or simply just <i>unbiased compressor</i> , if there ex $\mathbb{E}\left[\mathcal{C}(x)\right]=x, \mathbb{E}\left[\ \mathcal{C}(x)\ \right]$	zed) map $\mathcal{C}: \mathbb{R}^d o \mathbb{R}^d$ is an <i>unbiased compression operate</i> sts a constant $\omega \geq 0$ such that	or,
$\mathcal{C}: \mathbb{R}^d \to \mathbb{R}^d$ is a biased compression operator $0 < \alpha \leq 1$ such that	sts a constant $\omega \geq 0$ such that $-x\ ^2 \leq \omega \ x\ ^2$, $\forall x \in \mathbb{R}^d$. (9) $\mathbb{U}(\omega)$. Further, we say that a (possibly randomized) may or simply just <i>biased compressor</i> , if there exists a constant	(2) ap ant
The family of such operators will be denoted В строго больше, так как содержит компрессор, который не возникает в ре 2.2 Error feedback: what it is good for, and what we sti Несмещенность оказывается очень эффективным инструментом, облегча	by $\mathbb{B}(lpha)$. It is well known that, in a certain sense, the latt вультате масштабирования несмещенного компрессора. Biased эмпирически превосходят do not know вщим анализ распределенных методов первого порядка, использующих несмещенные ког	τ Unbiased.
богат и относительно хорошо изучен. Например, используя несмещенные Смещенные же наоборот сложны для изучения. Основная сложность в то $\mathbf{Example} \ 1 \ \ Consider \ n = d = 3 \ \ and \ \ defin$ $f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \ x\ _2^2 , \qquad f_2(x) = \langle b \rangle^2 .$	омпрессоры, мы знаем (список).	
where $t > 0$. Then $\nabla f_1(x^0) = \frac{t}{2}(-11, 9, 9), \qquad \nabla f_2(x^0) = \frac{t}{2}(-11, 9, 9)$	$\frac{t}{2}(9, -11, 9), \qquad \nabla f_3(x^0) = \frac{t}{2}(9, 9, -11).$ $f(x^0) = \frac{t}{2}(-11, 0, 0), \ \mathcal{C}(\nabla f_2(x^0)) = \frac{t}{2}(0, -11, 0)$	
Repeated application gives $x^{k} = \left(1\right)$	$f_i(x^0) = \left(1 + \frac{11\eta}{6}\right) x^0.$ $+ \frac{11\eta}{6} x^0.$	
	е $ntially\ fast\ to\ +\infty$. ей ошибок в общих условиях гетерогенных данных, опираясь только на смещенные компой ноде (в отличие от стохастического метода), то после T раундов коммуникации можно $\mathbb{E}ig[abla f(x) ^2ig]=\mathbb{O}(rac{G^{2/3}}{T^{2/3}}),$	
в предположении	$ \nabla f_i(x) ^2 \leq G^2, \forall x \in \mathbb{R}^d, i \in 1, 2, \dots, n$ Algorithm sCVX nCVX EF Stich et al. [2018] EF-SGD Stich and Karimireddy [2019]	DIST key limitation bounded gradients; sublinear rate in sCVX case single node only
	EF Ajalloeian and Stich [2020] SignSGD Karimireddy et al. [2019] EC-SGD Beznosikov et al. [2020]	\mathbf{x} single node only moment bound; single node only linear rate only if $\nabla f_i(x^\star) = 0 \ \forall i$ linear rate only using
	Gorbunov et al. [2020b] DoubleSqueeze Tang et al. [2020] Qsparse-SGD, CSER Basu et al. [2019], Xie et al. [2020] EC-SGD Koloskova et al. [2020]	an extra unbiased compressor bounded compression error; slow $O(1/T^{2/3})$ rate in nCVX case bounded gradients; slow $O(1/T^{1/2})$ rate in nCVX case bounded gradients; slow $O(1/T^{2/3})$ rate in nCVX case
function Teкущая ситуация с анализом смещенных компрессоров. 2.3 Summary of contributions • Новый error feedback механизм. Методы работают в условиях распре	Known results for first order methods using biased and nCVX= supports nonconvex functions, DIST =	works in the distributed regime. †decentralized
 Упростились условия, накладываемые на данные: ■ L_i-гладкость функций f_i ■ существование глобального минимума f^{inf} ∈ ℝ^d • Линейная скорость для функций Поляка-Лоясевича Main Results Раньше схема сжатия была построена на следующем способе. В каждой 	оде вычисляем градиент, сжимаем его с помощью biased метода и отправляем сжатый гр $x^{t+1} = x^t - rac{\gamma}{n} \sum_{i=1}^n C(abla f_i(x^t))$	радиент на сервер. Сервер агрегирует все посылки от n агентов и делает шаг.
