xticks(fontsize=fontsize)
48 _ = plt.yticks(fontsize=fontsize)

```



```

1 for m in methods:
2     print(m)
3     %timeit methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

GD fixed

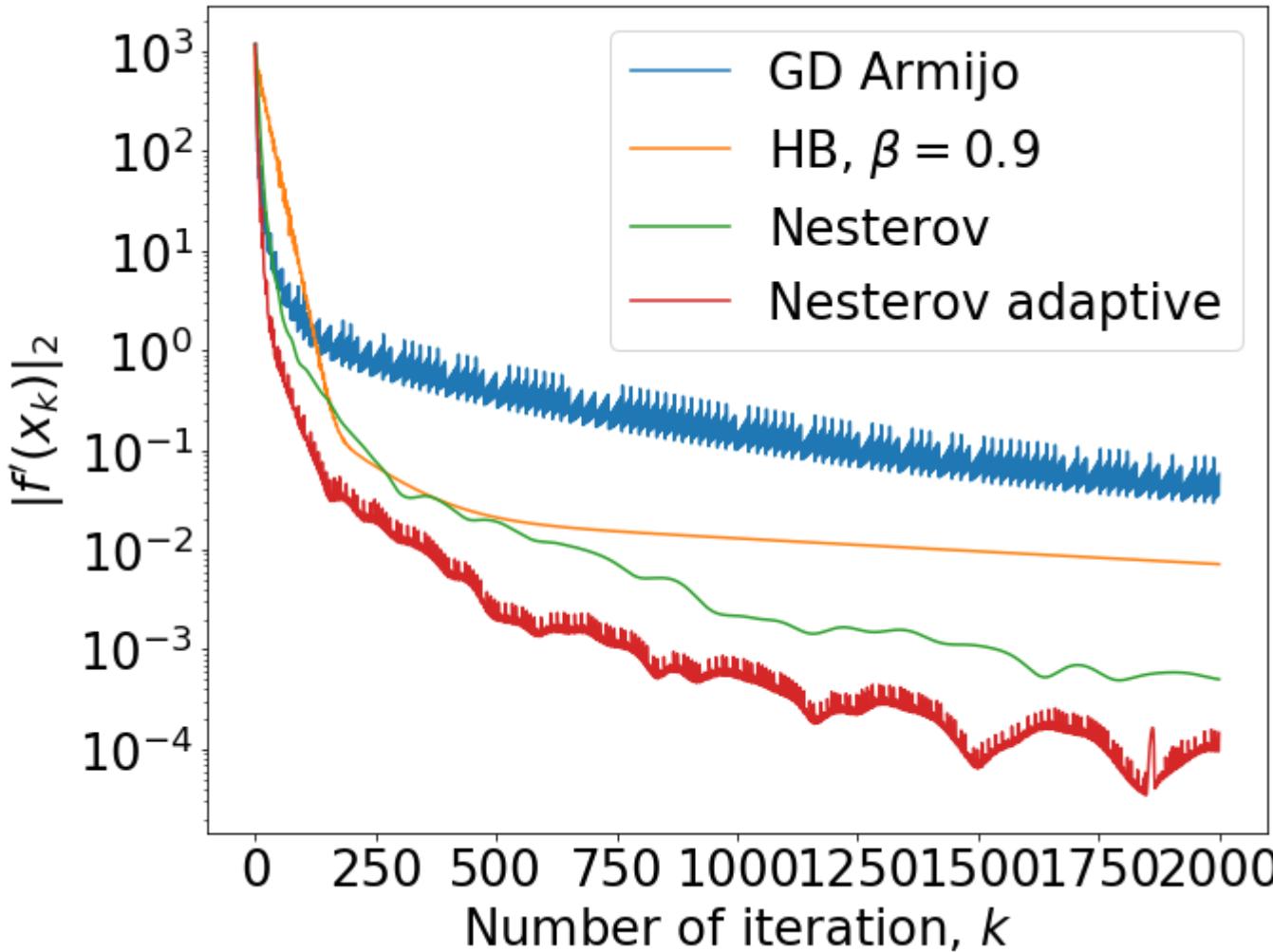
```

```

1 for m in methods:
2     _ = methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

