



ПРОЕКТНАЯ ДЕКЛАРАЦИЯ WWW.NOVOSTROYNN.RU.  
ЗАСТРОЙЩИК ООО "КРАСНАЯ ПОЛЯНА"



**ОДНУШКА 39 м<sup>2</sup>**  
**ТОЛЬКО В ИЮЛЕ 1,860 МЛН РУБ.**  
~~2,070 МЛН РУБ.~~



Информационное агентство INNOV.RU | Понедельник, 10 июля 2017 г. 13:51

Яндекс



**INNOV**  
РОССИЙСКИЙ БИЗНЕС ON-LINE

Иннов: электронный научный журнал

Главная страница журнала

Экономические науки

Технические науки

О журнале

Редакция

Общая лента

Выпуски

Опубликовать статью. Авторам

## Разработка вопросно-ответной системы с нейросетевым обучением на базе современных свободных технологий

### Design of Question Answering System Based on Neural Networks and Modern F Technologies

15.06.17 10:45

233

**Выходные сведения:** Науменко А.М., Шелудько С.Д., Юлдашев Р.Ю., Хлебников Н.О., Радугин В.Ю. Разработка вопросно-ответной системы с нейросетевым обучением на базе современных свободных технологий // Иннов: электронный научный журнал, 2017. №2 (31). URL: <http://www.innov.ru/science/tech/razrabotka-voprosno-otvetnoy-sistem/>

#### Авторы:

Науменко А.М.1, Шелудько С.Д.2, Юлдашев Р.Ю. 3, Хлебников Н.О. 4, Радугин В.Ю.5

1 студент 4-го курса бакалавриата по направлению «Информационные системы и технологии», ФГАОУ ВО Национальный исследовательский ядерный университет «МИФИ», Москва, Российская Федерация (115409, г. Москва, Каширское ш., 31), e-mail: [naumenko.mephi@gmail.com](mailto:naumenko.mephi@gmail.com).

2 студент 4-го курса бакалавриата по направлению «Информационные системы и технологии», ФГАОУ ВО Национальный исследовательский ядерный университет «МИФИ», Москва, Российская Федерация (115409, г. Москва, Каширское ш., 31), e-mail: [sheludko.serg@gmail.com](mailto:sheludko.serg@gmail.com).

3 студент 4-го курса бакалавриата по направлению «Информационные системы и технологии», ФГАОУ ВО Национальный исследовательский ядерный университет «МИФИ», Москва, Российская Федерация (115409, г. Москва, Каширское ш., 31), e-mail: [romanyuldashev@gmail.com](mailto:romanyuldashev@gmail.com).

4 студент 4-го курса бакалавриата по направлению «Информационные системы и технологии», ФГАОУ ВО Национальный исследовательский ядерный университет «МИФИ», Москва, Российская Федерация (115409, г. Москва, Каширское ш., 31), e-mail: [nikolay.khlebnikoff@gmail.com](mailto:nikolay.khlebnikoff@gmail.com).

5 к.т.н., доцент кафедры финансового мониторинга, ФГАОУ ВО Национальный исследовательский ядерный университет «МИФИ», Москва, Российская Федерация (115409, г. Москва, Каширское ш., 31), e-mail: [vyradugin@mephi.ru](mailto:vyradugin@mephi.ru).

#### Authors:

Naumenko A.M.1, Sheludko S.D.2, Uldashev R.Yu. 3, Hlebnikov N.O. 4, Radygin V.Yu.5

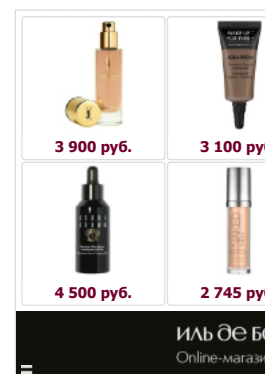
1 fourth year BA-student, specialty "Information systems and



#### Пластиковый погреб TINGARD!

TINGARD это надежно! Новая акция: «ПОТРЯСАЮЩАЯ ЦЕНА»! Подробнее об акции!

Яндекс.Директ



technologies”, National Research Nuclear University MEPhI (Moscow Engineering Physics Institute), Moscow, Russian Federation (115409, Moscow, Kashirskoe shosse, 31), e-mail: naumenko.mephi@gmail.com.

2 fourth year BA-student, specialty “Information systems and technologies”, National Research Nuclear University MEPhI (Moscow Engineering Physics Institute), Moscow, Russian Federation (115409, Moscow, Kashirskoe shosse, 31), e-mail: sheludko.serg@gmail.com.