Однако в минимуме градиент функции не обязан быть равен 0, а потому і	$n \geq 1$ с $(\nabla f_i(x^t))$ сходится к некоторму миник $\mathbb{E}[C(abla f_i(x^t)) - abla f_i(x^t) ^2] \leq (1-lpha) abla f_i(x^t) ^2$ равая часть может быть очень большой, поэтому метод как бы не сходится к некому $ ilde x$. оры, которые, как мы надеемся, будут сходиться к нулям. Поскольку ввиду (3) применени	
«легких для передачи» (некоторому объекту, который мы будем называть Определим рекрсивный способ пострения компрессора, который назовем	апример, градиентов), сгенерированная некоторым алгоритмом. Эта последовательность главным») оценок этих векторов, используя компрессор $\mathbb{C}\in\mathbb{B}(lpha)$. Давайте пройдем несомпрессором Маркова. $\mathbb{M}(v^0):=C(v^0)$ $\mathbb{M}(v^{t+1}):=\mathbb{M}(v^t)+C(v^{t+1}-\mathbb{M}(v^t)), t\geq 0$	
<pre>import numpy as np import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import torch from sklearn.model_selection import train_test_split import math</pre> <pre>Peaлизуем этот алгоритм. Рассмотрим Тор-k , который работает по сл</pre>	d	
где x - некоторый вектор, который мы хотим сжать, причем его координат Пример. Пусть у нас есть вектор $x=(-11,9,9)$. Тогда после применени def compressor(np_tensor, ratio=0.5): with torch.no_grad(): tensor = torch.from_numpy(np.asarray(np_tensor)) numel = tensor.numel() k = max(int(numel * ratio), 1)		
<pre>values, indexes = torch.topk(torch.abs(tensor), k=k) values = tensor[indexes] # np.append(values_, values) # np.append(indexes_, indexes) return tensor, indexes, values def get_compr_object(tensor, ratio=0.5): tensor, indexes, values = compressor(tensor, ratio) for idx in range(len(tensor)): if idx not in indexes: tensor[idx] = 0 return tensor.numpy()</pre>		
<pre>def topk(tensor, k): indexes = np.abs(tensor).argsort()[-k:] return indexes def other_get_compr(tensor, k=1): true_k = min(len(tensor), k) indexes = topk(tensor, true_k) for idx in range(len(tensor)): if idx not in indexes: tensor[idx] = 0 return tensor</pre> <pre> Пример работы метода Тор-К.</pre>		
vector = torch.tensor([9, -11, 2, 3, 4]) print(get_compr_object(vector, ratio=0.5)) [9 -11 0 0 0] Для оптимизации возьмем задачу минимизации эмпирического риска:	$\min_{w\in\mathbb{R}^d}rac{1}{n}\sum_{i=1}^n l(g(w,x_i),y_i)+rac{\lambda}{2}\ w\ _2^2.$ игмоидную функцию потерь $l(z,u)=\ln(1+\exp(-uz))$ (Важно: u должен принимать з	значения —1 или 1). Такую задачу еще называют догистической регрессией
	игмоидную функцию потерь $l(z,y)=\ln(1+\exp(-yz))$ (Важно: y должен принимать: $ abla f(w)=rac{1}{n}\sum_{i=1}^nrac{-x_iy_iexp(-\omega^T\cdot x_i\cdot y_i)}{1+exp(-\omega^T\cdot x_i\cdot y_i)}+\lambda\omega$ руками, что вообще говоря плохо. Поэтому попробуем вычислить его, основываясь на оце $\ abla f(x)- abla f(y)\ _2\leq L\ x-y\ _2$	
dataset = "mushrooms.txt" #файл должен лежать в той же деректории, что и потероок from sklearn.datasets import load_svmlight_file data = load_svmlight_file(dataset) X, y = data[0].toarray(), data[1] n, d = X.shape y = y * 2 - 3	_size=0.8, shuffle=True)	
<pre>X_train, X_test, y_train, y_test = train_test_split(X, y, train) def grad(x, y, w, lambd): n, d = x.shape gradient = np.zeros(d) for i in range(n): exp = np.exp(-(np.dot(w, x[i, :]) * y[i])) gradient += 1 / n * ((-x[i, :] * y[i] * exp) / (1 + exp)</pre>		
