

Ежегодная международная научно-практическая конференция
«РусКрипто'2020»

Подход к классификации последовательностей, сформированных алгоритмами сжатия и шифрования

Козачок Александр Васильевич,
д.т.н., сотрудник Академии ФСО России, г. Орёл
Спирин Андрей Андреевич,
сотрудник Академии ФСО России, г. Орёл

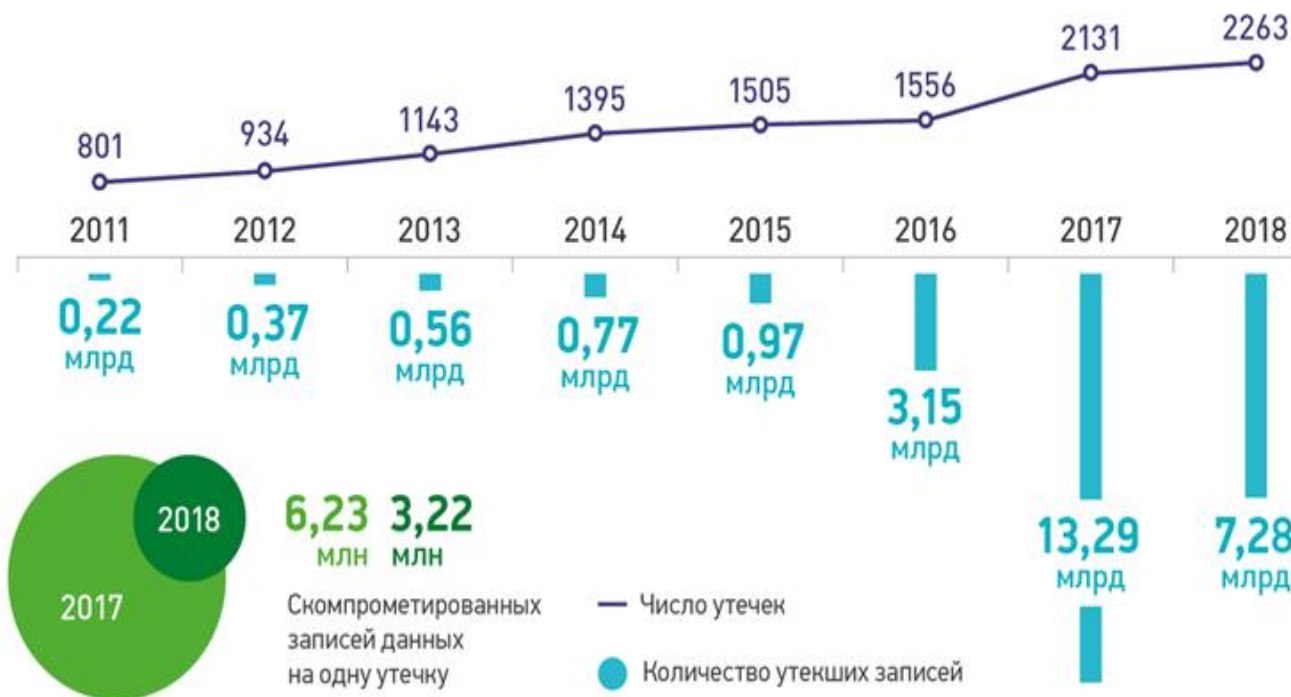


Рисунок 1. Число зарегистрированных утечек информации за 2006 – 2018 гг.