3 fourth year BA-student, specialty “Information systems and technologies”, National Research Nuclear University MEPhI (Moscow Engineering Physics Institute), Moscow, Russian Federation (115409, Moscow, Kashirskoe shosse, 31), e-mail: romanyuldashev@gmail.com.

4 fourth year BA-student, specialty “Information systems and technologies”, National Research Nuclear University MEPhI (Moscow Engineering Physics Institute), Moscow, Russian Federation (115409, Moscow, Kashirskoe shosse, 31), e-mail: nikolay.khlebnikoff@gmail.com.

5 Ph.D., assistant professor of dept. Financial Monitoring, National Research Nuclear University MEPhI (Moscow Engineering Physics Institute), Moscow, Russian Federation (115409, Moscow, Kashirskoe shosse, 31), e-mail: vyradygin@mephi.ru.

**Ключевые слова:** машинное обучение, вопросно-ответная система, синтаксическое дерево, word2vec, nltk, pymorphy2, python, нейросетевой анализ данных

**Keyword:** machine learning, question answering system, syntax tree, word2vec, nltk, pymorphy2, python, neural network analysis

**Аннотация:** Данная статья посвящена исследованию вопроса разработки автоматизированной вопросно-ответной системы на базе современных открытых технологий семантического сжатия текста. Выполнен подробный анализ существующих подходов к решению задачи извлечения информации из больших объемов текста. Показано отсутствие готовых решений. Исследован вопрос построения автоматизированного решения задачи часто задаваемых вопросов на основе больших массивов данных. Выполнен подробный анализ современных научных работ и технологий в данной области. Рассмотрены основные модели частично автоматизированного и полностью автоматизированного поиска ответов на вопросы. Выявлены недостатки применяемых на сегодняшний день подходов и подчеркнута ограниченность всех решений, имеющих доступ. Рассмотрены основные семантические модели представления текста. Подробно раскрыты идеи поиска на основе синтаксических деревьев и векторного описания вопросов. Исследованы технологии word2vec, NLTK, Python и их преимущества при реализации вопросно-ответной системы. Раскрыты основные идеи нейросетевого анализа. Приведено подробное описание разработанного программного продукта. Показаны примеры найденных решений для семантически сложных вопросов. Описан полученный полноценный программный продукт. Подчеркнута возможность применения отработанных технологий к решению аналогичных задач. Отмечено экономическое преимущество исследованных технологий, достигаемое использованием бесплатного программного обеспечения.

**Annotation:** This article is devoted to design of automatic question-answering system based on free technologies of semantic compression. A detailed analysis of modern solutions of knowledge extraction from big text data sets is performed. Lack of ready-made solutions is shown. The ways to design of automatic frequently answering questions system are investigated. Solutions for limited questions sets is discussed. A detailed analysis of scientific works in this field is performed. The main models of partially automatic and fully automatic semantic questions analysis are described. The disadvantages of all popular methods to semantic representation of text phrases are revealed. The Continuous Bag of Words (CBOW) and Skip-gram are discussed. The detailed description of designed question answering software is shown. Examples of automatically matched and semantically difficult real questions are given. The advantages of developed product are emphasized. The possibilities of using free software for solving similar tasks are discussed. In conclusion, the economic benefits of using free software are emphasized.

## Введение

Анализ информации – это одна из наиболее востребованных задач, возникающих сегодня во всех областях деятельности. Современные объемы данных и требуемая скорость их обработки побуждают людей всё чаще использовать средства анализа, базирующиеся на компьютерных технологиях. Некоторые аспекты данной области на сегодняшний день обеспечены соответствующим программным обеспечением. Это, например, статистический анализ чисел, распознавание печатного текста и т.д. К сожалению, до сих пор остаётся ряд широко востребованных в повседневной жизни задач, не имеющих качественного и общедоступного автоматизирующего программного обеспечения. В настоящее время задача семантического сжатия больших объемов текстовой информации по определению является актуальной.

одним из наиболее значимых подмножеств которой является задача поиска в тексте ответа на заданный вопрос (задача полнотекстового поиска).

Одним из решений задачи полнотекстового поиска можно считать современные поисковые системы. Например, исследованию особенностей семантического поиска с использованием технологий Google посвящены работы В. Мала (V. Mala) [2], В.Н. Пху (V.N. Phu) [3]. Построение вопросно-ответной системы (question-answering system – QA system) для Википедии рассмотрено в работе Ф. Аббас (F. Abbas) [4].

К сожалению, результаты полнотекстового поиска, полученные с помощью механизмов современных поисковых систем, неудовлетворительны для тематических задач поиска. Причинами этого являются использование речевых конструкций естественного языка (например, вопросительных слов), высокая частота встречаемости их в просматриваемом тексте и другие проблемы, связанные с отсутствием семантического анализа просматриваемых данных.

