# МИНИСТЕРСТВО ОБРАЗОВАНИЯ РЕСПУБЛИКИ БЕЛАРУСЬ БЕЛОРУССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ МЕХАНИКО-МАТЕМАТИЧЕСКИЙ ФАКУЛЬТЕТ

Кафедра дифференциальных уравнений и системного анализа

### ЖДАНОВИЧ

Евгения Сергеевна

## СТАТИСТИЧЕСКИЙ АНАЛИЗ ТЕКСТА И ЗАДАЧА ОПРЕДЕЛЕНИЯ АВТОРСТВА

Дипломная работа

Научный руководитель: кандидат физ.-мат. наук, доцент Е. М. Радыно

		доцен	т Е. М		
Дог	<b>у</b> щена к защите				
··	<u></u> »	_ 2018 г.			
Зав. кафедрой дифференциальных уравнений и системного анализа					
доктор физмат. наук, профессор В. И. Громак					

# ОГЛАВЛЕНИЕ

РЕФЕР	РЕФЕРАТ				
ABSTRA	ACT.		5		
РЭФЕР	ΔТ		6		
введе					
ГЛАВА	1. (	ОБРАБОТКА ЕСТЕССТВЕННОГО ЯЗЫКА МЕТОДАМИ МАШИННОГО ОБУЧІ	ЕНИЯ.8		
1.1	ВЕК	ТОРНАЯ МОДЕЛЬ ТЕКСТА	8		
1.2	ПРЕ	ДОБРАБОТКА ИСХОДНЫХ ДАННЫХ	8		
1.	2.1	Токенизация	9		
1.	2.2	Нормализация	10		
1.3	Пон	ИЖЕНИЕ РАЗМЕРНОСТИ	11		
1.	3.1	Алгоритм t-SNE	12		
1.	3.1.1	Математическое описание алгоритма	12		
1.	3.1.2	Физическая аналогия	14		
1.	3.2	Метод главных компонент	15		
1.	3.2.1	Описание метода	15		
1.	3.2.2	Сингулярное разложение	17		
ГЛАВА	2. (	СТАТИСТИЧЕСКИЙ АНАЛИЗ ТЕКСТА	18		
2.1	Изв	ЛЕЧЕНИЕ ПРИЗНАКОВ ИЗ ПРЕДОБРАБОТАННОГО ТЕКСТА	18		
2.	1.1	Мешок слов	18		
2.	1.2	Мешок термов	19		
2.	1.3	TF-IDF метод векторизации текстовых данных	19		
2.	1.4	N-граммы	20		
2.	1.5	Хэширование	21		
2.	1.6	Стоп-слова	22		
2.2	Выд	ЕЛЕНИЕ ХАРАКТЕРИСТИК	22		
ГЛАВА	<b>3.</b> 3	ЗАДАЧА ОПРЕДЕЛЕНИЯ АВТОРСТВА	24		
3.1	Пос	ТАНОВКА ЗАДАЧИ	24		
3.2		ОЛЬЗУЕМЫЕ ТЕХНОЛОГИИ			
3.3					
3.4		ЕНИЕ ЗАДАЧИ			
3.5	Критерии качества в задачах классификации		28		
3.6	Кла	ССИФИКАЦИЯ ТЕКСТОВ	31		

3.6.1	Матрица ошибок	31
	Динамика точности классификации	
	АЛИЗАЦИЯ	
	IE	
СПИСОК ИС	ПОЛЬЗОВАННОЙ ЛИТЕРАТУРЫ	44
приложені	ИЕ А. КОД ПРОГРАММЫ	45

### РЕФЕРАТ

В дипломной работе 45 страниц, 10 рисунков, 2 таблицы, 10 источников, одно приложение.

МАШИННОЕ ОБУЧЕНИЕ, СТАТИСТИЧЕСКИЙ АНАЛИЗ ТЕКСТА, ОПРЕДЕЛЕНИЕ АВТОРСТВА, АНАЛИЗ ДАННЫХ, ОБРАБОТКА ЕСТЕСТВЕННОГО ЯЗЫКА

В дипломной работе производится отбор статистических признаков текста, классификация текстов, принадлежащих различным авторам, и исследование динамики точности классификации в зависимости от длины текстовых фрагментов.

Задача решалась в программной среде Jupyter Notebook дистрибутива Anaconda, который позволяет сразу установить Python и необходимые библиотеки.

Для решения поставленной задачи использовались:

- Методы обработки естественного языка;
- Статистические характеристики текстов;
- Методы машинного обучения;
- Методы понижения размерности для возможности визуализации.

На основе полученной динамики изменения точности классификации в зависимости от длин текстовых фрагментов были сделаны соответствующие выводы об оптимальной длине текстов, используемых для обучения и тестирования моделей.

Дипломная работа выполнена автором самостоятельно.

### **ABSTRACT**

The project contains 45 pages, 10 pictures, 3 tables, 10 sources, one appendix.

MACHINE LEARNING, STATISTICAL TEXT ANALYSIS, AUTHORSHIP DETECTION, DATA ANALYSIS, NATURAL LANGUAGE PROCESSING

In this graduate work the selection of statistical text features, classification of texts belonging to different authors was made. Dynamics of the classification accuracy according to the text length was investigated.

The problem was solved in an interactive computing environment Jupyter Notebook belonging to the Anaconda distributive that allows Python installation with the necessary packages.

For the methods were used:

- Natural language processing methods;
- Generating statistical text characteristics;
- Machine learning;
- Dimensionality reduction for purpose of visualization.

The conclusion about the optimal text length that is used for model training and testing was made on the basis of the dynamics of the classification accuracy that depends on the length of text fragments.

The thesis project was done solely by the author.

### РЭФЕРАТ

У дыпломнай праце 45 старонак, 10 рысункаў, 2 табліцы, 10 крыніц, адзін дадатак.

МАШЫННАЕ НАВУЧАННЕ, СТАТЫСТЫЧНЫ АНАЛІЗ ТЭКСТА, ВЫЗНАЧЭННЕ АЎТАРСТВА, АНАЛІЗ ДАДЗЕНЫХ, АПРАЦОЎКА НАТУРАЛЬНАЙ МОВЫ

У дыпломнай праце робіцца адбор статыстычных прыкмет тэкста, класіфікацыя тэкстаў, якія належаць розным аўтарам, і даследаванне дынамікі акуратнасці класіфікацыі ў залежнасці ад даўжыні тэкставых фрагментаў.

Задача вырашалася ў праграмным асяроддзі Jupyter Notebook дыстрыбутыва Anaconda, які дазваляе адразу устанавіць Python і неабходныя бібліятэкі.

Для вырашэння пастаўленай задачы выкарыстоўваліся:

- Метады апрацоўкі натуральнай мовы;
- Статыстычныя характарыстыкі тэкстаў;
- Метады машиннага навучання;
- Метады паніжэння размернасці для магчымасці візуалізацыі.

На основе атрымленай дынамікі змены акуратнасці класіфікацыі ў залежнасці ад дліны тэкставых фрагментаў былі зроблемы адпаведныя вывады аб аптымальнай даўжыні тэкстаў, якія выкарыстоўваюцца для абучэння і тэста маделяў.

Дыпломная праца выканана аўтарам самастойна.

## **ВВЕДЕНИЕ**

В последнее время можно наблюдать тенденцию поиска и определения структур, которые характерны для текстов, принадлежащих различным авторам. Печатный текст делает невозможным проведение почерковедческой экспертизы, которая позволила бы определить принадлежность текста конкретному человеку. Поэтому встал вопрос о возможности определения авторского стиля при отсутствии возможности изучения почерка, опираясь на сам текст и его содержание. В связи с этим необходимо было выявить признаки, которые смогут решить задачу определения авторства для нерукописного текста. Были проведены исследования, В которых применялись статистические, формальноколичественные методы, позволяющие выявить характерные для авторов черты. Одними из первых данной задачей занимались Н.А. Морозов [1] и А.А. Марков Помимо выделялись и совершенствовались этого, новые методики определения авторского стиля.

Следует отметить, что все данные исследования ставили перед собой одну цель: установление авторства неизвестных текстов, полагаясь на имеющуюся выборку авторов. Для этого осуществлялось выделение признаков, на основе которых происходило обучение моделей, осуществляющих классификацию текстов. Но не исследовалась динамика точности классификации текстов при изменении длины классифицируемого текста. В рамках данной дипломной работы производился отбор признаков, частотный анализ текстов различной длины и обучение модели для получения описанной выше динамики.

### Г.ЛАВА 1

# ОБРАБОТКА ЕСТЕССТВЕННОГО ЯЗЫКА МЕТОДАМИ МАШИННОГО ОБУЧЕНИЯ

## 1.1 Векторная модель текста

Для решения различных задач, где в качестве исходных данных выступают тексты, требуется преобразование этих текстов к единому виду. Для этого используется векторная модель текста (vector space model), которая ставит текстам в соответствие вектор из общего для всей коллекции текстов векторного пространства. Данная модель используется в решении задач классификации, кластеризации документов, поиска в них.

Текст состоит не только из слов, но и из других элементов, таких как знаки пунктуации, числа, специальные обозначения. Все выше перечисленные единицы текста называются термами. Изначально указывается, каким образом определяется вес терма в тексте. Упорядочение весов всех термов есть вектор, который представляет текст в векторном пространстве. При этом в данном векторе может находиться информация о термах, которые отсутствуют в конкретном документе (тексте), но присутствуют во всем корпусе документов. Таким образом, получаем набор векторов одинакового размера, представляющих все тексты коллекции.

## 1.2 Предобработка исходных данных

Текст, представленный в привычном для нас виде, является непригодным для извлечения из него каких-либо признаков и построения моделей. Он представляет собой строку, содержащую пунктуационные символы, которые не несут

смысловой нагрузки; может состоять из слов, представленных во всевозможных своих формах. В связи с этим мы имеем большое разнообразие отдельных единиц текста и отсутствие единого представления, единой структуры, из-за чего исходный текст не позволяет выявить различные признаки, не подвергаясь предварительной обработке.

#### 1.2.1 Токенизация

К двум важным этапам предварительной обработки текста относятся токенизация и нормализация. Очевидно, что единицей текста является слово. Под токенизацией понимается разбиение текста на отдельные единицы-слова, которые носят название токенов. В результате токенизации исходная строка текста преобразовывается в список слов, из которых она состояла. Токенизация текста выполняется в несколько этапов:

- 1. Приведение исходного текста к нижнему регистру
- 2. Замена пунктуационных символов пробелами
- 3. Объявление слов отдельными токенами

Данный метод предобработки достаточно прост и понятен, но несмотря на это на каждом из этапов могут возникать различные нюансы, а именно:

- 1. Приведение аббревиатур к нижнему регистру может быть спутано с выражением эмоций. Например, «ООО» (Общество с ограниченной ответственностью) после обработки будет эквивалентно «ооо», которое, вероятнее всего, могло так же использоваться и для выражения удивления.
- 2. При замене всех знаков препинания на пробелы могут происходить потери исходных слов. Например, существуют сложные слова, которые пишутся через дефис (красно-синий, где-то). После удаления дефиса

мы получим два слова вместо исходного одного, который несет иную нагрузку. Текст, написанный в свободном стиле, может содержать смайлы, которые играют особую роль при семантическом анализе текста. При удалении знаков препинания пропадает возможность использования этого признака.

3. Имена собственные могут идти в парах. Например, названия отдельных стран, городов, организаций (Саудовская Аравия, Нижний Новгород). При токенизации они будут разделены на два отдельных слова. Но с точки зрения логики было бы полезно рассматривать их как одно целое. Аналогичный вопрос встает с сокращениями. Существуют принятые сокращения, которые при токенизации потеряются, и, более того, внесут ложную информацию в полученный набор.

Помимо этого можно выделить проблему, которая может возникнуть при токенизации текстов определенного типа. А именно, в китайском языке редко пользуются пробелами, вследствие чего текст представляет собой сплошной набор иероглифов, к которому нельзя так просто применить алгоритм токенизации. Помимо этого, в немецком языке часто используется практика образования сложных слов путем конкатенации нескольких. Для обработки подобных текстов требуется сегментация текста на слова.

### 1.2.2 Нормализация

После осуществления токенизации можно переходить к нормализации текста. Нормализация подразумевает под собой приведение слов к нормальной форме, где нормальная форма слова — это каноническая форма слова. Например, для существительных начальная форма слова — это форма единственного числа, стоящая в именительном падеже, для прилагательных — прилагательное мужского рода, единственного числа и стоящее в именительном падеже без предлога.

Данное преобразование не влечет особых потерь, так как конкретная форма слова редко обладает полезной информацией (смысл слова остается тем же). Часто встречаются задачи с небольшим объемом исходных данных, в связи с чем желательно уменьшить число признаков. Приводя же слова к начальной форме, мы уменьшаем количество уникальных слов.

При осуществлении нормализации можно выделить два подхода: стэмминг и Стэмминг лемматизация. занимается поиском основы слова, учитывая морфологию исходного слова. Таким образом, находя общую грамматических форм слова основу, отсекая суффиксы и окончания, стэмминг осуществляет морфологический разбор слова. Алгоритм стэмминга – это конкретный способ решения задачи поиска основы слов, а стэммер – конкретная реализация. Стэмминг обладает следующим недостатком: различные формы слова могут иметь разные основы. Поэтому после стэмминга результаты обработки этих слов будут различны. В связи с этим применяется иная техника – лемматизация. В отличие от стэмминга она работает на основе словаря, где и хранятся данные о словах и их начальных формах. Это позволяет избегать ошибок, описанных выше. В случае если слова нет в словаре, то происходит построение гипотезы о способе изменения слова.

## 1.3 Понижение размерности

Описанные ниже алгоритмы используются для понижения размерности, что позволяет эффективно визуализировать многомерные компоненты. Из многомерной переменной мы пытаемся получить новую переменную, которая будет лежать в двух- или трехмерном пространстве и сохранять закономерности исходной переменной.

### 1.3.1 Алгоритм t-SNE

Название алгоритма есть сокращение с полного варианта t-distributed stochastic neighbor embedding. На русский язык можно его попробовать перевести как «t-распределенное внедрение соседей». Данный метод относится к методам множественного обучения признаков. Классическое представление алгоритма SNE было изложено в 2002 году Ровейсом и Джеффри Хинтоном, а расширение t-SNE представлено в 2008 году Джеффри Хинтоном и Лоуренсом ван дер Маатеном.

### 1.3.1.1 Математическое описание алгоритма

Будем строить биективное отображение, то есть каждая точка будет представлять один экземпляр (строку) исходного набора текстов. Чтобы визуализация могла отображать разделение на классы, мы хотим, чтобы точки, принадлежащие одному классу, располагались рядом, что и покажет формула (1.1):

$$p_{j|i} = \frac{\exp\left(\frac{-|x_i - x_j|^2}{2\sigma_i^2}\right)}{\sum_{k \neq j}^{n} \exp\left(\frac{-|x_i - x_k|^2}{2\sigma_i^2}\right)}$$
(1.1)

Формула (1.1), основанная на гауссовском распределении вокруг точки  $x_i$  с заданной дисперсией  $\sigma_i$  определяет условное сходство двух точек. Дисперсии вычисляются таким образом, чтобы точки, расположенные в областях с малой плотностью, имели большую дисперсию, чем точки, расположенные в областях с большой плотностью. Для этого используется оценка перплексии (1.2). Обычно эта оценка используется для сравнения вероятностных моделей, при этом низкое значение перплексии означает, что распределение вероятностей (закон, описывающий область значений случайной величины и вероятности их исхода) хорошо работают на этапе предсказания. В нашем случае можно ее

интерпретировать как сглаженную оценку эффективного количества «соседей» для конкретной точки.

$$Perp(P_i) = 2^{H(P_i)} \tag{1.2}$$

Для этого вычисляется энтропия Шеннона в битах (1.3):

$$H(P_i) = -\sum_{j} p_{j|i} \log_2 p_{j|i}$$
(1.3)

Сходство двух точек будет определяться через формулу (1.4) условного сходства двух точек:

$$p_{j|i} = \frac{p_{j|i} + p_{i|j}}{2N} \tag{1.4}$$

На основе данной формулы вычисляется матрица сходства для исходных данных. Данная матрица является постоянной.