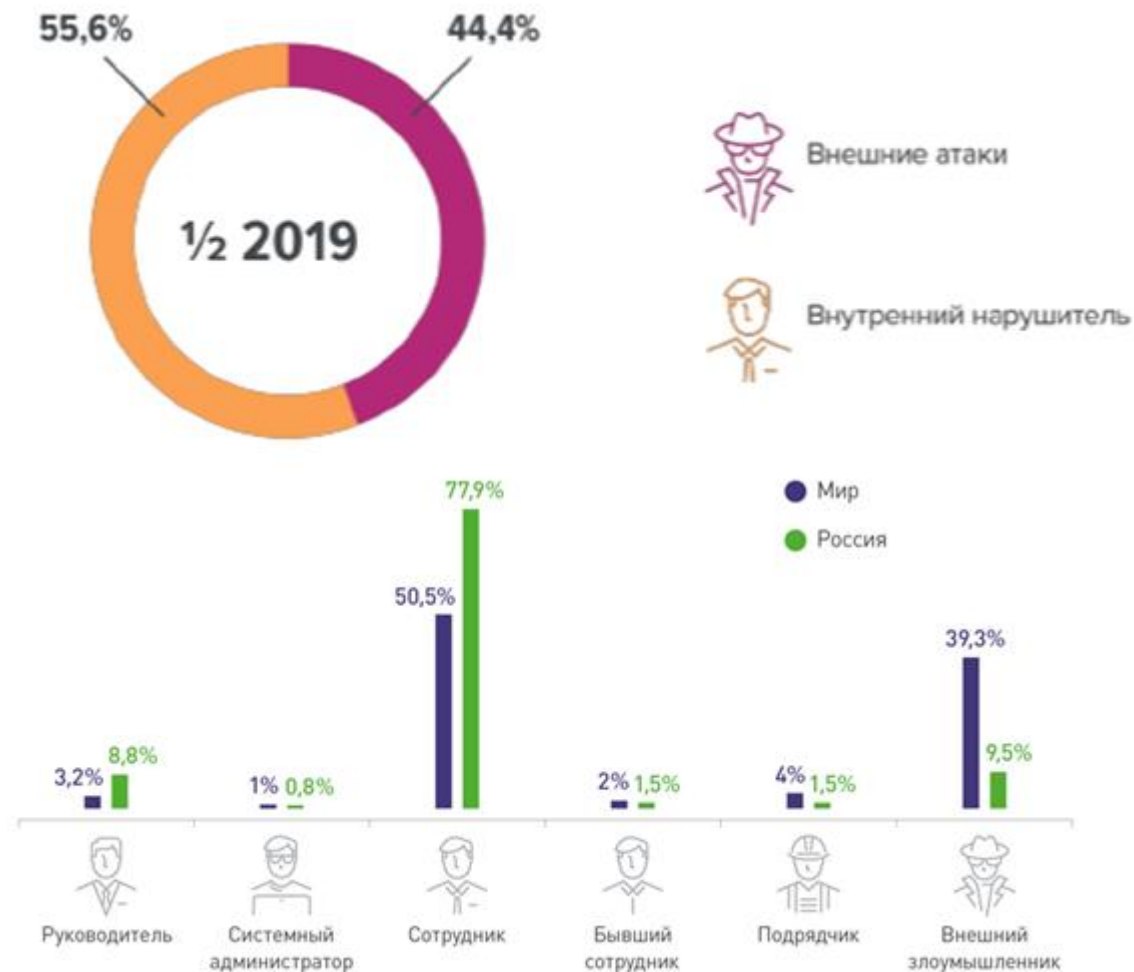


Рисунок 2. Соотношение нарушителей информационной безопасности

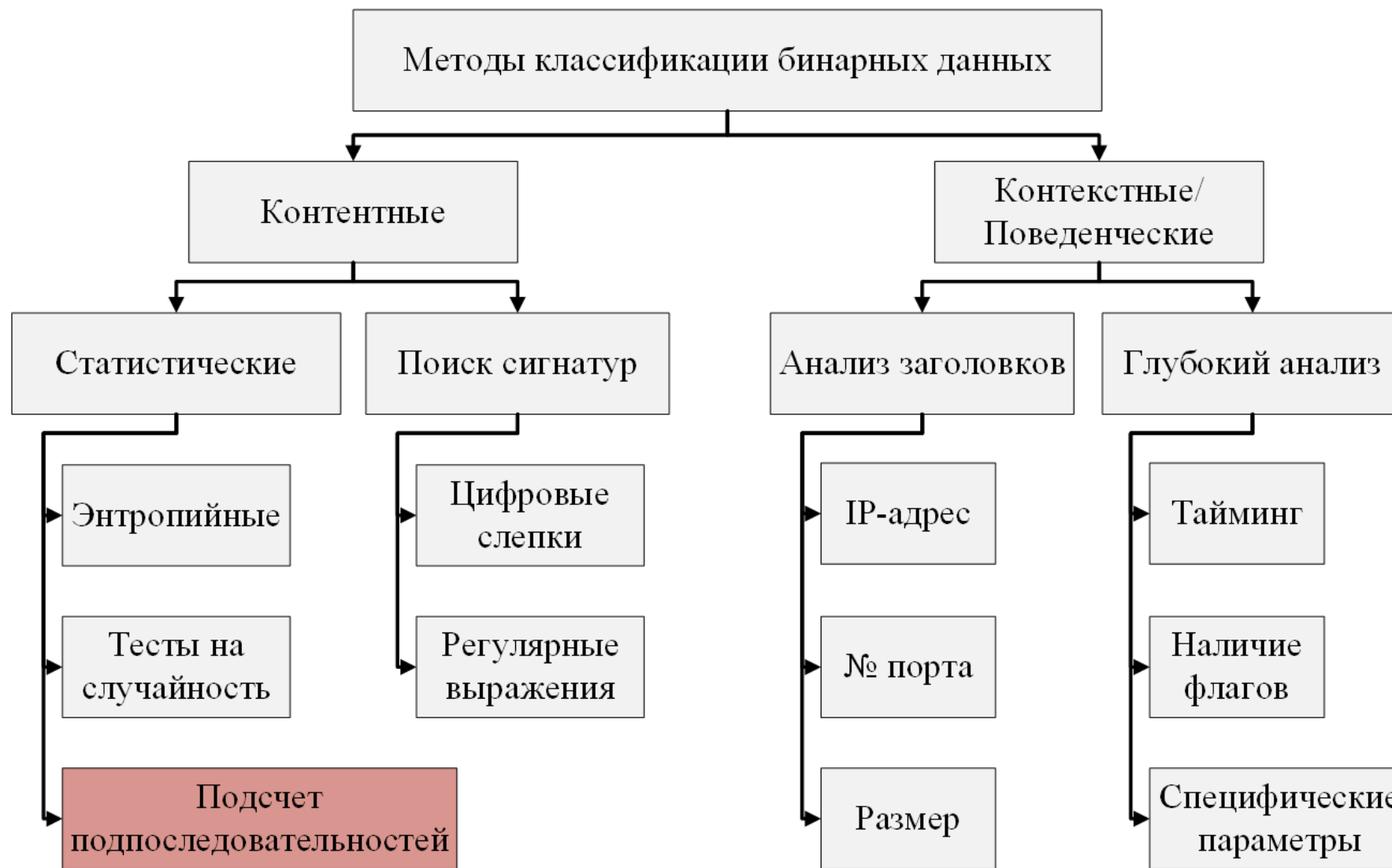
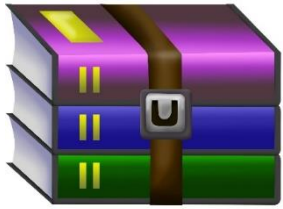
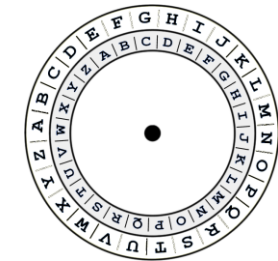


Рисунок 3. Методы, применяемые для классификации зашифрованных / сжатых данных

Авторы	Год	Задача	Признаки	Алгоритм	Результаты
Wright C., Monrose F., Masson G. M.	2005	Идентификация сетевых протоколов	Длина пакета, время жизни,	Скрытые Марковские Модели	HTTPS - 94.4%
Conti G.	2010	Классификация бинарных файлов	IP-адреса	(алгоритм Витерби)	HTTP - 96.7%
Alshammari R.	2011	Идентификация зашифрованного трафика	ip.flags, ip.len, ip.checksum, tcp.len*, tcp.flags*, oth.	C4.5, AdaBoost, team-based GA	SSH-71%
Dorfinger P., Panholzer G., John W.	2011	Идентификация зашифрованного трафика	Оценка энтропии полезной нагрузки	Критерий совпадения	eDonkey, Skype - 94% SMTP, HTTP, POP3, FTP -99%
Wang Y., Zhang Z., Guo L., Li S.	2011	Классификация трафика	Энтропийный подход, подсчет 4-битных подпоследовательностей	SVM	86.4% for text 79.8% enc
Zhang H., Papadopoulos C., Massey D.	2013	Идентификация зашифрованного трафика botnet	Энтропия потоков и пакетов	Метод Монте-Карло	flow (SSH,HTTPS) -95% packet (SSH,HTTPS) - 97%