В общем случае говорить об автоматизированном поиске ответов на вопросы нельзя, так как всесторонний поиск с использованием только компьютерных средств является AI-полной задачей [5], предполагающей разрабатываемый интеллект, сопоставимого по возможностям с человеческим. Тем не менее, данная задача может иметь определённые ограничения.

Альтернативой полноценного семантического поиска ответов на вопросы является подход, базирующийся на данных и наличии достаточной для выявления закономерностей подборки запросов к ним. Простейшим примером является решение задачи семантического поиска в виде набора часто задаваемых вопросов (frequently asked questions) исследуемой тематике. В данном случае есть два возможных направления автоматизации: решения с частичной исключением человеческого фактора и решения с полноценной компьютерной автоматизацией.

#### Частично автоматизированные подходы

Идея подхода с использованием часто задаваемых вопросов очень проста. При таком подходе пользователи осуществляют просмотр всего массива данных по искомой тематике, предлагается обратиться к краткой выдержке которой часто обращались другие посетители. Данная технология требует от пользователя времени на прочтение выявление среди них аналогичного своему, не гарантируя нахождения такового.

Решением задачи с частичной автоматизацией может быть подход на основе иерархических инструкций специалистом. Причем реализация данных инструкций может быть, как компьютеризованной, так и «бумажной» например, используется в системах телефонной поддержки клиентов банков или операторов сотовой связи. Предлагается выбрать тему своего вопроса из списка предложенных тем. Дальнейший поиск ответа осуществляется с оператором поддержки, имеющим инструкцию для решения часто возникающих ситуаций, заданную в виде дерева. Оператор, задавая пользователю вопросы и анализируя его ответы, продвигается по дереву сверху вниз и собирает сведения, полученные в конечном листе инструкции. Такая система эффективно отвечает на вопросы по организации и поддержке предполагают большие финансовые затраты на персонал.

#### Автоматизированные подходы

Автоматизированные подходы к задаче построения вопросно-ответной системы обычно также оперируют некоторыми массивами данных и подборкой заранее известных вопросов и ответов, позволяющей осуществить обучение с помощью работы Д. Бхардваджа (D. Bhardwaj) [6] рассматривается модель построения автоматизированной вопросно-ответной системы в формате FAQ, основанной на основе простейшего OR/AND-поиска и методов комбинаторики. Х. Баотьян (H. Baoti) используют для построения вопросно-ответной системы технологию нейронных сетей.

Особую сложность вопросу построения вопросно-ответной системы могут добавлять национальные особенности на сегодняшний день есть большой ряд узкоспециализированных работ, посвящённых разработке вопросно-ответных конкретных языковых групп. Например, в работе Силана А. (Saelan, A) [8] рассматривается построение вопросно-ответной системы на индонезийском языке. В работе Мегьюхота Х. (Meguehout, H) [9] показано построение вопросно-ответной системы на вьетнамском языке. Работы Медведа М. (Medved', M.) [10] и Фама С.Т. (Pham, S.T.) [11] посвящены исследованиям в данной области вьетнамского языков, соответственно.

Тем не менее, у большинства данных работ есть существенные недостатки. В основном это узко специализированные разработки, применение которых для задач других областей или других языковых групп невозможно. Существующие на сегодняшний день промышленные разработки являются преимущественно закрытыми проприетарными технологиями, недоступными для широкого использования. В какой-то мере, к готовым технологиям, используемым для поиска в тексте, можно отнести проект IBM Watson [12]. Данная разработка является дорогостоящим уникальным проектом, использование которого в построение вопросно-ответных систем современного интернет сообщества экономически обоснованным.

Таким образом, задача разработки полноценной вопросно-ответной системы является актуальной и востребованной. Исходя из востребованности данной тематики, в НИЯУ МИФИ была разработана вопросно-ответная система



себе свободные технологии семантической обработки текста, предоставляемые современными разработчиками, с актуальными алгоритмическими подходами обработки больших объемов текстовой информации

### Разработка вопросно-ответной системы на основе расширенной модели часто задаваемых вопросов

В основе разработанной системы лежит практика часто задаваемых вопросов в расширенном виде. В структуру вопросов по заданной тематике, которые могут задать пользователи. Каждому из них ставится в соответствие ответ, который могут соответствовать несколько вопросов. Таким образом, описанный подход сводит задачу нахождения ответа к поиску в базе, семантически близкого к заданному. Для того, чтобы решить эту задачу, формируются модели каждого вопроса из базы. В качестве моделей используются синтаксические деревья и геометрические представления рассматриваемой задачи сравнение подобных моделей является объективным показателем семантической близости.