В отличие от матрицы сходства исходных данных, матрица сходства для отображения зависит от точек отображения. Близость этих двух матриц будет нам доказывать, что похожие исходные точки отображаются в также похожие точки. При определении матрицы сходства для точек отображения (1.5) вместо гауссовского распределения используется распределение Стьюдента с одной степенью свободы или распределение Коши. Причина данной замены заключается в следующем: расстояние между парой точек в пространстве отображения, которые соответствуют паре среднеудаленных точек в исходном пространстве, должно быть намного больше, чем расстояние, которое можно получить при распределения. помощи гауссовского Α метод пытается воспроизвести одинаковые расстояния в обоих пространствах, И возникает проблема «скученности». Благодаря тяжелым «хвостам» распределения Стьюдента дисбаланс в распределении расстояний для соседей точек скомпенсирован (1.5).

$$q_{ij} = \sum_{k \neq i}^{f(|x_i - x_j|)} f(|x_i - x_k|)$$
, где  $f(y) = \frac{1}{1 + y^2}$  (1.5)

В ходе алгоритма происходит минимизация расстояния между матрицами сходства, осуществляемая за счет минимизации расстояния Кульбака-Лейблера между распределениями  $p_{ij}$  и  $q_{ij}$  (1.6):

$$KL(P||Q) = \sum_{i,j} p_{ij} log_{q_{ij}}^{p_{ij}}$$

$$(1.6)$$

Для минимизации данного расстояния применяется градиентный спуск. Под градиентом можно понимать сумму всех сил, которые приложены к точке отображения і. В формуле (1.7)  $u_{ij}$  обозначает единичный вектор, который идет от  $y_i$  к  $y_i$ .

$$\frac{\partial \text{ KL(P||Q)}}{\partial y_i} = 4 \sum_{j} (p_{ij} - q_{ij}) g(|x_i - x_j|) u_{ij}$$
, где  $g(y) = \frac{y}{1 + y^2}$  (1.7)

### 1.3.1.1 Физическая аналогия

Данный алгоритм можно интерпретировать с физической точки зрения следующим образом. Считаем, что все точки полученного отображения соединены между собой пружинами, обладающими разными жесткостями. Жесткости пружин зависят от величины  $p_{ij}$  -  $q_{ij}$ , которая обозначает разность сходства пары точек исходных данных и сходства пары точек отображения. Точки отображения притягиваются, если расстояние между точками данных малое, а между точками отображения большое. Отталкиваются в противном случае. Градиент является результирующей силой, которая действует на точку в пространстве отображения. Мы отпускаем систему, и она будет изменяться до

достижения равновесного состояния. Отображение на данном шаге и будет являться тем, к которому мы и стремимся.

#### 1.3.2 Метод главных компонент

Метод главных компонент (PCA – Principal Components Analysis) был изобретен Пирсоном в 1901 году. В нем происходит вычисление собственных векторов и собственных значений ковариационной матрицы исходных данных или находится сингулярное разложение матрицы.

### 1.3.2.1 Описание метода

Изначально мы работаем с матрицей X, которая имеет размерность  $I \times J$ , I обозначает количество образцов, J — количество независимых переменных. Матрица X раскладывается в произведение двух матриц T и P (1.8):

$$X = TP^t + E = \sum_{i}^{A} t_i p_i^t + E \tag{1.8}$$

Матрица Т представляет собой переменные, являющиеся комбинацией исходных переменных х (1.9), а именно:

$$t_i = p_{i_1} x_1 + \dots + p_{i_I} x_I$$
, где  $i = 1, \dots A$  (1.9)

Матрица Т, имеющая размерность  $I \times A$ , называется матрицей счетов. Эта матрица показывает структуру исходных данных. Она заключает в себе проекции исходных образцов (векторы  $x_1, ... x_I$  размерности J) на A-мерное подпространство главных компонент. Строки матрицы представляют собой координаты образцов в новой системе координат, а столбцы — проекции всех образцов на новую координатную ось. Если отобразить матрицу счетов на графике, то близкие точки будут являться схожими между собой, то есть обладать положительной

корреляцией; диаметрально противоположные точки будут обладать отрицательной корреляцией; а точки, расположенные под прямым углом, будут иметь нулевую корреляцию. Собственные значения матрицы  $T^tT$  показывают важность соответствующих им главных компонент.

Матрица Р перехода из исходного пространства (J-мерное) в пространство главных компонент (А-мерное), имеющая размерность  $J \times A$ , есть матрица нагрузок. График матрицы показывает зависимость переменных между собой. Строка матрицы нагрузок представляет собой проекцию переменных  $x_1 \dots x_J$  на соответствующую ось главных компонент. Столбец же матрицы есть проекция соответствующей переменной  $x_j$  на новую систему координат. Для матрицы Р верно (1.10):

$$P^t P = I (1.10)$$

Матрица E размерности  $I \times J$  есть матрица остатков. Визуально разложение можно представить следующим образом:

$$\mathbf{X}$$
 =  $\mathbf{t}_1$  + ... +  $\mathbf{t}_A$  +  $\mathbf{E}$ 

Рисунок 1.1 Сингулярное разложение

Главными компонентами будут являться новые переменные  $t_i$ . Количество главных компонент равно числу строк в матрице T и числу столбцов в матрице P, и равно A. Так как мы понижаем размерность, это число будет меньше количества образцов I и числа переменных J.

Понижая размерность, мы хотим выделить признаки, которые будут независимы друг от друга. Поэтому данный метод обладает свойством ортогональности. Если мы увеличиваем число главных компонент, то к матрице счетов Т добавляется столбец, соответствующий новой компоненте (новому направлению), при этом имеющаяся часть матрицы Т не изменяется. Матрица нагрузок Р аналогично матрице Т не перестраивается.

## 1.3.2.2 Сингулярное разложение

При разложении по сингулярным значениям (1.11) матрица X разлагается в произведение трех матриц:

$$X = USV^T \tag{1.11}$$

Матрица U образована ортонормированными собственными векторами  $u_i$  матрицы  $XX^t$ , которые соответствуют собственным значениям  $\lambda_i$ , то есть  $XX^tu_i=\lambda_iu_i$ . Матрица V образована ортонормированными собственными векторами  $v_i$  матрицы  $X^tX$ , а именно:  $X^tXv_i=\lambda_iv_i$ . Матрица S — положительно определенная диагональная матрица. Элементы матрицы S — сингулярные значения  $\sigma_1\geq\cdots\geq\sigma_i\geq0$ , где  $\sigma_i=\sqrt{\lambda_i}$ .

Метод главных компонент и сингулярное разложение связаны между собой согласно формуле (1.12):

$$T = US \text{ if } P = V \tag{1.12}$$

# ГЛАВА 2 СТАТИСТИЧЕСКИЙ АНАЛИЗ ТЕКСТА

### 2.1 Извлечение признаков из предобработанного текста

Для решения определенной задачи исходный текст либо предобработанный не является достаточным для обучения модели. Для более глубокого изучения исходного текста необходимо обратиться к методам, которые позволят выявить некоторые признаки, характерные для имеющихся данных.

### 2.1.1 Мешок слов

Название модели текста «мешок слов» есть перевод с английского «bag-of-words». Она была предложена в 1975 году Солтоном. Данная модель текста представляет собой суммативное единство слов, составляющих исходный текст. Единицы «мешка слов» — слова, каждое из которых имеет атрибут, а именно: количество встреч данного слова в тексте.

Важными особенностями данной модели является отсутствие учета порядка слов в документе и морфологических форм представления слов.

Описать работу «мешка слов» можно следующим образом:

- 1. Дана выборка, состоящая из n различных слов:  $w_1, \, \dots \, , \, w_n$ .
- 2. Текст кодируется при помощи п признаков, где і-тый признак обозначает долю вхождений слова w<sub>i</sub> среди всех вхождений слов в текст.

Обычно «мешок слов» применяется уже к предобработанному тексту, из которого были исключены все стоп-слова. Данные слова не несут в себе

смысловую нагрузку, поэтому часто доля их вхождений в текст не учитывается. Аналогично предлагается не учитывать в данной модели редкие слова, так как они, имея слишком малый вес, не смогут внести вклад в построенную модель.

### 2.1.2 Мешок термов

«Мешок термов» считается обобщением модели «мешок слов», названия происходит от английского «bag-of-terms». В отличие от мешка слов, его элементом является терм, который характеризуется частотой встречаемости в тексте. Под термом понимается символьное выражение объекта формальной модели (системы, языка). В качестве них могут использоваться всевозможные символьные выражения текста.

### 2.1.3 TF-IDF метод векторизации текстовых данных

Метод TF-IDF Vectorizer используется с целью назначения весов словам исходных текстов. Он основан на двух предположениях, а именно:

- 1. Слова, которые часто встречаются в тексте, важны для данного текста.
- 2. С другой же стороны, слово, которое редко встречается в других текстах (документах), оно важно для текущего. Иными словами, можно сказать, что данное слово может быть признаком, по которому можно идентифицировать данный документ среди остальных.

Данные предположения можно представить следующими формулами:

$$TDF(d, w) = n_{d\omega} \tag{2.1}$$

где  $n_{dw}$  обозначает долю вхождений слова w в текст (документ) d.

$$IDF(d, w) = log \frac{l}{n_{\omega}}$$
 (2.2)

Где 1 — общее количество текстов, а  $n_{\rm w}$  обозначает количество текстов, в которых встречалось слово w. Можно проанализировать, что если данное слово можно встретить в каждом тексте, то есть  $l=n_{\omega}$ , то тогда значение признака IDF(d,w)=log1=0. Это может свидетельствовать о том, что если слово встречается очень часто на всем корпусе документов, то оно вряд ли важно.

Тогда общая формула мне метода TF-IDF примет вид:

$$TDF - IDF(d, w) = n_{d\omega} \log \frac{l}{n_{\omega}}$$
 (2.3)

На основе данных формул можно сделать вывод, что результат выражения будет максимальным, если слово много раз встречается в тексте d, и число вхождений его в остальные документы минимально.

## 2.1.4 N-граммы

Достоинство n-грамм заключаем в том, что они позволяют учитывать порядок слов, в отличие от «мешка слов». Помимо этого, использование n-грамм расширяет признаковое пространство благодаря учету словосочетаний. Поэтому простые модели могут служить для поиска более сложных закономерностей по сравнению с «мешком слов».

N-граммы можно разделить на буквенные и словесные. Под словесными п-граммами понимаются наборы из n идущих подряд токенов. Среди n-грамм можно выделить:

- Униграммы: наборы, состоящие из одного токена;
- Биграммы: наборы, состоящие из двух подряд идущих токенов;

• Триграммы: наборы, состоящие из трех подряд идущих токенов и т.д.

Рассмотрим их построение на примере фразы «обработка естественного языка методами машинного обучения». Тогда:

- Униграммы: обработка, естественного, языка, методами, машинного, обучения;
- Биграммы: обработка естественного, естественного языка, языка методами, методами машинного, машинного обучения, и т.д.

В зависимости от выбора параметра n можно получить как огромное число признаков, так и свести всё к тому, что каждый текст будет сам по себе являться отдельным признаком (что будет означать подгонку под исходную выборку).

В буквенных п-граммах в качестве токенов рассматриваются буквы. Остальной процесс выделения униграмм, биграмм и т.д. аналогичен словесным биграммам.

Также можно обратить внимание на расширенную версию п-грамм, а именно k-skip-n-граммы. Это наборы из n токенов, где между соседними токенами должно быть не более чем k токенов. В качестве примера можно рассмотреть 1-skip-2-граммы. Тогда для приведенной ранее фразы получим следующий результат: обработка языка, естественного методами, языка машинного, методами обучения.

### 2.1.5 Хэширование

Перед рассмотрением данного метода следует привести определение хэшфункции. Поэтому h(x) - xэш-функция, которая принимает  $2^n$  возможных значений. Данный метод подразумевает замену всех слов x на ux хэши h(x) и использование этих хэшей как токенов. После этого можно применять описанные ранее методы, как «мешок слов», TF-IDF векторизатор u т.д.

Применение хэширования обладает рядом преимуществ. В первую очередь, это упрощение хранения модели. В данном случае значение хэша можно отождествить с индексом конкретного слова (токена). Например, при использовании «мешка слов» необходимо сохранять соответствие между словами и признаками. Также при использовании данного метода происходит сокращение числа признаков (возможность объединять исходные слова с одинаковыми хэшами).

#### 2.1.6 Стоп-слова

Стоп-слова — это слова, которые встречаются в большом объеме текстов и не несут особой смысловой нагрузки. Обычно к ним относят междометия, предлоги, частицы, некоторые местоимения, прилагательные и другие части речи. Поэтому чаще всего предобработанный текст очищают от стоп-слов.

В качестве примеров можно отметить следующие группы:

- Междометия: ах, ух, ну, уж, ой;
- Местоимения: я, мы, мой, вы, ваш;
- Вводные конструкции: скажем, допустим, например, в общем, на самом деле;
- Обобщения и неточные определения: всего, примерно, около, где-то, порядка и другие.

## 2.2 Выделение характеристик

Очевидно, что из текста можно выявить бесконечное число характеристик. Можно начать с подсчета числа предложений и дойти до количества запятых на абзац текста. Главный вопрос заключается в том, какие характеристики можно будет считать релевантными. Поэтому в данном вопросе можно опираться на

логику, и исходя из нее совершать отбор признаков. Признаки, извлекаемые из текста, можно, как минимум, разделить на три группы:

- 1. Признаки, основанные на знаках препинания;
- 2. Признаки, базирующиеся на словах как единицах текста;
- 3. Признаки, где единицей является отдельный символ буква.

Первая группа может исследовать распределение знаков препинания по тексту, взаимосвязь между знаками препинания и общим числом символов, слов в тексте.

Вторая группа дает возможность получить признаки, основанные на длинах слов, частотах употребления различных частей речи, длинах предложений, важности слов, «стоп-словах» либо уникальных словах.

Третья группа, например, позволит исследовать отношение между прописными и строчными буквами, строить биграммы и др.

### Г.ЛАВА 3

# ЗАДАЧА ОПРЕДЕЛЕНИЯ АВТОРСТВА

### 3.1 Постановка задачи

Задача определения авторства состоит в следующем: перед нами есть набор текстов (фрагменты произведений авторов) на русском или другом языке. Имеется неидентифицированный текст, но известно, что он принадлежит кому-то из перечисленных выше авторов. Необходимо определить автора данного текста. Будем исследовать динамику точности классификации текстов при изменении длины классифицируемого текста. Для решения данной задачи будут использованы методы, описанные выше, а также будут отбираться признаки и использоваться в сочетании с другими подходами для выявления сочетания, которое позволит получить наилучший результат.

## 3.2 Используемые технологии

Для реализации поставленной задачи в качестве программной среды использовался дистрибутив Anaconda, который позволяет сразу установить Python и необходимые библиотеки. Средой для выполнения задачи был выбран Jupyter Notebook. Основные пакеты, используемые в ходе исследования:

- Numpy, pandas, scikit-learn для анализа и обучения моделей;
- Matplotlib, seaborn для визуализации.

## 3.3 Построение тренировочного и тестового множеств

В сборе и построении данных, используемых для обучения, можно выделить следующие основные этапы:

- 1. Сбор текстов авторов;
- 2. Загрузка текстов;
- 3. Разбиение текстов на равные участки, предварительная обработка;
- 4. Конструирование тренировочного и текстового наборов.

Сначала производился сбор произведений для отобранных авторов и загрузка их из электронной библиотеки. Затем из текстов удалялись служебные символы, препятствующие дальнейшему анализу.

Все произведения для отдельного автора объединялись в единую строку, из которой затем нарезались строки определенной длины. Динамика точности классификации текстов исследовалась в зависимости из изменения длины данных строк. Во избежание пересечения тестового и тренировочного наборов, отрезки для тренировочного набора отбирались из первой половины строки, объединяющей произведения конкретного авторов, а тестового — из второй соответственно.

## 3.4 Решение задачи

Как говорилось ранее, изначально весь текст собирался по группам для каждого автора и подвергался предобработке, где удалялись служебные символы, не несущие никакой смысловой нагрузки. Так как нам важно исследовать динамику точности классификации текстов в зависимости от изменения длины текстов, генерировались множественные тестовые и тренировочные выборки для

различных длин текста L, где L принимало значения из множества {20000, 10000, 5000, 1000, 500, 200}.

Для обучения модели необходима векторизация текста. Для полученных текстов рассчитываются различные характеристики, которые затем представляются в качестве единого вектора. Длины векторов для всех элементов получаются равными между собой. Очевидно, что признаки могут быть разными, иметь разные единицы измерения и лежать в разных интервалах. Поэтому при построении вектора необходимо производить нормировку координат. Нормировка координат производится глобально и основывается на векторах, рассчитанных для тренировочной выборки. Производилось центрирование выборки: сдвигалась выборка так, чтобы средние значения признаков были равны нулю. Дисперсия чувствительна к масштабированию, то есть сильно зависит от порядков случайной величины. Поэтому рекомендуется стандартизировать признаки, если их единицы измерения сильно отличаются своими порядками. Стандартизация обычно требуется для обучающих алгоритмов: качество обучения заметно снижается, если отдельный признак не обладает нормальным распределением. Многие алгоритмы (например, метод опорных векторов с RBF ядром, L1 или L2 регуляризация линейных моделей) ожидают на входе центрированные признаки с центром в нуле и с одинаковым распределением. Если какой-то признак имеет дисперсию, которая на несколько порядков больше, чем дисперсия других признаков, она может доминировать над целевой функцией. В связи с этим модель не сможет обучаться на других признаках так, как предполагалось.

