

# Поиск аномалий

HBOS и ECOD

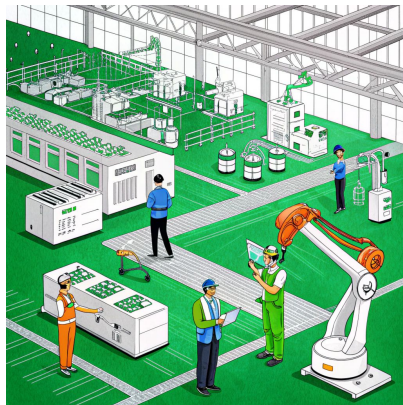
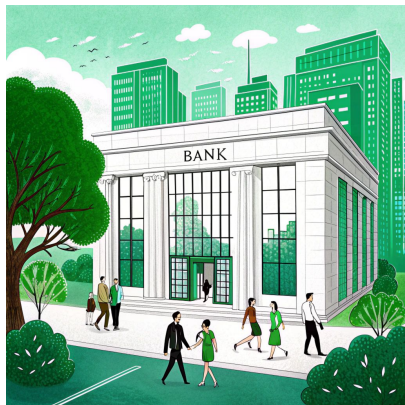
# Обо мне

- Старший специалист по машинному обучению
- Deep learning engineer
- NLP, CV, anomaly detection
- Open source contributor
- Выпускник и амбассадор Яндекс Практикума
- Выпускник DLS ФПМИ МФТИ