Синтаксическое дерево — это построенный по определенному алгоритму граф, узлами которого являются предложения. Ребрам, соединяющим узлы, соответствует их синтаксическая связь. В узле дерева таксономические единицы, отдельные слова предложения, или функциональные единицы, сочетания слов, которые перестают выполнять синтаксическую функцию. Существуют четыре основных алгоритма расположения узлов в предложении: грамматика Теньера, грамматика зависимостей, грамматика непосредственно составляющих предложения — это алгоритм, в котором в качестве вершины дерева выступает член предложения, не являющийся отношением ни к одной другой синтаксической единице. В соответствии с грамматикой Теньера вершиной синтаксического дерева является глагол-сказуемое. Кроме того, вводятся понятия актанты — функциональные единицы, обязательные для сказуемого и сирконстанты — необязательной (факультативной) функциональной единицы. Грамматика зависимостей — это алгоритм, при котором в узлах дерева располагаются таксономические единицы. Вершиной дерева является глагол-сказуемое. В случае составного глагола, все связи в дереве подчинительные. Грамматика непосредственно составляющих — это алгоритм, в ходе выполнения которого каждая грамматическая единица делится на две более простые. Деление происходит вплоть до выделения в качестве узла отдельного слова, каждому узлу соответствует грамматическая единица, среди которых все части речи, а также именная и глагольная составляющие.

Представление слова в векторном виде — сопоставление слова из словаря геометрическому вектору в пространстве слов. В пространстве слов понимается пространство конечной размерности  $N$ , равной количеству всех представленных слов. Задачей определения семантической близости между словами занимается дистрибутивная семантика. Увеличение размерности векторного пространства способствует повышению точности определения смысловой близости. Однако некоторая критическая размерность, превышая которую, модель не приносит заметного увеличения точности. Обычно устанавливается диапазон от 100 до 1000. Любой алгоритм построения векторного пространства для максимизации косинусного сходства между векторами семантически близких слов. Косинусное сходство определяется по формуле:

$$similarity = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

где  $A$  и  $B$  — вектора, расстояние между которыми вычисляется,  $\theta$  — угол между ними. Одним из наиболее известных семантических анализаторов на сегодняшний день является word2vec — программное средство для построения векторного пространства, разработанное компанией Google в 2013 году [14].

Word2vec основан на двухслойной нейронной сети прямого распространения, поэтому у пользователей есть возможность обучить сеть на собственных текстовых корпусах и, таким образом, получить наиболее подходящую для решения задачи векторную модель. Результаты обучения модели зависят от выбранной пользователем модельной архитектуры. Реализованы два алгоритма обучения: Continuous Bag of Words (CBOW) и Skip-gram.

При использовании архитектуры CBOW алгоритм предсказывает слово, исходя из его контекста, т.е. анализируя находящиеся по левую и правую стороны от данного. При этом результат работы алгоритма не зависит от порядка слов в контексте. Входным элементом в нейронную сеть выступает набор контекстных векторов  $w(t-k), \dots, w(t-1), w(t+1), \dots, w(t+k)$ , где  $w(t)$  — вектор предсказанного на основе контекста слова. Архитектура Skip-gram отличается тем, что предсказывает набор слов вокруг, основываясь на данном слове. Входным вектором выступает  $w(t)$ , а выходным множеством  $M = \{w(t-k), \dots, w(t-1), w(t+1), \dots, w(t+k)\}$ , где  $M$  — множество векторов. Каждое слово, соответствующее вектору в множестве  $M$ , характеризует слово, соответствующее входному вектору. Схема работы алгоритмов CBOW и Skip-gram представлена на рисунке 1.

Работу word2vec можно разделить на пять этапов. На первом этапе происходит статистическая обработка текстового корпуса, то есть для каждого слова рассчитывается количество вхождений его в исходный корпус.

На втором этапе происходит сортировка слов по частоте вхождения, а также, в целях оптимизации работы с данными, выделяются гапаксы — слова, встречающиеся редко в сравнении с другими словами текста. Результаты сортировки сохраняются в хеш-таблице.