w - веса x - данные return np.log(1 + np.exp(-1 * np.hstack(np.vstack(y)) * (return np.log(1 + np.exp(-1 * np.hstack(np.vstack(y))) * (return np.log(1 + np.exp(-1 * np.hstack(np.vstack(y)))) * (return np.exp(-1 * np.hstack(np.exp(-1 * np.h	.vstack(x) @ np.asarray(w))))	
Algorithm 1 EF21 (Single node) 1: Input: starting point $x^0 \in \mathbb{R}^d$, learning rate: 2: for $t = 0, 1, 2, \dots, T - 1$ do 3: $x^{t+1} = x^t - \gamma g^t$ 4: $g^{t+1} = g^t + \mathcal{C}(\nabla f(x^{t+1}) - g^t)$ 5: end for	e $\gamma > 0$, $g^0 = \mathcal{C}(\nabla f(x^0))$	
<pre># algo1 EF21 (single node) for logistic regression problem def ef21_singlenode(x, y, w_start, lr, iterations=1000, accura ''' w_start - starting point in R^d lr - learning rate (lr > 0) g_start = C(grad(f(x_start))) ''' w_temp = w_start g_start = other_get_compr(grad(x, y, w_start, lambd), k=1) history1 = [] # distance(g - g_start) history2 = [] # dist(gradient - 0) for _ in range(iterations): w = w_temp - lr * g_start</pre>	y=0.01, lambd=0.1):	
<pre>w = w_temp - lr * g_start gradient = grad(x, y, w, lambd) g = g_start + other_get_compr(gradient - g_start, k=1) cur_loss1 = np.linalg.norm(g - g_start, ord=2) history1.append(cur_loss1) cur_loss2 = np.linalg.norm(gradient) history2.append(cur_loss2) # if (cur_loss1 < accuracy): # return [history1, history2] w_temp = w g_start = g</pre>		
<pre>g_start = g return [history1, history2] # using Top-1 g_d = ef21_singlenode(x=X_train, y=y_train, w_start=np.zeros(X)) plt.plot(range(len(g_d[0][10:])), g_d[0][10:]) plt.title('Gradient Descent (EF21_single_node)') plt.xlabel('Iterations') plt.ylabel('\$ g_i - g_start_i \$') plt.grid() plt.show()</pre>	train.shape[1]), lr=0.000576, iterations=500, accuracy=0.001, lambd=0.1)	
0.10		
0.02 0.00 0 100 200 300 400 Iterations	500	
3.3 Compressed gradient descent using the Markov confidence $\mathbf{Algorithm~2}$ EF21 (Multiple nodes) 1: Input: starting point $x^0 \in \mathbb{R}^d$; $g_i^0 = \mathcal{C}(\nabla \mathbf{a})$ learning rate $\gamma > 0$; $g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0$ (known 2: for $t = 0, 1, 2, \ldots, T-1$ do	$f_i(x^0)$) for $i=1,\ldots,n$ (known by nodes and the master	er);
<pre>definit(self, x, y, lambd, number_machine, lr, iterat</pre>	числения	
x, y, w_start, lambd - параметры для инициализации машин number_machine - количество машин, на которых распределяем в lr - learning rate iterations - количество итераций ''' dimension = x[0].shape[1] self.x = x self.y = y self.w = np.zeros(dimension) self.machines = [EF21_Machine(self.x[i], self.y[i], self.y[i], self.g. self.g_start = np.sum(np.array([self.machines[i].get_g. self.number_machine = number_machine self.lr = lr self.iterations = iterations self.lambd = lambd def server_executes(self, accuracy): ''' accuracy - точность вычислений pаботаем на данных self.x, self.y ''' # history = [] for t in range(self.iterations): w_next = self.w - self.lr * self.g_start grads_from_machine = [self.machines[i].client_updaterine = self.g_start + np.sum(np.array(grads_from cur_loss1 = np.linalg.norm(g_next - self.g_start) history1.append(cur_loss1) # if history1[t] < accuracy: # return self.w, history1	<pre>f.w, dimension, lambd, lr) for i in range(number_machine)] start(x=self.x[i], y=self.y[i], w_start=self.w, lambd=lambd) for i in range e(self.x[i], self.y[i], w=w_next, lambd=self.lambd) for i in range(len(self.yelf.yelf.yelf.yelf.yelf.yelf.yelf.y</pre>	