# Аномалии

# Применение





# Свойства

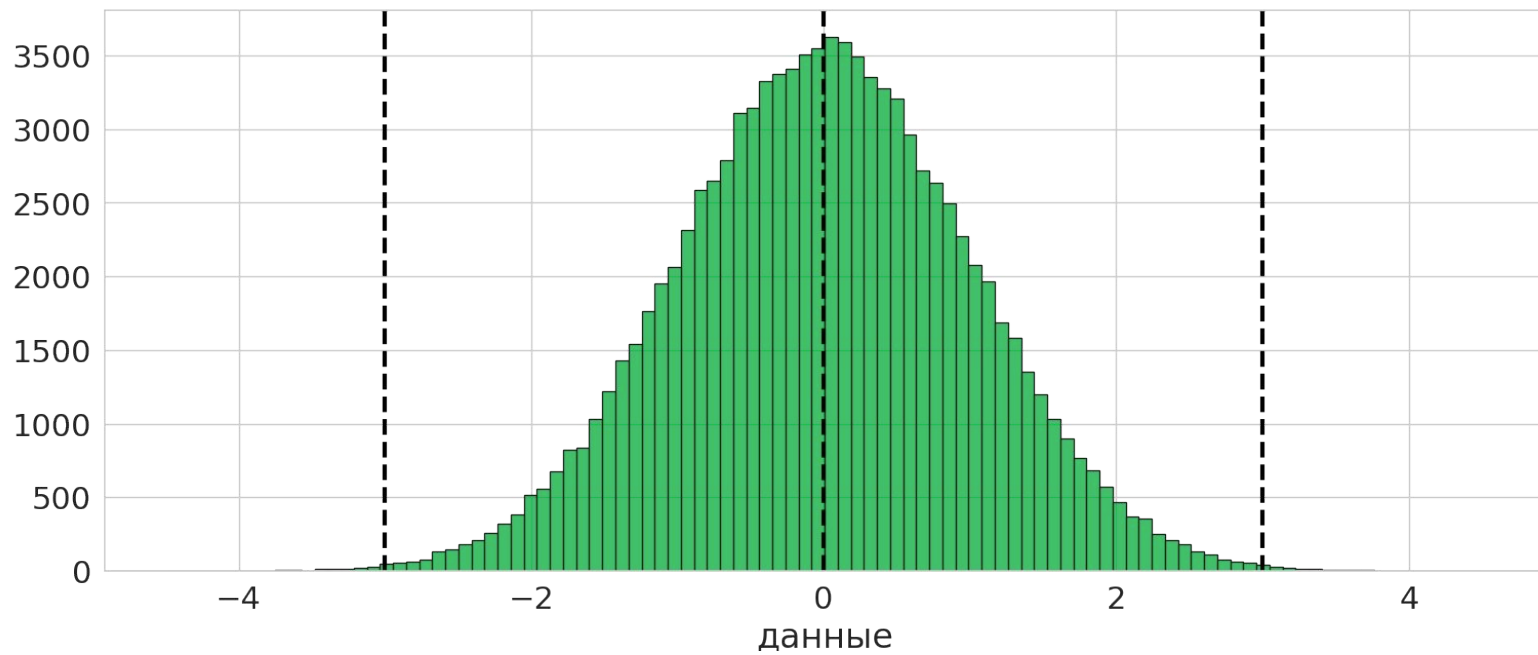


# Методы

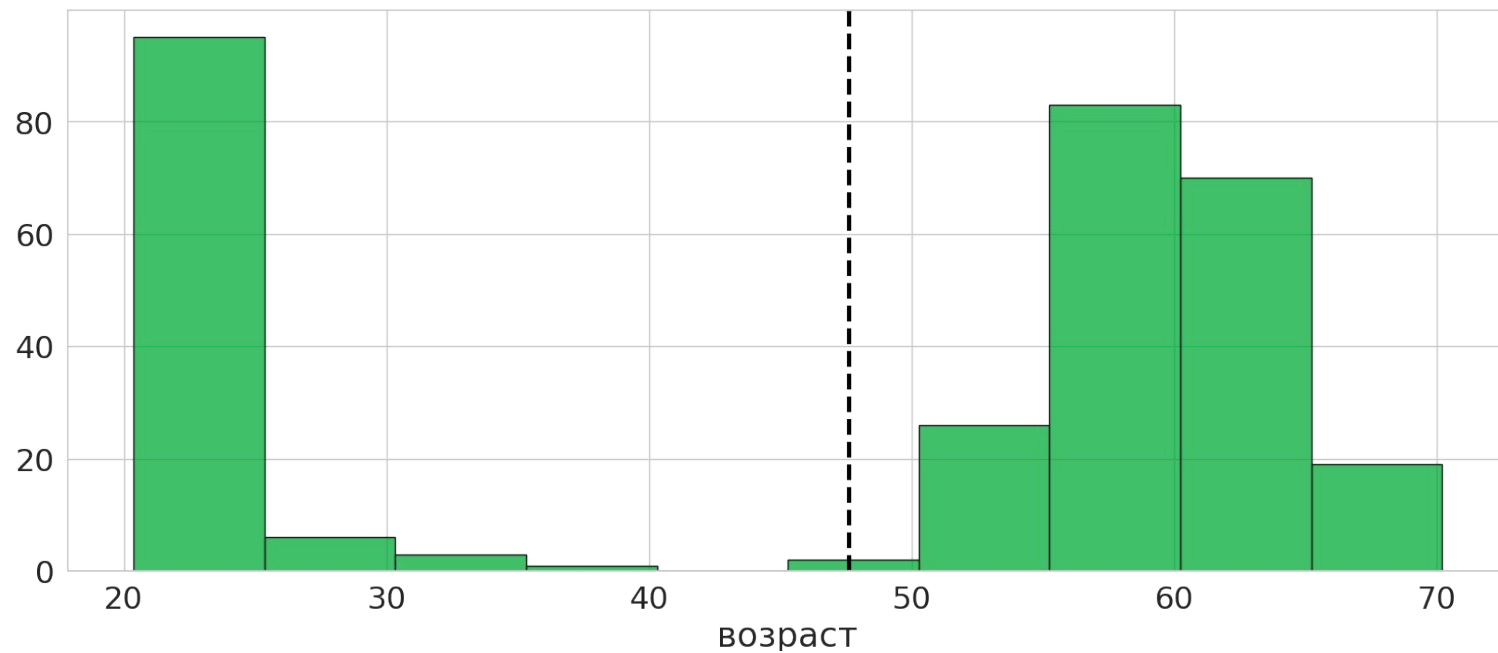
- [1] Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): a fast unsupervised anomaly detection algorithm. *KI-2012: Poster and Demo Track*, pages 59–63, 2012.
- [2] Zheng Li, Yue Zhao, Xiyang Hu, Nicola Botta, Cezar Ionescu, and H. George Chen. Ecod: unsupervised outlier detection using empirical cumulative distribution functions. *IEEE Transactions on Knowledge and Data Engineering*, 2022.

\$ pip install pyod

# Три сигмы



# Пионеры и пенсионеры





# Гистограмма в таблице

(20.4, 25.4]	95			
(25.4, 30.3]	6			
(30.3, 35.3]	3			
(35.3, 40.3]	1			
(40.3, 45.3]	0			
(45.3, 50.3]	2			
(50.3, 55.2]	26			
(55.2, 60.2]	83			
(60.2, 65.2]	70			
(65.2, 70.2]	19			