Для решения задачи определения авторства необходимо из огромного числа характеристик (признаков) выбрать те, которые будут использоваться для анализа текстов. На определенных этапах отбора признаков текст приводился к нижнему регистру, осуществлялась токенизация и лемматизация.

Вставал вопрос, какие признаки могли бы охарактеризовать стиль и манеру письма отдельного автора. В качестве основных были отобраны следующие статистические характеристики текста:

- Отношение количества заглавных букв к количеству строчных букв;
- Распределение различных знаков препинания по тексту. Производился подсчет вхождения каждого символа из заданного набора в текстовые строки.
- Распределение длин предложений. В качестве длины предложения бралось количество слов в нем. Для этого изначально проводилась токенизация и очистка предложений от знаков препинания. Полученные длины предложений подвергались нормализации;
- Распределение длин слов. Определялись длины слов, входящих в текст, и изучалось распределение подсчитанных длин.
- Определялась водность текста. Для этого рассчитывалось отношение количества слов после очистки текста от стоп-слов к количеству слов в исходном тексте. Разность единицы и данной величины и определяет водность текста. Например, текст, состоящий полностью из стоп-слов, будет иметь водность, равную единице; а текст, не содержащий стопслов, будет определяться нулевым значением водности. Считается, что неестественные тексты могут обладать повышенной водностью.
- Разнообразие речи. Под ней понимается отношение количества уникальных слов к общему количеству слов в тексте.
- Распределение частей речи в тексте. Для этого были отобраны наиболее часто встречающиеся части речи, осуществлялась нормализация текста, после чего производился подсчет встреч каждой части речи в предобработанном тексте.

Помимо этих признаков, рассчитывались буквенные биграммы для корпусов текстов. На основе тренировочной выборки определялись наиболее часто встречающиеся биграммы, их частоты использовались в качестве признаков. При тестировании подсчитывались частоты отобранных на этапе обучения наиболее часто встречающихся биграмм.

## 3.5 Критерии качества в задачах классификации

Для оценки качества моделей в задачах машинного обучения используются различные метрики. Для начала следует обратиться к понятию матрицы ошибок (confusion matrix). Выделяется два вида ошибок классификации: False Negative (FN) и False Positive (FP). Под ошибкой False Negative, или ошибкой второго рода, понимается неотнесение объекта к классу, которому он на самом деле принадлежит. Ошибка False Positive, или же ошибка первого рода, предполагает отнесение к данному классу объекта, который в реальности ему не принадлежит. Обозначим уответы модели, а у – истинная метка класса. Ошибки классификации для двух классов можно представить в следующей таблице:

	y = 1	y = 0
$\hat{y}=1$	True Positive (TP)	False Positive (FP)
$\hat{y} = 0$	False Negative (FN)	True Negative (TN)

Таблица 3.1 Ошибки классификации

Метрика ассигасу, или же аккуратность обозначает долю правильных ответов алгоритма и рассчитывается по формуле (3.1):

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{3.1}$$

Определим метрики точность (precision) по формуле (3.2) и полноту (recall) по формуле (3.3):

$$Precision = \frac{TP}{TP + FP} \tag{3.2}$$

$$Recall = \frac{TP}{TP + FN} \tag{3.3}$$

Точность определяется как отношение объектов, верно классифицированных как принадлежащих классу, по отношению ко всем объектам. Полнота показывает отношение верно положительно классифицированных объектов ко всем объектам положительного класса, которые были найдены алгоритмом. Последние две метрики не зависят от соотношения классов, что является преимуществом по сравнению с ассигасу, которая будет показывать нехорошие результаты для несбалансированных выборок.

Для определения точности классификации удобно иметь единую метрику, поэтому применяются различные комбинации метрик precision и recall. Чтобы можно было определять вес каждой метрики в итоговой, добавляется параметр  $\beta$ , который будет отвечать за вес метрики precision. Логично, что при метриках, близких к единице, итоговая метрика должна так же достигать максимума, и, наоборот, при малом значении какой-либо из метрик  $F_{\beta}$  должна стремиться к нулевому значению. В качестве примера можно рассмотреть F-меру (3.4) как среднее гармоническое этих двух метрик:

$$F_{\beta} = (1 + \beta^2) \frac{precision * recall}{(\beta^2 * precision) + recall}$$
(3.4)

Еще одной возможностью оценить качество модели является площадь (Area Under Curve) под кривой ошибок (Receiver Operating Characteristic curve) – ROC AUC. Кривая ошибок изображается в координатах True Positive Rate (TPR) и FPR (False Positive Rate), значения в точках которой вычисляются по формулам (3.5) и (3.6):

$$TPR = \frac{TP}{TP + FN} \tag{3.5}$$

$$FPR = \frac{FP}{FP + TN} \tag{3.6}$$

True Positive Rate эквивалентно полноте, а False Positive Rate отображает, какая доля объектов, не принадлежащих данному классу, была предсказана неверно. Площадь под полученной кривой будет являться метрикой, показывающей качество алгоритма. Логично, что наилучшим случаем являются значения TPR = 1 и FPR = 0, при котором площадь под кривой будет равна единице.

Следует обратить внимание на логистическую функцию потерь (logistic loss). Здесь используются ранее введенные обозначения:  $\hat{y}$  — ответы модели, у — истинная метка класса, — l определяет размер выборки. Она задается формулой (3.7):

$$logloss = -\frac{1}{l} \sum_{i=1}^{l} (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$
(3.7)

Данная метрика может использоваться в задачах максимизации точности предсказания путем штрафования за неверно предсказанные метки классов.

## 3.6 Классификация текстов

### 3.6.1 Матрица ошибок

Для оценки полученных классификаторов строилась матрица ошибок, таблице сопряженности. Представляя которая основана на распределение двух переменных, таблица сопряженности предназначена для исследования связи между ними. По одной оси располагаются оригинальные метки классов, по другой – предсказанные. Главная диагональ матрицы показывает количество верно предсказанных значений для каждого из классов. Неверно предсказанные элементы будут располагаться вне главной диагонали. Соответственно, диагональных сумма элементов есть все верно классифицированные объекты. Тогда общую точность классификации можно выразить как отношение суммы диагональных элементов (  $d_i$  ) к общему количеству элементов (N) формулой (3.8).

$$Overall\ Accuracy = \frac{\sum d_i}{N}$$
 (3.8)

Аналогичным образом можно посчитать точность определения реального (рассчитанного) класса: разделив число верно классифицированных объектов на общее количество объектов этого класса. Формула для первого класса будет иметь вид (3.9):

$$Producer's Accuracy = \frac{d_1}{a_{1i}}$$
 (3.9)

Матрица ошибок строилась для более глубокого анализа результатов классификации. На рисунках (3.1) – (3.6) приведены данные матрицы, построенные на основе классификаторов, обученных на текстах различной длины.

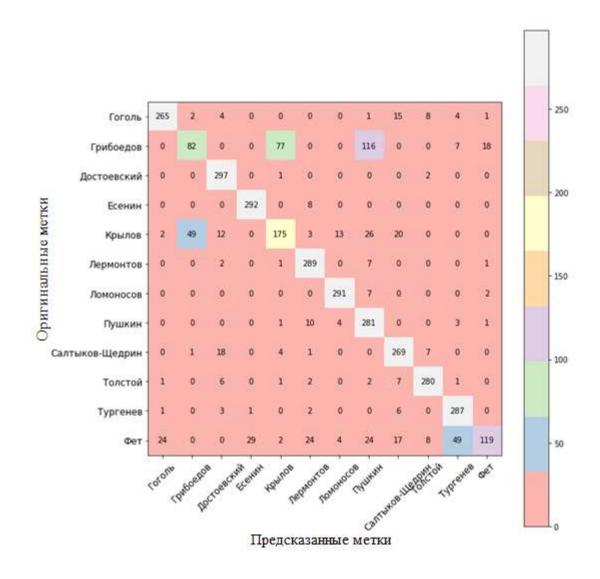


Рисунок 3.1 Матрица ошибок на текстах длиной 20000 символов

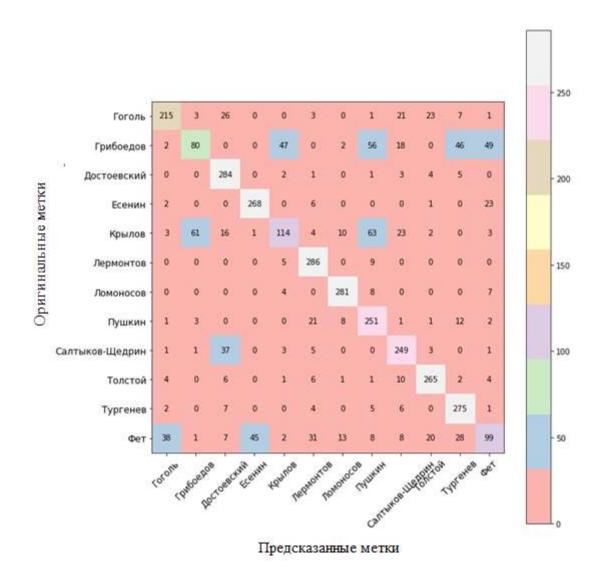


Рисунок 3.2 Матрица ошибок на текстах длиной 10000 символов

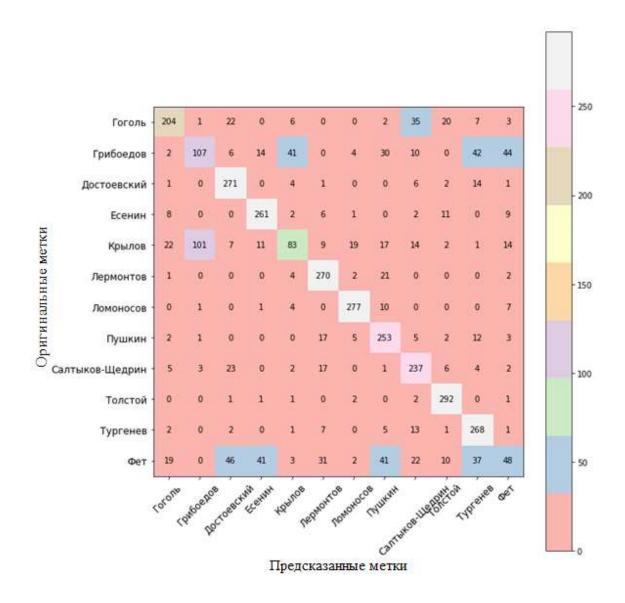


Рисунок 3.3 Матрица ошибок на текстах длиной 5000 символов

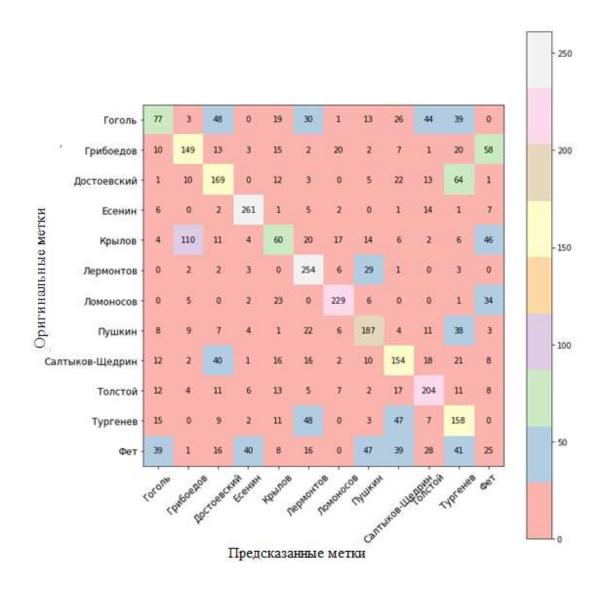


Рисунок 3.4 Матрица ошибок на текстах длиной 1000 символов

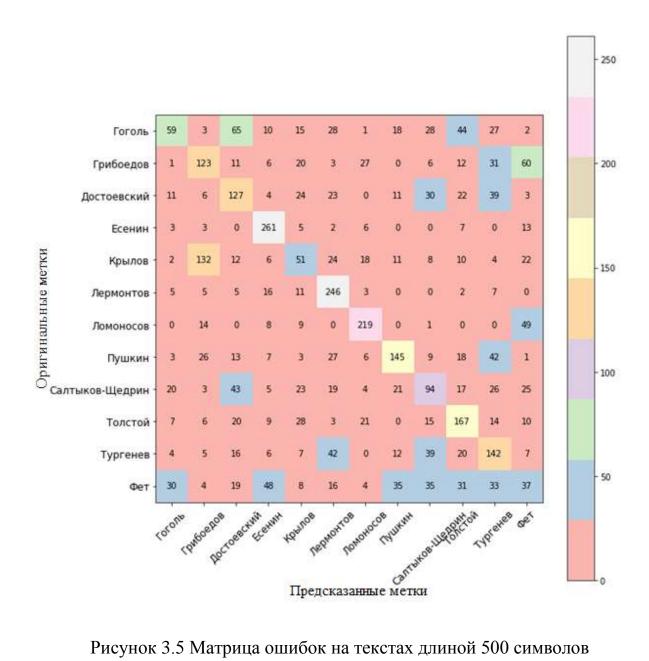


Рисунок 3.5 Матрица ошибок на текстах длиной 500 символов

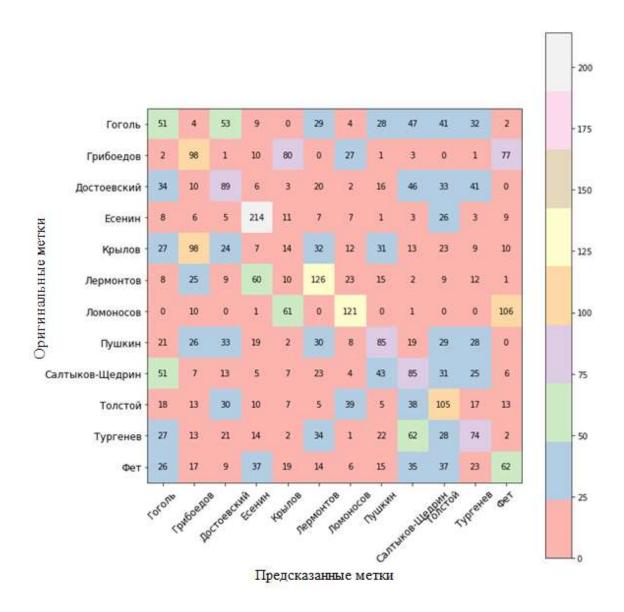


Рисунок 3.6 Матрица ошибок на текстах длиной 200 символов

# 3.6.2 Динамика точности классификации

Классификация текстов осуществлялась на основе векторов, построенных на статистических признаках, извлеченных из текстов, и посчитанных частотах биграмм. Для классификации использовалась различные алгоритмы: логистическая регрессия, метод k-ближайших соседей, метод опорных векторов с различными ядрами, стохастический градиентный спуск, стохастический градиентный спуск на векторизованных текстах методом TF-IDF. В качестве

основной метрики использовалась «аккуратность» (ассигасу). Наилучшие результаты показали логистическая регрессия и метод опорных векторов. Для достижения наилучшей точности происходил подбор параметров оптимальных параметров с использованием GridSearchCV.

Результаты работы алгоритмов представлены в табл.3.2.

	L = 20000	L = 10000	L = 5000	L = 1000	L = 500	L = 200
Logistic	0.813	0.751	0.725	0.536	0.467	0.317
Regression						
KNN	0.774	0.718	0.665	0.332	0.236	0.133
SVM	0.776	0.743	0.702	0.503	0.475	0.280
Tfidf +	0.709	0.721	0.733	0.606	0.534	0.332
SGDClassifier						
SGDClassifier	0.753	0.721	0.664	0.492	0.395	0.255

Таблица 3.2 Точность классификации (метрика accuracy)

Из данных, представленных в таблице, можно сделать вывод, что качество классификации значительно снижается при использовании для обучения текстов длины менее 5000 символов. При этом качество классификации при увеличении длины текстов с 5000 до 20000 символов повышается незначительно. Векторизация исходных текстов - достаточно трудоемкая задача. Поэтому при отсутствии достаточных ресурсов, можно не использовать тексты наибольшей длины, не ощущая при этом особых потерь качества.