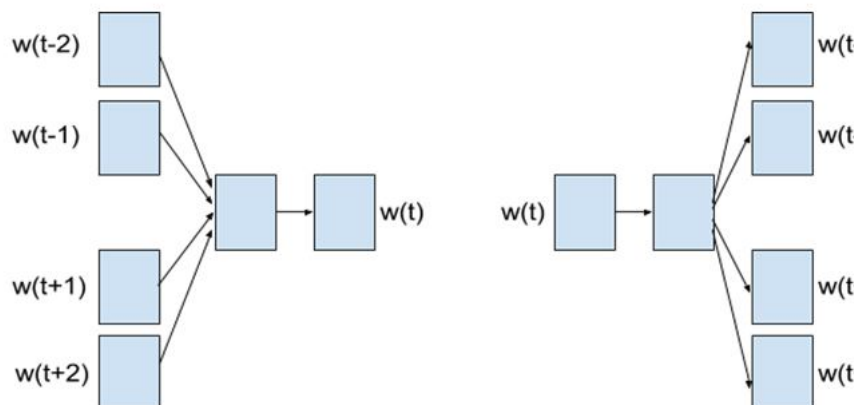


Рис. 1. Схема работы алгоритмов CBOW и Skip-gram

На третьем этапе для сжатия данных к полученной хеш-таблице применяется код Хаффмана — алгоритм оптимального кодирования. В результате применения данного алгоритма чаще встречающиеся слова кодируются меньшим кодом, а реже встречающиеся — большим.

Четвертый этап заключается в суб-сэмплировании самодостаточной выборки из текстового корпуса (например, абзаца). В ходе данного процесса из выборки удаляются наиболее часто встречающиеся слова, так как с ними не несется значимого смысла. Операция суб-сэмплирования применяется для уменьшения времени обучения модели.

На пятом этапе к получившейся выборке применяется один из алгоритмов обучения, рассмотренных выше: SBOV

Разработанная в НИЯУ МИФИ вопросно-ответная система включает в себя следующие элементы: управляющий модуль на Национальном корпусе русского языка векторная Skip-gram-модель в 300-мерном пространстве [15, 16], содержащая леммы; база ответов на вопросы, сами вопросы, а также их векторное представление.

Так как технология Word2vec применима как к отдельным словам, так и к текстам, в данном случае являющимся предложениями, то с её помощью на основе любой текстовой выборки можно построить соответствующее векторное представление. Каждое слово из предложения должно быть леммой. Векторную модель входят только леммы, то есть слова в словарной форме, а сама форма не имеет значения в векторном представлении. Некоторые части речи, такие как существительные, глаголы и прилагательные, в большинстве случаев являются значимыми, тогда как другие части речи, например, предлоги, союзы и местоимения не несут смысловой нагрузки. В анализе предложения к ним применяется фильтрация, в ходе которой из него исключаются так называемые «стоп-слова». Пример подготовки предложения к построению векторной модели изображен на рисунке 2.

Когда нужно оплачивать учебный отпуск работнику?



[ 'нужно' 'оплачивать' 'учебный' 'отпуск' 'работник' ]

Рис. 2. Пример анализа предложения

Управляющие скрипты написаны на языке Python. Лемматизация производится средствами библиотеки `rupt`. В настоящий момент способна обрабатывать до 100000 слов в секунду, при этом потребление оперативной памяти — 20 Мб [19]. Набор используемых стоп-слов взят из библиотеки Natural Language Toolkit [20], которая используется в связке с компьютерной лингвистикой и машинным обучением, и предназначена для обработки естественного языка.

Процесс реализации вопросно-ответной системы был разделён на несколько этапов. Так как система базируется на базе вопросов, то на первом этапе было осуществлено проектирование базы вопросов и ответов.

На втором этапе в базу были внесены векторные представления каждого вопроса, построенные по одному из алгоритмов. Затем этот же алгоритм применяется к вопросу, заданному пользователем системы. Таким образом, для определения семантической близости информации, на основе которой может делаться предположение о сходстве или различии заданного вопроса и вопросов из базы.

На заключительном этапе для каждого вопроса из базы и заданного пользователем вопроса были вычислены косинусы сходства. Результатом работы системы является ответ на тот вопрос из базы, косинусное сходство с которым заданного вопроса является наибольшим. Принцип работы системы изображена на рисунке 3.

Разработанная система оформлена в виде удалённого робота (бота), отвечающего на вопросы посредством чата. Для тестирования разработанной системы использовалась технология Microsoft Bot Framework Channel Emulator. Пример тестирования вопросно-ответной системы показан на рисунке 4.