# Плотность вероятности

(20.4, 25.4]	95	<b>95 / 305 / 5.0</b>		
(25.4, 30.3]	6	<b>6 / 305 / 5.0</b>		
(30.3, 35.3]	3	<b>3 / 305 / 5.0</b>		
(35.3, 40.3]	1	<b>1 / 305 / 5.0</b>		
(40.3, 45.3]	0	<b>0 / 305 / 5.0</b>		
(45.3, 50.3]	2	<b>2 / 305 / 5.0</b>		
(50.3, 55.2]	26	<b>26 / 305 / 5.0</b>		
(55.2, 60.2]	83	<b>83 / 305 / 5.0</b>		
(60.2, 65.2]	70	<b>70 / 305 / 5.0</b>		
(65.2, 70.2]	19	<b>19 / 305 / 5.0</b>		

# Плотность вероятности

(20.4, 25.4]	95	<b>0.063</b>		
(25.4, 30.3]	6	<b>0.004</b>		
(30.3, 35.3]	3	<b>0.002</b>		
(35.3, 40.3]	1	<b>0.001</b>		
(40.3, 45.3]	0	<b>0.000</b>		
(45.3, 50.3]	2	<b>0.001</b>		
(50.3, 55.2]	26	<b>0.017</b>		
(55.2, 60.2]	83	<b>0.055</b>		
(60.2, 65.2]	70	<b>0.046</b>		
(65.2, 70.2]	19	<b>0.013</b>		

# Перемножаем

(20.4, 25.4]	95	0.063	<del>0.063</del> * ...	
(25.4, 30.3]	6	0.004		
(30.3, 35.3]	3	0.002		
(35.3, 40.3]	1	0.001		
(40.3, 45.3]	0	0.000		
(45.3, 50.3]	2	0.001		
(50.3, 55.2]	26	0.017		
(55.2, 60.2]	83	0.055		
(60.2, 65.2]	70	0.046		
(65.2, 70.2]	19	0.013		

Перемножаем  $\log_2(x \ y) = \log_2(x) + \log_2(y)$

(20.4, 25.4]	95	0.063	<del>0.063</del> * ...	
(25.4, 30.3]	6	0.004		
(30.3, 35.3]	3	0.002		
(35.3, 40.3]	1	0.001		
(40.3, 45.3]	0	0.000		
(45.3, 50.3]	2	0.001		
(50.3, 55.2]	26	0.017		
(55.2, 60.2]	83	0.055		
(60.2, 65.2]	70	0.046		
(65.2, 70.2]	19	0.013		

# Складываем логарифмы

(20.4, 25.4]	95	0.063	$\log_2(0.063) + \dots$	
(25.4, 30.3]	6	0.004	$\log_2(0.004) + \dots$	
(30.3, 35.3]	3	0.002	$\log_2(0.002) + \dots$	
(35.3, 40.3]	1	0.001	$\log_2(0.001) + \dots$	
(40.3, 45.3]	0	0.000	$\log_2(0.000) + \dots$	
(45.3, 50.3]	2	0.001	$\log_2(0.001) + \dots$	
(50.3, 55.2]	26	0.017	$\log_2(0.017) + \dots$	
(55.2, 60.2]	83	0.055	$\log_2(0.055) + \dots$	
(60.2, 65.2]	70	0.046	$\log_2(0.046) + \dots$	
(65.2, 70.2]	19	0.013	$\log_2(0.013) + \dots$	

