



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

Internet Studies Lab, Department of Applied
Mathematics and Business Informatics

CLASSIFICATION IN ORANGE

1. KNN
2. SVM
3. LOGISTIC
REGRESSION

Анализ баз данных в публичном управлении
Кольцов С.Н.

Saint Petersburg, 1.10.2018



Задача классификации

Задача классификации в машинном обучении — это задача отнесения объекта к одному из заранее определенных классов на основании его формализованных признаков. Каждый из объектов в этой задаче представляется в виде вектора в N -мерном пространстве, каждое измерение в котором представляет собой описание одного из признаков объекта.

Для **обучения классификатора** необходимо иметь набор объектов, для которых заранее определены классы. Это множество называется обучающей выборкой, её разметка производится вручную, с привлечением специалистов в исследуемой области.

Например, у нас есть набор текстов, и у каждого текста есть оценка тональности. Алгоритм классификации может обучиться на этих текстах, и в дальнейшем, обученный алгоритм можно использовать для другого набора текстов. В этом случае, многомерное пространство признаков представляет собой матрица частот слов в текстах.

Другой пример, предположим есть таблица пациентов, с медицинскими показателями (виды болей, различные анализы) и диагноз, который был подтвержден. В этом случае можно обучить алгоритм распознавать диагноз у вновь поступивших пациентов.



ТИПЫ КЛАССИФИКАТОРОВ

Типичная задача статистического обучения – есть набор объектов с наблюдаемыми свойствами, и не наблюдаемыми свойствами. Нужно построить алгоритм, который бы позволял вычислить ненаблюдаемые свойства при помощи наблюдаемых, при этом хотелось бы что бы алгоритм ошибался не очень часто и не очень сильно.

Классификаторы основанные на таблице частот.

1. ZeroR (алгоритм строит таблицу частот и выбирает максимальную частоту).
2. OneR (Алгоритм строит таблицу частот и строит одно правило для каждой класса. Выбирает правило, которое дает минимальную ошибку. Это правило применяется для всего датасета)
3. **Naïve Bayesian**
4. Decision Tree (Алгоритм разбивает датасет на все меньшие куски данных, формируя тем самым дерево).

Классификаторы основанные на ковариационной матрице

1. Линейный дискриминационный анализ (Linear Discriminant Analysis)
2. **Логистическая регрессия (Logistic Regression)**

Классификатор основанный на функции сходства.

1. **Метод ближайших соседей (K Nearest Neighbors)**

Другие

1. Нейронные сети.
2. **Метод опорных векторов (Support Vector Machine)**

Меры качества классификаторов для бинарных классов

$$\text{Precision} = \frac{TP}{TP + FP}$$

Precision: число правильно предсказанных положительных значений деленных на число предсказанных классификатором положительных значений.

$$\text{Recall} = \frac{TP}{TP + FN}$$

Recall : число правильно предсказанных положительных значений деленных на число положительных ответов в данных.

	Predicted 1	Predicted 0
True 1	true positive	false negative
True 0	false positive	true negative

	Predicted 1	Predicted 0
True 1	TP	FN
True 0	FP	TN

	Predicted 1	Predicted 0
True 1	hits	misses
True 0	false alarms	correct rejections

	Predicted 1	Predicted 0
True 1	$P(pr1 tr1)$	$P(pr0 tr1)$
True 0	$P(pr1 tr0)$	$P(pr0 tr0)$



Confusion Matrix

TP – число правильно предсказанных положительных значений

FN – число неправильно предсказанных положительных значений

FN – число правильно предсказанных негативных значений

FP – число неправильно предсказанных негативных значений

Меры качества классификаторов для многих классов

Precision =

$$\frac{\sum_{i=1}^l tp_i}{\sum_{i=1}^l (tp_i + fp_i)}$$

Recall =

$$\frac{\sum_{i=1}^l tp_i}{\sum_{i=1}^l (tp_i + fn_i)}$$

F measure =

$$\frac{(\beta^2 + 1) Precision_{\mu} Recall_{\mu}}{\beta^2 Precision_{\mu} + Recall_{\mu}}$$

		OBTAINED CLASSES							Class	-2	-1	0	1	2	mean
TRUE CLASSES		-2	-1	0	1	2			Tp	5	5	5	5	5	
	-2	5	5	5	5	5			FP	=B9+B10+B11+B12					
	-1	5	5	5	5	5			Fn	20	20	20	20	20	
	0	5	5	5	5	5			Precision	0,20	0,20	0,20	0,20	0,20	0,40
	1	5	5	5	5	5			Recall	0,20	0,20	0,20	0,20	0,20	0,20
	2	5	5	5	5	5			F – measu	0,20	0,20	0,20	0,20	0,20	0,27

		OBTAINED CLASSES							Class	-2	-1	0	1	2	mean
TRUE CLASSES		-2	-1	0	1	2			Tp	5	5	5	5	5	
	-2	5	5	5	5	5			FP	20	20	=D8+D9+D11+D12			
	-1	5	5	5	5	5			Fn	20	20	20	20	20	
	0	5	5	5	5	5			Precision	0,20	0,20	0,20	0,20	0,20	0,40
	1	5	5	5	5	5			Recall	0,20	0,20	0,20	0,20	0,20	0,20
	2	5	5	5	5	5			F – measu	0,20	0,20	0,20	0,20	0,20	0,27

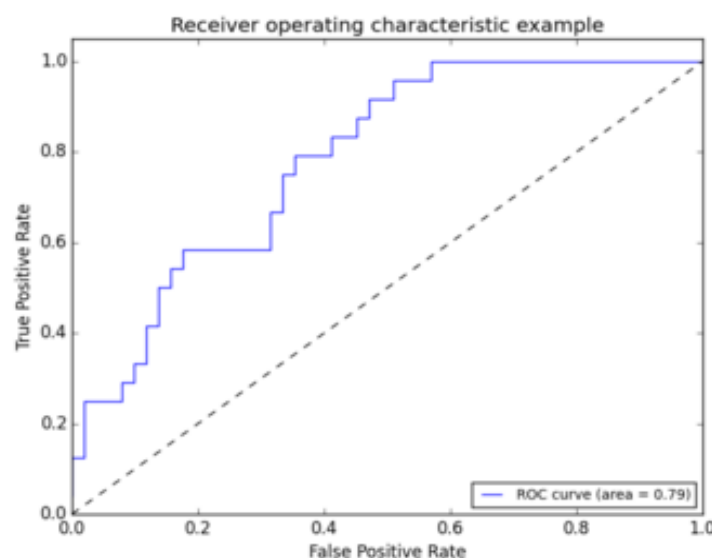
МЕРЫ КАЧЕСТВА КЛАССИФИКАТОРОВ ДЛЯ БИНАРНЫХ КЛАССОВ

F – measure

$$F - measure = \frac{(\beta^2 + 1) * precision * recall}{\beta^2 * precision + recall} \longleftrightarrow \frac{(\beta^2 + 1)tp}{(\beta^2 + 1)tp + \beta^2 fn + fp}$$

β – обычно берут равной 1. **F measure = 2 * (precision * recall) / (precision + recall)**

The F measure (F1, Fscore) можно интерпретировать как взвешенное среднее precision и recall. Если F1=1, то классификатор отработал на 100% и F1=0 тогда классификатор не справился с задачей.

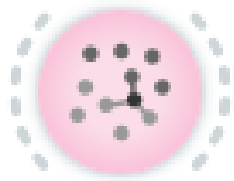


$$ROC = \frac{P(x|positive)}{P(x|negative)}$$

Рассчитывает отношение числа правильно распознанных случаев к числу не правильных. Процесс расчета таков: берутся данные, последовательно, и в них вычисляется это отношение. В какой то момент отношение становится константой.

AUC – интеграл под кривой.

ОБЗОР ВИДЖЕТОВ



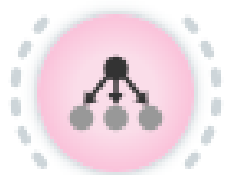
kNN

Виджет для работы с методом ближайших соседей (K Nearest Neighbors)



SVM

Виджет для работы с методом опорных векторов (SVN)



Naïve Bayes

Виджет для работы с моделью Naïve Bayes



Logistic Regression

Виджет для работы с моделью логистическая регрессия

ОБЗОР ВИДЖЕТОВ



Test & Score

Виджет для расчета метрик качества



Predictions

Виджет для расчета для расчета новых значений (предсказание)



Confusion Matrix

Виджет для расчета для расчета confusion matrix



ROC Analysis

Виджет для расчета для расчета ROC

ОБРАБОТКА ПРОПУЩЕННЫХ ЗНАЧЕНИЙ

in	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
1	1.0	0.0	3.0	1.0	22.00	1.0	0.0	45839.0000	S
2	2.0	1.0	1.0	0.0	38.00	1.0	0.0	71.2833	C
3	3.0	1.0	3.0	0.0	26.00	0.0	0.0	7.9250	S
4	4.0	1.0	1.0	0.0	35.00	1.0	0.0	53.1000	S
5	5.0	0.0	3.0	1.0	35.00	0.0	0.0	43228.0000	S
6	6.0	0.0	3.0	1.0	?	0.0	0.0	30529.0000	Q
7	7.0	0.0	1.0	1.0	54.00	0.0	0.0	51.8625	S
8	8.0	0.0	3.0	1.0	2.00	3.0	1.0		
9	9.0	1.0	3.0	0.0	27.00	0.0	2.0		
10	10.0	1.0	2.0	0.0	14.00	1.0	0.0		
11	11.0	1.0	3.0	0.0	4.00	1.0	1.0		
12	12.0	1.0	1.0	0.0	58.00	0.0	0.0		
13	13.0	0.0	3.0	1.0	20.00	0.0	0.0		
14	14.0	0.0	3.0	1.0	39.00	1.0	5.0		
15	15.0	0.0	3.0	0.0	14.00	0.0	0.0		
16	16.0	1.0	2.0	0.0	55.00	0.0	0.0		
17	17.0	0.0	3.0	1.0	2.00	4.0	1.0		
18	18.0	1.0	2.0	1.0	?	0.0	0.0		
19	19.0	0.0	3.0	0.0	31.00	1.0	0.0		
20	20.0	1.0	3.0	0.0	?	0.0	0.0		
21	21.0	0.0	2.0	1.0	35.00	0.0	0.0		
22	22.0	1.0	2.0	1.0	34.00	0.0	0.0		

Impute

Default Method

☐ Don't impute
☒ Average/Most frequent
☐ As a distinct value
☐ Model-based imputer (simple tree)
☐ Random values
☐ Remove instances with unknown values

Individual Attribute Settings

N

 PassengerId

C

 Survived

N

 Pclass

C

 Sex

N

 Age -> random

N

 SibSp

N

 Parch

N

 Fare

C

 Embarked

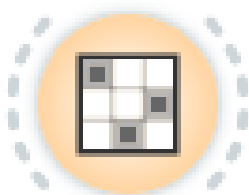
☐ Default (above)
☐ Don't impute
☐ Average/Most frequent
☐ As a distinct value
☐ Model-based imputer (simple tree)
☒ Random values
☐ Remove instances with unknown values
☐ Value

0,000

Restore All to Default

☒ Apply automatically

Apply



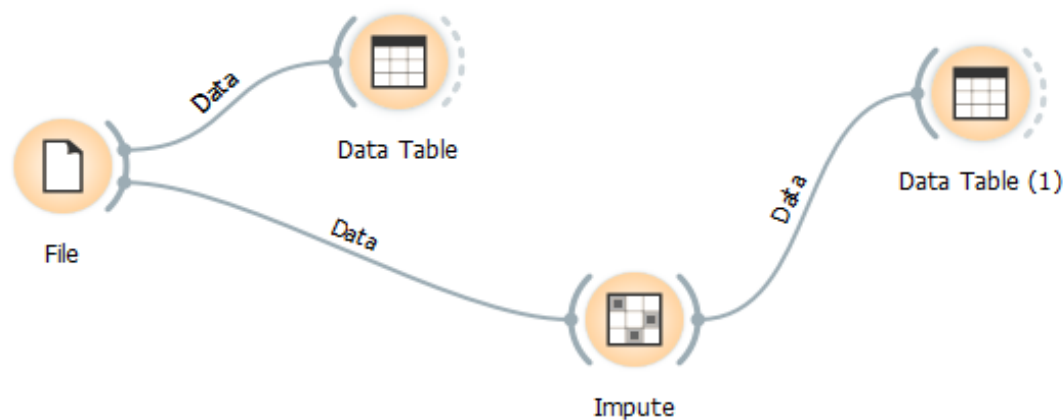
Impute

Виджет для заполнения пропущенных значений

ОБРАБОТКА ПРОПУЩЕННЫХ ЗНАЧЕНИЙ

in	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
1	1.0	0.0	3.0	1.0	22.00	1.0	0.0	45839.0000	S
2	2.0	1.0	1.0	0.0	38.00	1.0	0.0	71.2833	C
3	3.0	1.0	3.0	0.0	26.00	0.0	0.0	7.9250	S
4	4.0	1.0	1.0	0.0	35.00	1.0	0.0	53.1000	S
5	5.0	0.0	3.0	1.0	35.00	0.0	0.0	43228.0000	S
6	6.0	0.0	3.0	1.0	?	0.0	0.0	30529.0000	Q
7	7.0	0.0	1.0	1.0	54.00	0.0	0.0	51.8625	S
8	8.0	0.0	3.0	1.0	2.00	3.0	1.0	21.0750	S
9	9.0	1.0	3.0	0.0	27.00	0.0	2.0	11.1333	S
10	10.0	1.0	2.0	0.0	14.00	1.0	0.0	30.0708	C
11	11.0	1.0	3.0	0.0	4.00	1.0	1.0	43297.0000	S
12	12.0	1.0	1.0	0.0	58.00	0.0	0.0	26.5500	S
13	13.0	0.0	3.0	1.0	20.00	0.0	0.0	43228.0000	S
14	14.0	0.0	3.0	1.0	39.00	1.0	5.0	31.2750	S
15	15.0	0.0	3.0	0.0	14.00	0.0	0.0	15523.0000	S
16	16.0	1.0	2.0	0.0	55.00	0.0	0.0	16.0000	S
17	17.0	0.0	3.0	1.0	2.00	4.0	1.0	29.1250	Q
18	18.0	1.0	2.0	1.0	?	0.0	0.0	13.0000	S
19	19.0	0.0	3.0	0.0	31.00	1.0	0.0	18.0000	
20	20.0	1.0	3.0	0.0	?	0.0	0.0	7.2250	
21	21.0	0.0	2.0	1.0	35.00	0.0	0.0	26.0000	
22	22.0	1.0	2.0	1.0	34.00	0.0	0.0	13.0000	

Sex	Age	SibSp
1.0	22.00	1.0
0.0	38.00	1.0
0.0	26.00	0.0
0.0	35.00	0.0
1.0	35.00	0.0
1.0	36.00	0.0
1.0	54.00	0.0
1.0	2.00	3.0
0.0	27.00	0.0
0.0	14.00	1.0
0.0	4.00	1.0
0.0	58.00	0.0
1.0	20.00	0.0
1.0	39.00	1.0
0.0	14.00	0.0
0.0	55.00	0.0



МЕТОД БЛИЖАЙШИХ СОСЕДЕЙ (K NEAREST NEIGHBORS)

Все объекты расположены в многомерном пространстве

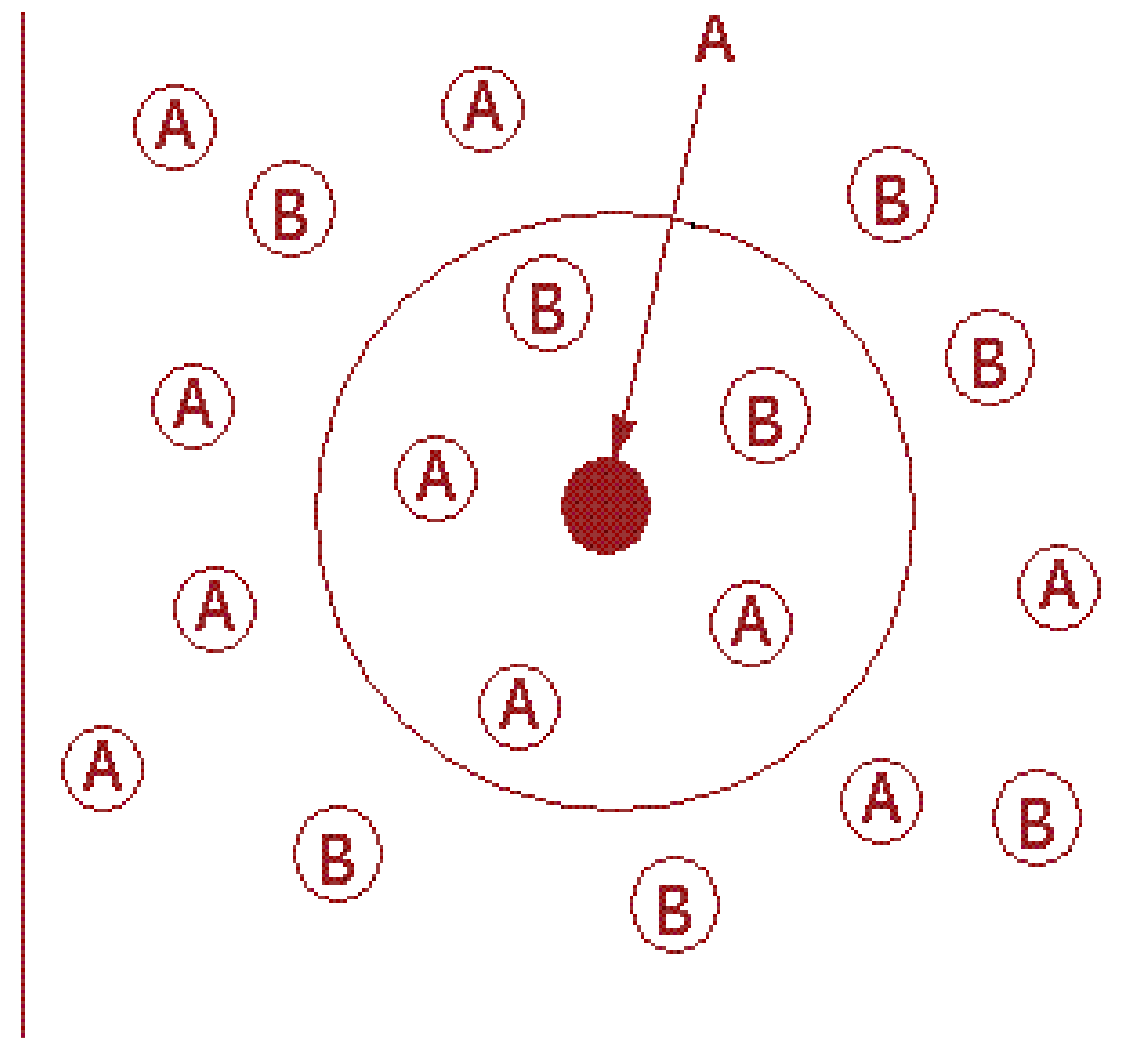
1. Задаем число k – количество ближайших соседей.