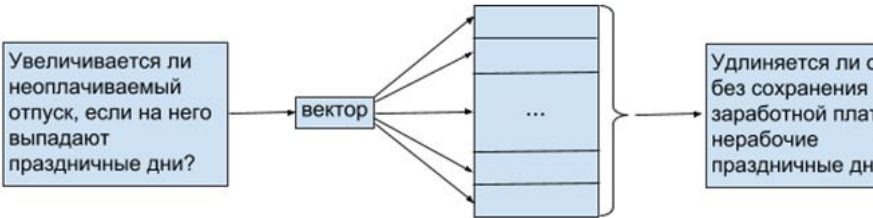


Рис 3. Схема принципа работы системы

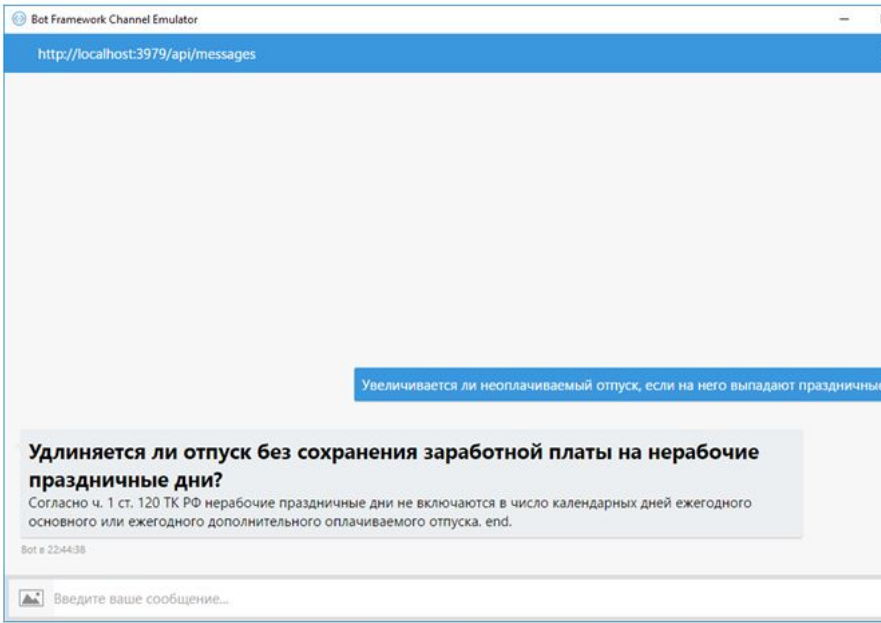


Рис 4. Пример работы системы

Тестирование системы выявило высокое качество поиска ответов, даже для вопросов семантически слож сопоставления с вопросами, для которых готовы ответы. Примеры заданных вопросов и поставленных им в соот базы FAQ показаны в таблице 1.

Примеры заданных вопросов и поставленных им в соответствие вопросов из базы FAQ

| Вопрос пользователя  | Вопрос в базе  |
|--|--|
| Должен ли сотрудник находиться на рабочем месте, если сейчас работы нет? | Обязан ли работник находиться во врем простоя на рабочем месте?                                  |
| Может ли сотрудник сменить банк, который выплачивает заработную плату?   | Имею ли я право поменять банк для выплат заработной платы?                                       |
| Включены ли в оклад надбавки за учёную степень?                          | Надбавки за учёную степень должны выплачиваться в виде премий или они уже включены в оклад?      |
| Может ли руководитель не проходить обучение по охране труда?             | Обязательно ли руководителю организации проходить обучение по охране труда?                      |
| Является ли договор без паспортных данных работника легитимным?          | Если в договоре отсутствуют паспортные данные работника, будет ли договор считаться заключённым? |
| Ограничен ли период действия ученического договора?                      | Установлен ли максимальный срок действия ученического договора?                                  |
| Можно ли уволить сотрудника, если тот подал                              | Может ли работник быть уволен :  |

|  |   |
|--|---|
| поддельный паспорт?  | предоставление поддельных документов при приеме на работу?  |
| Как индексируется зарплата сотрудника?                                       | Необходимо ли индексировать заработную плату работника?   |
| Предоставляется ли выходной в другой день, если в праздник был рабочий день? | Может ли быть предоставлен другой день отдыха работнику, работавшему в выходной или нерабочий праздничный день? |
| Должен ли работодатель давать сотрудникам отдыхать?                          | Обязан ли работодатель предоставлять работникам отпуска?  |

### Заключение

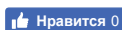
Разработанная в НИЯУ МИФИ система семантического анализа является полноценным продуктом и представляет эффективное решение задачи построения вопросно-ответных систем. Она способна найти применение во научных и бизнес задачах. Представленный подход к её реализации, основанный на технологиях word2vec, nlt, имеет значительные преимущества в сравнении с большинством используемых на настоящий момент разработок. Распространяются под свободной лицензией Apache 2.0. PyMorphy2 распространяется под свободной лицензией, данные технологии могут быть беспрепятственно использованы при разработке коммерческих продуктов, что обуславливает высокую экономическую эффективность созданного подхода. В итоге, можно отметить, что вопросно-ответных систем на основе нейросетей является перспективным направлением в практическом применении машинного обучения.