# Складываем логарифмы

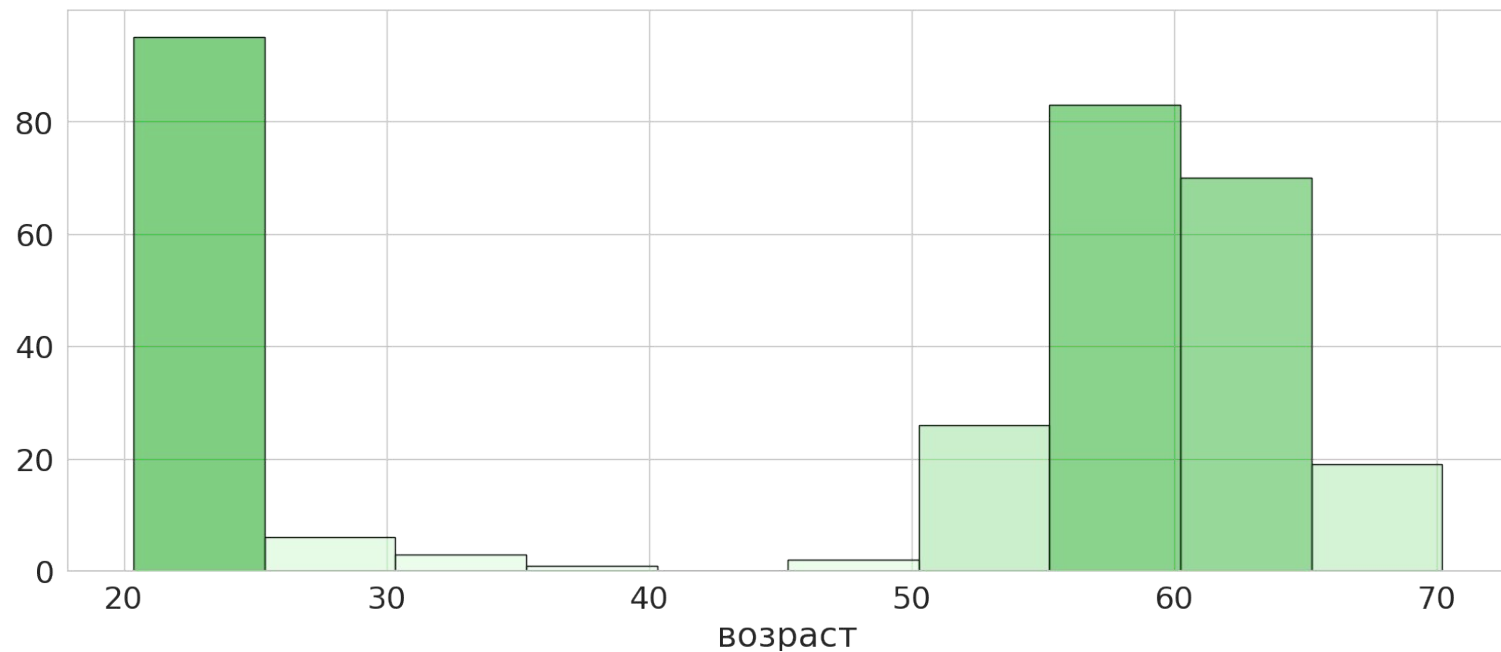
(20.4, 25.4]	95	0.063	<b>-2.621</b>	
(25.4, 30.3]	6	0.004	<b>-3.266</b>	
(30.3, 35.3]	3	0.002	<b>-3.294</b>	
(35.3, 40.3]	1	0.001	<b>-3.312</b>	
(40.3, 45.3]	0	0.000	<b>-3.322</b>	
(45.3, 50.3]	2	0.001	<b>-3.303</b>	
(50.3, 55.2]	26	0.017	<b>-3.094</b>	
(55.2, 60.2]	83	0.055	<b>-2.693</b>	
(60.2, 65.2]	70	0.046	<b>-2.775</b>	
(65.2, 70.2]	19	0.013	<b>-3.152</b>	