## 3.7 Визуализация

Мы получили длинные векторы, содержащие в себе множество признаков, вычисленных на созданной выборке. Было бы интересно посмотреть на эти данные и попытаться найти некоторые закономерности между ними. Если бы визуально полученные векторы разделялись на кластеры, соответствующие разным авторам, это бы означало, что признаки были выбрано достаточно успешно. Но возникает проблема, как можно визуализировать эти векторы, учитывая, что они имеют большую размерность. Данную проблему решают методы понижения размерности пространства. Следует отметить, что данные методы чувствительны к выбору единиц измерения. Величины, имеющие разброс, на порядок больший по сравнению с остальными признаками, вносили бы наибольший вклад при понижении размерности. Поэтому сначала производится нормировка выборки, после этого же применяется один из методов понижения размерности.

В рамках данной работы методы главных компонент и t-SNE применялись для получения наглядного представления полученных признаков. Векторы, которые состоят из статистических признаков и частот биграмм отображались на плоскости. Цвета точкам присваивались исходя из оригинальных меток (оригинальной принадлежности какому-либо автору), каждая точка обозначает конкретный текст, а именно вектор, который был получен после отображения длинного вектора (признаков текста) в пространство меньшей размерности. В связи с этим можно было наблюдать полученные распределения текстов.

Для текстов длины 20000 были получены графики, изображенные на рис.3.7, рис.3.8, рис.3.9. Каждая точка обладает цветом (номер), что означает ее принадлежность конкретному автору. Видно, что понижение размерности при помощи метода t-SNE для текстов длины 20000 помогает явно выделить разные

классы авторов. Первые два графика (рис.3.7, рис.3.8) были получены при понижении размерности до пространства  $\mathcal{R}^2$ , а третий график (рис.3.9) изображает компоненты в пространстве  $\mathcal{R}^3$ .

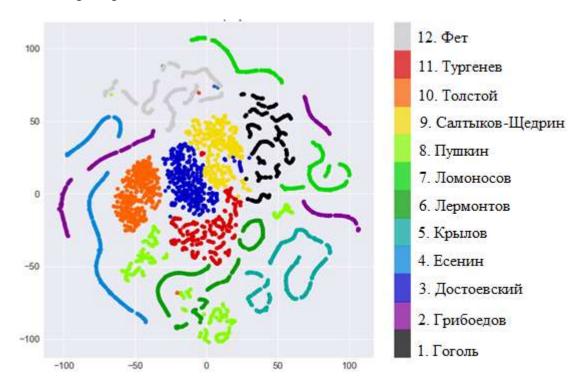


Рисунок 3.7 Алгоритм t-SNE для текстов длины 20000 символов тренировочного набора. Библиотека matplotlib.

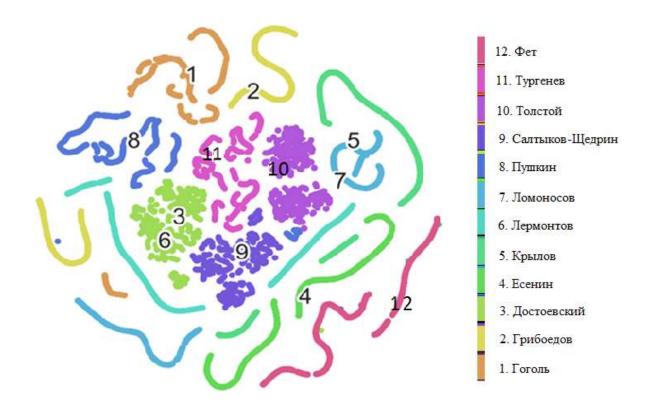


Рисунок 3.8 Алгоритм t-SNE для текстов длины 20000 символов тренировочного набора. Библиотека seaborn.

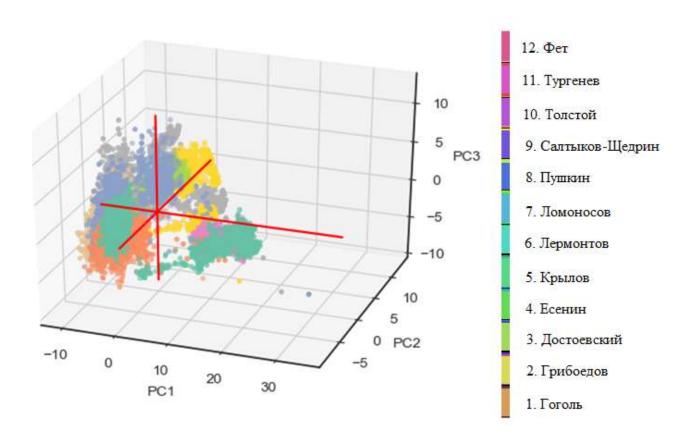


Рисунок 3.9 Метод главных компонент для текстов длины 20000 символов тренировочного набора. Библиотека seaborn.

### **ЗАКЛЮЧЕНИЕ**

В рамках дипломной работы решалась задача классификации текстов известных авторов. Помимо этого исследовалась зависимость точности классификации от длины текста, которая не изучалась в работах, описывающих проблематику и первые подходы к решению исходной задачи. Для решения данной проблемы были поставлены и решены задачи, включающие в себя:

- Подготовку текста к анализу;
- Исследование текста как последовательности символов;
- Поиск и выделение закономерностей, способных охарактеризовать текст;
- Получение векторного представления текста;
- Визуальное представление полученных признаков;
- Построение и обучение моделей, осуществляющих классификацию текстов;
- Сравнение построенных моделей для разных корпусов текстов и определение пороговых длин текстов, при которых точность классификации изменяется незначительно либо, наоборот, претерпевает спад или подъем точности классификации.

# СПИСОК ИСПОЛЬЗОВАННОЙ ЛИТЕРАТУРЫ

- Морозов Н.А. Лингвистические спектры: средство для отличения плагиатов от истинных произведений того или иного известного автора. Стилеметрический этюд. // Известия отд. русского языка и словестности Имп. Акад. наук, Т.ХХ, кн.4, 1915.
- 2. Марков А.А. Пример статистического исследования над текстом "Евгения Онегина", иллюстрирующий связь испытаний в цепь. // Известия Имп. Акад. наук, серия VI, Т.Х, N3, 1913, с.153.
- 3. Марков А.А. Об одном применении статистического метода. // Известия Имп. Акад. наук, серия VI, Т.Х, N4, 1916, с.239.
- 4. T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical Learning.

  Data Mining, Inference, and Prediction. 2nd Edition. Springer, 2013.
- 5. Мартыненко Г. Я . Основы стилеметрии . Л.: Изд-во ЛГУ, 1988 . 176 с.
- 6. Хмелев Д.В. Распознавание автора текста с использованием цепей А.А. Маркова //Вести. МГУ. Сер. 9. Филология. 2000. №2. С. 115-126.
- 7. Фоменко В.П., Фоменко Т.Г. Авторский инвариант русских литературных текстов // Методы количественного анализа текстов нарративных источников. М.: Ин-т истории СССР, 1983. С. 86-109.
- 8. Шитиков В. К., Мастицкий С. Э. Классификация, регрессия и другие алгоритмы Data Mining с использованием R. 2017.
- 9. От Нестора до Фонвизина. Новые методы определения авторства. М.: Издат. группа «Прогресс», 1994.

## приложение А.

# Код программы

Ниже приведен код алгоритма на Python, созданный в среде Jupyter Notebook.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import os
import glob
import nltk, re
import collections
import random
import operator
from functools import reduce
import time
from collections import Counter
from itertools import islice
from sklearn import preprocessing
                                                                        In [2]:
from nltk import word tokenize
from nltk.stem.wordnet import WordNetLemmatizer
from nltk.corpus import stopwords
import pymorphy2
morph = pymorphy2.MorphAnalyzer()
Import data
                                                                       In [11]:
FILE PATH = "E:/0UNI/7 сем/Диплом/Выборка/Русский/"
                                                                       In [12]:
folders = os.listdir(path=FILE PATH)
folders
```

```
Out[12]:
['Гоголь',
 'Грибоедов',
 'Достоевский',
 'Есенин',
 'Крылов',
 'Лермонтов',
 'Ломоносов',
 'Пушкин',
 'Салтыков-Щедрин',
 'Толстой',
 'Тургенев',
 'Фет']
                                                                          In [13]:
def read file(path):
    11 11 11
    Read txt file.
    :param path: path to the text-file
    :return: string
    with open(path, 'r') as myfile:
        data = myfile.read()
    return data
                                                                          In [14]:
def pure text(text):
    11 11 11
    Remove some punctuation
    norm text = text
    # Replace breaks with spaces
    for s in ['<br />', '\n', '&#46;']:
        norm_text = norm_text.replace(s, " ")
    # Replace punctuation with spaces
```

```
for char in ['\\', '/']:
            norm text = norm text.replace(char, " ")
        return norm text
                                                                            In [15]:
    def read files concat(folder path):
        Read all txt-files in the given folder path.
        :param folder path: path to a folder with txt-files.
        :return: list of strings.
        file names = [name for name in glob.glob(FILE PATH + folder path +'/*.*txt
')]
        read files = [read file(file name) for file name in file names]
        text = " ".join(read files)
        return text
                                                                            In [16]:
    def save_to_txt(file_name, data):
        thefile = open(file name, 'w')
        for item in data:
            thefile.write("%s\n" % item)
    Train, Test
```

#### **Authors encoding**

```
In [17]:
def create author dict(list names):
    Create a dictionary, where each author corresponds to the natural number
    author coded = {}
    ind = 1
    for folder in list names:
        author coded[folder] = ind
        ind += 1
    return author_coded
```

```
In [18]:
author coded = create author dict(folders)
author coded
                                                                        Out[18]:
{'Гоголь': 1,
 'Грибоедов': 2,
 'Достоевский': 3,
 'Есенин': 4,
 'Крылов': 5,
 'Лермонтов': 6,
 'Ломоносов': 7,
 'Пушкин': 8,
 'Салтыков-Щедрин': 9,
 'Толстой': 10,
 'Тургенев': 11,
 'Фет': 12}
                                                                        In [19]:
# author splitted = sum(author, [])
# y = [author coded[name] for name in author splitted]
For train, test dataset constructing
                                                                        In [20]:
def extract substrings(s, L, N):
    .. .. ..
    :param s: string
    :param n: len of the splitted strings
    :param N: number of strings to obtain
    chunks, chunk_size = len(s), L
    splitted = [ s[i:i+chunk_size] for i in range(0, chunks, chunk size) ]
    while len(splitted) < N:</pre>
        pos start = random.randint(0, len(s) - chunk size - 1)
        splitted.append(s[pos start:pos start + chunk size])
```

```
if len(splitted) > N:
        return splitted[:N]
    return splitted
                                                                       In [21]:
def string train test(s, L, N, M):
    :param s: string
    :param L: len
    :param N: number of strings for train
    :param M: number os strings for test
    length = len(s)
    train = extract substrings(s[: int(len(s) / 2)], L, N)
    test = extract substrings(s[int(len(s) / 2):], L, M)
    return train, test
                                                                       In [22]:
def get author texts(author folder, L, N, M):
    s = pure text(read files concat(author folder))
    train, test = string train test(s, L, N, M)
    return train, test
                                                                       In [23]:
def build train test sets(folders path, L, N, M):
    11 11 11
    :param folders path: path to all folders
    :param L: length of strings
    :param N: number of strings for the train set from one author
    :param M: number of strings for the test set from one author
    X train, y train = [], []
    X test, y test = [], []
    for folder in folders:
        cur_train, cur_test = get_author_texts(folder, L, N, M)
```

```
X_train.append(cur_train)
X_test.append(cur_test)
y_train.append([author_coded[folder]] * len(cur_train))
y_test.append([author_coded[folder]] * len(cur_test))
# y_train.append([folder] * len(cur_train))
# y_test.append([folder] * len(cur_test))
return sum(X_train, []), sum(y_train, []), sum(X_test, []), sum(y_test, [])
```

### **Text preprocessing**

#### **Feature Selection**

Посчитать распределение различных знаков препинания в тексте

```
In [24]:
def punctuation freq(lower text):
    Count the number of occurencies of punctuation characters in a text
    :param text: string
    :return: list with occurencies
    symbols = [',', '.', '?', '!', ':', ';', '(', '-', '"', """]
    punct dict = dict.fromkeys(symbols, 0) # ',.?!:;(-"\''
    if len(lower text):
        for symbol in symbols:
            punct dict[symbol] = lower text.count(symbol)
        punct dict['"'] /= 2
        punct dict["'"] /= 2
    else:
        return list(punct dict.values())
    return list(punct dict.values())
                                                                       In [25]:
punctuation freq('fewepofkwkm;-()ef')
                                                                       Out[25]:
[0, 0, 0, 0, 0, 1, 1, 0, 0.0, 0.0]
```

```
def remove punctuation(text):
         .....
        Strip punctuation/symbols from words
         :param text: string
         :return: modified string
        norm text = text
        symbols = [',', '.', '?', '!', ':', ';', '(', ')', '-', '"', "'"]
         # Replace punctuation with spaces
        for char in symbols:
             norm_text = norm_text.replace(char, ' ')
        return norm text
                                                                             In [27]:
    def word norm(word):
         11 11 11
        Convert word to its first form
        return morph.parse(word)[0].normal form
                                                                             In [28]:
    POS = ['NOUN', 'ADJF', 'ADJS', 'COMP', 'VERB', 'INFN', 'PRTF', 'PRTS', 'GRND', 'ADVB
','NPRO','PRED','PREP','CONJ','PRCL','INTJ']
                                                                             In [29]:
    def POS distribution(tokenized):
         mmm
         Распределение частей речи в тексте
         # normalized = [word norm(word) for word in tokenized]
         # Переходим к частям речи
        words POS = [morph.parse(word)[0].tag.POS for word in tokenized]
        POS distr = dict.fromkeys(POS, 0) # инициализация dict для всех частей реч
И
```

In [26]:

```
POS counted = dict(collections.Counter(words POS))
        for word in POS:
             if word in POS counted.keys():
                 POS distr[word] = POS counted[word]
        return list(POS distr.values())
                                                                            In [31]:
    def numpy_concatenate(a):
        Flatten the list with numbers
        return list(np.concatenate(a))
                                                                            In [32]:
    STOP WORDS = set(stopwords.words('russian'))
                                                                            In [33]:
    def sent length distribution(lower text):
        Sentence lengths distribution
         11 11 11
         # Количество предложений
        sentences = nltk.sent tokenize(lower text)
        sentence num = len(sentences)
         # FEATURE Среднее количество слов в предложениях
         # удаление пунктуации
        sentences no punct = [remove punctuation(sentence)] for sentence in sentence
es]
         # список из токенизированных предложений
        sent tokenized = [word tokenize(sent) for sent in sentences no punct]
         # длины предложений
         sentence lengths = [len(sent) for sent in sent tokenized]
         # интервалы для разбиения длин предложений
```

```
distr = list(np.histogram(sentence lengths, bins=bins)[0])
         return distr
                                                                            In [34]:
    def word length distribution(tokenized):
         Распределение длин слов
         # ДЛИНЫ СЛОВ
         word lengths = [len(word) for word in tokenized]
         # интервалы для разбиения длин слов
        bins=[1, 4, 7, 10, 100]
         distr = list(np.histogram(word lengths, bins=bins)[0])
         return distr
                                                                            In [35]:
    def capital to lower count(text):
         11 II II
        Number of Uppercase letters to the number of lowercase
         lowercase char num = sum([i.isalpha() and i.islower() for i in text]) # KO
личество строчных символов
         uppercase_char_num = sum([i.isalpha() and i.isupper() for i in text]) # κο
личество заглавных символов
         try:
             upper to lower = uppercase char num / lowercase char num # отношение з
аглавных букв к строчным
        except ZeroDivisionError:
             upper to lower = 0
         return upper_to_lower
```

bins=[1, 4, 7, 10, 15, 100]

```
In [36]:
    def non stopwords(normalized):
        11 11 11
        Водность текста -
        Разница между единицей и отношением «количество слов после очистки стоп-сл
ов/количество слов в исходном тексте».
        11 11 11
        # не стоп-слова
        filtered words = [word for word in normalized if word not in STOP WORDS]
        try:
            stop_words_to_full_length = 1 - len(filtered_words) / len(normalized)
        except ZeroDivisionError:
            return 0
        return stop words to full length
                                                                            In [37]:
    def lexical diversity(normalized):
         Unique words to the whole text length
         try:
            lexical div = len(set(normalized)) / len(normalized)
        except ZeroDivisionError:
            return 0
        return lexical_div
                                                                            In [38]:
    \# text1 = X test[0]
     # res1 = []
     # upper_to_lower1 = capital_to_lower_count(text1)
     # res1.append([upper to lower1])
    # lower text1 = str(text1).lower()
```