Khakpour A. R., Liu A. X.	2013	Идентификация бинарных, текстовых, зашифрованных данных	Энтропия	CART, SVM	88%
Hahn D., Apthorpe N., Feamster N.	2018	Идентификация зашифрованных и сжатых данных	Энтропия, Хи-квадрат	k - ближайших соседей, сверточные нейронные сети, сети прямого распространения	k-Nearest Neighbors 60.0% Feed Forward Neural Net 54.1% CNN 66.9%
Casino F., Choo K.-K. R., Patsakis C.	2019	Классификация зашифрованного/сжатого трафика	Тесты NIST: frequency block test, cum sums test, approximate entropy test.	HEDGE (High Entropy DistinGuishEr)	70.61%
Tang Z., Zeng X., Sheng Y.	2019	Классификация зашифрованного/сжатого трафика	Энтропия 4,8,16,24 битных подпоследовательностей в последовательностях, полученных из исходной с различным шагом (от 1 до 8)	SVM, RF	for traffic (ip-адреса, № портов, протокол) 97.9% video <70% Images, text <72% audio <66%



$$F : X \rightarrow Y \quad (1)$$



$$p(F(x_i) = y_j \mid i = j) \rightarrow 1 \quad (2)$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

$$Precision = \frac{TP}{TP + FP} \quad (4) \quad Recall = \frac{TP}{TP + FN} \quad (5)$$



где TP – количество объектов верно отнесенных к классу i , TN – количество объектов верно отнесенных к классу j ($j \neq i$), FP – количество ложных срабатываний (ошибка первого рода), FN – количество пропусков цели (ошибка второго рода).