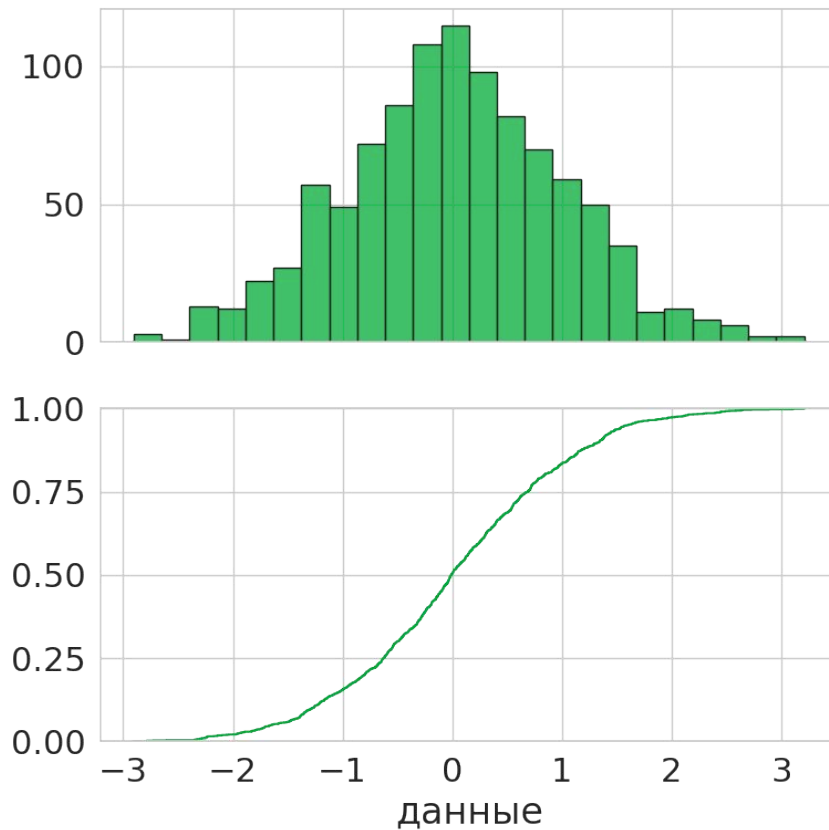
# Меняем знак

(20.4, 25.4]	<b>95</b>	0.063	-2.621	<b>2.621</b>
(25.4, 30.3]	<b>6</b>	0.004	-3.266	<b>3.266</b>
(30.3, 35.3]	<b>3</b>	0.002	-3.294	<b>3.294</b>
(35.3, 40.3]	<b>1</b>	0.001	-3.312	<b>3.312</b>
(40.3, 45.3]	<b>0</b>	0.000	-3.322	<b>3.322</b>
(45.3, 50.3]	<b>2</b>	0.001	-3.303	<b>3.303</b>
(50.3, 55.2]	<b>26</b>	0.017	-3.094	<b>3.094</b>
(55.2, 60.2]	<b>83</b>	0.055	-2.693	<b>2.693</b>
(60.2, 65.2]	<b>70</b>	0.046	-2.775	<b>2.775</b>
(65.2, 70.2]	<b>19</b>	0.013	-3.152	<b>3.152</b>

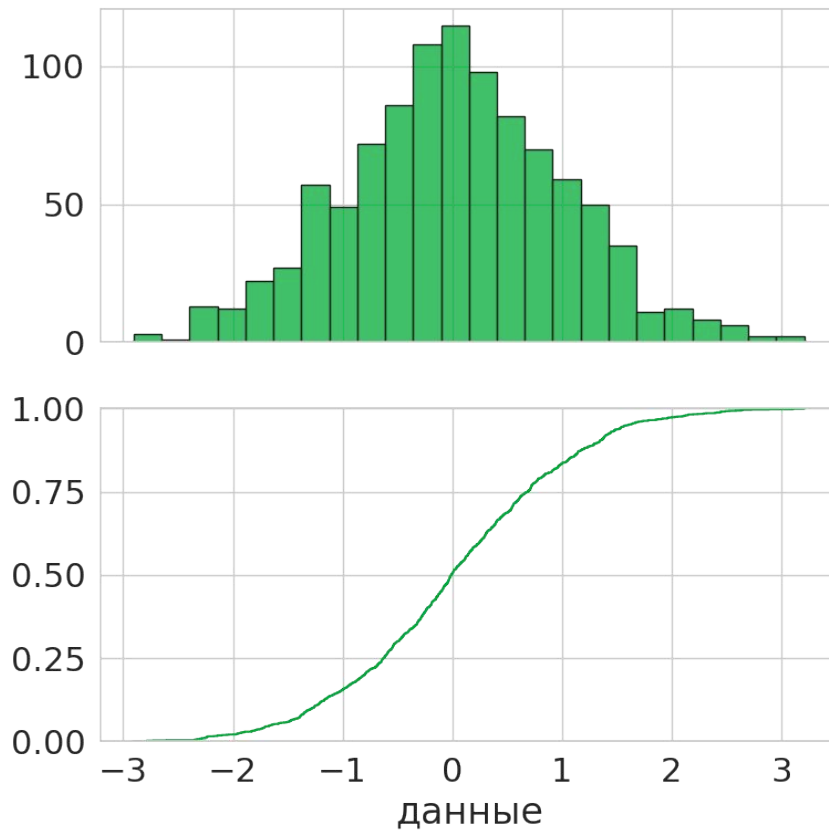
# Histogram-Based Outlier Score [1]