### Библиографический список

1. Ceglarek, D.: Semantic Compression for Text Document Processing. // Proceedings of Transactions on Computational Linguistics, Springer, Heidelberg, 2014. – C20–48.
2. Mala V., Lobiyal D.K. Semantic and keyword based web techniques in information retrieval // Proceedings of Computing and Automation (ICCA), International Conference, 2016 – C23–26.
3. Phu, V.N., Chau, V.T.N., Dat, N.D., Tran, V.T.N., Nguyen, T.A. A valences-totaling model for English sentiment classification and Information Systems, 2017 – C1–58.
4. Abbas, F., Malik, M.K., Rashid, M.U., Zafar, R. WikiQA - A question answering system on Wikipedia using freebase, // Proceedings of Sixth International Conference on Innovative Computing Technology (INTECH), 2016 – C185–193.
5. Raymond E.S., The New Hacker's Dictionary — MIT Press, 1996. —547 c.
6. Bhardwaj, D., Pakray, P., Bentham, J., Saha, S., Gelbukh, A. Question answering system for frequently asked questions. Proceedings, Vol. 1749, 2016 – C1–5.
7. Hu, B., Lu, Z., Li, H., Chen, Q. Convolutional neural network architectures for matching natural language sentences. Information Processing Systems, № 3, 2014 – C2042-2050.
8. Saetan, A., Purwarianti, A., Widyantoro, D.H. Question analysis for Indonesian comparative question // Journal of Indonesian Series, Vol. 801, Iss. 1, 2017 – C1–6.
9. Meguehout, H., Bouhadada, T., Laskri, M.T. Semantic role labeling for Arabic language using case-based reasoning. // Journal of Speech Technology № 2, 2017, – C1-10
10. Medved', M., Horák, A. AQA: Automatic question answering system for Czech // Lecture Notes in Computer Science (Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), Vol. 9924, 2016, – C270-278.
11. Pham, S.T., Nguyen, D.T. A Computational and Inferential Method for Analyzing the Semantics of Phrase and Sentence. Question Answering System Model (VietQASM) // Proceedings of Asia Modelling Symposium 2015 – Asia 9th International Mathematical Modelling and Computer Simulation, 2016, – C107-112.
12. Fan, J., Kalyanpur, A., Gondek, D.C., Ferrucci, D.A., Automatic knowledge extraction from documents // IBM Journal of Research and Development, Vol. 56, Iss. 3-4, 2012 – C5:1–5:10.
13. Касевич В.Б. Структура предложения. Элементы общей лингвистики. — М.: Наука, 1977. —183 с.
14. Mikolov T., Sutskever I., Chen K., Corrado G., Dean J. Distributed Representations of Words and Phrases and their Compositions // Proceedings of the 26th International Conference on Neural Information Processing Systems, USA, 2013 – C3111–3119.
15. Kutuzov A., Andreev I. Texts in, Meaning Out: Neural Language Models in Semantic Similarity Task for Russian // Proceedings of the 2015 Conference. Moscow, Russia, 2015 – C143–154.
16. <http://ling.go.mail.ru> – официальный сайт проекта RusVectores (дата последнего обращения 22.05.2017).
17. Manning C.D., Raghavan P., Schütze H. An Introduction to Information Retrieval, Cambridge University Press, Cambridge, 2008. —547 c.
18. Korobov M.: Morphological Analyzer and Generator for Russian and Ukrainian Languages // Proceedings of International Conference on Analysis of Images, Social Networks and Texts, 2016 – C107-112.
19. <http://pymorphy2.readthedocs.io> – официальный сайт проекта pymorphy2 (дата последнего обращения 22.05.2017).
20. <http://www.nltk.org/> – официальный сайт проекта Natural Language Toolkit (дата последнего обращения 22.05.2017).
21. <https://docs.microsoft.com/en-us/bot-framework/cognitive-services/bot-intelligence-overview> – обзор технологии Microsoft Bot Channel Emulator (дата последнего обращения 22.05.2017).