2. Ищем k объектов с минимальным расстоянием до нашего нового объекта. Используем меру для расчета расстояний.

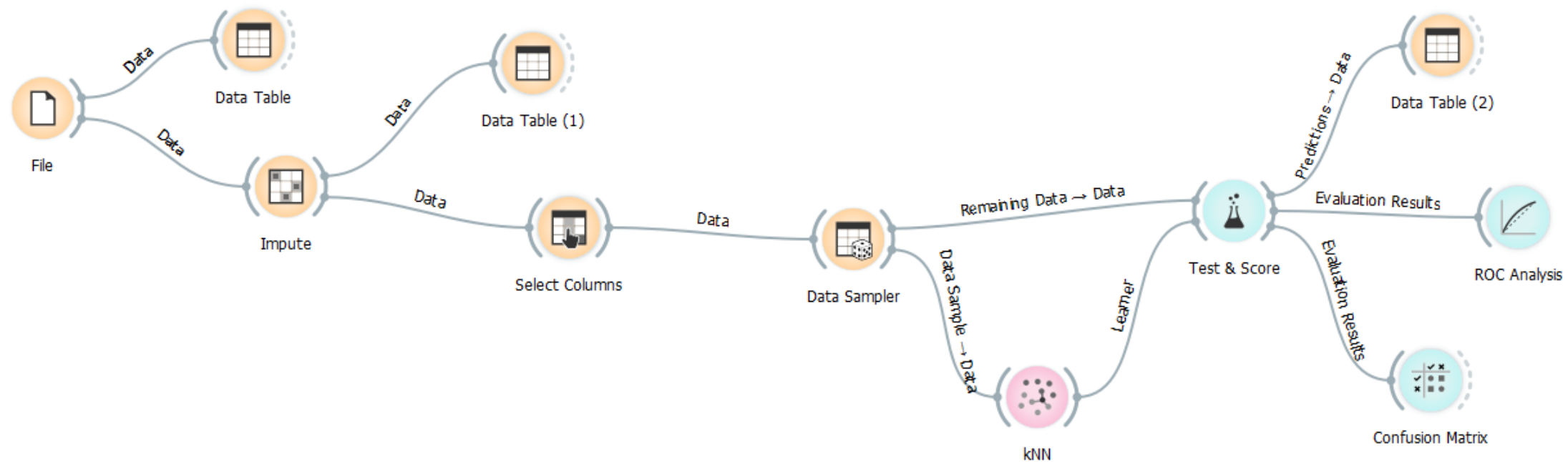
3.1 Простое невзвешенное голосование. Считаем сколько объектов с классами присутствует внутри заданного расстояния. Например, если число объектов с классом А большинство, то новый объект относится к классу А.

3.2. Взвешенное голосование

В такой ситуации учитывается также и расстояние до новой записи. Чем меньше расстояние, тем более значимый вклад вносит голос.



КЛАССИФИКАЦИЯ ТАБЛИЧНЫХ ДАННЫХ

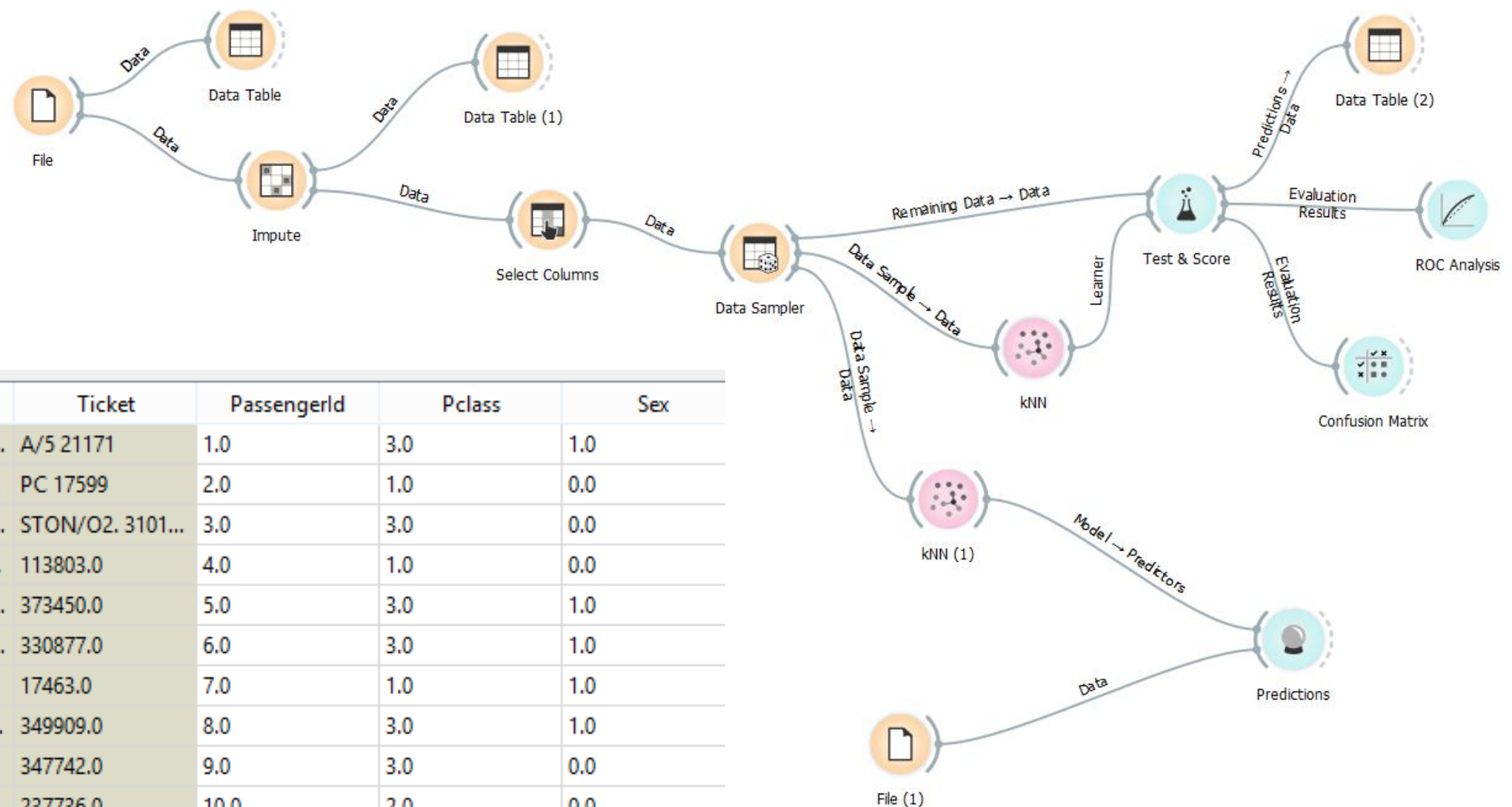


	kNN	kNN (0.0)	kNN (1.0)	Fold	PassengerId	Age	Pclass	Sex	SibSp
870		0.600	0.400	10	362.0	29.00	2.0	1.0	1.0
871		0.800	0.200	10	839.0	32.00	3.0	1.0	0.0
873		0.800	0.200	10	526.0	40.50	3.0	1.0	0.0
874		0.800	0.200	10	574.0	2.00	3.0	0.0	0.0
875		0.600	0.400	10	493.0	55.00	1.0	1.0	0.0
877		0.800	0.200	10	268.0	25.00	3.0	1.0	1.0
878		0.800	0.200	10	198.0	42.00	3.0	1.0	0.0
879		1.000	0.000	10	822.0	27.00	3.0	1.0	0.0
880		1.000	0.000	10	818.0	31.00	2.0	1.0	1.0
881		0.800	0.200	10	24.0	28.00	1.0	1.0	0.0
882		0.600	0.400	10	708.0	42.00	1.0	1.0	0.0
883		0.600	0.400	10	234.0	5.00	3.0	0.0	4.0
886		0.800	0.200	10	859.0	24.00	3.0	0.0	0.0

		Predicted		
		0.0	1.0	Σ
Actual	0.0	422	128	550
	1.0	235	125	360
Σ		657	253	910

Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
kNN	0.591	0.601	0.584	0.584	0.601

ПРЕДСКАЗАНИЯ С ИСПОЛЬЗОВАНИЕМ КЛАССИФИКАТОРА KNN



	kNN	Name	Ticket	PassengerId	Pclass	Sex
1	<u>1.00 : 0.00 → 0.0</u>	Braund, Mr. Ow...	A/5 21171	1.0	3.0	1.0
2	<u>0.60 : 0.40 → 0.0</u>	Cumings, Mrs. ...	PC 17599	2.0	1.0	0.0
3	<u>0.00 : 1.00 → 1.0</u>	Heikkinen, Miss...	STON/O2. 3101...	3.0	3.0	0.0
4	<u>0.40 : 0.60 → 1.0</u>	Futrelle, Mrs. Ja...	113803.0	4.0	1.0	0.0
5	<u>1.00 : 0.00 → 0.0</u>	Allen, Mr. Willia...	373450.0	5.0	3.0	1.0
6	<u>0.60 : 0.40 → 0.0</u>	Moran, Mr. Jam...	330877.0	6.0	3.0	1.0
7	<u>0.60 : 0.40 → 0.0</u>	McCarthy, Mr. ...	17463.0	7.0	1.0	1.0
8	<u>0.60 : 0.40 → 0.0</u>	Palsson, Master...	349909.0	8.0	3.0	1.0
9	<u>0.20 : 0.80 → 1.0</u>	Johnson, Mrs. ...	347742.0	9.0	3.0	0.0
10	<u>0.40 : 0.60 → 1.0</u>	Nasser, Mrs. Ni...	237736.0	10.0	2.0	0.0
11	<u>1.00 : 0.00 → 0.0</u>	Sandstrom, Mis...	PP 9549	11.0	3.0	0.0
12	<u>0.60 : 0.40 → 0.0</u>	Bonnell, Miss. E...	113783.0	12.0	1.0	0.0



ПРЕДСКАЗАНИЯ АКЦИЙ ГАЗПРОМА

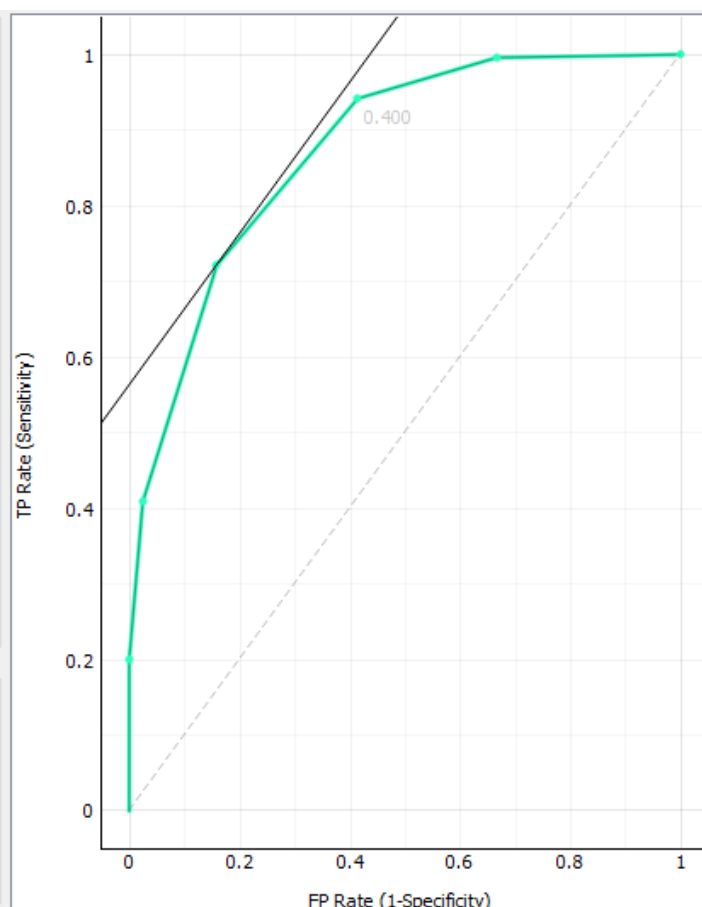
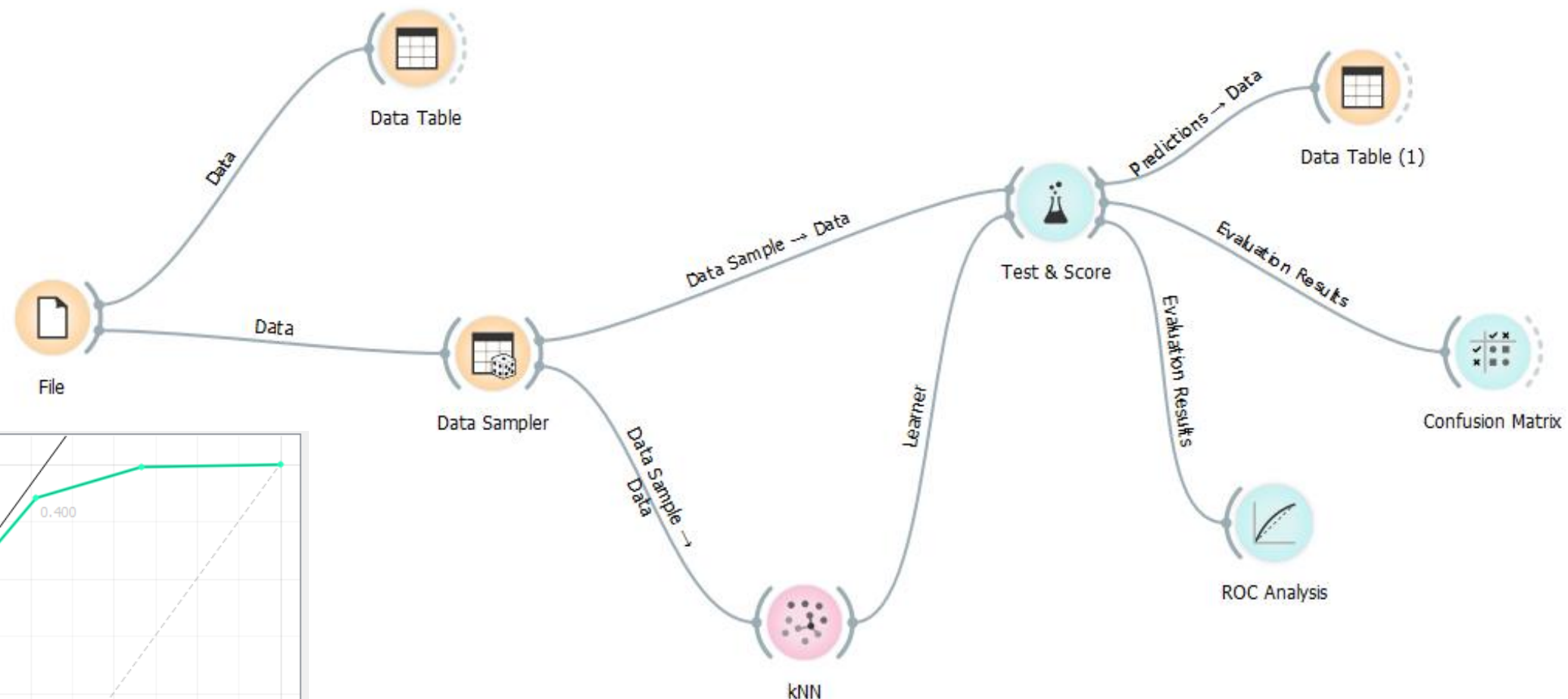
<DATE>	<TIME>	<OPEN>	<HIGH>	<LOW>	<CLOSE>	difference
20160104	110000	135.89	136.65	135.62	135.96	0.07
20160104	120000	135.92	135.92	134.42	134.67	-1.25
20160104	130000	134.71	135.24	134.42	135.15	0.44
20160104	140000	135.12	135.22	134.7	135.22	0.1
20160104	150000	135.23	135.49	134.88	135.11	-0.12
20160104	160000	135.11	135.35	134.94	134.97	-0.14
20160104	170000	134.97	135.3	134.8	135.13	0.16
20160104	180000	135.13	135.4	134.96	135.15	0.02
20160104	190000	135.14	135.45	134.9	134.91	-0.23
20160105	110000	134.85	135.35	134.85	135.06	0.21
20160105	120000	135.06	136.16	135	135.93	0.87
20160105	130000	135.92	135.94	135.1	135.31	-0.61
20160105	140000	135.3	135.3	134.61	135.07	-0.23
20160105	150000	135.08	135.23	134.76	135.22	0.14
20160105	160000	135.21	135.75	135.01	135.64	0.43
20160105	170000	135.64	137.28	135.53	136.74	1.1
20160105	180000	136.71	137.15	136.18	136.19	-0.52
20160105	190000	136.23	136.58	136.08	136.45	0.22
20160106	110000	136.35	136.93	136.17	136.51	0.16
20160106	120000	136.51	136.72	136	136.46	-0.05

Задача заключается в том ,что бы построить на основе KNN алгоритм предсказания котировок ценных бумаг.

Решение: будем использовать пакет Orange.

1. Возьмем исходный датасет.
2. Обучим классификатора на наших данных.
3. Посмотрим как классификатор предсказывает на наших данных котировки акций Газпрома.

ПРЕДСКАЗАНИЯ АКЦИЙ ГАЗПРОМА



		Predicted			Σ
		down	flat	up	
Actual	down	363	0	137	500
	flat	4	0	6	10
	up	95	0	505	600
	Σ	462	0	648	1110

Evaluation Results

Method	AUC	CA	F1	Precision	Recall
kNN	0.864	0.782	0.777	0.775	0.782

ВЛИЯНИЕ ЧИСЛА СОСЕДЕЙ НА КЛАССИФИКАЦИЮ

kNN

?