Рисунок 4. Функциональная модель системы классификации последовательностей, сформированных алгоритмами сжатия и шифрования



Data: $P: |P|=Q, S: |S| = 2^{N-1}$
Result: $F_{Q,E}$

```

1  $F_{Q,E} \leftarrow \langle \rangle$ 
2 for  $p \in P$  do
3    $M_p \leftarrow \text{Len}(p)$ 
4   for  $s \in S$  do
5      $n_s \leftarrow \text{Count}(p,s)$ 
6      $f_{p,s} \leftarrow \frac{n_s}{M_p - N_s + 1}$ 
7      $F_{Q,E} \leftarrow F_{Q,E} \cup \langle f_{p,s}, y_i \rangle$ 
8 return  $F_{Q,E}$ 
    
```

a)

Data: ПСП p , классификатор $\langle K \rangle, \langle V \rangle$
Result: Класс y ПСП p

```

1  $F_{Q,V} \leftarrow \langle \rangle$ 
2  $State \leftarrow \langle \rangle$ 
3  $M_p \leftarrow \text{Len}(p)$ 
4 for  $v \in V$  do
5    $N_v \leftarrow \text{Len}(v)$ 
6    $n_v \leftarrow \text{Count}(p,v)$ 
7    $f_{p,v} = \frac{n_v}{M_p - N_v + 1}$ 
8    $F_{Q,V} = F_{Q,V} \cup f_{p,v}$ 
9  $State \leftarrow \text{Next}(k)$ 
10 while  $State[7] \neq \text{True}$  do
11   if  $f_{p,State[2]} \geq State[3]$  then
12      $State \leftarrow \text{NextRight}(State)$ 
13   else
14      $State \leftarrow \text{NextLeft}(State)$ 
15  $y_p \leftarrow State[4]$ 
16 return  $y_p$ 
    
```

б)

Рисунок 5. Алгоритмы выделения признаков (а) и классификации (б).



Исходные данные	Алгоритм преобразования	Метка класса	Выборка, кол-во файлов	Размер файла в выборке, Кб
Текст	AES(CBC)	0	2000	600
	3DES(CBC)	0	2000	600
	Camellia(CBC)	0	2000	600
	RC4(CBC)	0	2000	600
	Кузнечик(ECB)	0	2000	600
Видео-файлы	ZIP	1	2000	600
	RAR	1	2000	600
	7Z	1	2000	600
	XZ	1	2000	600
	GZ	1	2000	600
	BZ2	1	2000	600