```
# punct_occur_distr1 = punctuation_freq(lower_text1)
# sent len distr1 = sent length distribution(lower text1)
# pure text1 = remove punctuation(lower text1)
# tokenized1 = word tokenize(pure text1)
# pos distr1 = POS distribution(tokenized1)
# word len distr1 = word length distribution(tokenized1)
# normalized1 = [word norm(word) for word in tokenized1]
# lexical divers1 = lexical diversity(normalized1)
# nonstop to volume1 = non stopwords(normalized1)
                                                                      In [39]:
def text features(input text):
    text = input text
    res = []
    # FEATURE Отношение количества заглавных букв к строчным
    upper to lower = capital to lower count(text)
    res.append([upper to lower])
    # Приведение к нижнему регистру
    lower text = str(text).lower()
    # FEATURE Pacпределение знаков препинания (',.?!:;(-"\'') по тексту
    punct occur distr = punctuation freq(lower text)
    res.append(punct occur distr)
    # FEATURE Распределение длин предложений
    sent len distr = sent length distribution(lower text)
    res.append(sent len distr)
    # Токенизация
    pure text = remove punctuation(lower text)
    tokenized = word tokenize(pure text)
```

```
# FEATURE Распределение частей речи в тексте
        pos_distr = POS_distribution(tokenized)
        res.append(pos distr)
         # FEATURE Распределение длин слов
        word len distr = word length distribution(tokenized)
        res.append(word len distr)
        # Нормализация
        normalized = [word norm(word) for word in tokenized]
         \#_{\_\_\_} FEATURE_{\_\_\_} Разнообразие речи: Количество уникальных слов ко всему объ
ему текста
        lexical divers = lexical diversity(normalized)
        res.append([lexical_divers])
         # FEATURE Водность речи
        nonstop_to_volume = non_stopwords(normalized)
        res.append([nonstop to volume])
        return numpy concatenate(res)
    Bigrams
                                                                           In [40]:
    ALPHABET RU = 'абвгдеёжзийклмнопрстуфхцчшщъыьэюя'
                                                                           In [41]:
    russian bigrams = []
    for i in ALPHABET RU:
        for j in ALPHABET RU:
            russian bigrams.append(i+j)
                                                                           In [42]:
    print(russian bigrams[:10])
```

```
print('Количество биграм: ', len(russian bigrams))
     ['aa', 'aб', 'aв', 'ar', 'aд', 'ae', 'aë', 'aж', 'aз', 'aи']
    Количество биграм: 1089
                                                                             In [43]:
    def count bigrams(text):
         text modified = pure text(str(text).lower())
         bigrams = Counter(x+y for x, y in zip(*[text_modified[i:] for i in range(2
)]) if x.isalpha() and y.isalpha())
        return bigrams
                                                                             In [44]:
    count_bigrams('CTATИCTИЧЕСКИЙ АНАЛИЗ ТЕКСТА')
                                                                             Out[44]:
    Counter({ 'an': 1,}
              'ан': 1,
              'ar': 1,
              'ек': 1,
              'ec': 1,
              'из': 1,
              'ий': 1,
              'ис': 1,
              'ич': 1,
              'ки': 1,
              'KC': 1,
              'ли': 1,
              'на': 1,
              'CK': 1,
              'CT': 3,
              'та': 2,
              're': 1,
              'ти': 2,
              'че': 1})
```

#### **Construct TRAIN/TEST**

#### L = 20000

```
In [45]:
   X train, y train, X test, y test = build train test sets(folders, 20000, 1000,
300)
                                                                         In [212]:
   print('Количество строк в трейне: ', len(X train))
   print('Длина строки: ', len(X train[0]))
   print('Количество строк в test: ', len(X test))
   Количество строк в трейне: 12000
   Длина строки: 20000
   Количество строк в test: 3600
   Apply features to train/test
                                                                         In [235]:
   start_time = time.time()
   text features(X train[0])
   print("--- %s seconds ---" % (time.time() - start time))
   --- 6.630380392074585 seconds ---
                                                                         In [291]:
   start time = time.time()
   X train features = []
   i = 0
   for text in X train:
       X train features.append(text features(text))
       i += 1
       if i % 1000 == 0:
           print('Выполнено ', i, ' итераций.')
   print("--- %s seconds ---" % (time.time() - start time))
   Выполнено 1000 итераций.
   Выполнено 2000 итераций.
   Выполнено 3000 итераций.
   Выполнено 4000 итераций.
   Выполнено 5000 итераций.
   Выполнено 6000 итераций.
```

```
Выполнено 7000 итераций.
Выполнено 8000 итераций.
Выполнено 9000 итераций.
Выполнено 10000 итераций.
Выполнено 11000 итераций.
Выполнено 12000 итераций.
--- 51494.5682182312 seconds ---
```

#### **Common function**

```
def apply features(data, print iterations = 1000):
        start time = time.time()
        data features = []
        i = 0
        for text in data:
            data features.append(text features(text))
            i += 1
            if i % print iterations == 0:
                print('Выполнено ', i, ' итераций.')
        print("--- %s seconds ---" % (time.time() - start time))
        return data features
                                                                            In [ ]:
    X train 20000, y train 20000, X test 20000, y test 20000 = build train test se
ts(folders, 20000, 1000, 300)
    X_train_10000, y_train_10000, X_test_10000, y_test_10000 = build_train_test_se
ts(folders, 10000, 1000, 300)
    X train 5000, y train 5000, X test 5000, y test 5000 = build train test sets(f
olders, 5000, 1000, 300)
    X train 1000, y train 1000, X test 1000, y test 1000 = build train test sets(f
olders, 1000, 1000, 300)
                                                                          In [566]:
    X test features = apply features(X test, 400)
    Выполнено 400 итераций.
    Выполнено 800 итераций.
    Выполнено 1200 итераций.
```

In [46]:

```
Выполнено 1600 итераций.
    Выполнено 2000 итераций.
    Выполнено 2400 итераций.
    Выполнено 2800 итераций.
    Выполнено 3200 итераций.
    Выполнено 3600 итераций.
    --- 22765.061748743057 seconds ---
                                                                        In [513]:
    X_test_features_n500 = apply_features(X_test[:500], 100)
    --- 2678.3511624336243 seconds ---
                                                                        In [516]:
    save_to_txt('X_test_features_20000_n500.txt' , X_test_features_n500)
                                                                        In [567]:
    save to txt('X test features 20000.txt', X test features)
    DOWNLOAD FEATURES FROM TEXT FILE
                                                                         In [47]:
    TXT FILE PATH = 'C:/Users/Yauheniya/Documents/Diploma'
                                                                         In [48]:
    import pickle
                                                                         In [50]:
    def read features file(path):
        list of lists = []
        with open(path) as f:
            for line in f:
                line = line.replace('[', '')
                line = line.replace(']', '')
                inner list = [float(elt.strip()) for elt in line.split(',')]
                list of lists.append(inner list)
        return list of lists
                                                                         In [51]:
    X test features 1000 = read features file('C:/Users/Yauheniya/Documents/Diplom
a\\X test features 1000.txt')
```

```
X train features 1000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X train features 1000.txt')
    X test features 5000 = read features file('C:/Users/Yauheniya/Documents/Diplom
a\\X test features 5000.txt')
    X train features 5000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X train features 5000.txt')
    X test features 10000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X test features 10000.txt')
    X train features 10000 = read features file('C:/Users/Yauheniya/Documents/Dipl
oma\\X train features 10000.txt')
    X test features 20000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X test features 20000.txt')
    X train features 20000 = read features file('C:/Users/Yauheniya/Documents/Dipl
oma\\X train features 20000.txt')
                                                                           In [56]:
    def read bigrams file(path):
        with open(path, 'r') as myfile:
             data=myfile.read().replace('\n', '')
        data = data.replace(']','')
        data = data.split('[')
        list lines = []
        data = [i for i in data if i]
        for line in data:
             inner list = [float(elt.strip()) for elt in line.split(' ') if elt]
             list lines.append(inner list)
        return list lines
                                                                           In [57]:
     # X test bigrams 1000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X test bigrams 1000.txt')
     # X train bigrams 1000 = read features file('C:/Users/Yauheniya/Documents/Dipl
oma\\X train bigrams 1000.txt')
```

```
# X test bigrams 5000 = read features file('C:/Users/Yauheniya/Documents/Diplo
ma\\X test bigrams 5000.txt')
     # X_train_bigrams_5000 = read features file('C:/Users/Yauheniya/Documents/Dipl
oma\\X train bigrams 5000.txt')
     # X test bigrams 10000 = read features file('C:/Users/Yauheniya/Documents/Dipl
oma\\X test bigrams 10000.txt')
     # X train bigrams 10000 = read features file('C:/Users/Yauheniya/Documents/Dip
loma\\X train bigrams 10000.txt')
    X test bigrams 20000 = read bigrams file('C:/Users/Yauheniya/Documents/Diploma
\\X test bigrams 20000.txt')
    X train bigrams 20000 = read bigrams file('C:/Users/Yauheniya/Documents/Diplom
a\\X train bigrams 20000.txt')
    L = 10000
                                                                          In [58]:
    X train 10000, y train 10000, X test 10000, y test 10000 = build train test se
ts(folders, 10000, 1000, 300)
                                                                         In [615]:
    X train 10000, y train 10000, X test 10000, y test 10000 = build train test se
ts(folders, 10000, 1000, 300)
    print('Apply features to train: ')
    X train 10000 features = apply features (X train 10000, 1000)
    print('Apply features to test: ')
    X test 10000 features = apply features(X test 10000, 300)
    Apply features to train:
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
    Выполнено 3000 итераций.
    Выполнено 4000 итераций.
    Выполнено 5000 итераций.
    Выполнено 6000 итераций.
    Выполнено 7000 итераций.
    Выполнено 8000 итераций.
```

```
Выполнено 9000 итераций.
    Выполнено 10000 итераций.
    Выполнено 11000 итераций.
    Выполнено 12000 итераций.
    --- 71220.50283145905 seconds ---
    Apply features to test:
    Выполнено 300 итераций.
    Выполнено 600 итераций.
    Выполнено 900 итераций.
    Выполнено 1200 итераций.
    Выполнено 1500 итераций.
    Выполнено 1800 итераций.
    Выполнено 2100 итераций.
    Выполнено 2400 итераций.
    Выполнено 2700 итераций.
    Выполнено 3000 итераций.
    Выполнено 3300 итераций.
    Выполнено 3600 итераций.
    --- 12608.028089284897 seconds ---
                                                                        In [616]:
    save to txt('X test features 10000.txt', X test 10000 features)
    save to txt('X train features 10000.txt', X train 10000 features)
    L = 5000
                                                                         In [59]:
    X train 5000, y train 5000, X test 5000, y test 5000 = build train test sets(f
olders, 5000, 1000, 300)
                                                                        In [617]:
    X_train_5000, y_train_5000, X_test_5000, y_test_5000 = build_train_test_sets(f
olders, 5000, 1000, 300)
    X train 5000 features = apply features(X train 5000, 1000)
    X_test_5000_features = apply_features(X_test_5000, 300)
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
```

```
Выполнено 3000 итераций.
    Выполнено 4000 итераций.
    Выполнено 5000 итераций.
    Выполнено 6000 итераций.
    Выполнено 7000 итераций.
    Выполнено 8000 итераций.
    Выполнено 9000 итераций.
    Выполнено 10000 итераций.
    Выполнено 11000 итераций.
    Выполнено 12000 итераций.
    --- 16093.26831483841 seconds ---
    Выполнено 300 итераций.
    Выполнено 600 итераций.
    Выполнено 900 итераций.
    Выполнено 1200 итераций.
    Выполнено 1500 итераций.
    Выполнено 1800 итераций.
    Выполнено 2100 итераций.
    Выполнено 2400 итераций.
    Выполнено 2700 итераций.
    Выполнено 3000 итераций.
    Выполнено 3300 итераций.
    Выполнено 3600 итераций.
    --- 4962.587314367294 seconds ---
                                                                       In [620]:
    save to txt('X test features 5000.txt', X test 5000 features)
    save to txt('X train features 5000.txt', X train 5000 features)
    L = 1000
                                                                        In [60]:
    X train 1000, y train 1000, X test 1000, y test 1000 = build train test sets(f
olders, 1000, 1000, 300)
                                                                       In [619]:
```

```
X train 1000, y train 1000, X test 1000, y test 1000 = build train test sets(f
olders, 1000, 1000, 300)
    X train 1000 features = apply features (X train 1000, 1000)
    X_test_1000_features = apply_features(X test 1000, 300)
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
    Выполнено 3000 итераций.
    Выполнено 4000 итераций.
    Выполнено 5000 итераций.
    Выполнено 6000 итераций.
    Выполнено 7000 итераций.
    Выполнено 8000 итераций.
    Выполнено 9000 итераций.
    Выполнено 10000 итераций.
    Выполнено 11000 итераций.
    Выполнено 12000 итераций.
    --- 3074.2685482501984 seconds ---
    Выполнено 300 итераций.
    Выполнено 600 итераций.
    Выполнено 900 итераций.
    Выполнено 1200 итераций.
    Выполнено 1500 итераций.
    Выполнено 1800 итераций.
    Выполнено 2100 итераций.
    Выполнено 2400 итераций.
    Выполнено 2700 итераций.
    Выполнено 3000 итераций.
    Выполнено 3300 итераций.
    Выполнено 3600 итераций.
    --- 893.1362199783325 seconds ---
                                                                       In [621]:
    save to txt('X test features 1000.txt', X test 1000 features)
```

save to txt('X train features 1000.txt', X train 1000 features)

### Add L 500, 200, 100

```
In [53]:
    X test features 500 = read features file('C:/Users/Yauheniya/Documents/Diploma
\\X test features 500.txt')
    X train features 200 = read features file('C:/Users/Yauheniya/Documents/Diplom
a\\X train features 200.txt')
    X train features 100 = read features file('C:/Users/Yauheniya/Documents/Diplom
a\\X train features 100.txt')
    X train 500, y train 500, X test 500, y test 500 = build train test sets(folde
rs, 500, 1000, 300)
    X train 200, y train 200, X test 200, y test 200 = build train test sets(folde
rs, 200, 1000, 300)
    X train 100, y train 100, X test 100, y test 100 = build train test sets(folde
rs, 100, 1000, 300)
    print('Apply features to train 500: ')
    X train features 500 = apply features (X train 500, 1000)
    save to txt('X train features 500.txt', X train features 500)
    print('Apply features to test 200: ')
    X test features 200 = apply features (X test 200, 300)
    save to txt('X test features 200.txt', X test features 200)
    print('Apply features to test 100: ')
    X test features 100 = apply features(X test 100, 300)
    save to txt('X test features 100.txt', X test features 100)
    Apply features to train 500:
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
    Выполнено 3000 итераций.
    Выполнено 4000 итераций.
    Выполнено 5000 итераций.
    Выполнено 6000 итераций.
```

```
Выполнено 7000 итераций.
```

Выполнено 8000 итераций.

Выполнено 9000 итераций.

Выполнено 10000 итераций.

Выполнено 11000 итераций.

Выполнено 12000 итераций.

--- 1906.083199262619 seconds ---

Apply features to test 200:

Выполнено 300 итераций.

Выполнено 600 итераций.

Выполнено 900 итераций.

Выполнено 1200 итераций.

Выполнено 1500 итераций.

Выполнено 1800 итераций.

Выполнено 2100 итераций.

Выполнено 2400 итераций.

Выполнено 2700 итераций.

Выполнено 3000 итераций.

Выполнено 3300 итераций.

Выполнено 3600 итераций.

--- 238.18416953086853 seconds ---

Apply features to test 100:

Выполнено 300 итераций.

Выполнено 600 итераций.

Выполнено 900 итераций.

Выполнено 1200 итераций.

Выполнено 1500 итераций.

Выполнено 1800 итераций.

Выполнено 2100 итераций.

Выполнено 2400 итераций.

Выполнено 2700 итераций.

Выполнено 3000 итераций.

Выполнено 3300 итераций.