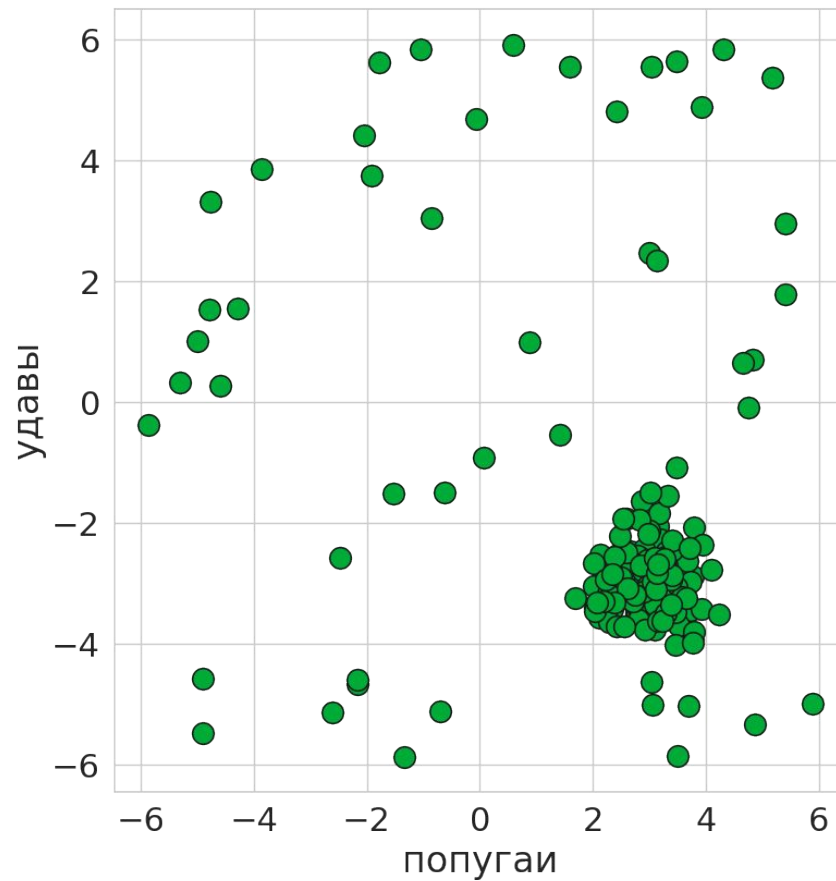
# Выборочная эмпирическая функция распределения

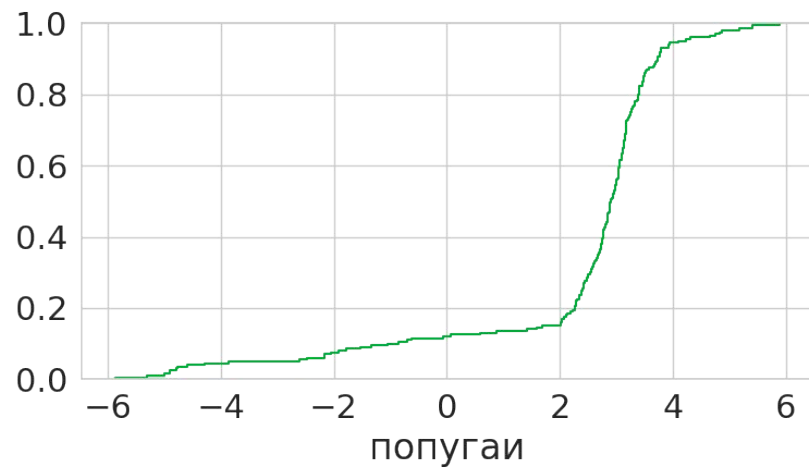
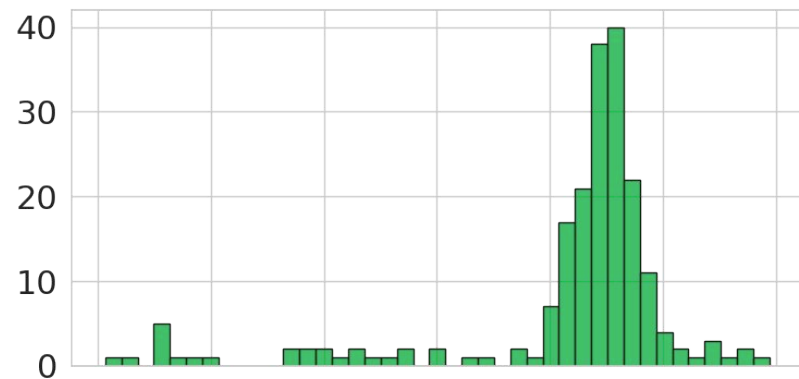