Таблица 2. Исходные данные для проведения исследования

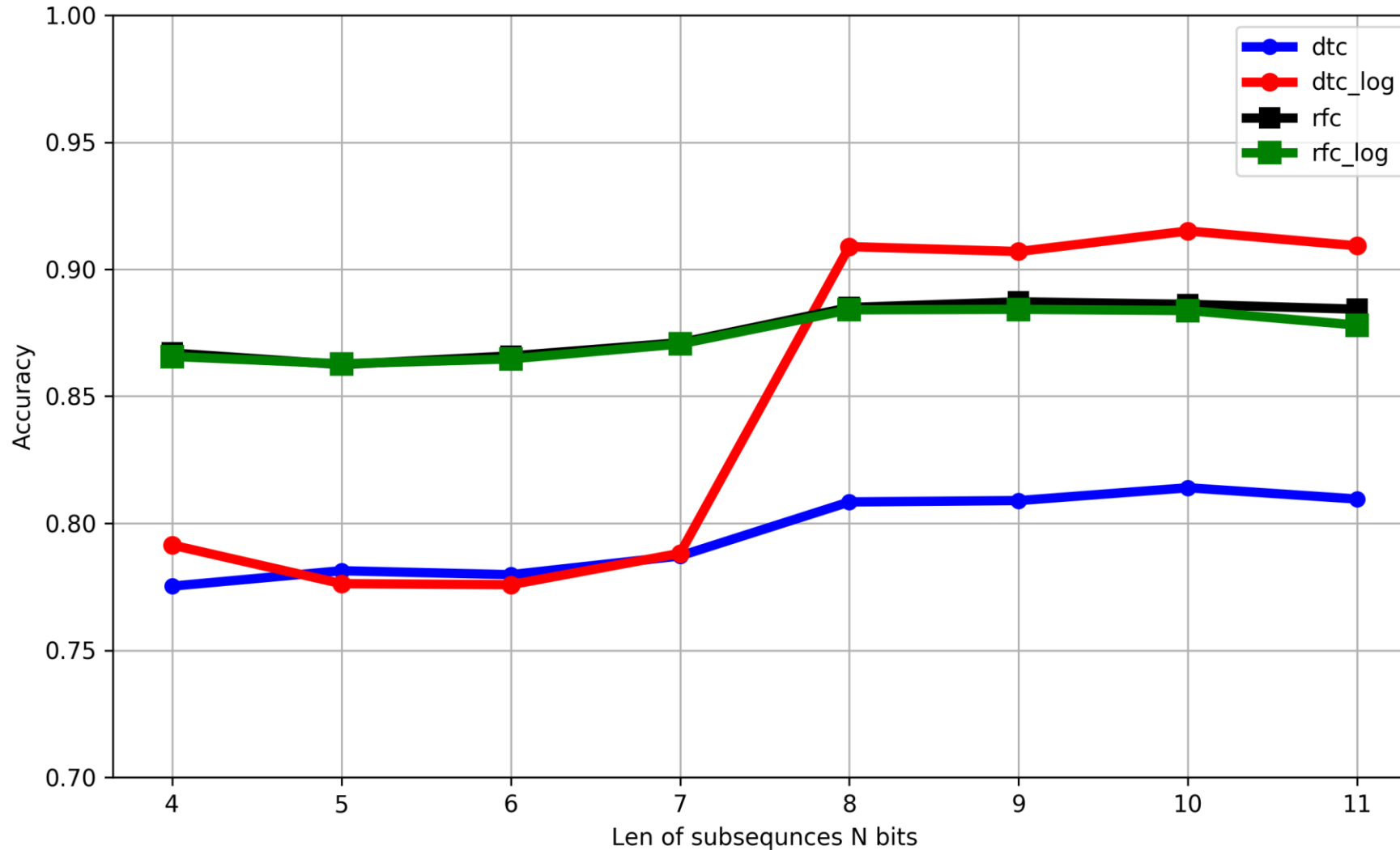


Рисунок 6. Оценка зависимости точности (Accuracy) классификации от длины подпоследовательностей (признаков)