```
Выполнено 3600 итераций.
    --- 126.03307580947876 seconds ---
                                                                         In [193]:
    X train 500, y train 500, X test 500, y test 500 = build train test sets(folde
rs, 500, 1000, 300)
    print('Apply features to train 500: ')
    X train 500 features = apply features(X train 500, 1000)
    print('Apply features to test 500: ')
    X test 500 features = apply features(X test 500, 300)
    save to txt('X test features 500.txt', X test 500 features)
    X train 200, y train 200, X test 200, y test 200 = build train test sets(folde
rs, 200, 1000, 300)
    print('Apply features to train 200: ')
    X train 200 features = apply features(X train 200, 1000)
    print('Apply features to test 200: ')
    X test 200 features = apply features(X test 200, 300)
    save to txt('X train features 200.txt', X train 200 features)
    X train 100, y train 100, X test 100, y test 100 = build train test sets(folde
rs, 100, 1000, 300)
    print('Apply features to train 100: ')
    X train 100 features = apply features(X train 100, 1000)
    print('Apply features to test 100: ')
    X test 100 features = apply features(X test 100, 300)
    save to txt('X train features 100.txt', X train 100 features)
    Apply features to train 500:
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
    Выполнено 3000 итераций.
    Выполнено 4000 итераций.
    Выполнено 5000 итераций.
    Выполнено 6000 итераций.
    Выполнено 7000 итераций.
```

```
Выполнено 8000 итераций.
```

Выполнено 9000 итераций.

Выполнено 10000 итераций.

Выполнено 11000 итераций.

Выполнено 12000 итераций.

--- 1637.513708114624 seconds ---

Apply features to test 500:

Выполнено 300 итераций.

Выполнено 600 итераций.

Выполнено 900 итераций.

Выполнено 1200 итераций.

Выполнено 1500 итераций.

Выполнено 1800 итераций.

Выполнено 2100 итераций.

Выполнено 2400 итераций.

Выполнено 2700 итераций.

Выполнено 3000 итераций.

Выполнено 3300 итераций.

Выполнено 3600 итераций.

--- 502.9348421096802 seconds ---

Apply features to train 200:

Выполнено 1000 итераций.

Выполнено 2000 итераций.

Выполнено 3000 итераций.

Выполнено 4000 итераций.

Выполнено 5000 итераций.

Выполнено 6000 итераций.

Выполнено 7000 итераций.

Выполнено 8000 итераций.

Выполнено 9000 итераций.

Выполнено 10000 итераций.

Выполнено 11000 итераций.

Выполнено 12000 итераций.

```
--- 679.055750131607 seconds ---
    Apply features to test 200:
    Выполнено 300 итераций.
    Выполнено 600 итераций.
               900 итераций.
    Выполнено
    Выполнено 1200 итераций.
    Выполнено 1500 итераций.
    Выполнено 1800 итераций.
    Выполнено 2100 итераций.
    Выполнено 2400 итераций.
    Выполнено 2700 итераций.
    Выполнено
               3000 итераций.
    Выполнено
               3300 итераций.
    Выполнено 3600 итераций.
    --- 19852.91144132614 seconds ---
    Apply features to train 100:
    Выполнено 1000 итераций.
    Выполнено 2000 итераций.
Выполнено 3000 итераций.
Выполнено 4000 итераций.
Выполнено 5000 итераций.
Выполнено 6000 итераций.
Выполнено 7000 итераций.
Выполнено 8000 итераций.
Выполнено 9000 итераций.
Выполнено 10000 итераций.
Выполнено 11000 итераций.
Выполнено 12000 итераций.
--- 375.43145871162415 seconds ---
Apply features to test 100:
Выполнено 300 итераций.
Выполнено 600 итераций.
Выполнено 900 итераций.
Выполнено 1200 итераций.
Выполнено 1500 итераций.
Выполнено 1800 итераций.
Выполнено 2100 итераций.
Выполнено 2400 итераций.
Выполнено 2700 итераций.
Выполнено 3000 итераций.
Выполнено 3300 итераций.
Выполнено 3600 итераций.
--- 110.28601455688477 seconds ---
```

#### Count bigrams for train/test

0)

```
In [253]:
start time = time.time()
text bigrams = []
for text in X train[:1]:
    text_bigrams.append(count_bigrams(text))
print("--- %s seconds ---" % (time.time() - start time))
--- 0.8690032958984375 seconds ---
                                                                       In [61]:
def count bigrams distr(data, label, max freq):
    total bigrams = Counter()
    for text in data:
        total_bigrams.update(count_bigrams(text))
    most common bigrams freq = total bigrams.most common(200)
    most common bigrams = [key for key, val in most common bigrams freq]
    all freqs = []
    for text in data:
        cur frequencies = []
        cur bigrams = count bigrams(text)
        for bigram in most_common_bigrams:
            cur frequencies.append(cur bigrams[bigram])
        all freqs.append(cur frequencies)
    if label == 'train':
        largest freq = np.max(all freqs)
        return largest freq, all freqs / largest freq
    else:
        return max_freq, all_freqs / max_freq
                                                                      In [277]:
start time = time.time()
bigrams largest freq, X train bigrams = count bigrams distr(X train, 'train',
```

```
print("--- %s seconds ---" % (time.time() - start time))
    --- 608.1697566509247 seconds ---
                                                                          In [280]:
    start time = time.time()
    bigrams largest freq, X test bigrams = count bigrams distr(X test, 'test', big
rams largest freq)
    print("--- %s seconds ---" % (time.time() - start time))
    --- 178.92333936691284 seconds ---
                                                                          In [294]:
    save to txt('X train bigrams 20000.txt', X train bigrams)
                                                                          In [295]:
    start time = time.time()
    save to txt('X test bigrams 20000.txt', X test bigrams)
    print("--- %s seconds ---" % (time.time() - start time))
    --- 25.53329849243164 seconds ---
                                                                          In [296]:
    start time = time.time()
    save to txt('X train features 20000.txt', X train features)
    print("--- %s seconds ---" % (time.time() - start time))
    --- 0.7941687107086182 seconds ---
    Bigrams with scaling
                                                                           In [62]:
    def count_total_bigrams(X, list_bigrams = russian_bigrams):
        total bigrams = dict.fromkeys(list bigrams, 0)
        for text in X:
            total bigrams.update(count bigrams(text))
        total bigrams1 = dict.fromkeys(list bigrams, 0)
        for bigram in list bigrams:
            total bigrams1[bigram] = total bigrams[bigram]
        return total bigrams1
```

```
In [187]:
count total bigrams(['пипиор', 'неуаан'], ['аб', 'бв', 'пи', 'не'])
                                                                      Out[187]:
{'аб': 0, 'бв': 0, 'не': 1, 'пи': 2}
                                                                       In [63]:
def separate bigrams(data, list bigrams):
    Return frequencies only for bigrams mentioned in list bigrams
    train bigrams freq = []
    for text in data:
        cur frequencies = []
        cur bigrams = count bigrams(text)
        for bigram in list bigrams:
            cur frequencies.append(cur bigrams[bigram])
        train bigrams freq.append(cur frequencies)
    return train bigrams freq
                                                                       In [64]:
def dataset bigrams(X train, X test):
    11 11 11
    Count bigrams frequencies for train and test
    .....
    train bigrams = count total bigrams(X train)
    most common bigrams freq = Counter(train bigrams).most common(200)
    train list bigrams = [key for key, val in most common bigrams freq]
    # train bigrams freq = [val for key, val in most common bigrams freq]
    train bigrams freq = separate bigrams(X train, train list bigrams)
    test bigrams freq = separate bigrams(X test, train list bigrams)
#
      train list bigrams = [key for key, val in most common bigrams freq]
      test bigrams = count total bigrams (X test, train list bigrams)
      test bigrams freq = [val for key, val in list(test bigrams.items())]
```

```
return train bigrams freq, test bigrams freq
                                                                           In [581:
    X train bigrams freq 20000, X test bigrams freq 20000 = dataset bigrams(X trai
n 20000, X test 20000)
    X train bigrams freq 10000, X test bigrams freq 10000 = dataset bigrams(X trai
n 10000, X test 10000)
    X train bigrams freq 5000, X test bigrams freq 5000 = dataset bigrams(X train
5000, X test 5000)
    X train bigrams freq 1000, X test bigrams freq 1000 = dataset bigrams(X train
1000, X test 1000)
                                                                           In [59]:
    X train bigrams freq 500, X test bigrams freq 500 = dataset bigrams (X train 50
0, X test 500)
    X train bigrams freq 200, X test bigrams freq 200 = dataset bigrams (X train 20
0, X test 200)
    X train bigrams freq 100, X test bigrams freq 100 = dataset bigrams (X train 10
0, X test 100)
```

# Нормализация

An alternative standardization is scaling features to lie between a given minimum and maximum value, often between zero and one, or so that the maximum absolute value of each feature is scaled to unit size. This can be achieved using MinMaxScaler or MaxAbsScaler, respectively.

The motivation to use this scaling include robustness to very small standard deviations of features and preserving zero entries in sparse data.

#### Vectorized features normalize

```
In []:
# min_max_scaler = preprocessing.MinMaxScaler()
# X_train_minmax = min_max_scaler.fit_transform(X_train)
# X_test_minmax = min_max_scaler.transform(X_test)
In [596]:
min_max_scaler_20000 = preprocessing.MinMaxScaler()
X_train_features_norm = min_max_scaler_20000.fit_transform(X_train_features)
X test features norm = min_max_scaler_20000.transform(X_test_features)
```

## Bigrams normalize

```
In [199]:
    def bigrams normalize (X train bigrams freq, X test bigrams freq):
        Нормализация биграммных частот
        max value = np.max(X train bigrams freq)
        return X train bigrams freq / max value, X test bigrams freq / max value
                                                                          In [486]:
    X train bigrams freq, X test bigrams freq = bigrams normalize(X train bigrams
freq, X test bigrams freq)
    Dataset with L = 20000
                                                                          In [524]:
    print('Количество строк в X train features: ', len(X train features))
    print('Количество строк в X test features n500: ', len(X test features n500))
    Количество строк в X train features: 12000
    Количество строк в X test features n500:
                                                                          In [523]:
    print('Количество строк в X train bigrams freq: ', len(X train bigrams freq))
    print('Количество строк в X_test_bigrams_freq: ', len(X_test_bigrams_freq))
    Количество строк в X train bigrams freq: 12000
    Количество строк в X test bigrams freq: 3600
                                                                          In [200]:
    def concat features bigrams(features, bigrams):
        .....
        Concatenates two matrices (multi-dim-arrays): with features and bigrams fr
equencies
        concat features = []
        for i, j in zip(features, bigrams):
            concat features.append(list(i) + list(j))
        return concat features
                                                                          In [568]:
```

```
X train 20000 concat = concat features bigrams(X train features, X train bigra
ms freq)
    X test 20000 concat = concat features bigrams(X test features, X test bigrams
freq)
                                                                          In [603]:
    X train 20000 norm = concat features bigrams(X train features norm, X train bi
grams freq)
    X test 20000 norm = concat features bigrams(X test features norm, X test bigra
ms freq)
    Train
                                                                             In [1]:
    from sklearn.svm import SVC
    from sklearn import metrics
    from sklearn.linear_model import LogisticRegression
    from sklearn.naive bayes import GaussianNB
    from sklearn.ensemble import GradientBoostingClassifier
    from sklearn.tree import DecisionTreeClassifier
    from sklearn.metrics import accuracy score, roc auc score
    from sklearn.model selection import validation curve
    from sklearn.neighbors import KNeighborsClassifier
                                                                             In [2]:
    from sklearn.model selection import RandomizedSearchCV
    from sklearn.model selection import GridSearchCV
                                                                             In [ ]:
    y_train = y_train_20000
    y test = y test 20000
                                                                          In [569]:
    model = SVC(kernel='linear')
    model.fit(X_train_20000_concat, y_train)
    print(model)
    print(metrics.classification report(y test, model.predict(X test 20000 concat)
) )
```

# print(metrics.confusion matrix(expected, predicted))

SVC(C=1.0, cache\_size=200, class\_weight=None, coef0=0.0,
 decision\_function\_shape='ovr', degree=3, gamma='auto', kernel='linear',
 max\_iter=-1, probability=False, random\_state=None, shrinking=True,
 tol=0.001, verbose=False)

	precision	recall	f1-score	support
1	0.52	0.48	0.50	300
2	0.26	0.24	0.25	300
3	0.66	0.96	0.78	300
4	1.00	0.97	0.98	300
5	1.00	0.17	0.29	300
6	0.72	0.94	0.81	300
7	0.96	0.86	0.91	300
8	0.58	0.89	0.70	300
9	0.79	0.73	0.76	300
10	0.79	0.91	0.85	300
11	0.58	0.91	0.70	300
12	0.07	0.01	0.01	300
avg / total	0.66	0.67	0.63	3600

In [570]:

print(metrics.confusion matrix(y test, model.predict(X test 20000 concat))) [[145 0 27 0 39 26 57 0] [ 23 0 11 0 113 2] 0 288 1] [ 0 291 0] 0 209 0 283 0] 0 258 0 268 0] [ 1 [ 10 0 68 0 220 0] 0 ] 0 12 3 274 

```
0 4 0 0 8 0 272
 [ 0 0 16 0
                                             01
 [101 0 20
               0
                   0 52
                         0 81
                                 0 33 11
                                               211
                                                                   In [585]:
svc = SVC(kernel='rbf').fit(X train 20000 concat, y train)
svc.score(X test 20000 concat, y test)
                                                                   Out[585]:
0.08333333333333333
                                                                   In [586]:
svc = SVC(kernel='linear').fit(X_train_20000_concat, y_train)
svc.score(X test 20000 concat, y test)
                                                                   Out[586]:
0.6733333333333333
                                                                   In [584]:
svc = SVC(kernel='sigmoid').fit(X train 20000 concat, y train)
svc.score(X test 20000 concat, y test)
                                                                   Out[584]:
0.08333333333333333
                                                                     In [ ]:
# X train features, X train bigrams freq
                                                                   In [587]:
svc = SVC(kernel='linear').fit(X train features, y train)
svc.score(X test features, y test)
                                                                   Out[587]:
0.6741666666666667
                                                                   In [588]:
svc = SVC(kernel='linear').fit(X train bigrams freq, y train)
svc.score(X test bigrams freq, y test)
                                                                   Out[588]:
0.679722222222222
                                                                   In [598]:
svc = SVC(kernel='linear').fit(X train features norm, y train)
svc.score(X test features norm, y test)
                                                                   Out[598]:
```

```
0.7066666666666667
                                                                      In [604]:
svc = SVC(kernel='linear').fit(X train 20000 norm, y train)
svc.score(X test 20000 norm, y test)
                                                                      Out[604]:
0.7436111111111111
LOGRegression
                                                                      In [211]:
from sklearn.linear model import LogisticRegression
                                                                      In [212]:
model log = LogisticRegression()
model log.fit(X train 20000 norm, y train)
# predicted = model log.predict(X test 20000 norm)
model log.score(X test 20000 norm, y test)
# print(metrics.classification report(y test, predicted))
Naive Bayes
                                                                      In [210]:
from sklearn.naive bayes import GaussianNB
                                                                      In [612]:
model gauss = GaussianNB()
model gauss.fit(X train 20000 norm, y train)
model gauss.score(X test 20000 norm, y test)
                                                                      Out[612]:
0.5966666666666667
Not normalized data
                                                                      In [647]:
svc = SVC(kernel='linear').fit(X train 20000 concat, y train)
svc.score(X test 20000 concat, y test)
                                                                      Out[647]:
0.6733333333333333
Common function
                                                                      In [208]:
```

```
def count score (X train features, X test features, X train bigrams freq, X tes
t bigrams freq):
         # объединяем признаки и биграммы
        X train concat = concat features bigrams(X train features, X train bigrams
freq)
        X test concat = concat features bigrams(X test features, X test bigrams fr
eq)
         # скалируем
        scaler = preprocessing.StandardScaler()
        X train standard = scaler.fit transform(X train concat)
        X test standard = scaler.fit transform(X test concat)
         # svc = SVC(kernel='linear').fit(X train standard, y train)
        svc = LogisticRegression().fit(X train standard, y train)
        return svc.score(X test standard, y test)
                                                                          In [704]:
    from sklearn.ensemble import GradientBoostingClassifier as gbc
                                                                          In [214]:
     # log
    count score(X train features 20000, X test features 20000, X train bigrams fre
q, X test bigrams freq)
                                                                          Out[214]:
    0.77222222222223
                                                                          In [674]:
     # логистическая без стандартизация
    count score(X train features, X test features, X train bigrams freq, X test bi
grams freq)
                                                                          Out[674]:
    0.64111111111111111
                                                                          In [665]:
     # без нормализации
     # count score(X train features, X test features, X train bigrams freq, X test
bigrams freq)
```