x, y, w_start, lambd - параметры для инициализации машин number_machine - количество машин, на которых распределяем в lr - learning rate iterations - количество итераций ''' dimension = x[0].shape[1]	f.w, dimension, lambd, lr) for i in range(number_machine)] start(x=self.x[i], y=self.y[i], w_start=self.w, lambd=lambd) for i in range e(self.x[i], self.y[i], w=w_next, lambd=self.lambd) for i in range(len(self.machine), axis=0) / self.number_machine втасет на 5 частей. Они будут представлять из себя 5 устройств.	
x, y, w_start, lambd — параметры для инициализации машин number_machine — количество машии, на которых распределяем в lr — learning rate iterations — количество итераций dimension = x[0].shape[1] self.x = x self.y = y self.w = pp.zeros(dimension) self.machines = [EF21_Machine(self.x[i], self.y[i], se self.g_start = np.sum(np.array([self.machines[i].get_c self.number_machine = number_machine self.lr = lr self.iterations = iterations self.lambd = lambd def server_executes(self, accuracy): """ accuracy — точность вычислений pa6oraem на данных self.x, self.y """ # history = [] history1 = [] for t in range(self.iterations): w_next = self.w — self.lr * self.g_start grads_from_machine = [self.machines[i].client_updatates g_next = self.g_start + np.sum(np.array(grads_from_machine = [self.machines[i].client_updatates g_next = self.g_start + np.sum(np.array(grads_from_machine = [self.machines[i].client_updatates g_next = self.g_start + np.sum(np.array(grads_from_machine = [self.machines[i].client_updatates g_next = self.g_start = np.sinalg.norm(g_next — self.g_start) history1.append(cur_loss1) # if history1(t) < accuracy; # return self.w, history1 self.w = w_next	f.w, dimension, lambd, lr) for i in range (number_machine) start(x=self.x[i], y=self.y[i], w_start=self.w, lambd=lambd) for i in range (eself.x[i], self.y[i], w=w_next, lambd=self.lambd) for i in range (len(self-machine), axis=0) / self.number_machine for i in range (len(self-machine), axis=0) / self-machine for i in range (len(self-machine), axis=0) / self	
x, y, w_start, lambd — параметры для инициализации машии number_machine — количество машии, на которых распределяем в lr — learning rate iterations — количество итераций dimension = x(0].shape(1)	f.w, dimension, lambd, lr) for i in range(number_machine)] start(x=self.x(i), y=self.y(i), w_start=self.w, lambd=lambd) for i in range (self.x(i), self.y(i), w=w_next, lambd=self.lambd) for i in range (len(self-machine), axis=0) / self.number_machine $0 \leq \gamma \leq \left(L + \widetilde{L}\sqrt{\frac{\beta}{\theta}}\right)^{-1},$ 0.1. $t(i)$ for i in range(NUM_MACHINE)) $-1)$	
x, y, w_start, lambd — параметры для инициализация машии number_machine — моничество машия, на моторых распределем я lr — learning rate iterations — моничество иторыций dimension = x[0].shape[1]	f.w, dimension, lambd, lr) for i in range(number_machine)] start(x=self.x(i), y=self.y(i), w_start=self.w, lambd=lambd) for i in range (self.x(i), self.y(i), w=w_next, lambd=self.lambd) for i in range(len(self-machine), axis=0) / self.number_machine $0 \le \gamma \le \left(L + \widetilde{L}\sqrt{\frac{\beta}{\theta}}\right)^{-1},$ 0.1. $t(i)$ for i in range(NUM_MACHINE)) $-1)$	