1 figsize = (10, 8)
2 fontsize = 26
3 plt.figure(figsize=figsize)
4 for m in methods:
5     plt.semilogy([np.linalg.norm(grad(x)) for x in methods[m].get_convergence()], label=m)
6 plt.legend(fontsize=fontsize, loc="best")
7 plt.xlabel("Number of iteration, $k$", fontsize=fontsize)
8 plt.ylabel(r"$\| f'(x_k) \|_2$", fontsize=fontsize)
9 plt.xticks(fontsize=fontsize)
10 _ = plt.yticks(fontsize=fontsize)

```



```

1 for m in methods:
2     print(m)
3     %timeit methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

GD Armijo
402 ms ± 2.8 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
HB, $\beta = 0.9$
43.4 ms ± 949 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)
Nesterov
52.5 ms ± 1.15 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)
Nesterov adaptive
209 ms ± 7.13 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```

▼ Эксперимент на неквадратичной задаче

$$f(w) = \frac{1}{2} \|w\|_2^2 + C \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle x_i, w \rangle)) \rightarrow \min_w$$

```

1 import jax
2 import jax.numpy as jnp
3 from jax.config import config
4 config.update("jax_enable_x64", True)
5 import sklearn.datasets as skldata
6 n = 300
7 m = 1000
8
9 X, y = skldata.make_classification(n_classes=2, n_features=n, n_samples=m, n_informative=n//3, random_state=42)
10 C = 1
11
12 @jax.jit
13 def f(w):
14     return jnp.linalg.norm(w)**2 / 2 + C * jnp.mean(jnp.logaddexp(jnp.zeros(X.shape[0]), -y * (X @ w)))
15