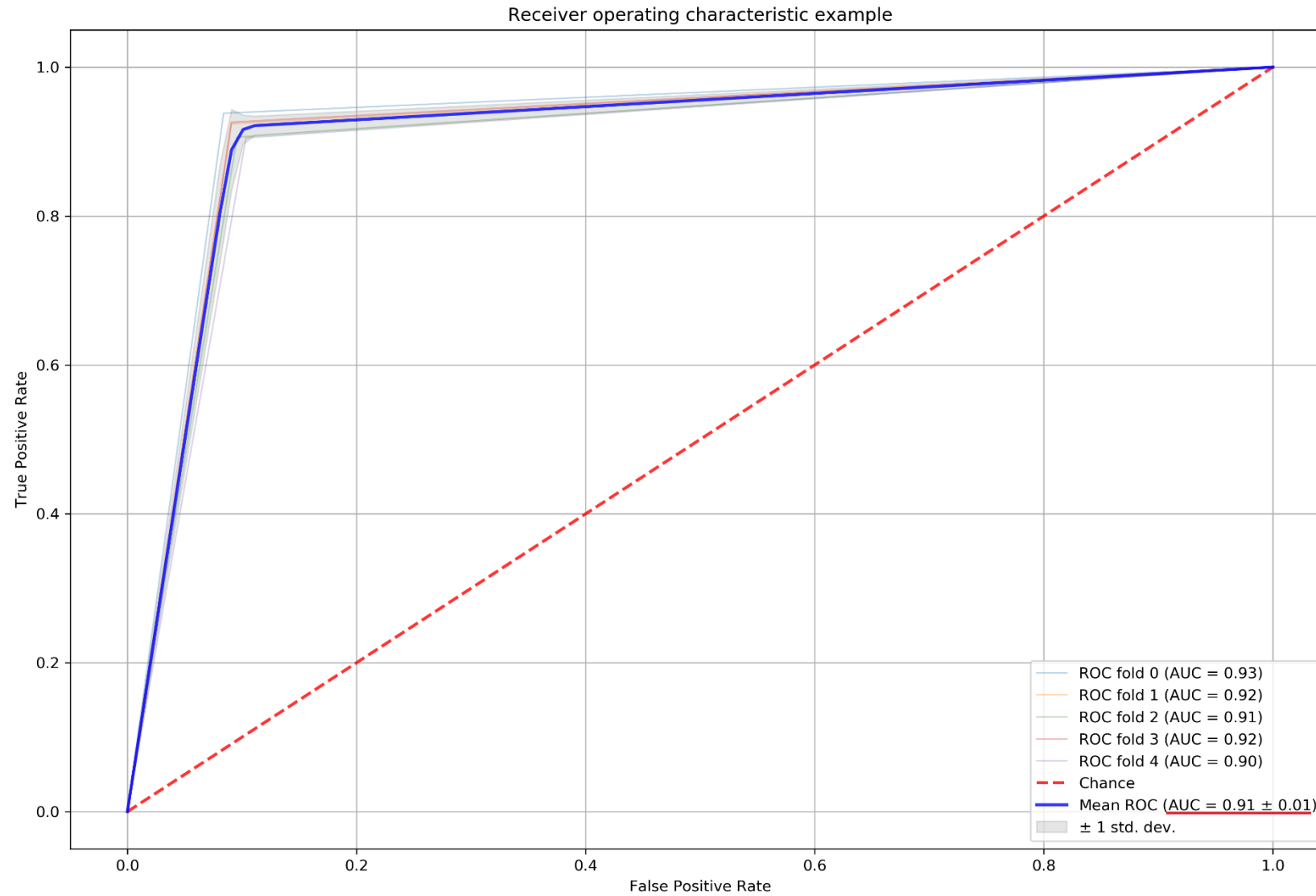


Рисунок 7. ROC кривая классификации зашифрованных/сжатых данных

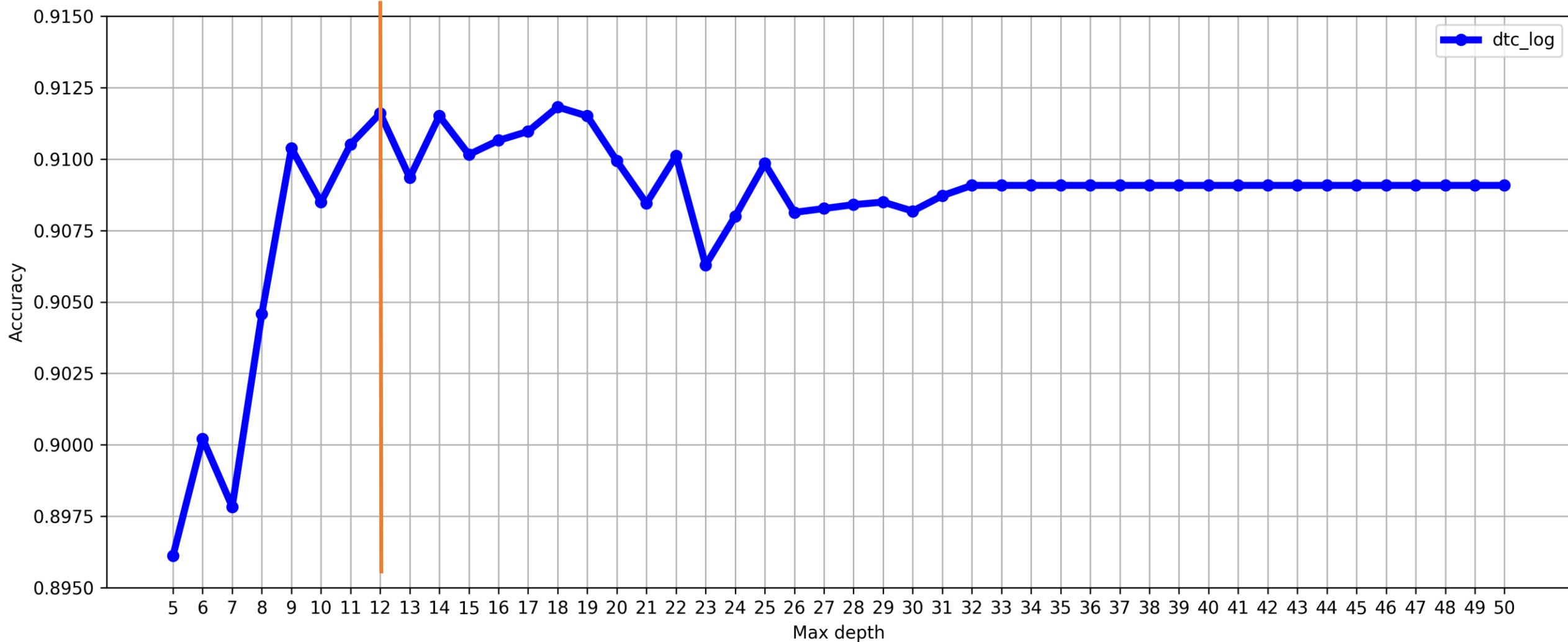


Рисунок 8. Оценка зависимости точности (Accuracy) классификации от максимальной глубины дерева решений