```
Out[665]:
    0.67333333333333333
                                                                          In [701]:
    count score (X train features, X test features, X train bigrams freq, X test bi
grams freq)
                                                                          Out[701]:
    0.667222222222223
                                                                          In [667]:
    count score (X train features, X test features, X train bigrams freq, X test bi
grams freq)
                                                                          Out[667]:
    0.7475
                                                                          In [654]:
    count score (X train 10000 features, X test 10000 features, X train bigrams fre
q 10000, X test bigrams freq 10000)
                                                                          Out[654]:
    0.6966666666666667
                                                                          In [640]:
    count_score(X_train_features, X_test_features, X_train_bigrams freq, X test bi
grams freq)
                                                                          Out[640]:
    0.74527777777778
                                                                          In [641]:
    # c MinMaxScaler 0.707
    count score(X train 10000 features, X test 10000 features, X train bigrams fre
q 10000, X test bigrams freq 10000)
                                                                          Out[641]:
    0.6994444444444444
                                                                          In [642]:
     # c MinMaxScaler 0.6325
    count_score(X_train_5000_features, X_test_5000_features, X_train bigrams freq
5000, X test bigrams freq 5000)
                                                                          Out[642]:
```

0.6936111111111111

```
In [643]:
     # c MinMaxScaler 0.47916
    count score(X train 1000 features, X test 1000 features, X train bigrams freq
1000, X test bigrams freq 1000)
                                                                          Out[643]:
    0.50777777777778
    Confusion Matrix
                                                                          In [81]:
    import seaborn as sns
                                                                          In [82]:
     # cm = metrics.confusion matrix(y test, predictions)
    def draw comfusion matrix sns(cm):
        plt.figure(figsize=(9,9))
        sns.heatmap(cm, annot=True, fmt=".3f", linewidths=.5, square = True, cmap
= 'Blues r');
        plt.ylabel('Actual label');
        plt.xlabel('Predicted label');
        all sample title = 'Accuracy Score: {0}'.format(score)
        plt.title(all sample title, size = 15);
                                                                          In [83]:
    def draw cm matplot(cm):
        plt.figure(figsize=(10,10))
        plt.imshow(cm, interpolation='nearest', cmap='Pastel1')
        plt.title('Confusion matrix', size = 15)
        plt.colorbar()
        tick marks = np.arange(len(folders))
        plt.xticks(tick marks, folders, rotation=45, size = len(folders))
        plt.yticks(tick marks, folders, size = len(folders))
        plt.tight layout()
        plt.ylabel('Actual label', size = 15)
        plt.xlabel('Predicted label', size = 15)
        width, height = cm.shape
```

```
for x in range(width):
       for y in range(height):
           plt.annotate(str(cm[x][y]), xy=(y, x),
           horizontalalignment='center',
           verticalalignment='center')
                                                                In [185]:
lr = LogisticRegression().fit(X train scale 20000, y train 20000)
cm 20000 = metrics.confusion matrix(y test, lr.predict(X test scale 20000))
                                                                In [186]:
lr.score(X test scale 20000, y test 20000)
                                                                Out[186]:
0.8130555555555555
                                                                In [187]:
draw cm matplot(cm 20000)
                                                                In [190]:
lr = LogisticRegression().fit(X train scale 10000, y train 10000)
print('Score = ', lr.score(X test scale 10000, y test 10000))
cm 10000 = metrics.confusion matrix(y test, lr.predict(X test scale 10000))
draw cm matplot(cm 10000)
In [191]:
lr = LogisticRegression().fit(X train scale 5000, y train 5000)
print('Score = ', lr.score(X test scale 5000, y test 5000))
cm 5000 = metrics.confusion matrix(y test, lr.predict(X test scale 5000))
draw cm matplot(cm 5000)
In [192]:
lr = LogisticRegression().fit(X train scale 1000, y train 1000)
print('Score = ', lr.score(X test scale 1000, y test 1000))
cm 1000 = metrics.confusion matrix(y test, lr.predict(X test scale 1000))
```

```
draw cm matplot(cm 1000)
Score = 0.53527777777777
                                                                    In [207]:
lr = LogisticRegression().fit(X train scale 500, y train 500)
print('Score = ', lr.score(X test scale 500, y test 500))
cm 500 = metrics.confusion matrix(y test, lr.predict(X test scale 500))
draw cm matplot(cm 500)
Score = 0.46416666666666667
                                                                    In [208]:
lr = LogisticRegression().fit(X train scale 200, y train 200)
print('Score = ', lr.score(X test scale 200, y test 200))
cm 200 = metrics.confusion matrix(y test, lr.predict(X test scale 200))
draw cm matplot(cm 200)
Score = 0.312222222222223
                                                                    In [209]:
lr = LogisticRegression().fit(X train scale 100, y train 100)
print('Score = ', lr.score(X test scale 100, y test 100))
cm 100 = metrics.confusion matrix(y test, lr.predict(X test scale 100))
draw cm matplot(cm 100)
Score = 0.25
Tune logistic regression
                                                                    In [212]:
def tune param logistic(X train, X test, y train, y test):
    start time = time.time()
    C param range = [0.0001, 0.001, 0.01, 0.1, 1, 10]
    acc table = pd.DataFrame(columns = ['C parameter', 'Accuracy'])
```

### Out[214]:

	C_para meter	Acc uracy
0	0.0001	0.7 81389
1	0.0010	0.8 04167
2	0.0100	0.8 01667
3	0.1000	0.8 11111
4	1.0000	0.8 13056

	C_para meter	Acc uracy
5	10.0000	0.8

In [215]:

tune\_param\_logistic(X\_train\_scale\_10000, X\_test\_scale\_10000, y\_train, y\_test)
--- 144.0762186050415 seconds ---

Out[215]:

	C_para meter	Acc uracy
0	0.0001	0.7 26111
1	0.0010	0.7 50556
2	0.0100	0.7 45833
3	0.1000	0.7 48611
4	1.0000	0.7 40833
5	10.0000	0.7 33333

In [216]:

```
tune_param_logistic(X_train_scale_5000, X_test_scale_5000, y_train, y_test)
--- 151.885760307312 seconds ---
```

Out[216]:

	C_para meter	Acc uracy
0	0.0001	0.7 05278
1	0.0010	0.7 25278
2	0.0100	0.7 16111
3	0.1000	0.7 10278
4	1.0000	0.7 14167
5	10.0000	0.6 95833

In [217]:

tune\_param\_logistic(X\_train\_scale\_1000, X\_test\_scale\_1000, y\_train, y\_test)
--- 188.87707138061523 seconds ---

Out[217]:

	C_para meter	Acc uracy
0	0.0001	0.5 06389
1	0.0010	0.5

	C_para meter	Acc uracy
		18889
2	0.0100	0.5 20556
3	0.1000	0.5 36389
4	1.0000	0.5 35278
5	10.0000	0.5 34444

In [218]:

tune\_param\_logistic(X\_train\_scale\_500, X\_test\_scale\_500, y\_train, y\_test)
--- 183.8874614238739 seconds ---

Out[218]:

	C_para meter	Acc uracy
0	0.0001	0.4 41389
1	0.0010	0.4 58333
2	0.0100	0.4 625

	C_para meter	Acc uracy
3	0.1000	0.4 66944
4	1.0000	0.4 64167
5	10.0000	0.4 62778

In [219]:

tune\_param\_logistic(X\_train\_scale\_200, X\_test\_scale\_200, y\_train, y\_test)
--- 202.77686429023743 seconds ---

Out[219]:

	C_para meter	Acc uracy
0	0.0001	0.2 56944
1	0.0010	0.2 76667
2	0.0100	0.2 92222
3	0.1000	0.3 06389
4	1.0000	0.3 12222

	C_para meter	Acc uracy
5	10.0000	0.3 12222

#### k-NN

```
In [84]:
knc = KNeighborsClassifier(n neighbors = 10, metric = "cityblock")
knc.fit(X train scale 20000, y train)
knc.score(X test scale 20000, y test)
                                                                     Out[84]:
0.7686111111111111
                                                                     In [89]:
knc = KNeighborsClassifier(n_neighbors = 12, metric = "euclidean")
knc.fit(X train scale 20000, y train)
knc.score(X test scale 20000, y test)
                                                                     Out[89]:
0.77333333333333333
                                                                     In [88]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X_train_scale_20000, y_train)
knc.score(X test scale 20000, y test)
                                                                     Out[88]:
0.7744444444444445
                                                                    In [172]:
params = {"n_neighbors": range(3,50,10)} # "metric": ["euclidean", "cityblock"
model = KNeighborsClassifier()
grid = GridSearchCV(model, params)
start = time.time()
grid.fit(X train scale 20000, y train)
```

```
# evaluate the best grid searched model on the testing data
print("[INFO] grid search took {:.2f} seconds".format(time.time() - start))
acc = grid.score(X test scale 20000, y test)
print("[INFO] grid search accuracy: {:.2f}%".format(acc * 100))
print("[INFO] grid search best parameters: {}".format(grid.best params ))
[INFO] grid search took 2614.65 seconds
[INFO] grid search accuracy: 72.83%
[INFO] grid search best parameters: {'n neighbors': 3}
                                                                     In [90]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X train scale 10000, y train)
knc.score(X_test_scale_10000, y_test)
                                                                     Out[90]:
0.718055555555556
                                                                     In [91]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X train scale 5000, y train 5000)
knc.score(X test scale 5000, y test 5000)
                                                                     Out[91]:
0.664722222222222
                                                                     In [92]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X train scale 1000, y train 1000)
knc.score(X test scale 1000, y test 1000)
                                                                     Out [92]:
0.33222222222222
                                                                     In [95]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "cityblock")
knc.fit(X train scale 1000, y train 1000)
knc.score(X test scale 1000, y test 1000)
                                                                     Out[95]:
0.304444444444446
                                                                     In [96]:
```

```
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X train scale 500, y train 500)
knc.score(X test scale 500, y test 500)
                                                                    Out[96]:
0.2361111111111111
                                                                     In [97]:
knc = KNeighborsClassifier(n neighbors = 11, metric = "euclidean")
knc.fit(X train scale 200, y train 200)
knc.score(X test scale 200, y test 200)
                                                                    Out [97]:
0.1327777777777777
SVM
                                                                   In [104]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.0001)
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                   Out[104]:
0.7705555555555555
                                                                   In [105]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.001)
svm.fit(X train scale 10000, y train)
svm.score(X test scale 10000, y test)
                                                                   Out[105]:
0.743055555555556
                                                                   In [106]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.001)
svm.fit(X train scale 5000, y train)
svm.score(X_test_scale_5000, y_test)
                                                                   Out[106]:
0.7016666666666667
                                                                   In [107]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.001)
```

```
svm.fit(X train scale 1000, y train)
svm.score(X test scale 1000, y test)
                                                                   Out[107]:
0.5025
                                                                   In [108]:
svm = SVC(kernel='linear')
svm.fit(X train scale 500, y train)
svm.score(X test scale 500, y test)
                                                                   Out[108]:
0.46194444444444444
                                                                   In [109]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.001)
svm.fit(X train scale 500, y train)
svm.score(X test scale 500, y test)
                                                                   Out[109]:
0.474722222222222
                                                                   In [110]:
svm = SVC(kernel='linear', C=0.0001, gamma = 0.001)
svm.fit(X train scale 200, y train 200)
svm.score(X test scale 200, y test 200)
                                                                   Out[110]:
0.279722222222222
                                                                   In [177]:
svm = SVC()
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                   Out[177]:
0.705277777777778
                                                                   In [185]:
def svc param selection(X, y, nfolds=3):
    Cs = [0.001, 0.01, 0.1, 1, 10]
    gammas = [0.001, 0.01, 0.1, 1]
    param grid = {'C': Cs, 'gamma' : gammas}
```

```
grid search = GridSearchCV(SVC(kernel='linear'), param grid, cv=nfolds)
    grid search.fit(X, y)
    grid search.best params
    return grid search.best params
                                                                    In [186]:
svc param selection(X train scale 20000, y train, nfolds=3)
                                                                    Out[186]:
{'C': 0.01, 'gamma': 0.001}
                                                                    In [190]:
svm = SVC(kernel = 'linear', C=0.01, gamma = 0.001)
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                    Out[190]:
0.749444444444445
Best params = 0.776944
                                                                    In [193]:
svm = SVC(kernel = 'linear', C=0.001, gamma = 0.0001)
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                    Out[193]:
0.7769444444444444
                                                                    In [205]:
svm = SVC(kernel = 'linear', C=100, gamma = 0.0001)
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                    Out[205]:
0.72
                                                                    In [206]:
def svc param selection1(X, y, nfolds=3):
    Cs = [0.001, 0.01]
    gammas = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1]
    kernels = ['linear']
```

```
param grid = {'C': Cs, 'gamma' : gammas}
    grid search = GridSearchCV(SVC(), param grid, cv=nfolds)
    grid search.fit(X, y)
    grid search.best params
    return grid search.best params
                                                                   In [207]:
svc param selection1(X train scale 20000, y train, nfolds=3)
                                                                   Out[207]:
{'C': 0.01, 'gamma': 0.01}
                                                                   In [208]:
svm = SVC(kernel = 'linear', C=0.01, gamma = 0.01)
svm.fit(X train scale 20000, y train)
svm.score(X test scale 20000, y test)
                                                                   Out[208]:
0.749444444444445
                                                                      In [ ]:
svc param selection(X train scale 10000, y train, nfolds=3)
                                                                      In [ ]:
svc_param_selection(X_train_scale_5000, y train, nfolds=3)
                                                                      In [ ]:
svc param selection(X train scale 1000, y train, nfolds=3)
GradientBoostingClassifier
                                                                   In [111]:
gbc = GradientBoostingClassifier()
gbc.fit(X train scale 20000, y train)
gbc.score(X test scale 20000, y test)
                                                                   Out[111]:
0.6913888888888889
                                                                      In [ ]:
param test1 = {'n estimators':range(20,81,10)}
gsearch1 = GridSearchCV(
    estimator = GradientBoostingClassifier(
```

```
learning rate=0.1, min samples split=500, min samples leaf=50, max depth
=8, max features='sqrt', subsample=0.8, random state=10),
        param grid = param test1, scoring='roc auc', n jobs=4, cv=5)
    gsearch1.fit(X train scale 20000,y train)
                                                                           In [ ]:
    def gbc param selection(X, y, nfolds=5):
        estimators = range (20, 81, 10)
        param grid = {'n estimators': estimators}
        grid search = GridSearchCV(GradientBoostingClassifier(learning rate=0.1, m
in samples split=500,
                                                               min samples leaf=50,
max depth=8,
                                                               max features='sqrt',
subsample=0.8,random state=10), param grid, cv=nfolds)
        grid search.fit(X, y)
        grid search.best params
        return grid search.best params , grid search.best params , grid search.bes
t_score
                                                                           In [ ]:
    gsearch1.grid scores , gsearch1.best params , gsearch1.best score
    DecisionTreeClassifier
                                                                         In [180]:
    dt = DecisionTreeClassifier()
    dt.fit(X train scale 20000, y train)
    dt.score(X test scale 20000, y test)
                                                                         Out[180]:
    0.49
    Понижение размерности
                                                                         In [140]:
    from sklearn import decomposition
                                                                         In [141]:
    colors = ['#990033', '#CC00FF', '#330066', '#00FF33', '#FF9900',
```

```
'#3366FF', '#333300', '#999999', '#336666', '#FF33CC', '#009900', '#
FF9999'1
                                                                          In [142]:
                                                                           In [215]:
    # That's an impressive list of imports.
    import numpy as np
    from numpy import linalg
    from numpy.linalg import norm
    from scipy.spatial.distance import squareform, pdist
    # We import sklearn.
    import sklearn
    from sklearn.manifold import TSNE
    from sklearn.datasets import load digits
    from sklearn.preprocessing import scale
    # We'll hack a bit with the t-SNE code in sklearn 0.15.2.
    from sklearn.metrics.pairwise import pairwise distances
    from sklearn.manifold.t_sne import (_joint_probabilities,
    kl divergence)
    # from sklearn.utils.extmath import ravel
    # Random state.
    RS = 20150101
    # We'll use matplotlib for graphics.
    import matplotlib.pyplot as plt
    import matplotlib.patheffects as PathEffects
    import matplotlib
    %matplotlib inline
     # We import seaborn to make nice plots.
    import seaborn as sns
```

```
sns.set style('darkgrid')
sns.set palette('muted')
sns.set context("notebook", font scale=1.5,
rc={"lines.linewidth": 2.5})
                                                                      In [216]:
def scatter(x, colors):
    # We choose a color palette with seaborn.
    palette = np.array(sns.color palette("hls", 13))
    # We create a scatter plot.
    f = plt.figure(figsize=(8, 8))
    ax = plt.subplot(aspect='equal')
    sc = ax.scatter(x[:,0], x[:,1], lw=0, s=40,
                    c=palette[colors.astype(np.int)])
    plt.xlim(-25, 25)
    plt.ylim(-25, 25)
    ax.axis('off')
    ax.axis('tight')
    # We add the labels for each digit.
    txts = []
    for i in range (10):
        # Position of each label.
        xtext, ytext = np.median(x[colors == i, :], axis=0)
        txt = ax.text(xtext, ytext, str(i), fontsize=24)
        txt.set path effects([
            PathEffects.Stroke(linewidth=5, foreground="w"),
            PathEffects.Normal()])
        txts.append(txt)
    return f, ax, sc, txts
                                                                      In [217]:
y = np.hstack(y train)
```