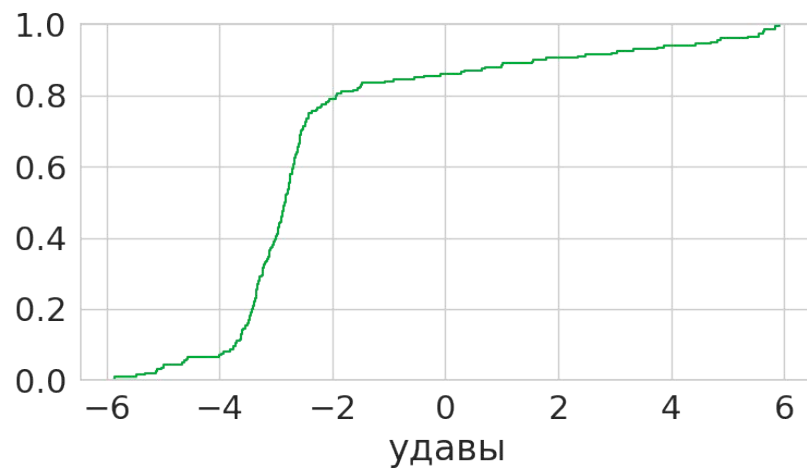
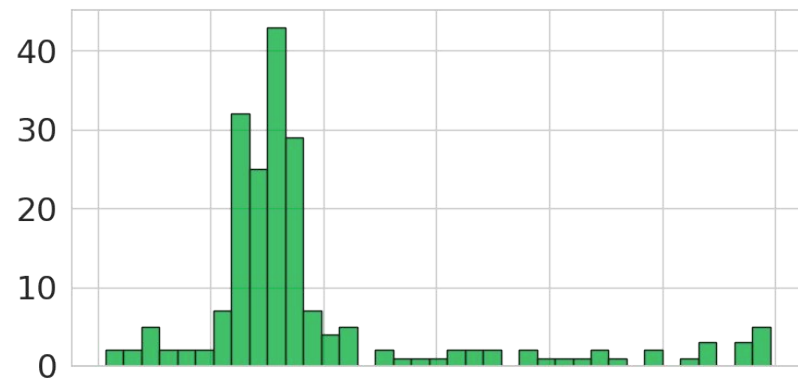
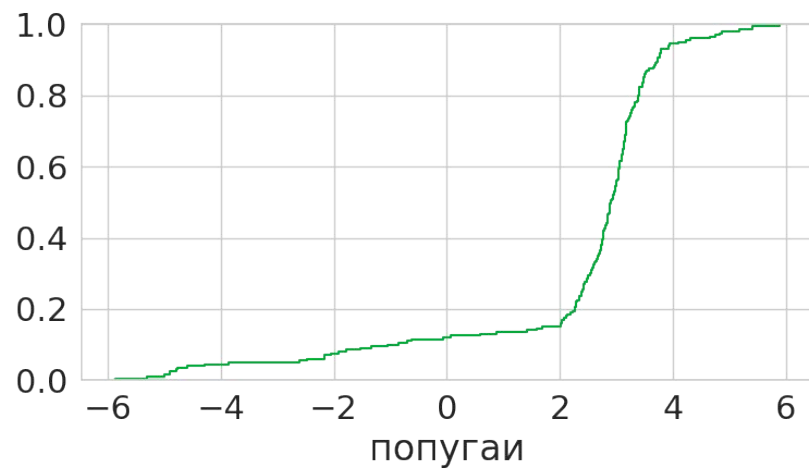
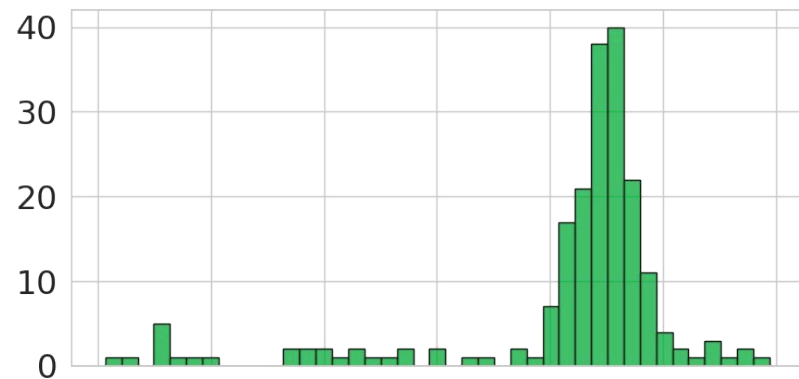
```
>>> x = -1  
>>> (df['данные'] <= x).mean()  
np.float64(0.155)
```



Данные









# ЭФР и 1 - ЭФР

Удавы	ЭФР	1 - ЭФР			
-2.87	0.48	0.52			
-2.83	0.51	0.50			
-3.34	0.26	0.75			
-2.88	0.48	0.53			
...	...	...			
-0.55	0.85	0.16			
1.53	0.90	0.11			
1.01	0.89	0.12			
4.81	0.96	0.05			

# Negative log probs

Удавы	ЭФР	1 - ЭФР	- log(ЭФР)	-log(1 - ЭФР)	
-2.87	0.48	0.52	<b>0.73 + ...</b>	<b>0.64 + ...</b>	
-2.83	0.51	0.50	<b>0.67 + ...</b>	<b>0.70 + ...</b>	
-3.34	0.26	0.75	<b>1.37 + ...</b>	<b>0.29 + ...</b>	
-2.88	0.48	0.53	<b>0.74 + ...</b>	<b>0.63 + ...</b>	
...	...	...	...	...	
-0.55	0.85	0.16	<b>0.16 + ...</b>	<b>1.86 + ...</b>	
1.53	0.90	0.11	<b>0.11 + ...</b>	<b>2.21 + ...</b>	
1.01	0.89	0.12	<b>0.12 + ...</b>	<b>2.16 + ...</b>	
4.81	0.96	0.05	<b>0.05 + ...</b>	<b>3.00 + ...</b>	