x, y, w_start, lambd - параметры для инперацизацию мощном number_machine - можностом мощно, на моторых растроромием I ir - learning rate iterations - мошчество итераций dimension = x[0].shape[1]	f.w, dimension, lambd, lr) for i in range(number_machine)] start(x=self.x(i), y=self.y(i), w_start=self.w, lambd=lambd) for i in range (self.x(i), self.y(i), w=w_next, lambd=self.lambd) for i in range(len(self-machine), axis=0) / self.number_machine $0 \le \gamma \le \left(L + \widetilde{L}\sqrt{\frac{\beta}{\theta}}\right)^{-1},$ 0.1. $t(i)$ for i in range(NUM_MACHINE)) $-1)$	
x, y, w_start, lambd - параметры для изисправизацию машли number_sachine - моличество миции, из моторых растрежением in rearning sate iterations - моличество итераций disension = x(0).shap(i) self.x = x self.y = y self.w = np.seros(dimension) self.machines = [EF21_Xanhine(self.x[i], self.y[i], s	tracer wa 5 vacreA. Own Sygyr speggragases is a cost of Sycrophysical content of the cost	
$x_1, y_1, y_2, y_3, y_4, y_4, y_4, y_4, y_4, y_4, y_4, y_4$	the end of the state of the st	f.montiren);]
The contract of the property	better. $0 \leq r \leq \left(k + \frac{T_{i}(\overline{S})}{T_{i}}\right)^{1},$ where S is the conjugate of the state of	f.montiren);]
The property control of apparents of immeriations the material process of the control of the co	For dissipation, lands in for it is approximately stated with present yield victorized to landwisered for it is concerned with present yield victorized to landwisered for it is concerned with present of the concerned with present of the concerned with the con	f.montiren);]
The state of the	the decides and the property of the state o	f.montiren);]
The control is a second property of the control is an interest to the control of	the decides and the property of the state o	f.montiren);]
The property of the property of the companion of the property	The contribution of the c	Constituent
For the state of the property of the company common property of the company of t	Fig. 6 statements, leaves, the form is reasonable to the content of the content	Constituent
The probability and the companies an employ pulpoparies of the probability and the pr	Fig. 6 statements, leaves, the form is reasonable to the content of the content	Constituent
The companion was a constrained with the control of the control o	Fig. 6 statements, leaves, the form is reasonable to the content of the content	Constituent
The processor of the pr	The denotion of the control of the control of product t and t anothing t and t	Constituent
The state of the	The contract of the contract	Constituent
The state of the	The state of the	Constituent
The contract of the contract	The contract of the contract	Constituent
Set of the control of	to zero. $a^{(j)} = 0$ The state of the st	Constituent
Set of the control of	For all and a second points of the consequence of the transfer of the transfer of the consequence of the transfer of the transfer of the consequence of the transfer of the t	Constituent
The control of the co	For considering the first is a stage across another constraints of the constraints of th	Excitator : 1
Compute starting point 20 Eggs and 10 and 20 Eggs and	For consider, such that for a large action makes on the constraint of the constrain	Excitator : 1
Algorithm 3 FF21 (Multiple nodes) 1. Imput: statistical beautiful and statistical statist	For consider, such that for a large action makes on the constraint of the constrain	Excitator : 1
Against a service would be a service of the service	For consider, such that for a large action makes on the constraint of the constrain	Excitator : 1