×

Name

kNN

Neighbors

Number of neighbors:

50

Metric:

Eudidean

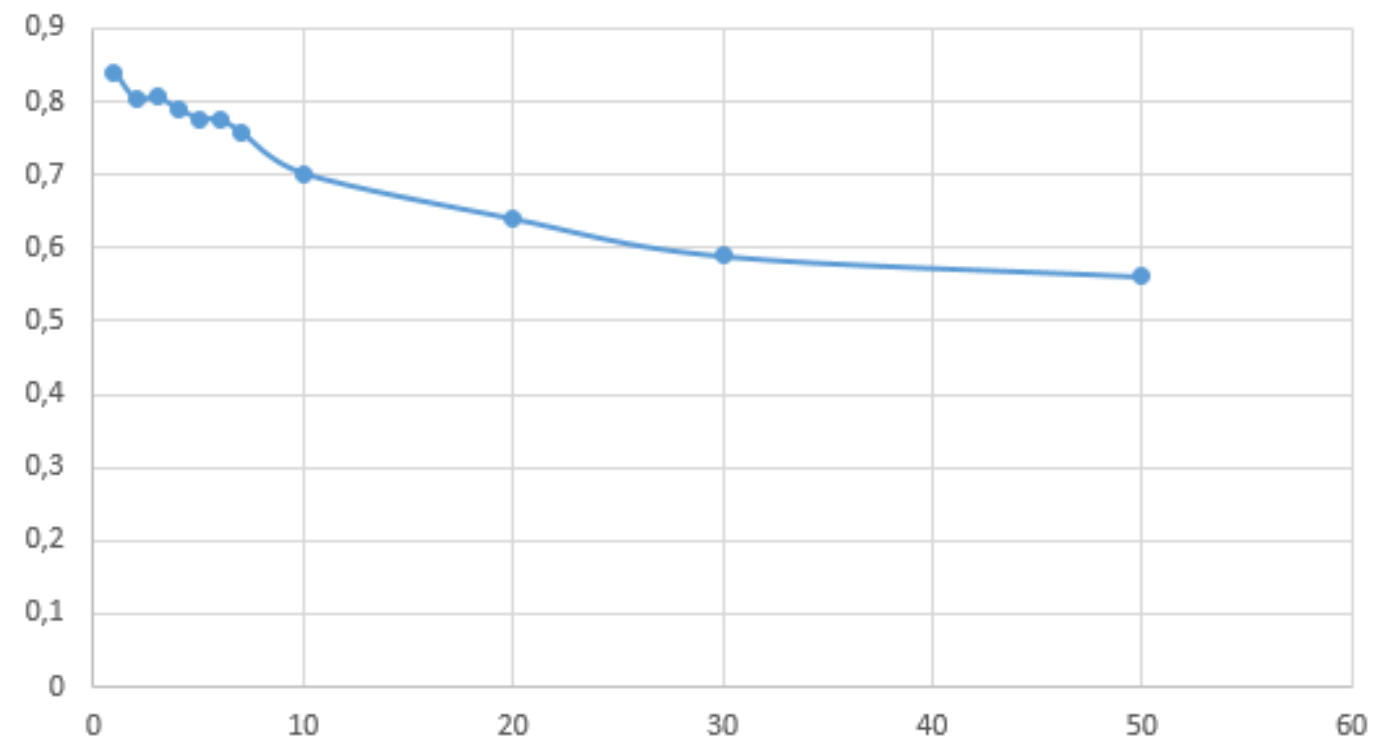
Weight:

Uniform

☒

Apply Automatically

?



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
kNN	0.625	0.561	0.559	0.556	0.561

ВЛИЯНИЕ ТИПА РАССТОЯНИЯ НА КЛАССИФИКАЦИЮ

kNN ? x

Name
kNN

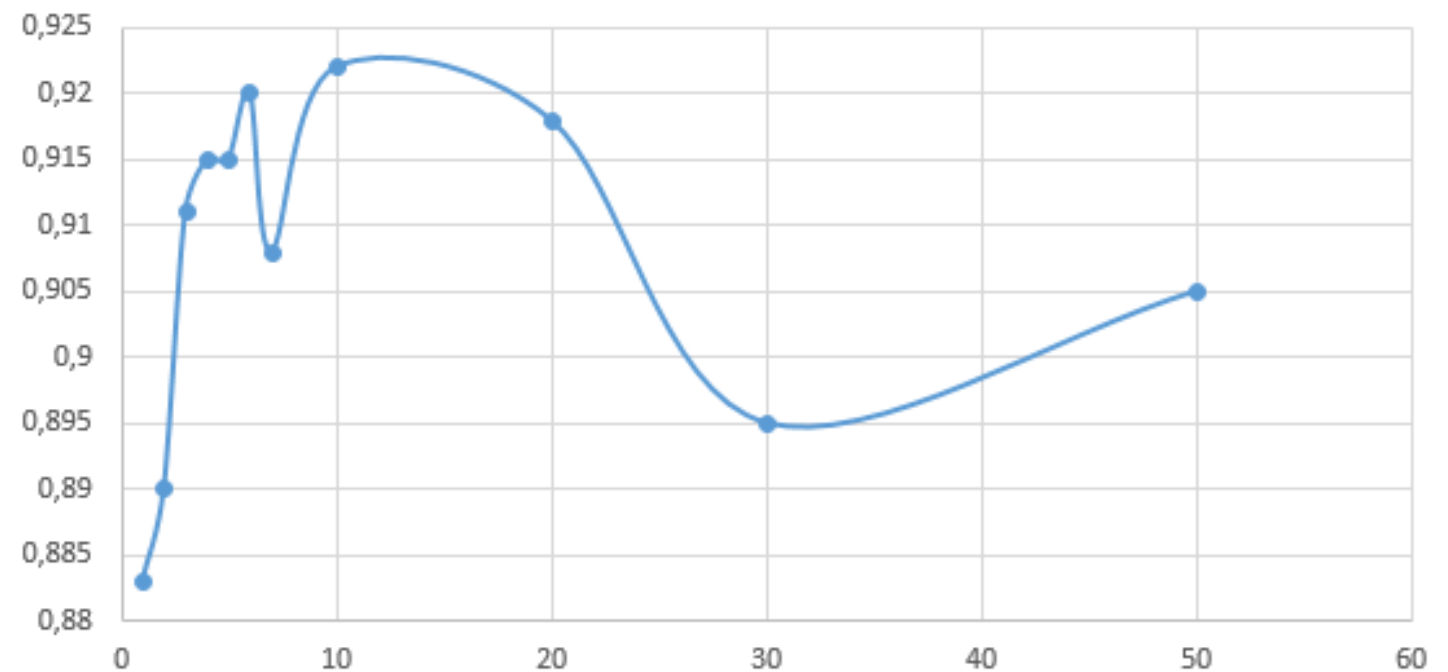
Neighbors
Number of neighbors: 50

Metric: Mahalanobis

Weight: Uniform

☒ Apply Automatically

? | [icon]



kNN ? x

Name
kNN

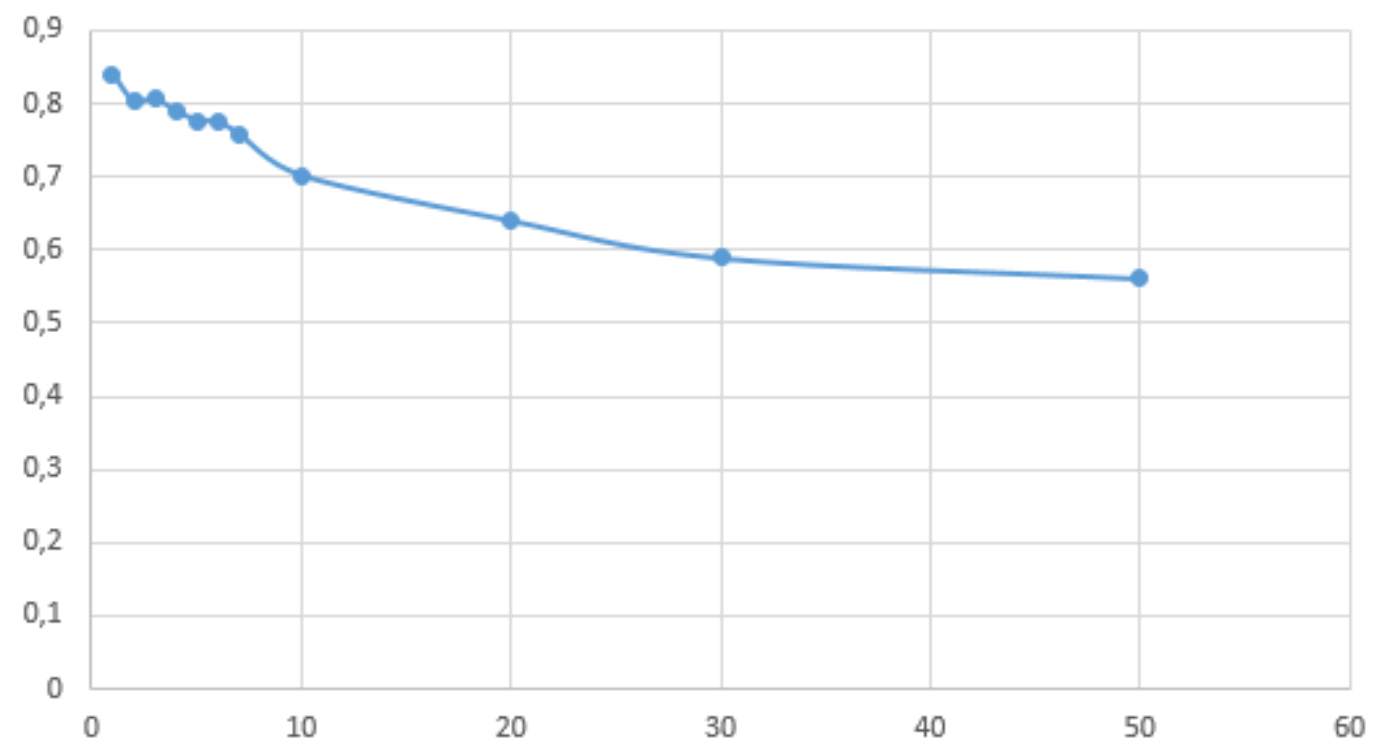
Neighbors
Number of neighbors: 50

Metric: Euclidean

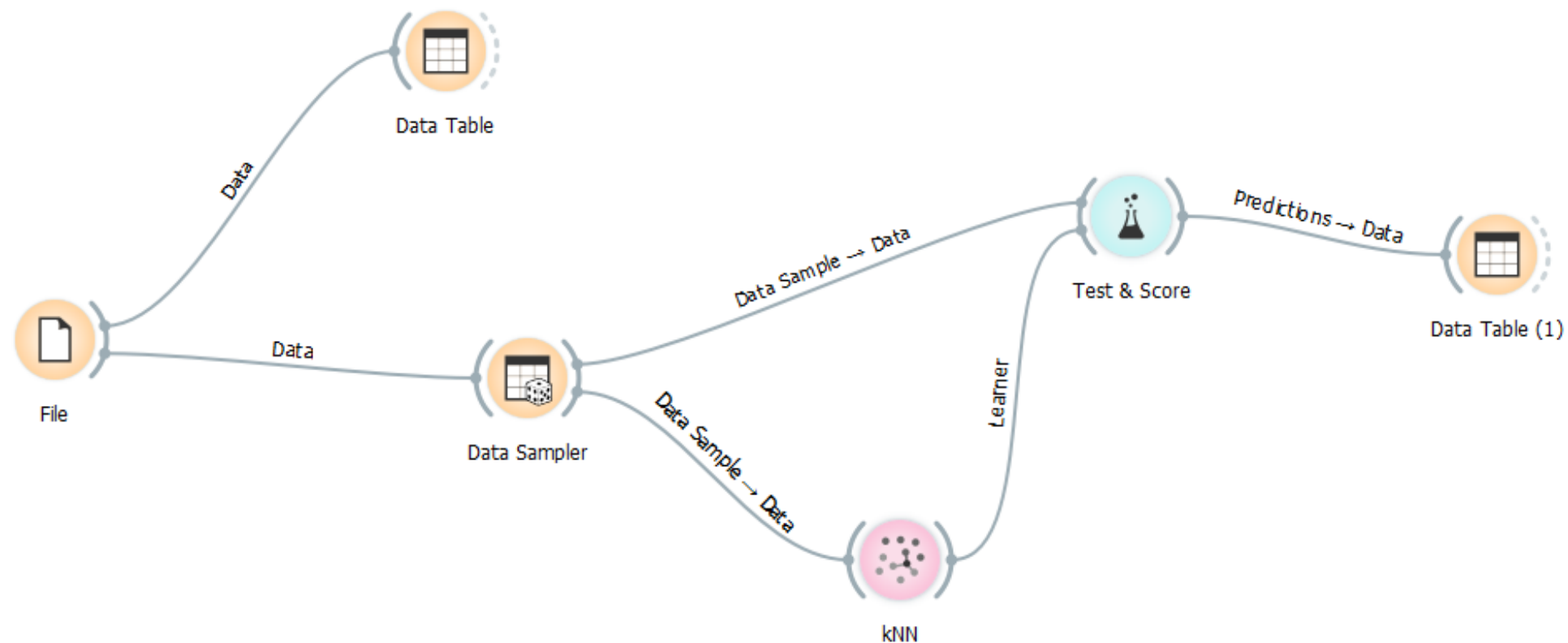
Weight: Uniform

☒ Apply Automatically

? | [icon]



НЕФТЬ - РУБЛЬ



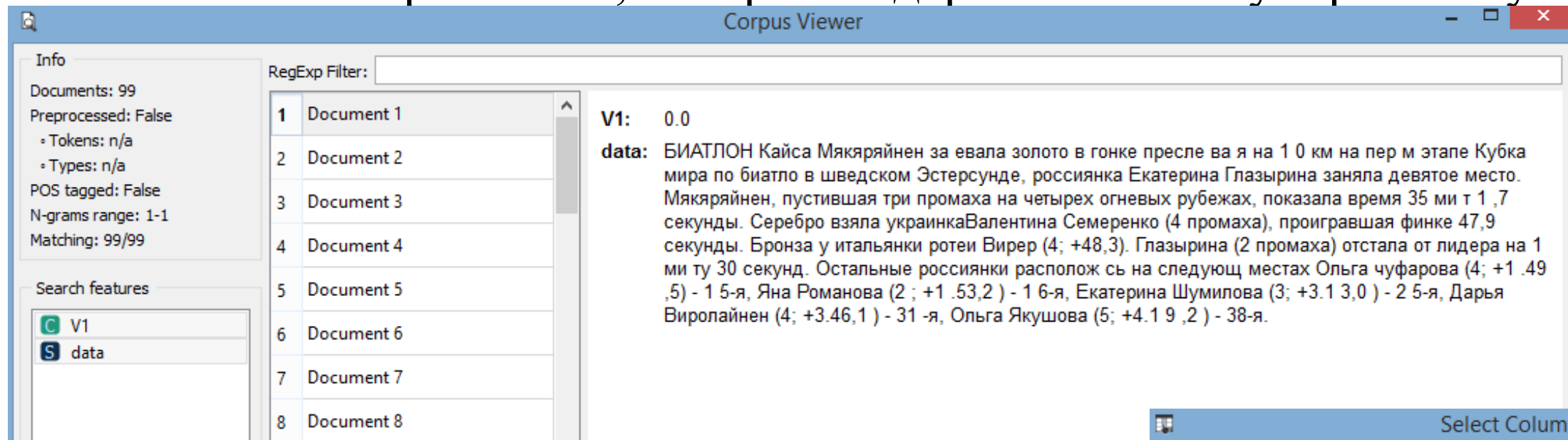
	RubUsd	Date	Brent	UsdRub
1	0.030202	41641.0	107.65	33.1100
2	0.030120	41642.0	106.65	33.2000
3	0.030130	41645.0	106.67	33.1900
4	0.030175	41646.0	107.01	33.1400
5	0.030166	41647.0	106.83	33.1500
6	0.030111	41648.0	106.00	33.2100
7	0.030312	41649.0	106.83	32.9900
8	0.030312	41651.0	106.92	32.9900
9	0.030048	41652.0	105.59	33.2800
10	0.030030	41653.0	105.25	33.3000
11	0.029940	41654.0	105.97	33.4000
12	0.029958	41655.0	105.48	33.3800
13	0.029771	41656.0	106.32	33.5900
14	0.029621	41659.0	106.20	33.7600

Evaluation Results

Method	MSE	RMSE	MAE	R2
kNN	0.000	0.001	0.001	0.964

КЛАССИФИКАЦИЯ РУССКОЯЗЫЧНЫХ ТЕКСТОВ

Возьмем набор текстов, которые содержат тональную разметку.