```

```

16 # def grad(w):
17 #     denom = scspec.expit(-y * x.dot(w))
18 #     return w - C * x.T.dot(y * denom) / x.shape[0]
19
20 autograd_f = jax.jit(jax.grad(f))
21 x0 = jnp.ones(n)
22 print("Initial function value = {}".format(f(x0)))
23 print("Initial gradient norm = {}".format(jnp.linalg.norm(autograd_f(x0)))))

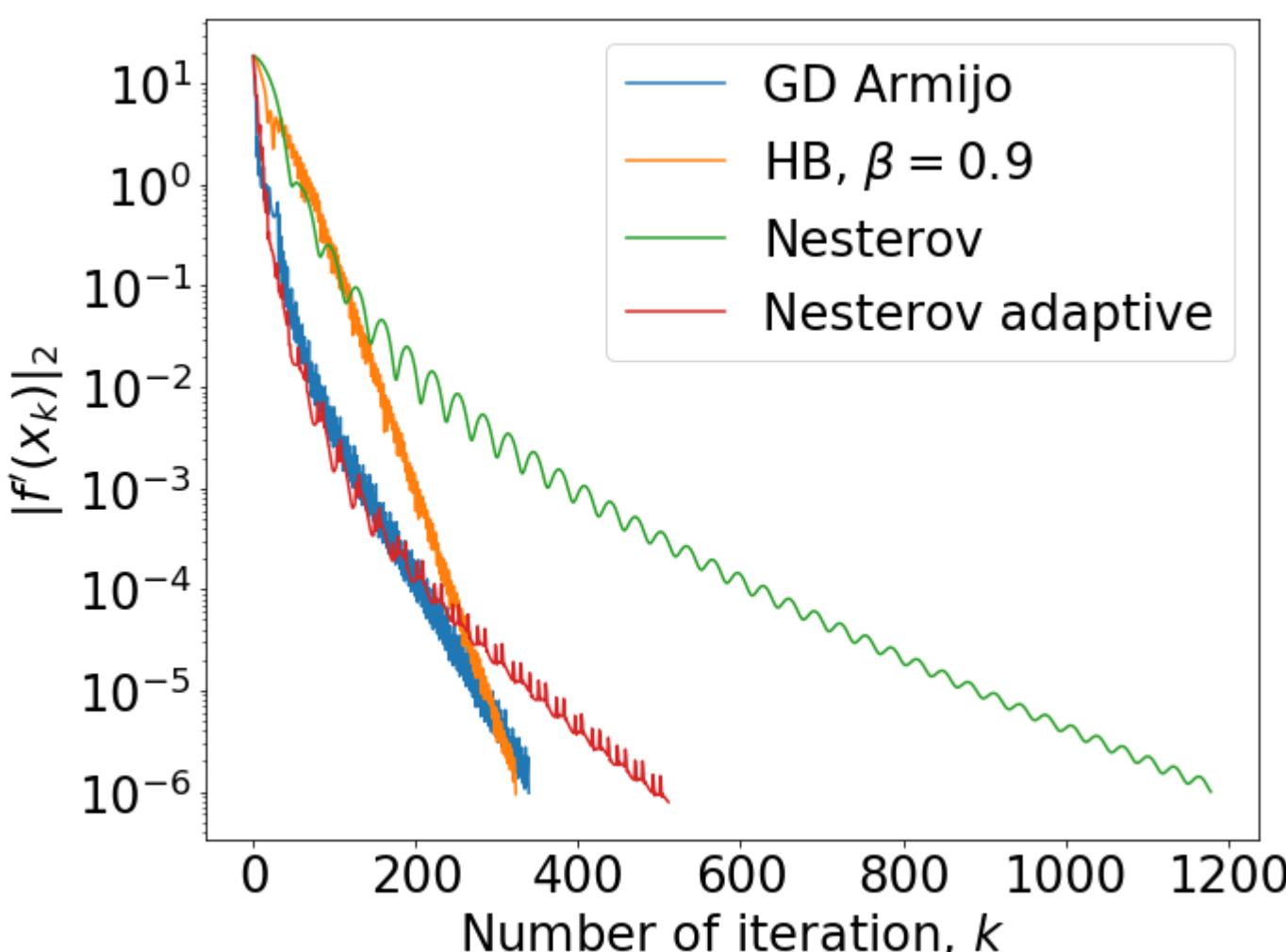
WARNING:absl:No GPU/TPU found, falling back to CPU. (Set TF_CPP_MIN_LOG_LEVEL=0 and rerun for more info.)
Initial function value = 162.63079545367793
Initial gradient norm = 18.723119978580637

1 beta_test = 0.9
2 L_trial = (1 + C*100)
3 methods = {
4     "GD Armijo": fo.GradientDescent(f, autograd_f,
5         ss.Backtracking("Armijo", rho=0.5, beta=0.01, init_alpha=1.)),
6     r"HB, $\beta = {}$".format(beta_test): HeavyBall(f, autograd_f, ss.ConstantStepSize(1 / L_trial), beta=beta_test),
7     "Nesterov": fo.AcceleratedGD(f, autograd_f, ss.ConstantStepSize(1 / L_trial)),
8     "Nesterov adaptive": fo.AcceleratedGD(f, autograd_f, ss.Backtracking(rule_type="Lipschitz", rho=0.5, init_alpha=1))
9 }
10 max_iter = 2000
11 tol = 1e-6

1 for m in methods:
2     _ = methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

1 figsize = (10, 8)
2 fontsize = 26
3 plt.figure(figsize=figsize)
4 for m in methods:
5     plt.semilogy([np.linalg.norm(autograd_f(x)) for x in methods[m].get_convergence()], label=m)
6 plt.legend(fontsize=fontsize, loc="best")
7 plt.xlabel("Number of iteration, $k$", fontsize=fontsize)
8 plt.ylabel(r"$\| f'(x_k) \|_2^2$", fontsize=fontsize)
9 plt.xticks(fontsize=fontsize)
10 _ = plt.yticks(fontsize=fontsize)

```



```

1 for m in methods:
2     print(m)
3     %timeit methods[m].solve(x0=x0, max_iter=max_iter, tol=tol)

GD Armijo
3.9 s ± 259 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
HB, $\beta = 0.9$
379 ms ± 8.48 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
Nesterov
2.45 s ± 108 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```

Nesterov adaptive
9.77 s ± 1.21 s per loop (mean ± std. dev. of 7 runs, 1 loop each)

Выходы

- Нижние оценки для методов первого порядка
- Метод тяжёлого шарика
- Ускоренный метод Нестерова
- Поиск константы Липшица градиента L