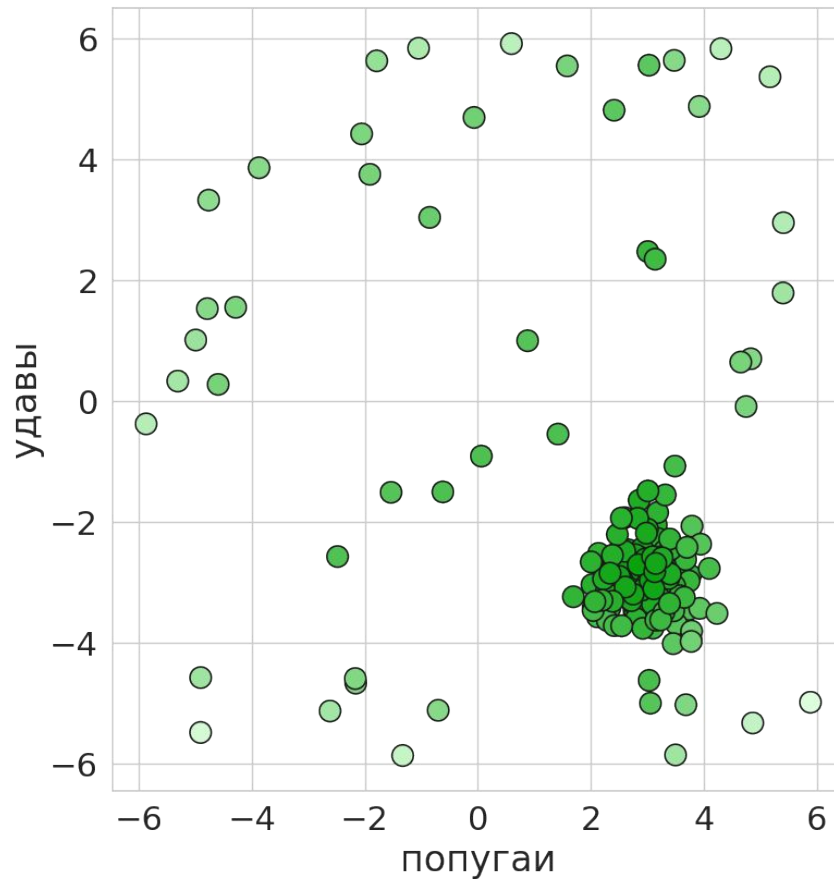
# Negative log probs

Удавы	ЭФР	1 - ЭФР	- log(ЭФР)	-log(1 - ЭФР)	
-2.87	0.48	0.52	<b>1.03</b>	<b>1.99</b>	
-2.83	0.51	0.50	<b>1.34</b>	<b>1.42</b>	
-3.34	0.26	0.75	<b>1.61</b>	<b>1.78</b>	
-2.88	0.48	0.53	<b>0.84</b>	<b>2.99</b>	
...	...	...	...	...	
-0.55	0.85	0.16	<b>2.13</b>	<b>2.01</b>	
1.53	0.90	0.11	<b>3.62</b>	<b>2.23</b>	
1.01	0.89	0.12	<b>4.32</b>	<b>2.17</b>	
4.81	0.96	0.05	<b>1.39</b>	<b>3.29</b>	

# ECOD

Удавы	ЭФР	1 - ЭФР	- log(ЭФР)	-log(1 - ЭФР)	max
-2.87	0.48	0.52	1.03	1.99	<b>1.99</b>
-2.83	0.51	0.50	1.34	1.42	<b>1.42</b>
-3.34	0.26	0.75	1.61	1.78	<b>1.78</b>
-2.88	0.48	0.53	0.84	2.99	<b>2.99</b>
...	...	...	...	...	...
-0.55	0.85	0.16	2.13	2.01	<b>2.13</b>
1.53	0.90	0.11	3.62	2.23	<b>3.62</b>
1.01	0.89	0.12	4.32	2.17	<b>4.32</b>
4.81	0.96	0.05	1.39	3.29	<b>3.29</b>

## Empirical Cumulative Outlier Detection [2]



# Резюме

Вопросы?