Corpus Viewer

Info

Documents: 99
Preprocessed: False
• Tokens: n/a
• Types: n/a
POS tagged: False
N-grams range: 1-1
Matching: 99/99

Search features

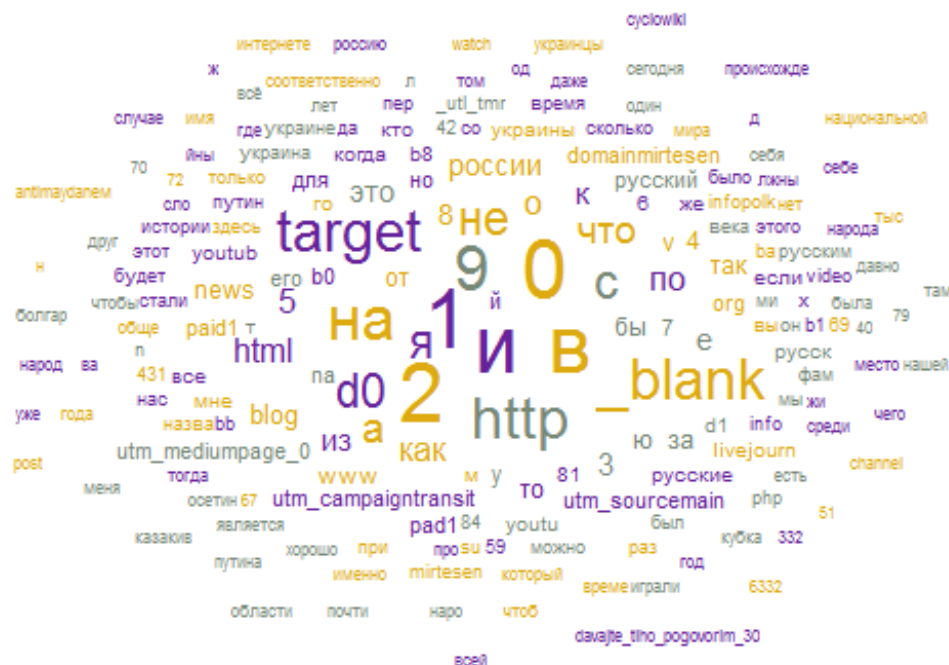
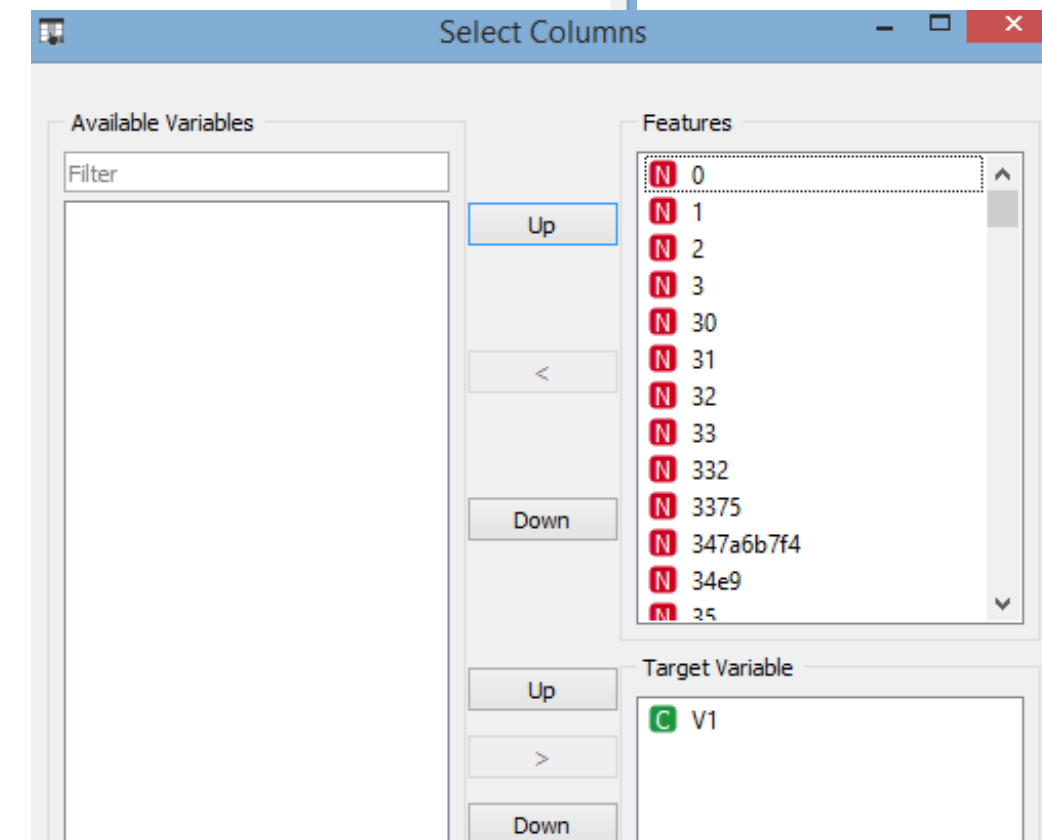
V1
data

RegExp Filter:

Document	Text
1	Document 1
2	Document 2
3	Document 3
4	Document 4
5	Document 5
6	Document 6
7	Document 7
8	Document 8

V1: 0.0

data: БИАТЛОН Кайса Мякярйнен за евала золото в гонке пресле ва я на 1 0 км на пер м этапе Кубка мира по биатло в шведском Эстерсунде, россиянка Екатерина Глазырина заняла девятое место. Мякярйнен, пустившая три промаха на четырех огневых рубежах, показала время 35 ми т 1,7 секунды. Серебро взяла украинкаВалентина Семеренко (4 промаха), проигравшая финке 47,9 секунды. Бронза у итальянки ротеи Вирер (4; +48,3). Глазырина (2 промаха) отстала от лидера на 1 ми ту 30 секунд. Остальные россиянки располож сь на следующ местях Ольга чуфарова (4; +1 .49 ,5) - 1 5-я, Яна Романова (2 ; +1 .53,2) - 1 6-я, Екатерина Шумилова (3; +3.1 3,0) - 2 5-я, Дарья Виролайнен (4; +3.46,1) - 31 -я, Ольга Якушова (5; +4.1 9 ,2) - 38-я.

Select Columns

Available Variables

Filter

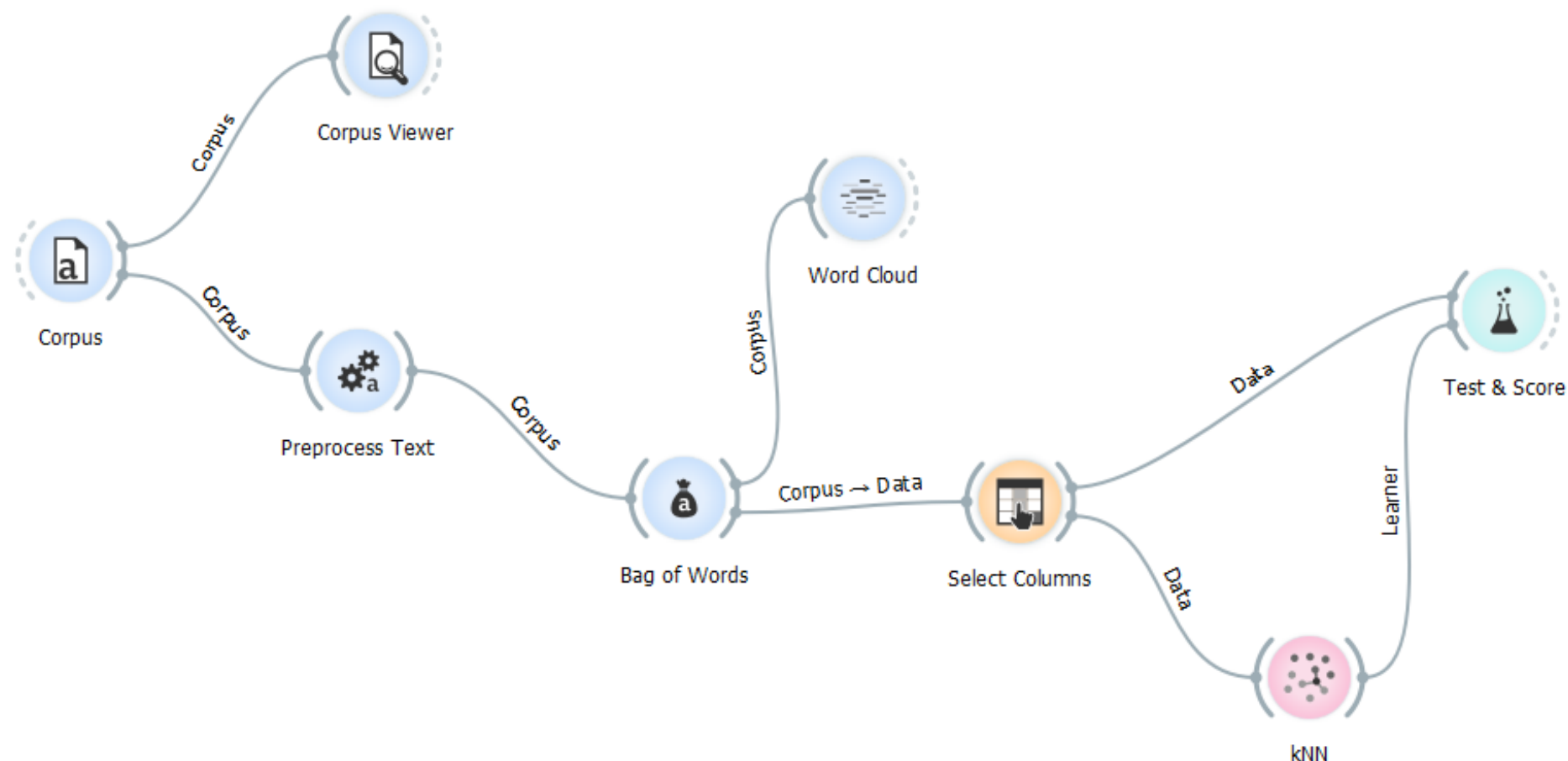
Features

Feature
N 0
N 1
N 2
N 3
N 30
N 31
N 32
N 33
N 332
N 3375
N 347a6b7f4
N 34e9
N 35

Target Variable

V1

ЭФФЕКТ УДАЛЕНИЯ СТОП СЛОВ НА КЛАССИФИКАЦИЮ



Результат классификации без очистки

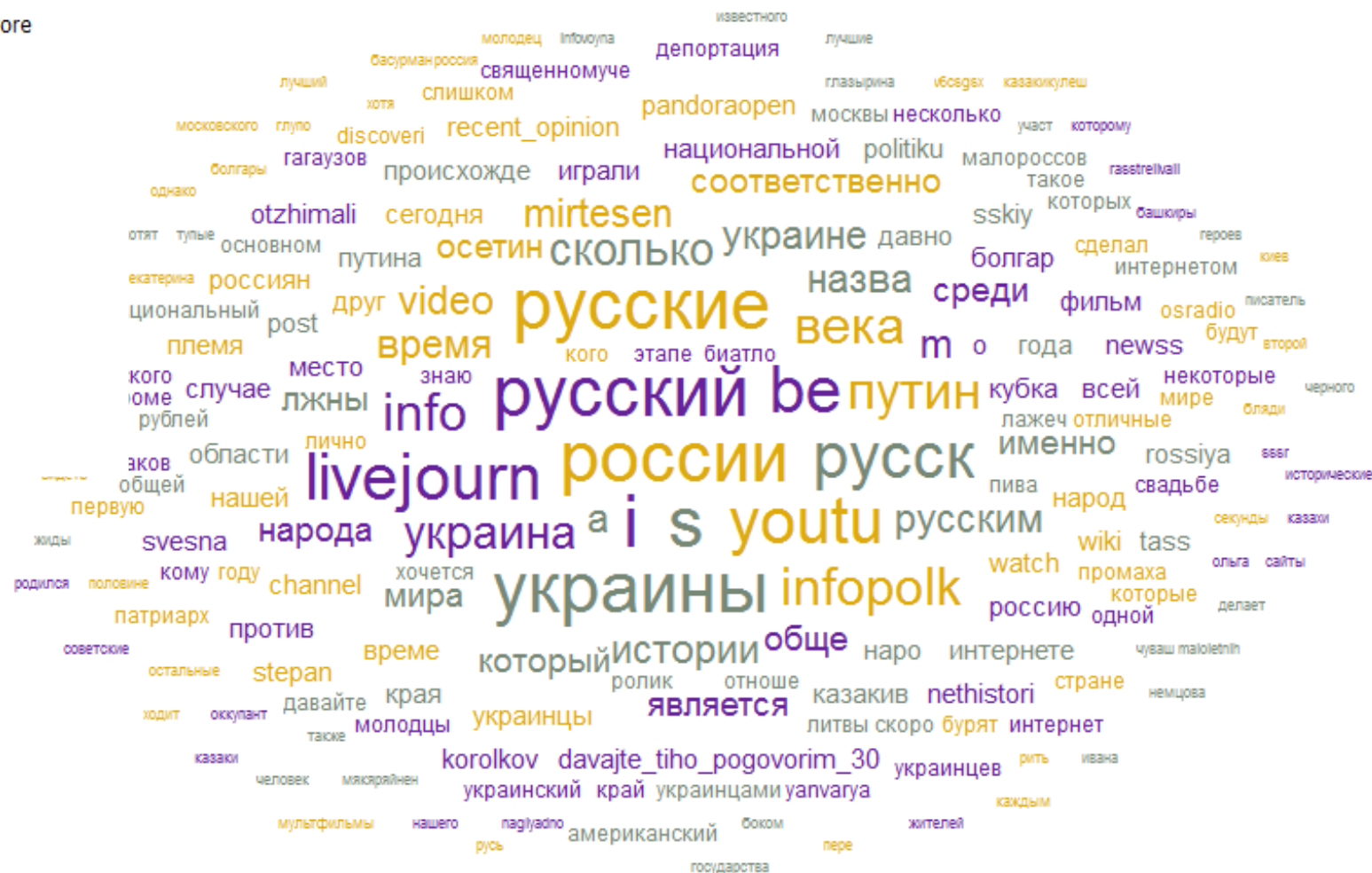
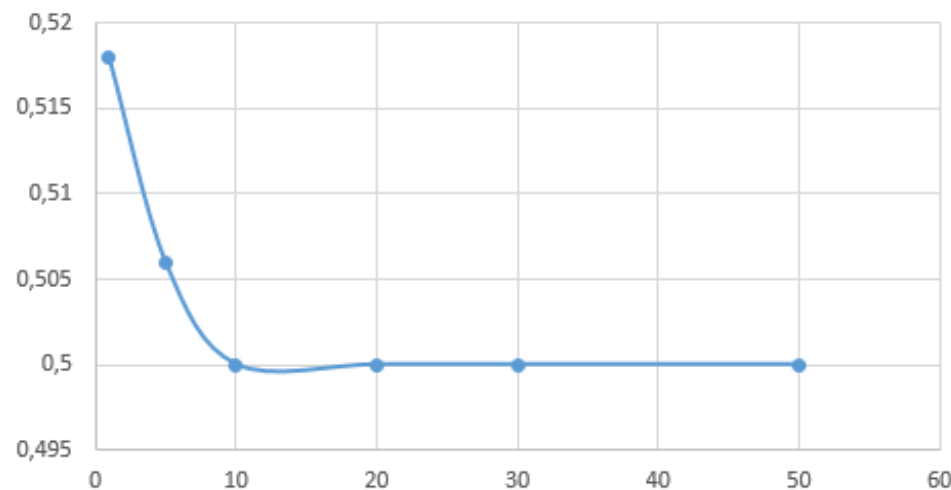
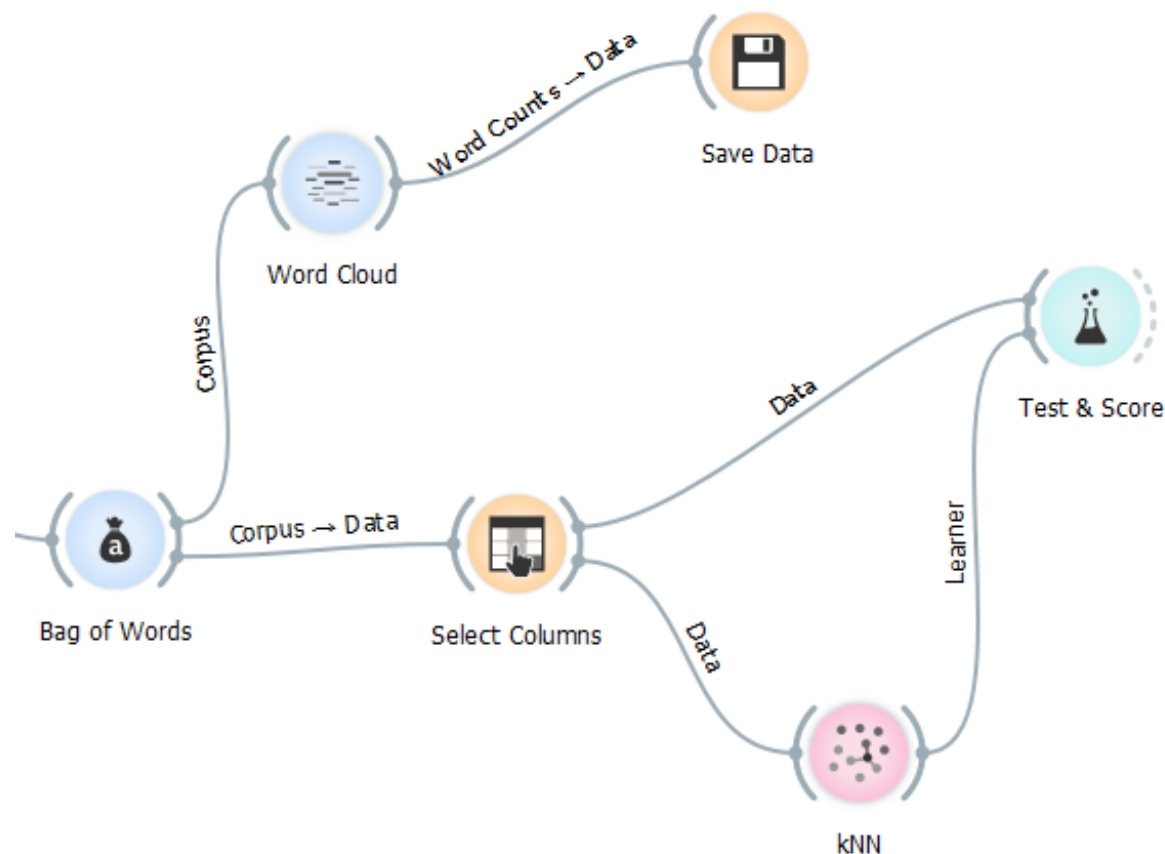
Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
kNN	0.514	0.506	0.351	0.626	0.506

Результат классификации после удаления стоп слов

Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
kNN	0.539	0.500	0.401	0.500	0.500

ЭФФЕКТ ВЫБОРА ЧИСЛА СОСЕДЕЙ

Используется очищенный датасет и эвклидово расстояние



КЛАССИФИКАЦИЯ АНГЛОЯЗЫЧНЫХ ТЕКСТОВ

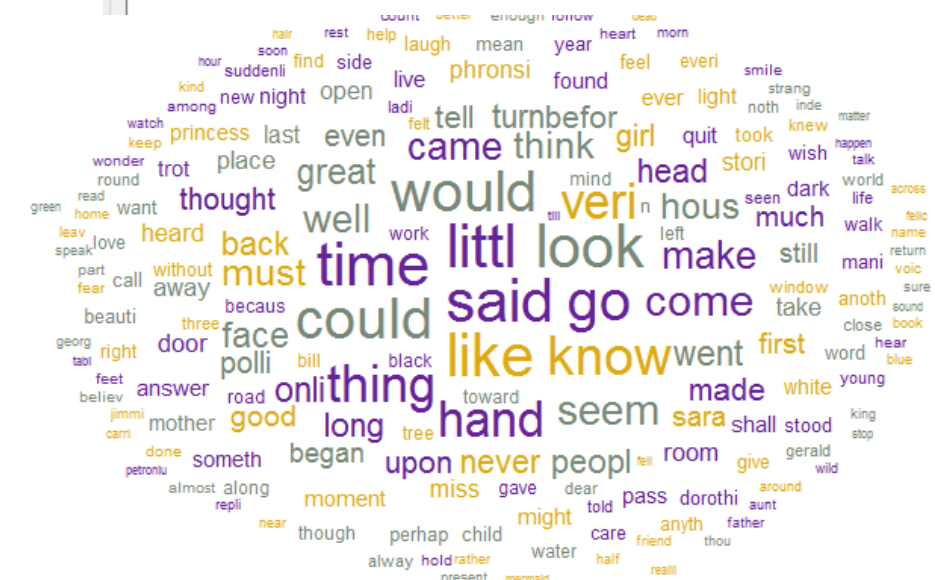
1	Document 1	children:	children
2	Document 2	the house Jim says he rum ; and as he	has lived rough and I'll raise Cain Your doctor hisself said
3	Document 3	spoke he reeled a little and caught himself	one glass wouldn't hurt me I'll give you a golden guinea for
4	Document 4	with one hand against the wall Are you	a noggin Jim He was growing more and more excited and
5	Document 5	hurt? cried I Rum he repeated I must get	this alarmed me for my father who was very low that day
6	Document 6	away from here Rum! Rum! I ran to fetch it	and needed quiet; besides I was reassured by the doctor's
7	Document 7	but I was quite unsteadied by all that had	words now quoted to me and rather offended by the offer of
8	Document 8	fallen out and I broke one glass and fouled	a bribe I want none of your money said I but what you owe
9	Document 9	the tap and while I was still getting in my	my father I'll get you one glass and no more When I
10	Document 10	own way I heard a loud fall in the parlour	brought it to him he seized it greedily and drank it out Aye
		and running in beheld the captain lying full	aye said he that's some better sure enough And now
		length upon the floor At the same instant my	matey did that doctor say how long I was to lie here in this
		mother alarmed by the cries and fighting	old berth? A week at least said I Thunder! he cried A week!
		came running downstairs to help me	I can't do that; they'd have the black spot on me by then
		Between us we raised his head He was	The lubbers is going about to get the wind of me this
		breathing very loud and hard but his eyes	blessed moment; lubbers as couldn't keep what they got
		were closed and his face a horrible colour	
		Dear deary me cried my mother what a	