### References

1. Ceglarek, D.: Semantic Compression for Text Document Processing in Proceedings of Transactions on Computational Linguistics, Springer, Heidelberg, 2014. – pp. 20–48.
2. Mala V., Lobiya D.K. Semantic and keyword based web techniques in information retrieval in Proceedings of Computational Intelligence and Automation (ICCA), International Conference, 2016 – pp. 23–26.
3. Phu, V.N., Chau, V.T.N., Dat, N.D., Tran, V.T.N., Nguyen, T.A. A valences-totaling model for English sentiment classification in Proceedings of International Conference on Information Systems, 2017 – pp. 1–58.
4. Abbas, F., Malik, M.K., Rashid, M.U., Zafar, R. WikiQA - A question answering system on Wikipedia using freebase, in Proceedings of Sixth International Conference on Innovative Computing Technology (INTECH), 2016 – pp. 185–193.
5. Raymond E.S., The New Hacker's Dictionary — MIT Press, 1996. —547 p.
6. Bhardwaj, D., Pakray, P., Benthani, J., Saha, S., Gelbukh, A. Question answering system for frequently asked questions in Proceedings of International Conference on Information Systems, Vol. 1749, 2016 – pp. 1–5.
7. Hu, B., Lu, Z., Li, H., Chen, Q. Convolutional neural network architectures for matching natural language sentences in Proceedings of International Conference on Information Processing Systems, # 3, 2014 – pp. 2042-2050.
8. Saelan, A., Purwarianti, A., Widyantoro, D.H. Question analysis for Indonesian comparative question in Journal of Computational Linguistics, Vol. 801, Iss. 1, 2017 – pp. 1–6.
9. Meguehou, H., Bouhadada, T., Laskri, M.T. Semantic role labeling for Arabic language using case-based reasoning in Proceedings of International Journal of Speech Technology # 2, 2017, – pp. 1-10
10. Medved', M., Horák, A. AQA: Automatic question answering system for Czech in Lecture Notes in Computer Science (Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), Vol. 9924, 2016, – pp. 270-278.
11. Pham, S.T., Nguyen, D.T. A Computational and Inferential Method for Analyzing the Semantics of Phrase and Sentence in Proceedings of Question Answering System Model (VietQASM) in Proceedings of Asia Modelling Symposium 2015 – Asia 9th International Conference on Mathematical Modelling and Computer Simulation, 2016, – pp. 107-112.
12. Fan, J., Kalyanpur, A., Gondek, D.C., Ferrucci, D.A., Automatic knowledge extraction from documents in IBM Journal of Research and Development, Vol. 56, Iss. 3-4, 2012 – pp. 5:1–5:10.
13. Kasevich V.B. Structure of Sentence. The Elements of Common Linguistic. — Moscow, Science, 1977. —183 p.
14. Mikolov T., Sutskever I., Chen K., Corrado G., Dean J. Distributed Representations of Words and Phrases and their Compositions in Proceedings of the 26th International Conference on Neural Information Processing Systems, USA, 2013 – pp. 3111–3119
15. Kutuzov A., Andreev I. Texts in, Meaning Out: Neural Language Models in Semantic Similarity Task for Russian Dialog 2015 Conference. Moscow, Russia, 2015 – pp. 143–154.
16. <http://ling.go.mail.ru> – official site of RusVectores project (last access date 22.05.2017).
17. Manning C.D., Raghavan P., Schütze H. An Introduction to Information Retrieval, Cambridge University Press, Cambridge, 2008. —547 p.
18. Korobov M.: Morphological Analyzer and Generator for Russian and Ukrainian Languages in Proceedings of International Conference on Analysis of Images, Social Networks and Texts, 2016. <http://pymorphy2.readthedocs.io> – official site of pymorphy2 project (last access date 22.05.2017).
19. <http://www.nltk.org/> – official site of Natural Language Toolkit project (last access date 22.05.2017).
21. <https://docs.microsoft.com/en-us/bot-framework/cognitive-services/bot-intelligence-overview> – review of Microsoft Bot Framework technology (last access date 22.05.2017).



архив: [2013](#) [2012](#) [2011](#) [1999-2011](#) [новости ИТ](#) [гость портала 2013](#) [тема недели 2013](#) [поздравления](#)

иль де ботэ  
Online-магазин



3 900 руб.



999 руб.



999 руб.



1 099 руб.



1 940 руб.



1 940 руб.



Реклама на INNOV.RU | Партнеры | История компании | О компании | Услуги | Создать сайт | Стена памяти | Поиск

© 1996-2017 INNOV.RU (Иннов.ру) - информационное агентство, ООО «Иннов».

[\\* - правила пользования](#)

Свидетельство Управления Федеральной службы по надзору в сфере связи, информационных технологий и массовых коммуникаций по Нижегородской области ИА № ТУ 52-0604 от 29 февраля 2012 г.

ISSN: 2414-5122

Веб-студия «INNOV» - продвижение и [разработка сайта](#)