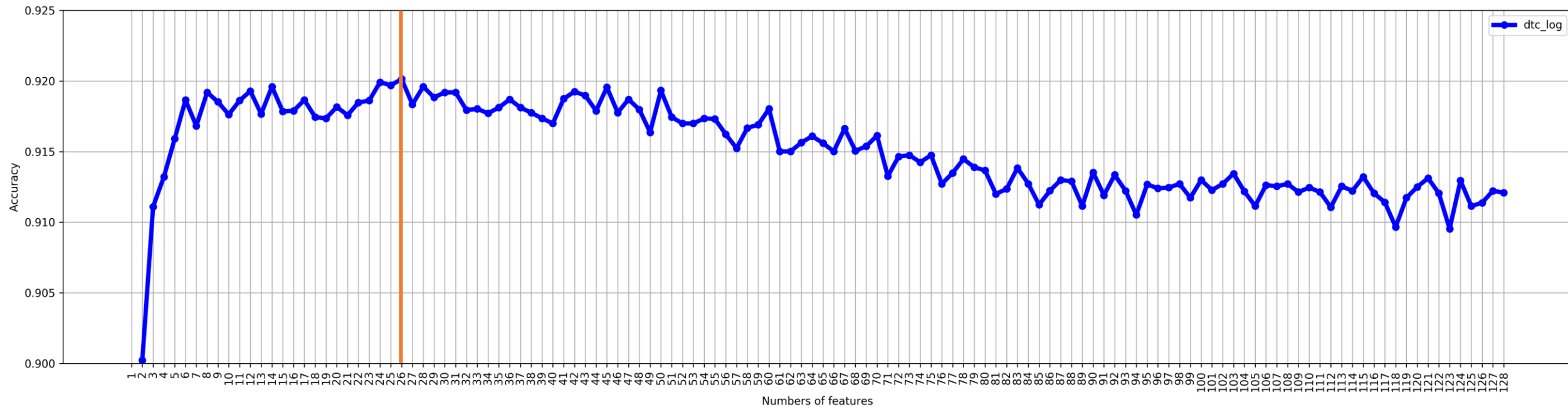
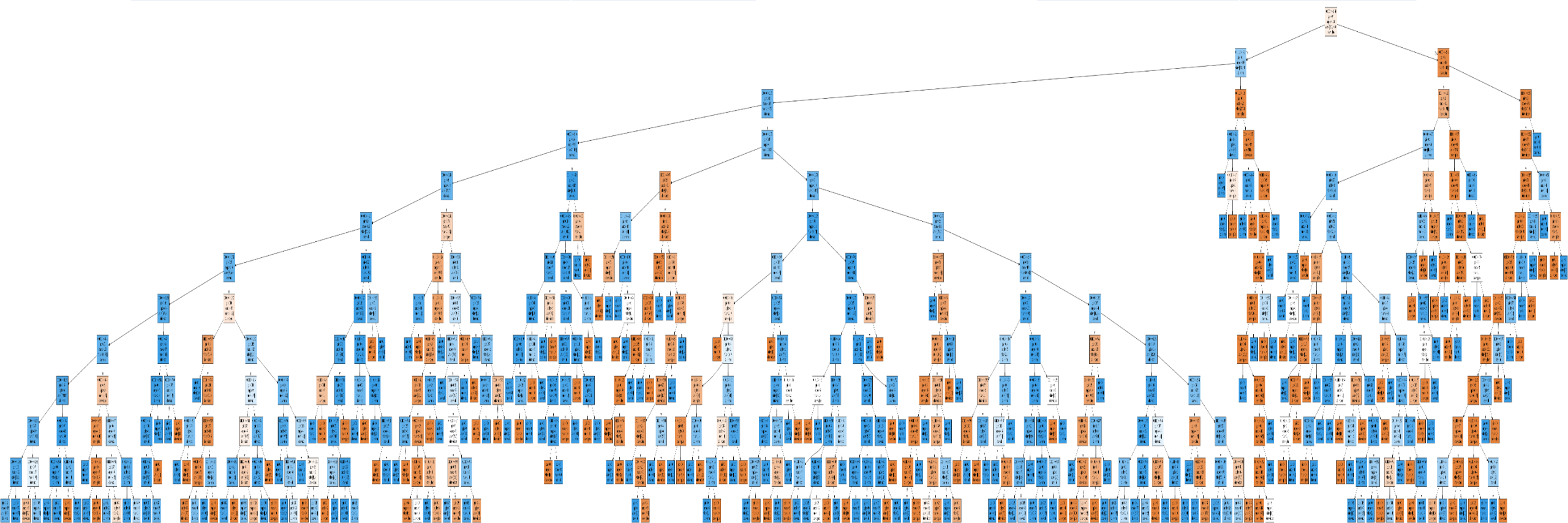