Method	AUC	CA	F1	Precision	Recall
kNN	0.878	0.762	0.760	0.774	0.762

Удаление стоп слов



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
kNN	0.896	0.787	0.786	0.797	0.787



ЭФФЕКТ ВЫБОРА РАССТОЯНИЯ НА ТЕКСТОВУЮ КЛАССИФИКАЦИЮ

kNN ? x

Name
kNN

Neighbors
Number of neighbors: 1

Metric: Euclidean

Weight: Uniform

☐ Apply

? | [icon]



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
kNN	0.925	0.925	0.925	0.926	0.925

kNN ? x

Name
kNN

Neighbors
Number of neighbors: 1

Metric: Manhattan

Weight: Uniform

☐ Apply

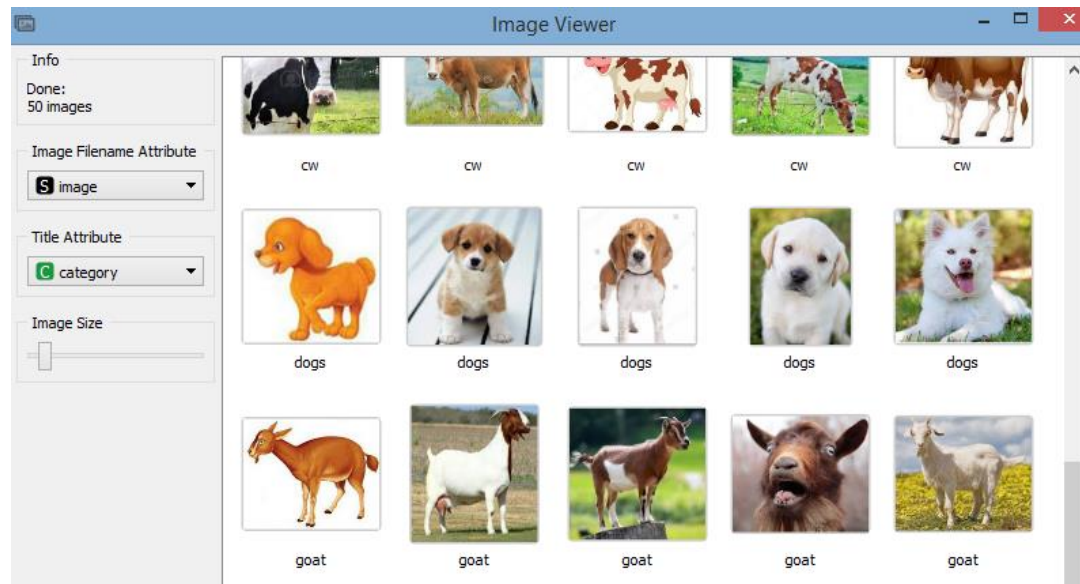
? | [icon]



Evaluation Results

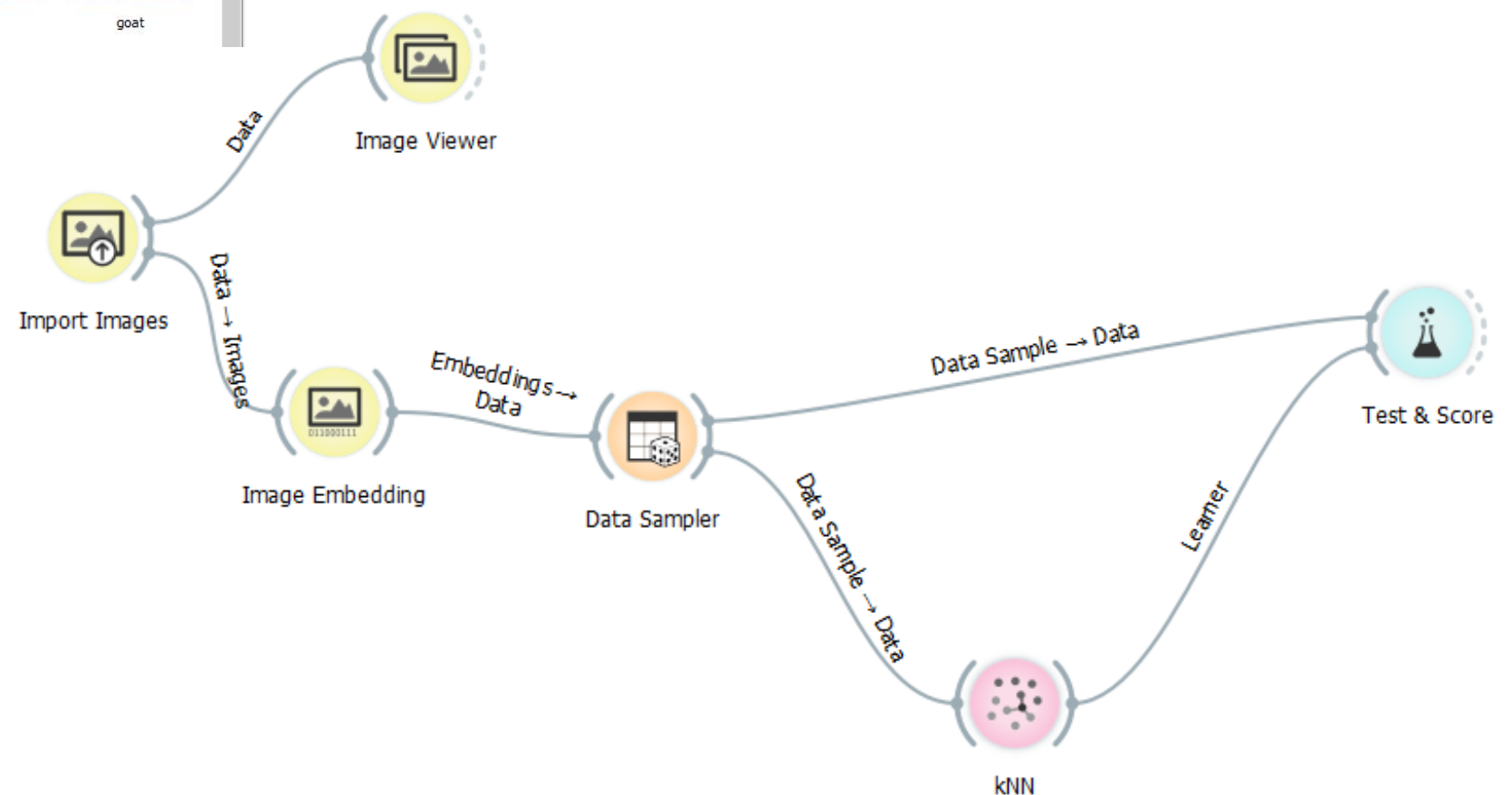
Method	AUC	CA	F1	Precision	Recall
kNN	0.673	0.673	0.637	0.786	0.673

КЛАССИФИКАЦИЯ ИЗОБРАЖЕНИЙ

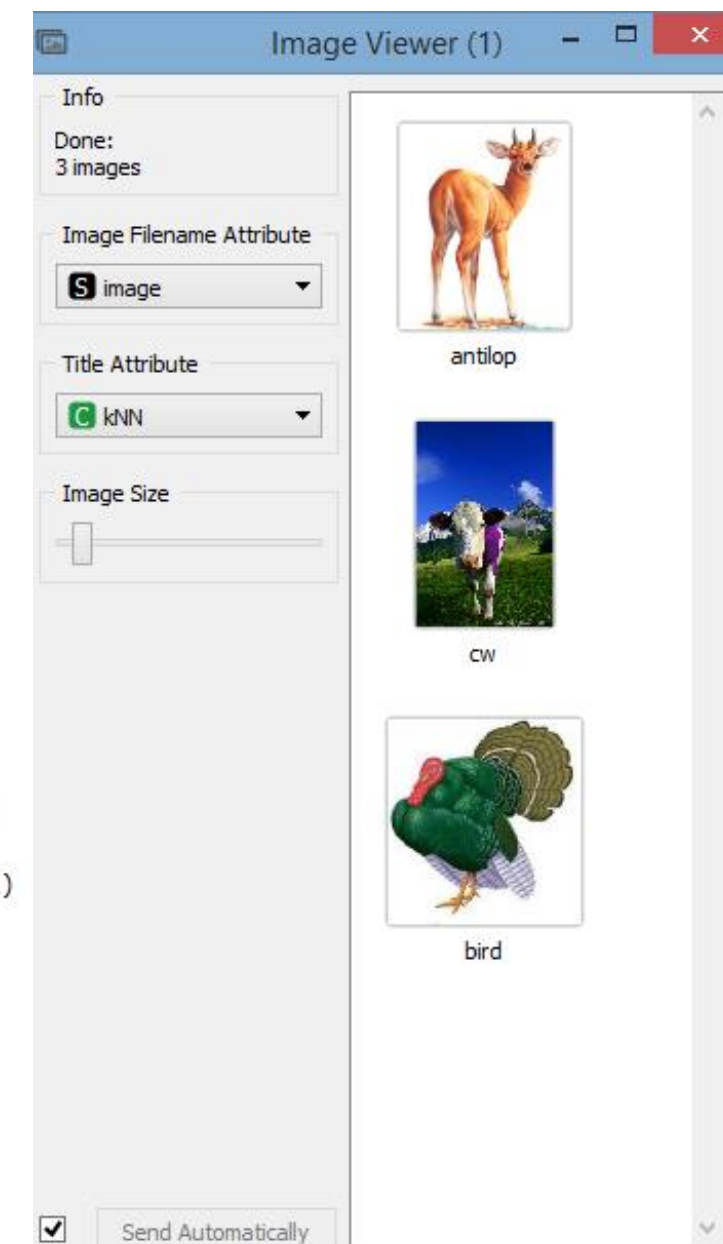
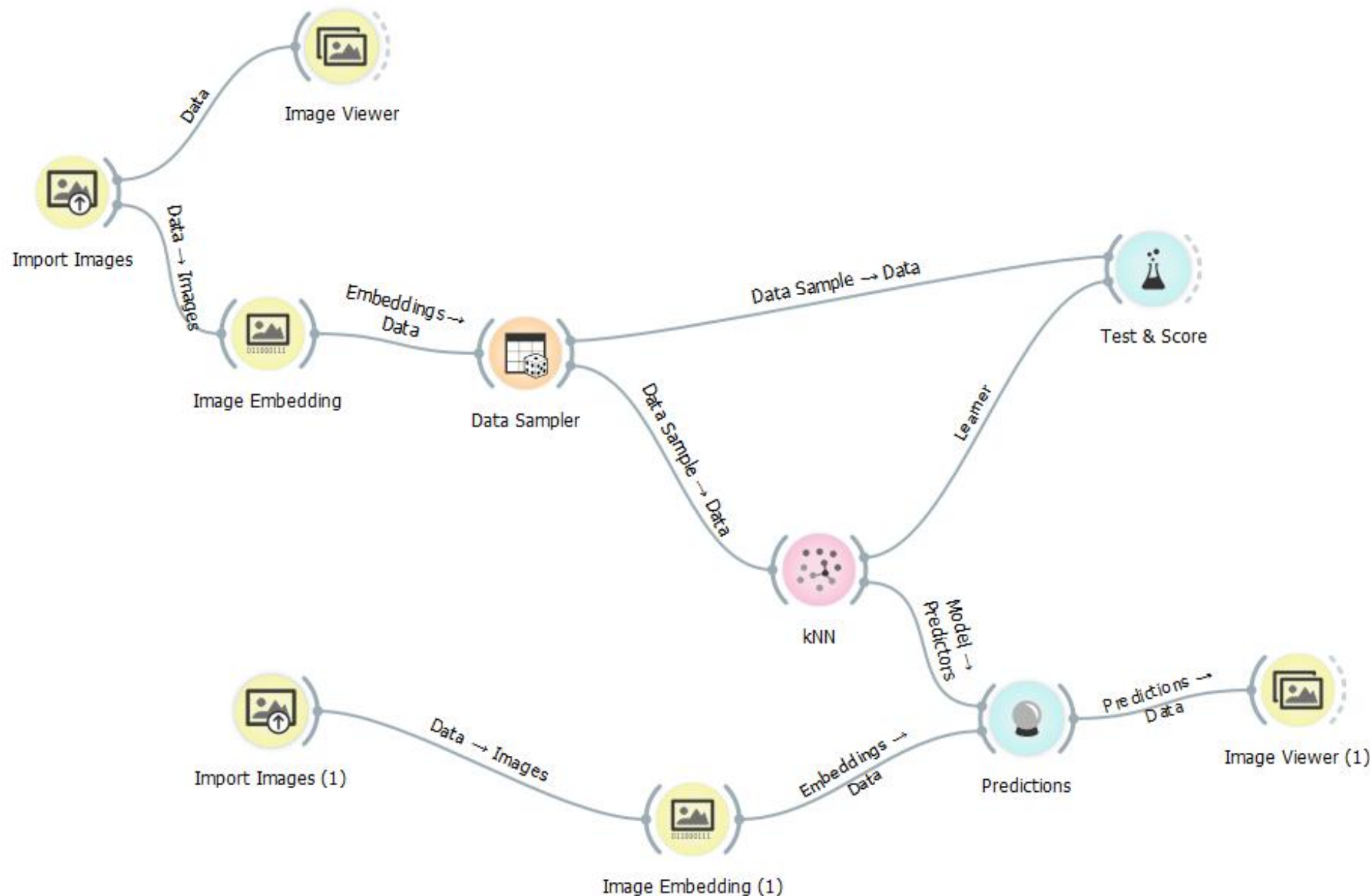


Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
kNN	0.777	0.625	0.592	0.597	0.625

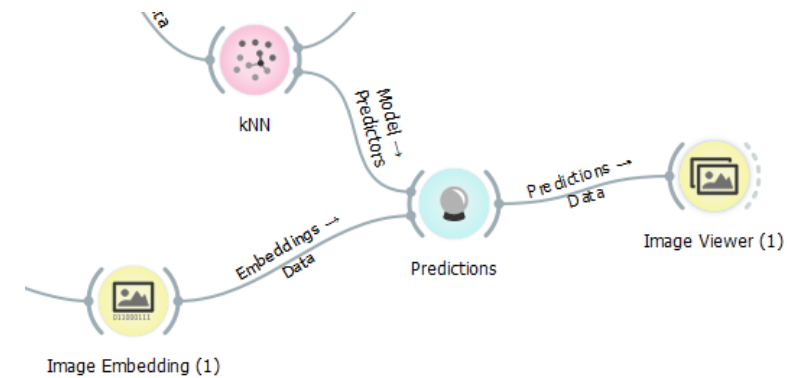
Изображения разложены по категориям. Каждая категория это отдельный каталог.



КЛАССИФИКАЦИЯ ИЗОБРАЖЕНИЙ



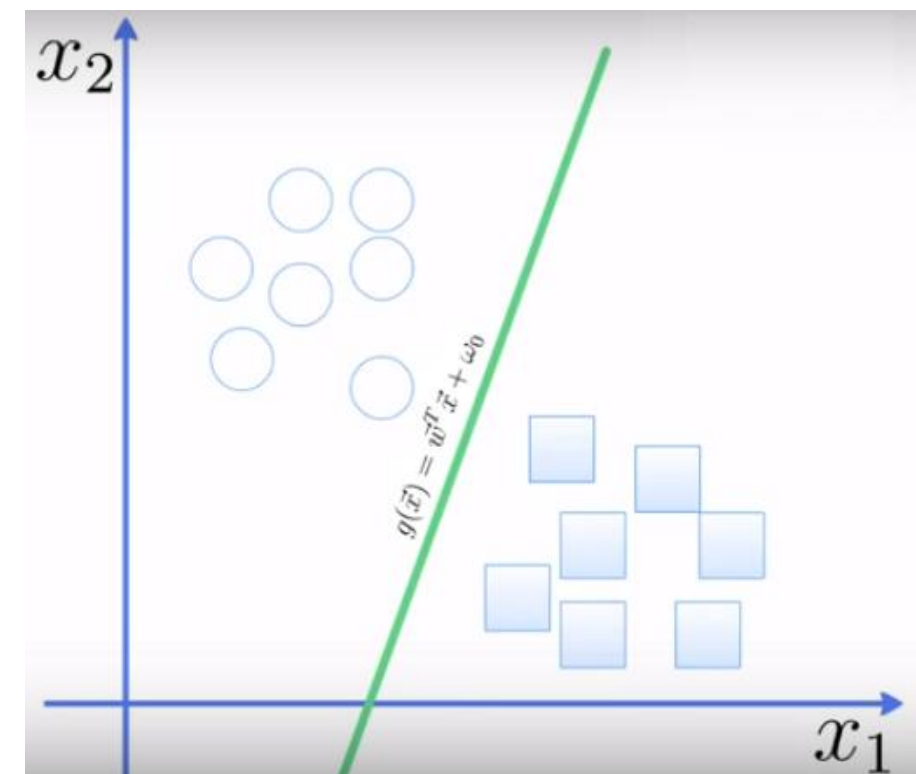
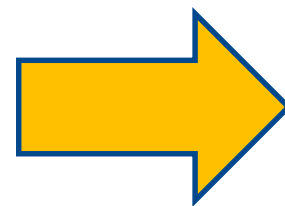
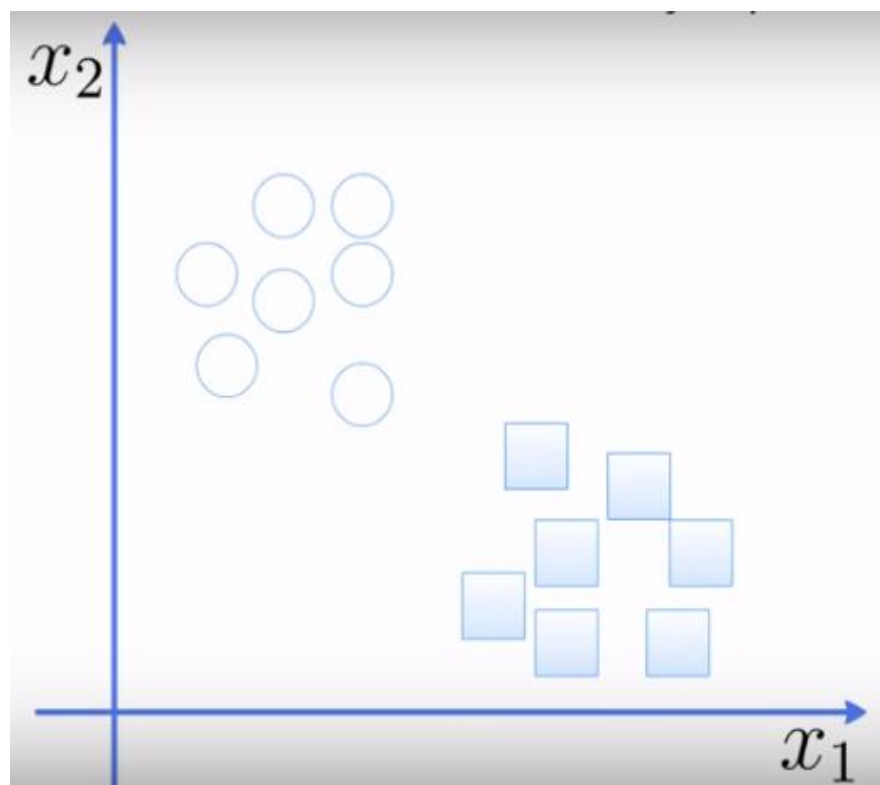
КЛАССИФИКАЦИЯ ИЗОБРАЖЕНИЙ



МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM)

Каждый объект данных (например, документ, котировки ценных бумаг или компании) представлен как вектор в P мерном пространстве (последовательность чисел). Пусть у нас есть тестовая коллекция, в которой есть набор объектов (features) и есть набор классов. Математическая задача обучения заключается в том что бы найти функцию, которая адекватно сопоставляла объекты и классы, то есть найти такую функцию, которая эффективно разделяла бы объекты в пространстве features.