```
In [692]:
    digits_proj = TSNE(random_state=RS).fit_transform(X_train_10000_concat)
    scatter(digits proj, y)
    plt.savefig('digits tsne-generated.png', dpi=120)
                                                                          In [693]:
    proj = TSNE(random_state=RS).fit_transform(X_train_features_norm)
    scatter(proj, y)
    plt.savefig('X train features norm.png', dpi=120)
                                                                          In [694]:
    proj = TSNE(random state=RS).fit transform(X train bigrams freq 10000 norm)
    scatter(proj, y)
    plt.savefig('Bigrams_freq_10000_norm.png', dpi=120)
                                                                          In [695]:
    proj = TSNE(random_state=RS).fit_transform(X_train_10000_features)
    scatter(proj, y)
    plt.savefig('Features 10000 no norm.png', dpi=120)
                                                                          In [218]:
    def count_tsne(X_train_features, X_test_features, X_train_bigrams_freq, X_test
_bigrams_freq, name_pic):
        # объединяем признаки и биграммы
        X_train_concat = concat_features_bigrams(X_train_features, X_train_bigrams
freq)
        X test concat = concat features bigrams(X test features, X test bigrams fr
eq)
        # скалируем
        scaler = preprocessing.StandardScaler()
        X train standard = scaler.fit transform(X train concat)
        X_test_standard = scaler.fit_transform(X_test_concat)
        proj = TSNE(random state=RS).fit transform(X train standard)
        scatter(proj, np.hstack(y_train))
        plt.savefig(name_pic + '_train_stand.png', dpi=120)
```

```
proj = TSNE(random state=RS).fit transform(X test standard)
        scatter(proj, np.hstack(y test))
        plt.savefig(name pic + ' test stand.png', dpi=120)
                                                                          In [697]:
    count tsne(X train features, X test features, X train bigrams freq, X test big
rams freq, 'L 20000 concat')
                                                                          In [219]:
    count tsne(X train features 1000, X test features 1000, X train bigrams freq 1
000, X test bigrams freq 1000, 'L 1000 concat')
    PCA
                                                                          In [220]:
    from sklearn import decomposition
                                                                          In [223]:
    colors = ['#990033', '#CC00FF', '#330066', '#00FF33', '#FF9900',
               '#3366FF', '#333300', '#999999', '#336666', '#FF33CC', '#009900', '#
FF9999']
                                                                          In [716]:
     # Прогоним встроенный в sklearn PCA
    pca = decomposition.PCA(n components=2)
    pca.fit(X train features norm)
    X pca = pca.transform(X train features norm)
    colors = ['#990033', '#CC00FF', '#330066', '#00FF33', '#FF9900',
               '#3366FF', '#333300', '#999999', '#336666', '#FF33CC', '#009900', '#
FF9999'1
     # И нарисуем получившиеся точки в нашем новом пространстве
    for i, color, author in zip(range(0, 13), colors, author coded.keys()):
        plt.plot(X pca[y == i+1, 0], X pca[y == i+1, 1], color, label=str(author))
    plt.legend(loc=0);
                                                                          In [222]:
     import matplotlib.pyplot as plt
```

```
import seaborn as sns; sns.set(style='white')
%matplotlib inline
from sklearn import decomposition
from sklearn import datasets
from mpl toolkits.mplot3d import Axes3D
                                                                      In [236]:
pca = decomposition.PCA(n components=3)
XX = pca.fit transform(X train scale)
XX[y == 1, 2].mean()
                                                                      Out[236]:
2.8819675671186733
                                                                      In [230]:
def plot pca 3d(X in, y):
    pca = decomposition.PCA(n components=3)
    X = pca.fit transform(X in)
    print('Projecting %d-dimensional data to 3D')
    # Заведём красивую трёхмерную картинку
    fig = plt.figure(1, figsize=(6, 5))
    plt.clf()
    ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azim=134)
    plt.cla()
    for name, label in zip(colors, range(1,13)):
        ax.text3D(X[y == label, 0].mean(),
                  X[y == label, 1].mean() + 1.5,
                  X[y == label, 2].mean(), name,
                  horizontalalignment='center',
                  bbox=dict(alpha=.5, edgecolor='w', facecolor='w'))
    # Поменяем порядок цветов меток, чтобы они соответствовали правильному
```

```
\# y_{clr} = np.choose(y, [1, 2, 0]).astype(np.float)
    ax.scatter(X[:, 0], X[:, 1], X[:, 2], c=y, cmap=plt.cm.spectral)
    ax.w xaxis.set ticklabels([])
    ax.w yaxis.set ticklabels([])
    ax.w zaxis.set ticklabels([])
                                                                      In [238]:
plot pca 3d(X train scale, y train)
Projecting %d-dimensional data to 3D
                                                                      In [221]:
def plot pca(X, y):
    pca = decomposition.PCA(n components=2)
    X reduced = pca.fit transform(X)
    print('Projecting %d-dimensional data to 2D' % X.shape[1])
    plt.figure(figsize=(12,10))
    plt.scatter(X reduced[:, 0], X reduced[:, 1], c=y,
                edgecolor='none', alpha=0.7, s=40,
                cmap=plt.cm.get cmap('nipy spectral', 12))
    plt.colorbar()
    plt.title('Authors. PCA projection')
                                                                      In [732]:
plot pca(X test features norm, y test)
Projecting 38-dimensional data to 2D
                                                                      In [734]:
from sklearn.manifold import TSNE
tsne1 = TSNE(random state=17)
X tsne = tsne1.fit transform(X_train_features_norm)
plt.figure(figsize=(12,10))
```

```
plt.scatter(X tsne[:, 0], X tsne[:, 1], c=y train,
                 edgecolor='none', alpha=0.7, s=40,
                cmap=plt.cm.get cmap('nipy spectral', 12))
    plt.colorbar()
    plt.title('Authors. t-SNE projection')
                                                                           Out[734]:
    Text(0.5,1,'Authors. t-SNE projection')
                                                                           In [735]:
    pca = decomposition.PCA(n components=2)
    pca.fit(X train features norm)
    X pca = pca.transform(X train features norm)
    pca = decomposition.PCA(n components=2)
    pca.fit(X test features norm)
    X test pca = pca.transform(X test features norm)
                                                                           In [749]:
    def count pca(data):
        pca = decomposition.PCA(n components=2)
        pca.fit(data)
        X pca = pca.transform(data)
        return X pca
                                                                           In [739]:
    from sklearn.tree import DecisionTreeClassifier
    from sklearn.metrics import accuracy score, roc auc score
                                                                           In [745]:
    logr = LogisticRegression().fit(X pca, y train)
    logr.score(X test pca, y test)
                                                                           Out[745]:
    0.225833333333333333
                                                                           In [226]:
    def scaler (X train features, X test features, X train bigrams freq, X test big
rams freq):
```

```
# объединяем признаки и биграммы
        X train concat = concat features bigrams(X train features, X train bigrams
freq)
        X test concat = concat features bigrams(X test features, X test bigrams fr
eq)
        # скалируем
        scaler = preprocessing.StandardScaler()
        X train standard = scaler.fit transform(X train concat)
        X test standard = scaler.fit transform(X test concat)
        # svc = SVC(kernel='linear').fit(X train standard, y train)
        return X train standard, X test standard
                                                                          In [228]:
    X train scale, X test scale = scaler(X train features 20000, X test features 2
0000, X train bigrams freq, X test bigrams freq)
                                                                          In [128]:
    X train scale 20000, X test scale 20000 = scaler(X train features 20000, X tes
t features 20000, X train bigrams freq 20000, X test bigrams freq 20000)
    X train scale 10000, X test scale 10000 = scaler(X train features 10000, X tes
t features 10000, X train bigrams freq 10000, X test bigrams freq 10000)
    X train scale 5000, X test scale 5000 = scaler(X train features 5000, X test f
eatures 5000, X train bigrams freq 5000, X test bigrams freq 5000)
    X train scale 1000, X test scale 1000 = scaler(X train features 1000, X test f
eatures 1000, X train bigrams freq 1000, X test bigrams freq 1000)
                                                                          In [129]:
    X train scale 500, X test scale 500 = scaler(X train features 500, X test feat
ures 500, X train bigrams freq 500, X test bigrams freq 500)
    X train scale 200, X test scale 200 = scaler(X train features 200, X test feat
ures 200, X train bigrams freq 200, X test bigrams freq 200)
    X train scale 100, X test scale 100 = scaler(X train features 100, X test feat
ures 100, X train bigrams freq 100, X test bigrams freq 100)
                                                                          In [748]:
    logr = LogisticRegression().fit(X train_scale, y_train)
    logr.score(X test scale, y test)
                                                                          Out[748]:
```

```
0.774444444444445
                                                                          In [750]:
    logr = LogisticRegression().fit(count pca(X train scale), y train)
    logr.score(count pca(X test scale), y test)
                                                                          Out[750]:
    0.284722222222222
                                                                          In [751]:
    svc = SVC(kernel='linear').fit(count pca(X train scale), y train)
    svc.score(count pca(X test scale), y test)
                                                                          Out[751]:
    0.298055555555556
                                                                          In [752]:
    logr = SVC(kernel='linear').fit(X train scale, y train)
    logr.score(X test scale, y test)
                                                                          Out[752]:
    0.7475
                                                                          In [755]:
    def plot decision regions(X,y,classifier,test idx=None,resolution=0.02):
        # Initialise the marker types and colors
        markers = ('s','x','o','^','v', '!', '*', '-', '.', '$', '#', '@')
        colors = ('#990033', '#CC00FF', '#330066', '#00FF33', '#FF9900',
               '#3366FF', '#333300', '#999999', '#336666', '#FF33CC', '#009900', '#
FF9999')
        color Map = ListedColormap(colors[:len(np.unique(y))]) #we take the color
mapping correspoding to the
                                                                 #amount of classes
in the target data
        # Parameters for the graph and decision surface
        x1 min = X[:,0].min() - 1
        x1 max = X[:,0].max() + 1
        x2 min = X[:,1].min() - 1
```

x2 max = X[:,1].max() + 1

```
xx1, xx2 = np.meshgrid(np.arange(x1 min,x1 max,resolution),
                            np.arange(x2 min,x2 max,resolution))
    Z = classifier.predict(np.array([xx1.ravel(),xx2.ravel()]).T)
    Z = Z.reshape(xx1.shape)
    plt.contour(xx1,xx2,Z,alpha=0.4,cmap = color Map)
    plt.xlim(xx1.min(),xx1.max())
    plt.ylim(xx2.min(),xx2.max())
    # Plot samples
    X_test, Y_test = X[test_idx,:], y[test_idx]
    for idx, cl in enumerate(np.unique(y)):
        plt.scatter(x = X[y == cl, 0], y = X[y == cl, 1],
                    alpha = 0.8, c = color Map(idx),
                    marker = markers[idx], label = cl
                   )
                                                                       In [759]:
X combined sepal standard = np.vstack((X train scale, X test scale))
Y combined sepal = np.hstack((y train, y test))
                                                                       In [761]:
from matplotlib.colors import ListedColormap
                                                                       In [764]:
from sklearn.linear model import LogisticRegression
from sklearn.metrics import accuracy score
from sklearn.learning curve import validation curve
C \text{ param range} = [0.001, 0.01, 0.1, 1, 10, 100]
j = 0
for i in C_param_range:
```

```
# Apply logistic regression model to training data
   lr = LogisticRegression(penalty = '12', C = i,random state = 0)
   lr.fit(X train scale, y train)
   print ('Score = ', lr.score(X test scale, y test), ' . C = ', i)
Score = 0.766388888888888 . C = 0.01
Score = 0.775555555555555 . C = 0.1
Score = 0.77444444444444 . C = 1
Score = 0.77277777777777 . C = 10
Score = 0.7658333333333333 . C = 100
                                                                 In [765]:
tsne1 = TSNE(random state=17)
X train tsne = tsne1.fit transform(X train features norm)
X test tsne = tsne1.fit transform(X test features norm)
Создание конвейерной обработки
                                                                 In [112]:
from sklearn.pipeline import Pipeline
                                                                 In [113]:
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature extraction.text import TfidfTransformer
from sklearn.naive bayes import MultinomialNB
                                                                 In [114]:
text clf = Pipeline([('vect', CountVectorizer()),
                    ('tfidf', TfidfTransformer()),
                    ('clf', MultinomialNB()),
])
                                                                 In [116]:
start time = time.time()
text_clf = text_clf.fit(X_train_20000, y_train)
```

```
predicted = text clf.predict(X test 20000)
    print('Acc: ', np.mean(predicted == y test))
    print("--- %s seconds ---" % (time.time() - start time))
    Acc: 0.68833333333333334
    --- 170.39190793037415 seconds ---
                                                                           In [117]:
    from sklearn.linear model import SGDClassifier
                                                                           In [118]:
    text clf = Pipeline([('vect', CountVectorizer()),
                          ('tfidf', TfidfTransformer()),
                          ('clf', SGDClassifier(loss='hinge', penalty='12',
                                                alpha=1e-3, n iter=5, random state=
42)),
    1)
    _ = text_clf.fit(X_train_20000, y_train)
    predicted = text_clf.predict(X test 20000)
    np.mean(predicted == y test)
    C:\Users\Yauheniya\Anaconda3\lib\site-packages\sklearn\linear model\stochastic
gradient.py:117: DeprecationWarning: n iter parameter is deprecated in 0.19 and wi
11 be removed in 0.21. Use max iter and tol instead.
      DeprecationWarning)
                                                                           Out[118]:
    0.7091666666666666
                                                                           In [123]:
    text clf = Pipeline([('vect', CountVectorizer()),
                          ('tfidf', TfidfTransformer()),
                          ('clf', SGDClassifier(loss='hinge', penalty='12',
                                                alpha=1e-3, n iter=5, random state=
42)),
    1)
    _ = text_clf.fit(X_train_200, y_train)
    predicted = text clf.predict(X test 200)
    np.mean(predicted == y test)
```

C:\Users\Yauheniya\Anaconda3\lib\site-packages\sklearn\linear\_model\stochastic \_gradient.py:117: DeprecationWarning: n\_iter parameter is deprecated in 0.19 and wi ll be removed in 0.21. Use max iter and tol instead.

```
DeprecationWarning)
```

```
Out[123]:
    0.3319444444444443
                                                                          In [129]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=1000, rand
om state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[129]:
    0.73222222222222
                                                                          In [130]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=1000, rand
om state=42)
    clf.fit(X train scale 10000, y train)
    clf.score(X test scale 10000, y test)
                                                                          Out[130]:
    0.72083333333333333
                                                                          In [131]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=1000, rand
om state=42)
    clf.fit(X train scale 5000, y train)
    clf.score(X test scale 5000, y test)
                                                                          Out[131]:
    0.6636111111111112
                                                                          In [137]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=50, random
state=42)
    clf.fit(X train scale 1000, y train)
    clf.score(X test scale 1000, y test)
                                                                          Out[137]:
    0.4655555555555556
```

```
In [136]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=50, random
state=42)
    clf.fit(X train scale 500, y train)
    clf.score(X test scale 500, y test)
                                                                          Out[136]:
    0.394722222222222
                                                                          In [135]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=50, random
state=42)
    clf.fit(X train scale 200, y train)
    clf.score(X test scale 200, y test)
                                                                          Out[135]:
    0.2547222222222224
                                                                          In [124]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-3, max iter=500, random
state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[124]:
    0.74777777777778
                                                                          In [238]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-3, max iter=3, random s
tate=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X_test_scale_20000, y_test)
                                                                          Out[238]:
    0.75888888888888888
                                                                          In [125]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=5, random
state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[125]:
```

```
0.725
                                                                          In [241]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-10, max iter=5, random
state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[241]:
    0.766388888888888
                                                                          In [244]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-25, max iter=3, random
state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[244]:
    0.744444444444445
                                                                          In [243]:
    clf = SGDClassifier(loss='log', penalty='12', alpha=1e-25, max iter=5, random
state=42)
    clf.fit(X train scale 20000, y train)
    clf.score(X test scale 20000, y test)
                                                                          Out[243]:
    0.76277777777778
                                                                          In [246]:
    text clf = Pipeline([('vect', CountVectorizer()),
                          ('tfidf', TfidfTransformer()),
                          ('clf', SVC(kernel = 'linear', C=0.01, gamma = 0.01)),
    1)
    = text clf.fit(X train, y train)
    predicted = text clf.predict(X test)
    np.mean(predicted == y test)
                                                                          Out[246]:
```

0.5311111111111111