Рисунок 9. Оценка зависимости точности (Accuracy) классификации от числа признаков

Гиперпараметр классификатора	Значение
Длина подпоследовательности	8 бит
Максимальная глубина дерева	12
Количество признаков	26



Мера	Значение
Accuracy	0,919
Precision	0,911
Recall	0,926

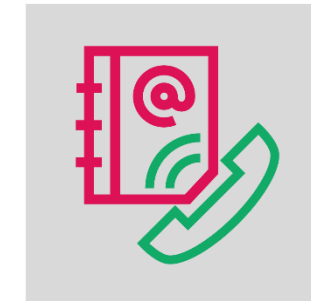


Контактная информация

Электронная почта: a.kozachok@academ.msk.rsnet.ru

Телефон: +7 920 284-57-43

Сайт: www.academ.msk.rsnet.ru



Подход к классификации последовательностей, сформированных алгоритмами сжатия и шифрования

Козачок Александр Васильевич,
д.т.н., сотрудник Академии ФСО России, г. Орёл
Спирин Андрей Андреевич,
сотрудник Академии ФСО России, г. Орёл

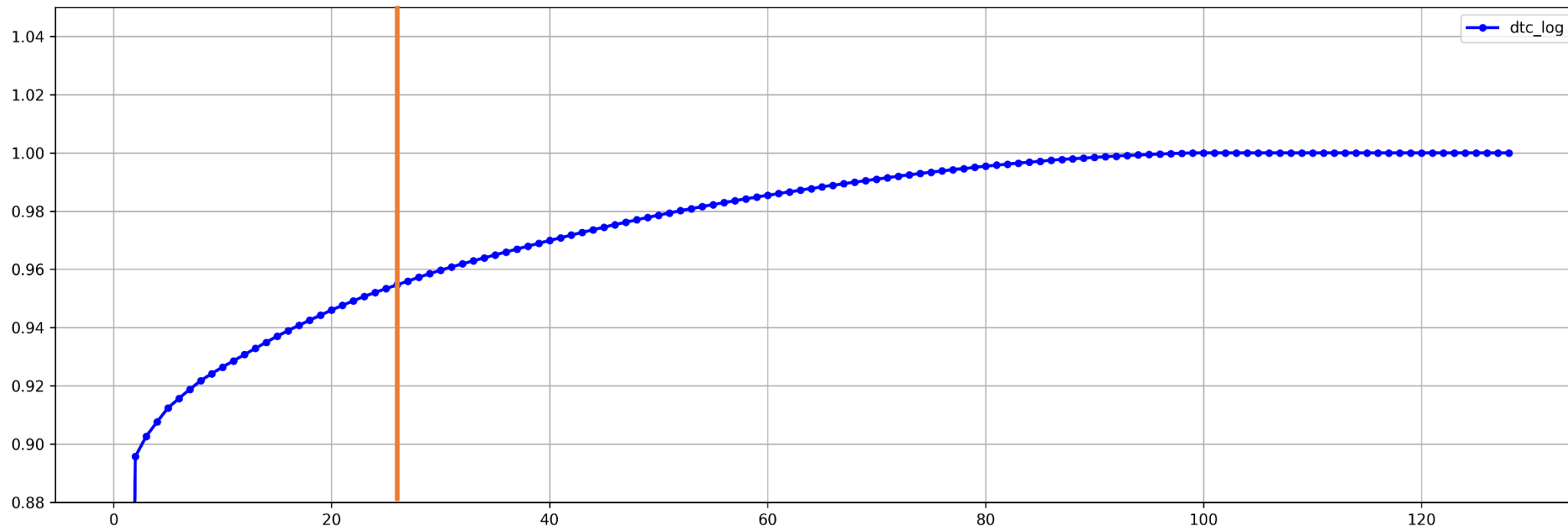


Рисунок 10. Оценка зависимости кумулятивной значимости признаков от их числа