Рассмотрим пример на плоскости: У нас есть два класса с двумя features (x_1 , x_2). Нужно найти прямую линию, которая оптимально разделяла два класса.

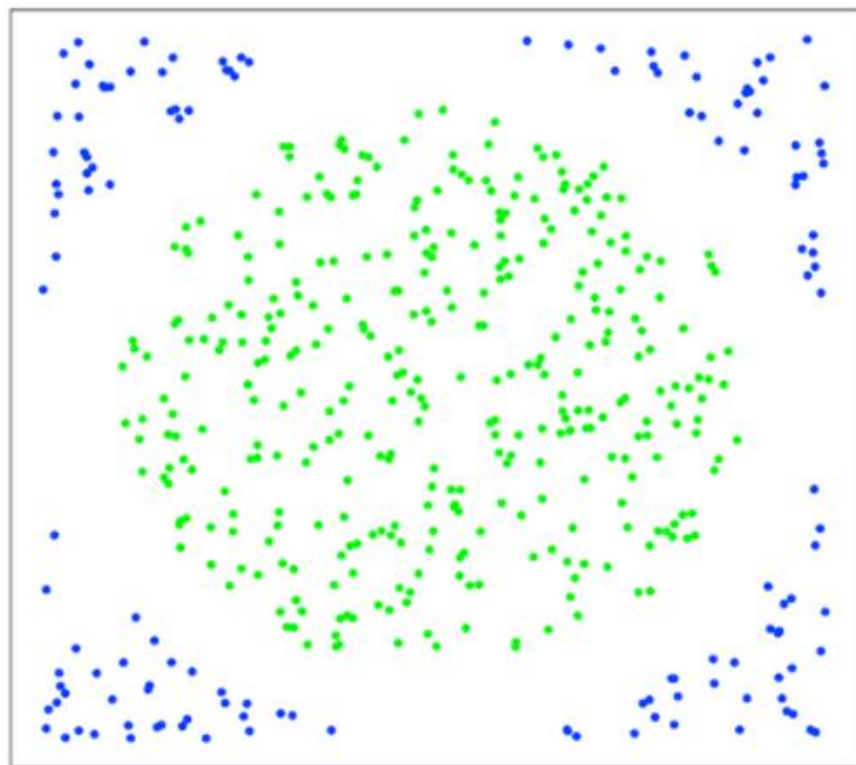


МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM)

Нахождение уравнения плоскости является стандартной задачей квадратичного программирования и решается с помощью множителей Лагранжа. Собственно в этом заключается процесс обучения.

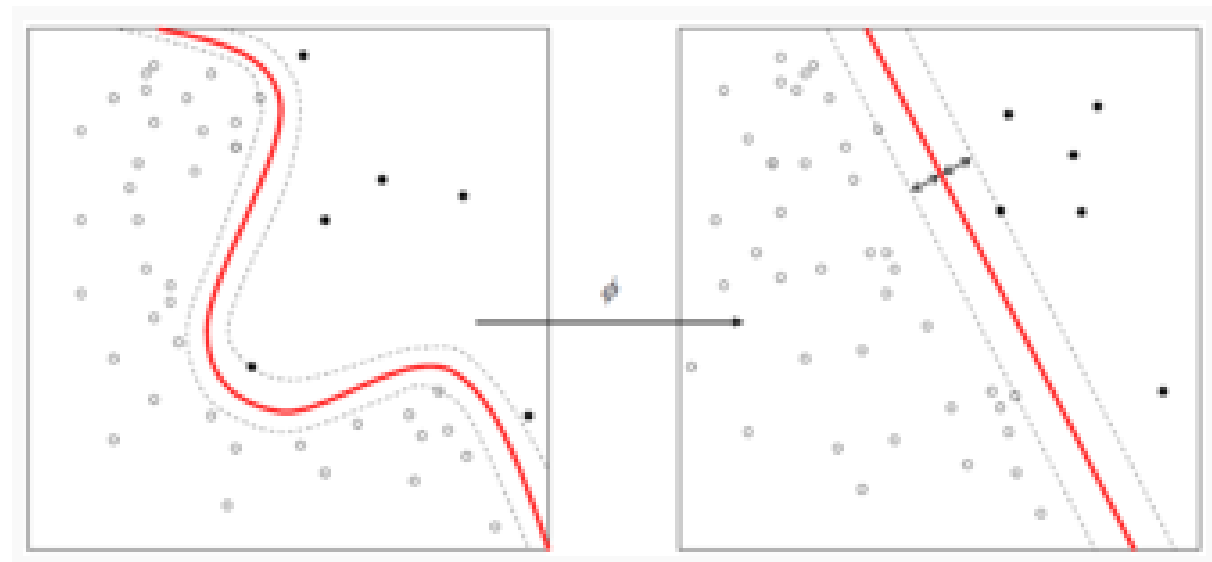
Как только плоскость найдена, берем новый объект и смотрим где он расположен относительно плоскости. Если справа, то принадлежит одному классу, если слева, то наш объект принадлежит другому классу.

Однако, как правило на практике встречаются случаи когда объекты расположены, так что на плоскости невозможно провести разделяющую прямую. В этом случае плоскость вкладывается в пространство большей размерности. При вложении плоскость трансформируется таким образом, что бы появилась возможность провести разделяющую плоскость.



Демонстрация подобного преобразования

<https://youtu.be/3liCbRZPrZA>



МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM)

SVM

Name: SVM

SVM Type

☒ SVM Cost (C): 1,00

☐ v-SVM Regression loss epsilon (ε): 0,10

Regression cost (C): 1,00

Complexity bound (v): 0,50

Kernel

☐ Linear Kernel: $\tanh(g \cdot x \cdot y + c)$

☐ Polynomial g: auto

☐ RBF c: 0,00

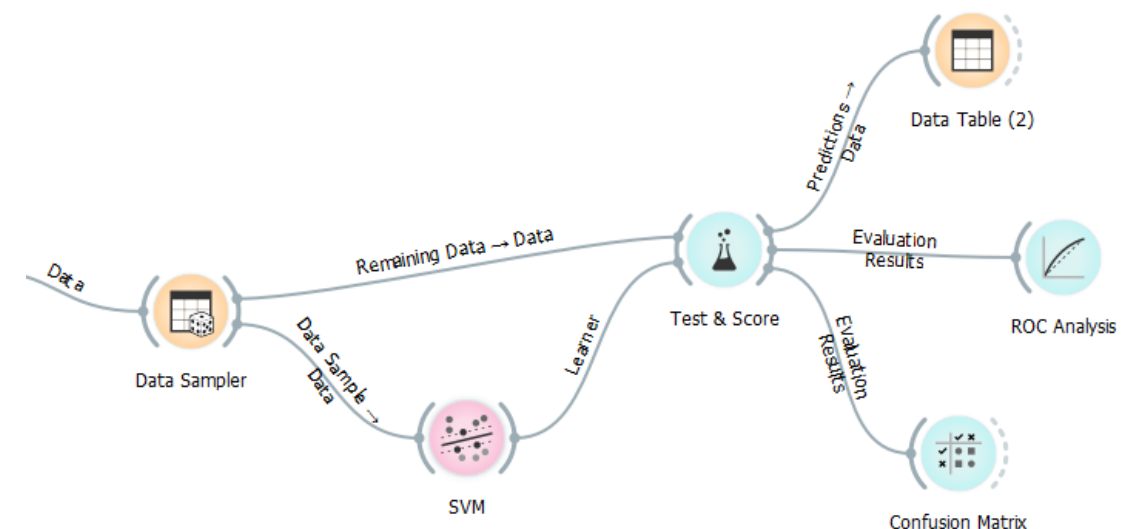
☒ Sigmoid

Optimization Parameters

Numerical tolerance: 0,0010

☒ Iteration limit: 100

☒ Apply Automatically



Типы kernels (ядер):

1. Линейное ядро.

2. Полиномиальное ядро.

$$k(x, y) = (\alpha x^T y + c)^d$$

3. Radial basis function kernel

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

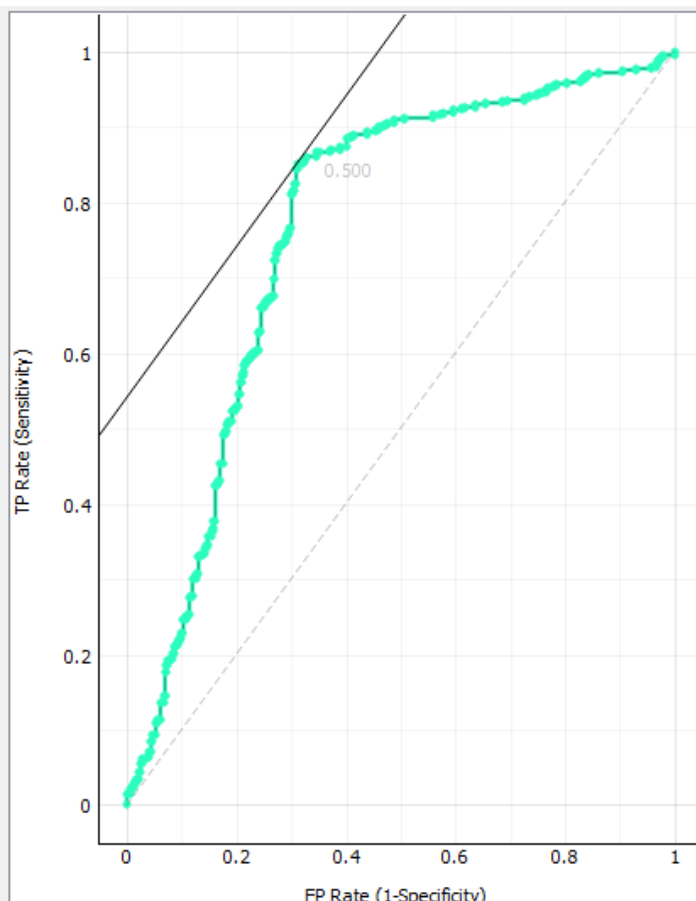
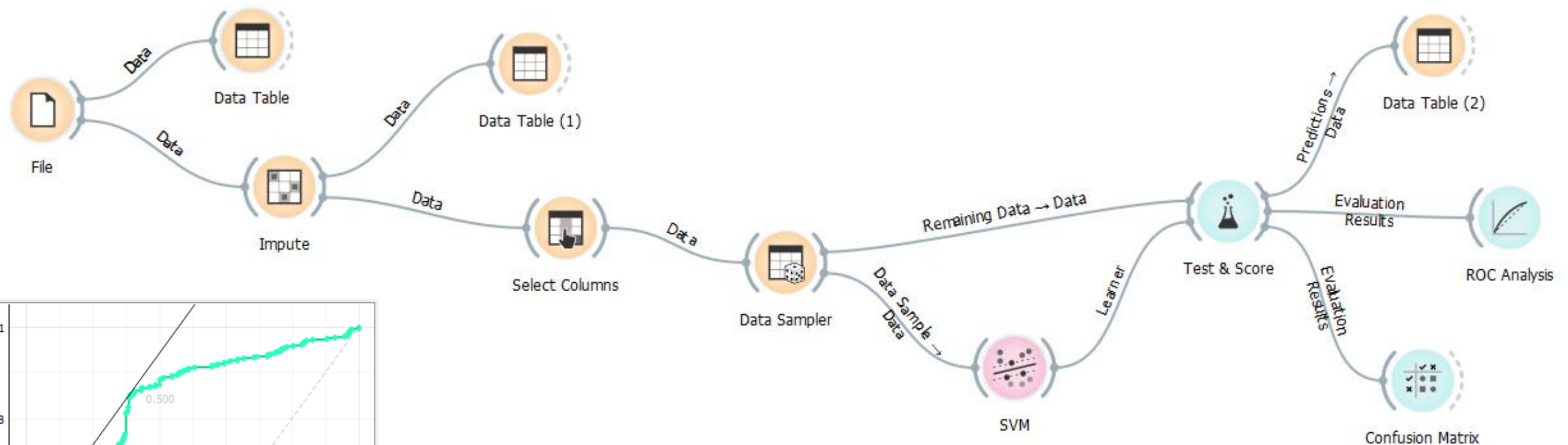
4. Hyperbolic Tangent (Sigmoid) Kernel

$$k(x, y) = \tanh(\alpha x^T y + c)$$

c - const (default =0)

d - степень ядра

SVM - ТАБЛИЧНЫЕ ДАННЫЕ



	Predicted		Σ
	0.0	1.0	
0.0	473	77	550
1.0	120	240	360
Σ	593	317	910

Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.770	0.784	0.781	0.782	0.784

SVM – ТАБЛИЧНЫЕ ДАННЫЕ – ЭФФЕКТ ВЫБОРА ЯДРА

Kernel

☒ Linear Kernel: $x \cdot y$

☐ Polynomial

☐ RBF

☐ Sigmoid



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
SVM	0.776	0.780	0.776	0.779	0.780

Kernel

☐ Linear Kernel: $(g x \cdot y + c)^d$

☒ Polynomial

☐ RBF

☐ Sigmoid

g:

c:

d:



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
SVM	0.803	0.795	0.789	0.796	0.795

Kernel

☐ Linear Kernel: $\tanh(g x \cdot y + c)$

☐ Polynomial

☐ RBF

☒ Sigmoid

g:

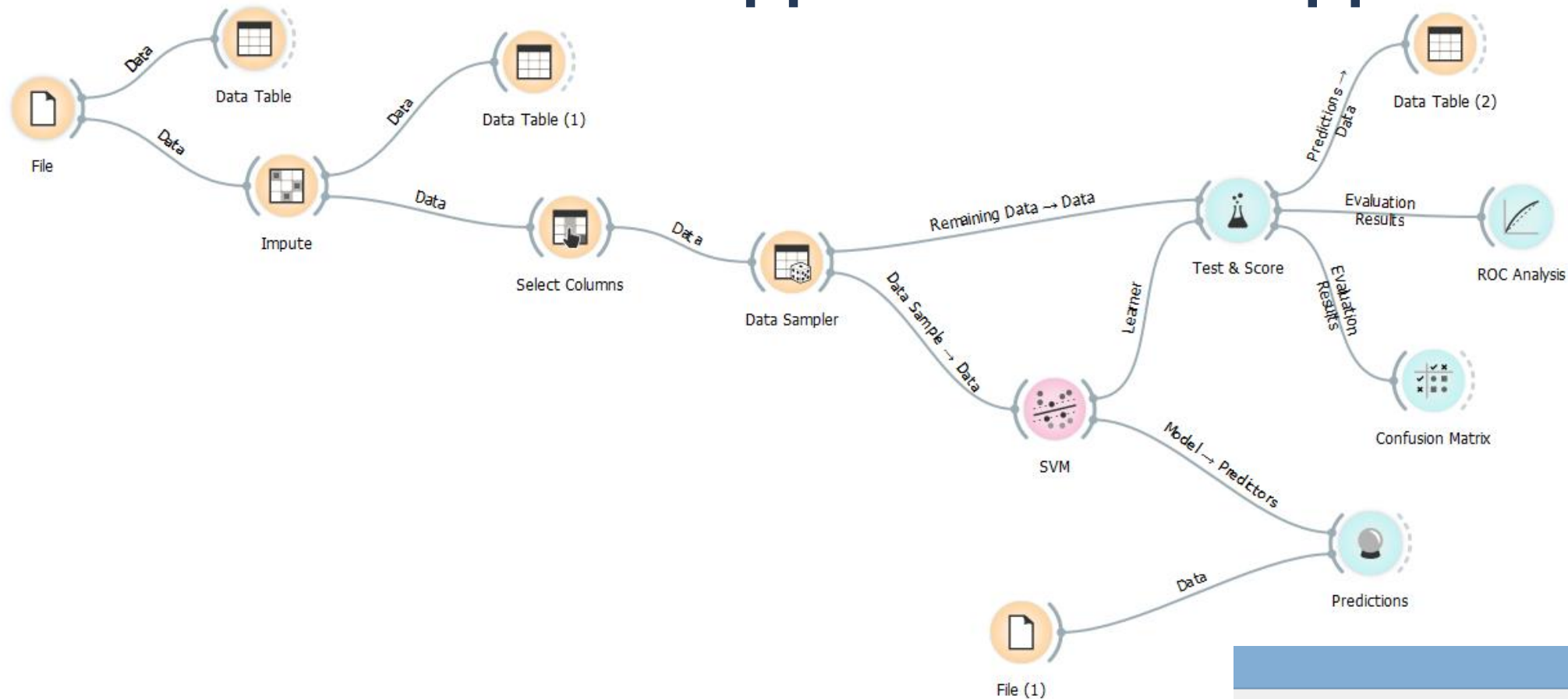
c:



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
SVM	0.770	0.784	0.781	0.782	0.784

SVM – ТАБЛИЧНЫЕ ДАННЫЕ - ПРЕДСКАЗАНИЕ

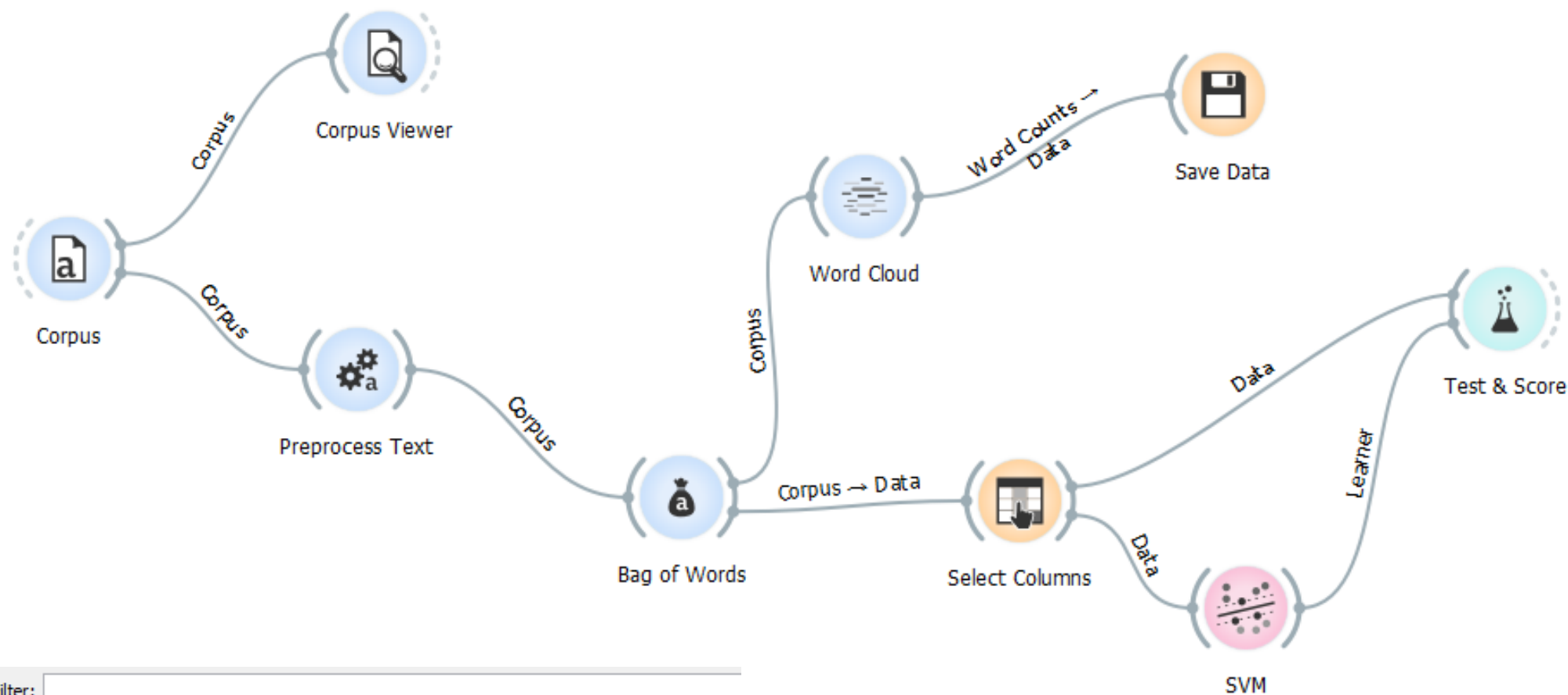


	Name	Ticket	PassengerId	Pclass	Sex	Age
1	Braund, Mr. Ow...	A/5 21171	1.0	3.0	1.0	22.0
2	Cumings, Mrs. ...	PC 17599	2.0	1.0	0.0	38.0
3	Heikkinen, Miss...	STON/O2. 3101...	3.0	3.0	0.0	26.0
4	Futrelle, Mrs. Ja...	113803.0	4.0	1.0	0.0	35.0
5	Allen, Mr. Willia...	373450.0	5.0	3.0	1.0	35.0
6	Moran, Mr. Jam...	330877.0	6.0	3.0	1.0	?
7	McCarthy, Mr. ...	17463.0	7.0	1.0	1.0	54.0
8	Palsson, Master...	349909.0	8.0	3.0	1.0	2.0
9	Johnson, Mrs. ...	347742.0	9.0	3.0	0.0	27.0
10	Nasser, Mrs. Ni...	237736.0	10.0	2.0	0.0	14.0
11	Sandstrom, Mis...	PP 9549	11.0	3.0	0.0	4.0
12	Bonnell, Miss. E...	113783.0	12.0	1.0	0.0	58.0
13	Saunderscock, ...	A/5. 2151	13.0	3.0	1.0	20.0
14	Andersson, Mr. ...	347082.0	14.0	3.0	1.0	39.0



Predictions				
	SVM	Name	Ticket	PassengerId
1	0.78 : 0.22 → 0.0	Braund, Mr. Ow...	A/5 21171	1.0
2	0.34 : 0.66 → 1.0	Cumings, Mrs. ...	PC 17599	2.0
3	0.33 : 0.67 → 1.0	Heikkinen, Miss...	STON/O2. 3101...	3.0
4	0.42 : 0.58 → 1.0	Futrelle, Mrs. Ja...	113803.0	4.0
5	0.63 : 0.37 → 0.0	Allen, Mr. Willia...	373450.0	5.0
6	0.67 : 0.33 → 0.0	Moran, Mr. Jam...	330877.0	6.0
7	0.38 : 0.62 → 1.0	McCarthy, Mr. ...	17463.0	7.0
8	0.59 : 0.41 → 0.0	Palsson, Master...	349909.0	8.0
9	0.44 : 0.56 → 1.0	Johnson, Mrs. ...	347742.0	9.0
10	0.41 : 0.59 → 1.0	Nasser, Mrs. Ni...	237736.0	10.0
11	0.69 : 0.31 → 0.0	Sandstrom, Mis...	PP 9549	11.0
12	0.26 : 0.74 → 1.0	Bonnell, Miss. E...	113783.0	12.0
13	0.62 : 0.38 → 0.0	Saunderscock, ...	A/5. 2151	13.0
14	0.12 : 0.88 → 1.0	Andersson, Mr. ...	347082.0	14.0

SVM – ТЕКСТОВЫЕ ДАННЫЕ



RegExp Filter:

	Document
1	Document 1
2	Document 2
3	Document 3
4	Document 4
5	Document 5
6	Document 6
7	Document 7
8	Document 8

children:

the house Jim says he rum ; and as he spoke he reeled a little and caught himself with one hand against the wall Are you hurt? cried I Rum he repeated I must get away from here Rum! Rum! I ran to fetch it but I was quite unsteadied by all that had fallen out and I broke one glass and fouled the tap and while I was still getting in my own way I heard a loud fall in the parlour and running in beheld the captain lying full length upon the floor At the same instant my mother alarmed by the cries and fighting came running downstairs to help me

Evaluation Results

Method	AUC	CA	F1	Precision	Recall
SVM	0.998	0.965	0.965	0.965	0.965

SVM – ТЕКСТОВЫЕ ДАННЫЕ – ЭФФЕКТ ВЫБОРА ЯДРА

Kernel

☒ Linear Kernel: $x \cdot y$

☐ Polynomial

☐ RBF

☐ Sigmoid



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.998	0.965	0.965	0.965	0.965

Kernel

☐ Linear Kernel: $(g \cdot x \cdot y + c)^d$

☒ Polynomial

☐ RBF

☐ Sigmoid

g:

c:

d:



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.033	0.500	0.333	0.250	0.500

Kernel

☐ Linear Kernel: $\tanh(g \cdot x \cdot y + c)$

☐ Polynomial

☐ RBF

☒ Sigmoid

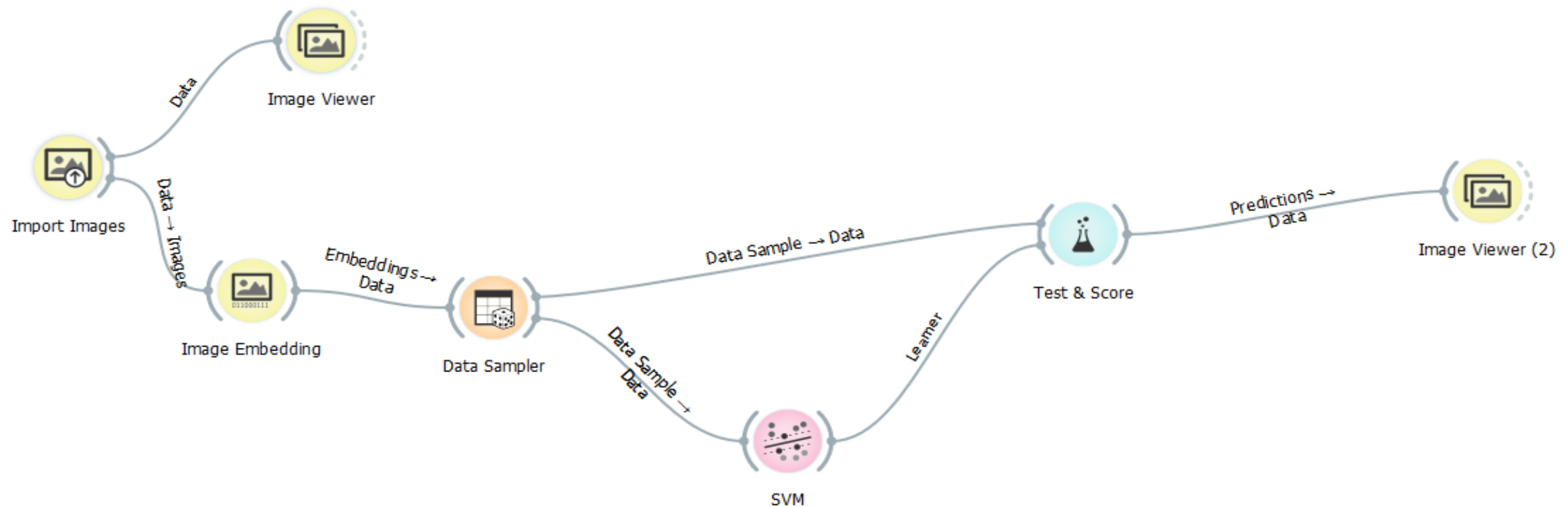
g:

c:



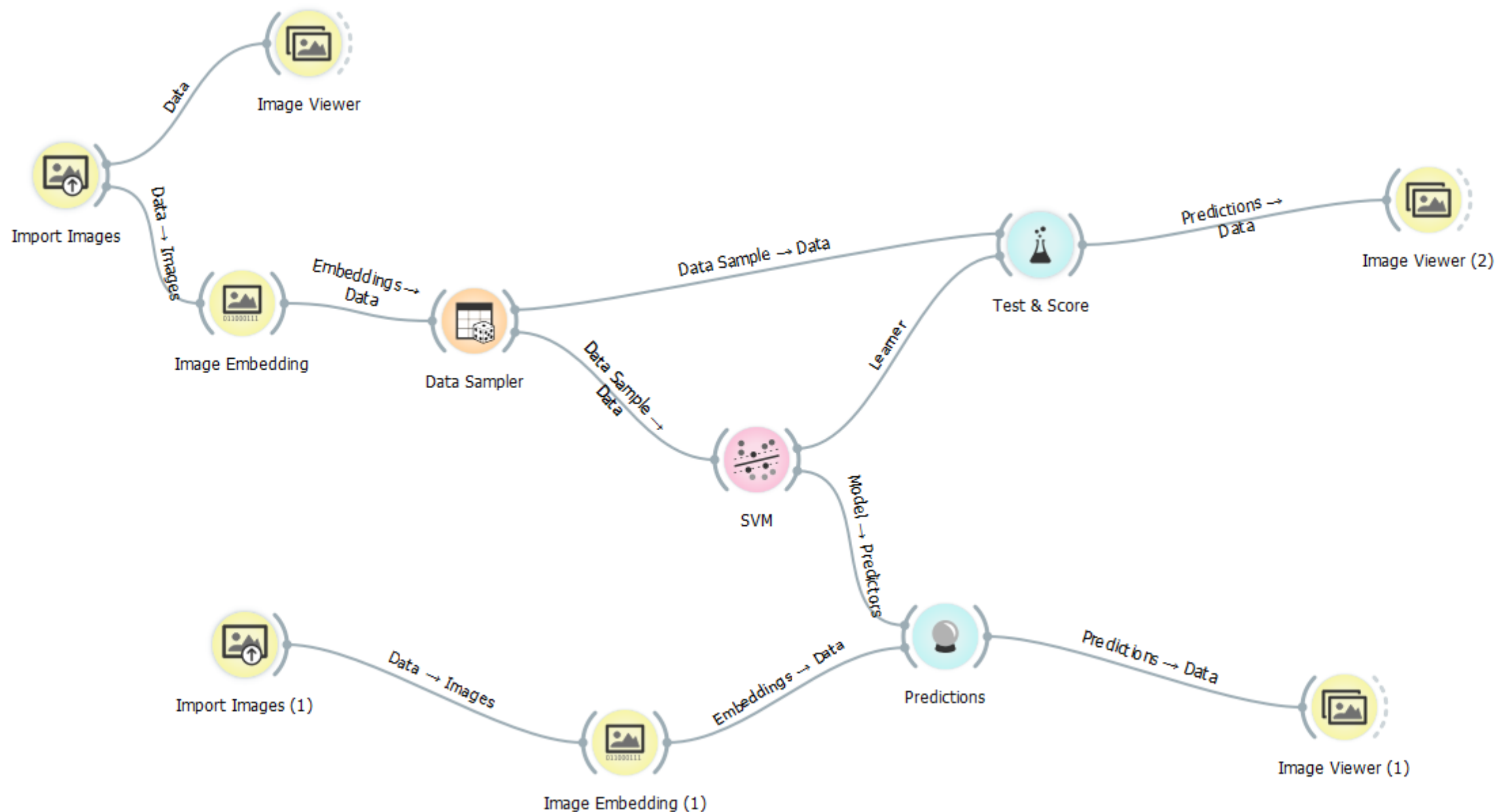
Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.953	0.625	0.565	0.778	0.625

МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM) ИЗОБРАЖЕНИЯ



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.852	0.803	0.770	0.833	0.803

МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM) ИЗОБРАЖЕНИЯ - ПРЕДСКАЗАНИЯ



МЕТОД ОПОРНЫХ ВЕКТОРОВ (SVM)

ЭФФЕКТ IMAGE EMBEDDING

Kernel

☒ Linear Kernel: $x \cdot y$

☐ Polynomial

☐ RBF

☐ Sigmoid



Method	AUC	CA	F1	Precision	Recall
SVM	0.999	0.917	0.914	0.926	0.917

Kernel

☐ Linear Kernel: $(g \cdot x \cdot y + c)^d$

☒ Polynomial

☐ RBF

☐ Sigmoid

g:

c:

d:



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.272	0.240	0.236	0.612	0.240

Kernel

☐ Linear Kernel: $\tanh(g \cdot x \cdot y + c)$

☐ Polynomial

☐ RBF

☒ Sigmoid

g:

c:



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
SVM	0.956	0.883	0.880	0.897	0.883

ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ В ЗАДАЧАХ КЛАССИФИКАЦИИ

Логистическая регрессия – это разновидность множественной регрессии, общее назначение которой состоит в анализе связи между несколькими независимыми переменными (называемыми также регрессорами или предикторами) и зависимой переменной. Бинарная логистическая регрессия, как следует из названия, применяется в случае, когда зависимая переменная является бинарной (т.е. может принимать только два значения). Иными словами, с помощью логистической регрессии можно оценивать вероятность того, что событие наступит для конкретного испытуемого (больной/здоровый, возврат кредита/дефолт и т.д.).

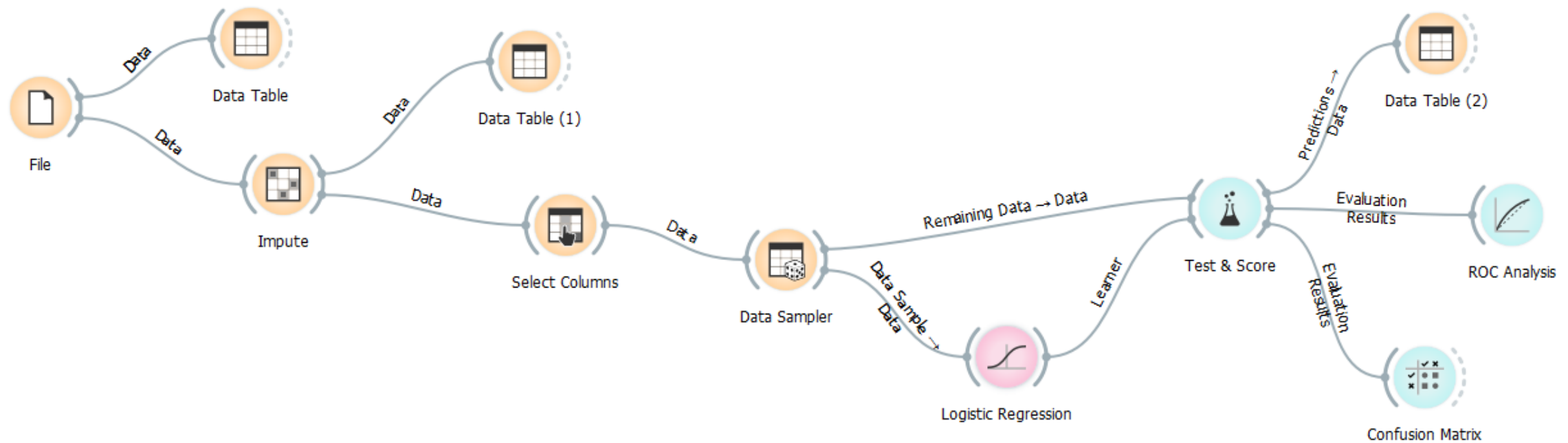
$$y = a + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Для решения проблемы задача регрессии может быть сформулирована иначе: вместо предсказания бинарной переменной, мы предсказываем непрерывную переменную со значениями на отрезке $[0,1]$ при любых значениях независимых переменных. Это достигается применением следующего регрессионного уравнения (логит-преобразование):

$$p = \frac{1}{1 + e^{-y}}$$

где P – вероятность того, что произойдет интересующее событие; e – основание натуральных логарифмов $2,71\dots$; y – стандартное уравнение регрессии.

LOGISTIC REGRESSION



Logistic Regressi... ? x

Name
Logistic Regression

Regularization type: Ridge (L2) Lasso (L1) Ridge (L2)

Strength:
Weak Strong
C=1

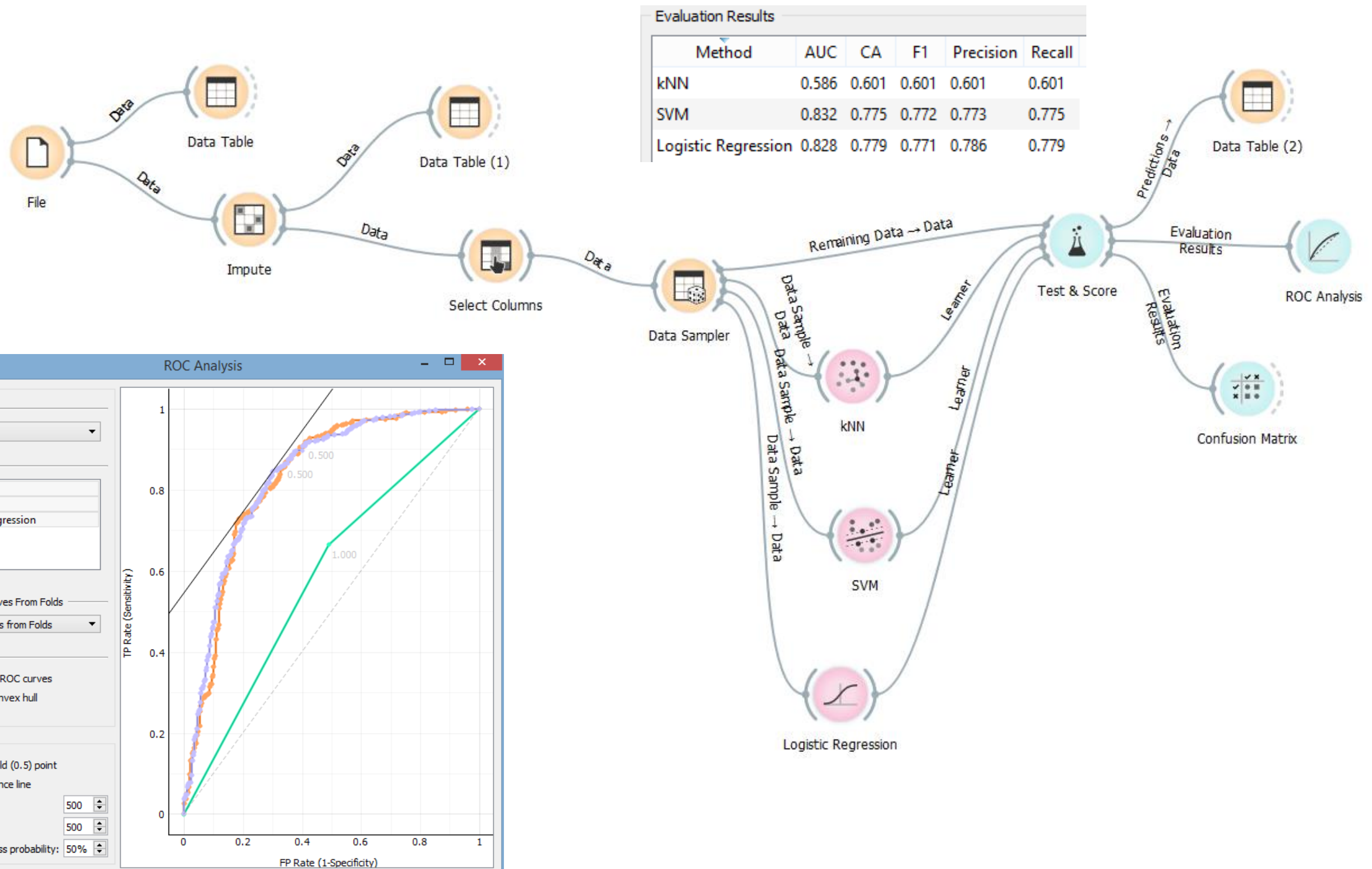
☒ Apply Automatically

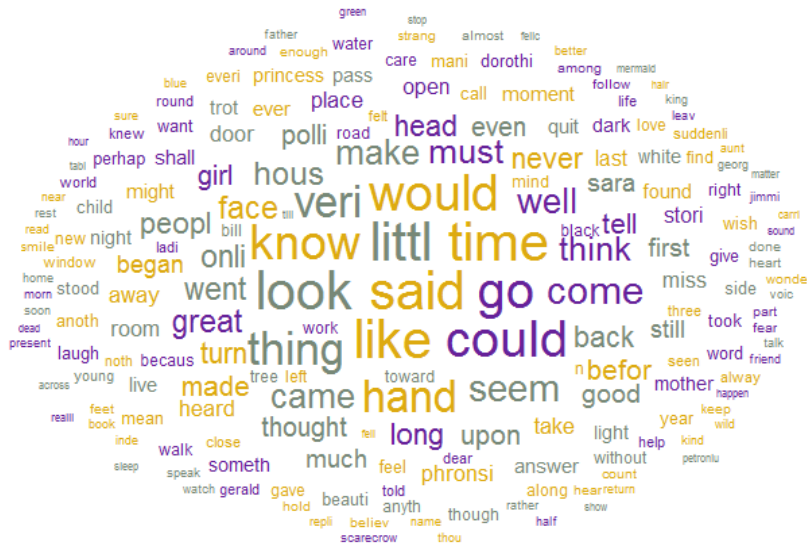
? [Icon]

	Name	Ticket	Cabin	PassengerId	Survived
1	Braund, Mr. Ow...	A/5 21171	?	1.0	0.0
2	Cumings, Mrs. ...	PC 17599	C85	2.0	1.0
3	Heikkinen, Miss...	STON/O2. 3101...	?	3.0	1.0
4	Futrelle, Mrs. Ja...	113803.0	C123	4.0	1.0
5	Allen, Mr. Willia...	373450.0	?	5.0	0.0
6	Moran, Mr. Jam...	330877.0	?	6.0	0.0
7	McCarthy, Mr. ...	17463.0	E46	7.0	0.0
8	Palsson, Master...	349909.0	?	8.0	0.0
9	Johnson, Mrs. ...	347742.0	?	9.0	1.0
10	Nasser, Mrs. Ni...	237736.0	?	10.0	1.0
11	Sandstrom, Mis...	PP 9549	G6	11.0	1.0
12	Bonnell, Miss. E...	113783.0	C103	12.0	1.0
13	Saunderscock, ...	A/5. 2151	?	13.0	0.0

Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.828	0.779	0.771	0.786	0.779

СРАВНЕНИЕ РАБОТЫ ТРЕХ КЛАССИФИКАТОРОВ





Evaluation Results						
Method	AUC	CA	F1	Precision	Recall	
Logistic Regression	0.997	0.956	0.956	0.956	0.956	

LOGISTIC REGRESSION – ЭФФЕКТ РЕГУЛЯРИЗАЦИИ

Name

Logistic Regression

Regularization type: Lasso (L1)

Strength:

Weak Strong

C=1



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.968	0.915	0.915	0.915	0.915

Name

Logistic Regression

Regularization type: Ridge (L2)

Strength:

Weak Strong

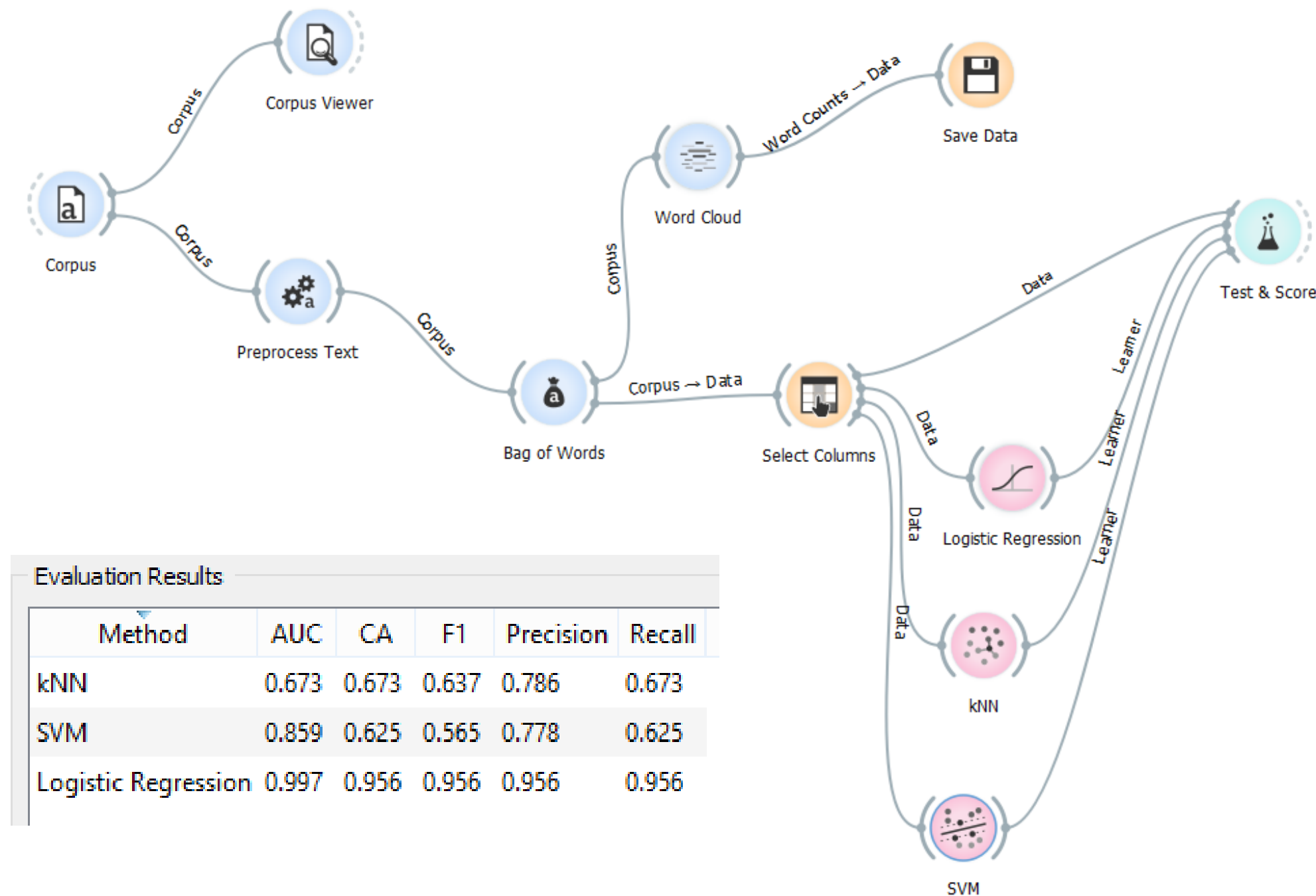
C=1



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.997	0.956	0.956	0.956	0.956

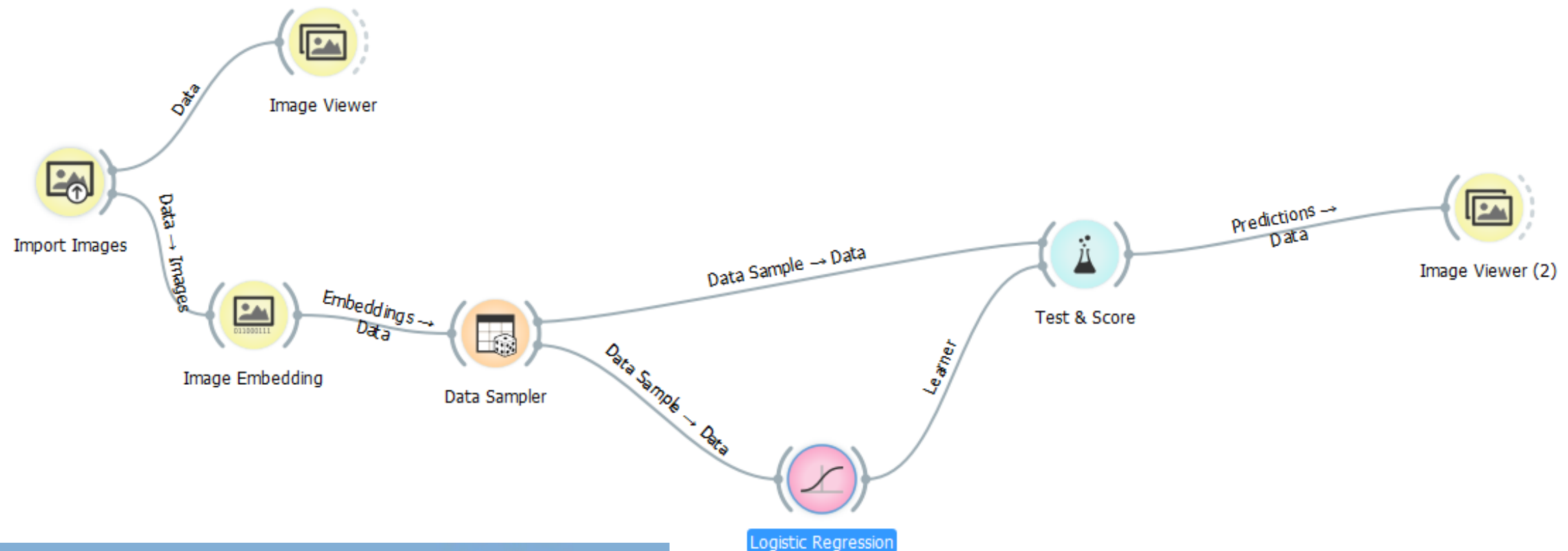
LOGISTIC REGRESSION – KNN – SVM



Evaluation Results

Method	AUC	CA	F1	Precision	Recall
kNN	0.673	0.673	0.637	0.786	0.673
SVM	0.859	0.625	0.565	0.778	0.625
Logistic Regression	0.997	0.956	0.956	0.956	0.956

LOGISTIC REGRESSION – IMAGE CLASSIFICATION



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.997	0.940	0.936	0.947	0.940

LOGISTIC REGRESSION – ЭФФЕКТ РЕГУЛЯРИЗАЦИИ

Logistic Regression... ? x

Name
Logistic Regression

Regularization type: Ridge (L2) ▼

Strength:
Weak ——— Strong
C=1



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.997	0.940	0.936	0.947	0.940

Name
Logistic Regression

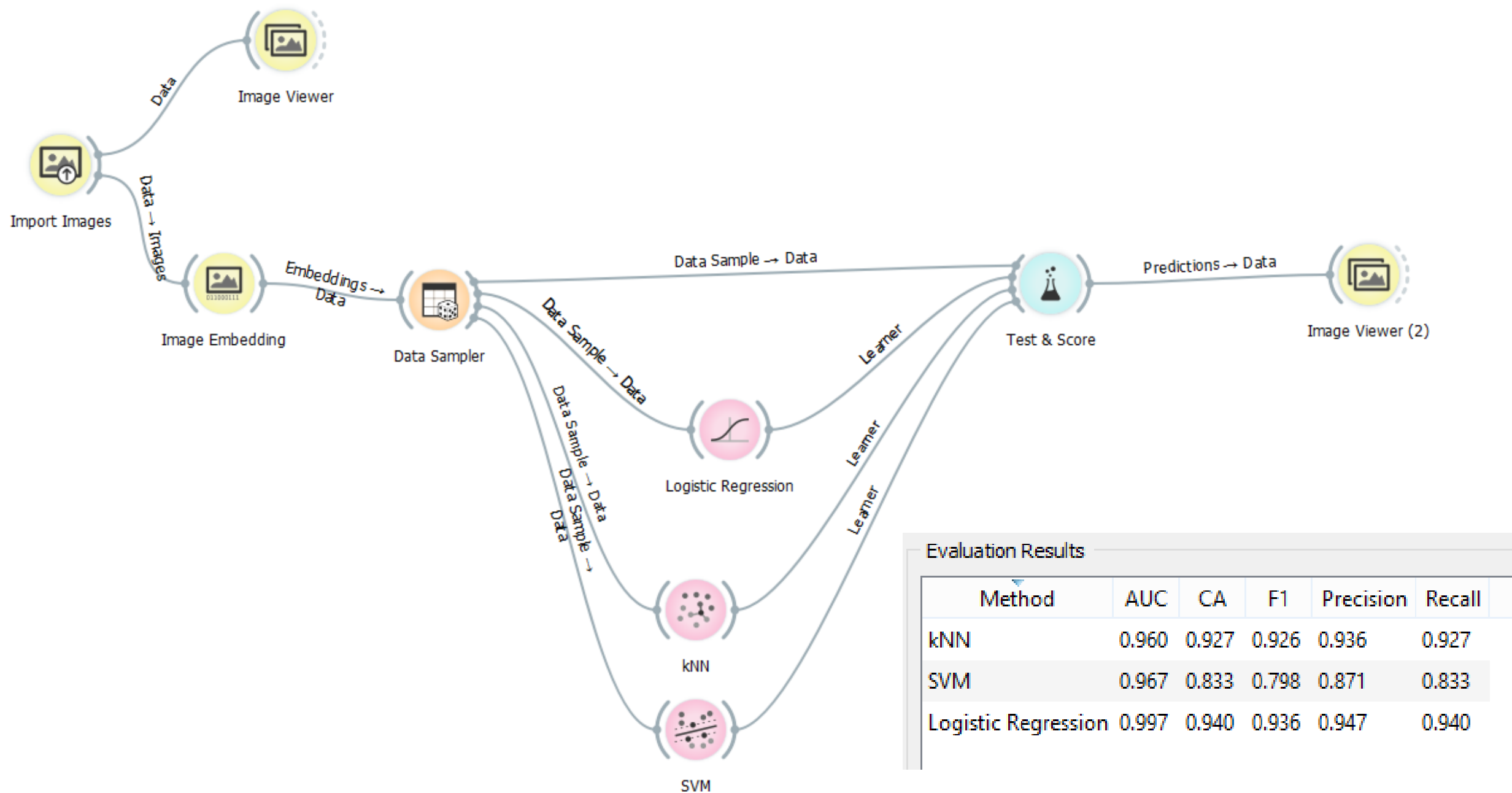
Regularization type: Lasso (L1) ▼

Strength:
Weak ——— Strong
C=1



Evaluation Results					
Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.974	0.807	0.788	0.820	0.807

LOGISTIC REGRESSION – KNN - SVM





НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

<https://linis.hse.ru/>

Phone: +7 (911) 981 9165

Email: skoltsov@hse.ru