

IV межвузовская конференция  
*Бизнес-аналитика. Использование аналитической  
платформы Loginot (Deductor) в учебном процессе вуза*



## **Опыт использования серверных компонентов аналитической платформы Deductor Enterprise**

при решении учебных и практических задач  
на кафедре «Информатика и программное обеспечение»  
Брянского государственного технического университета

*Подвесовский Александр Георгиевич*  
*заведующий кафедрой, к.т.н., доцент*  
*apodv@tu-bryansk.ru*

*Лазерев Дмитрий Григорьевич*  
*к.т.н., доцент*  
*lagerevdg@mail.ru*

*Бабурин Артем Николаевич*  
*ассистент*  
*ababurin@bk.ru*

г. Москва, 27 июня 2017 г.



# Содержание

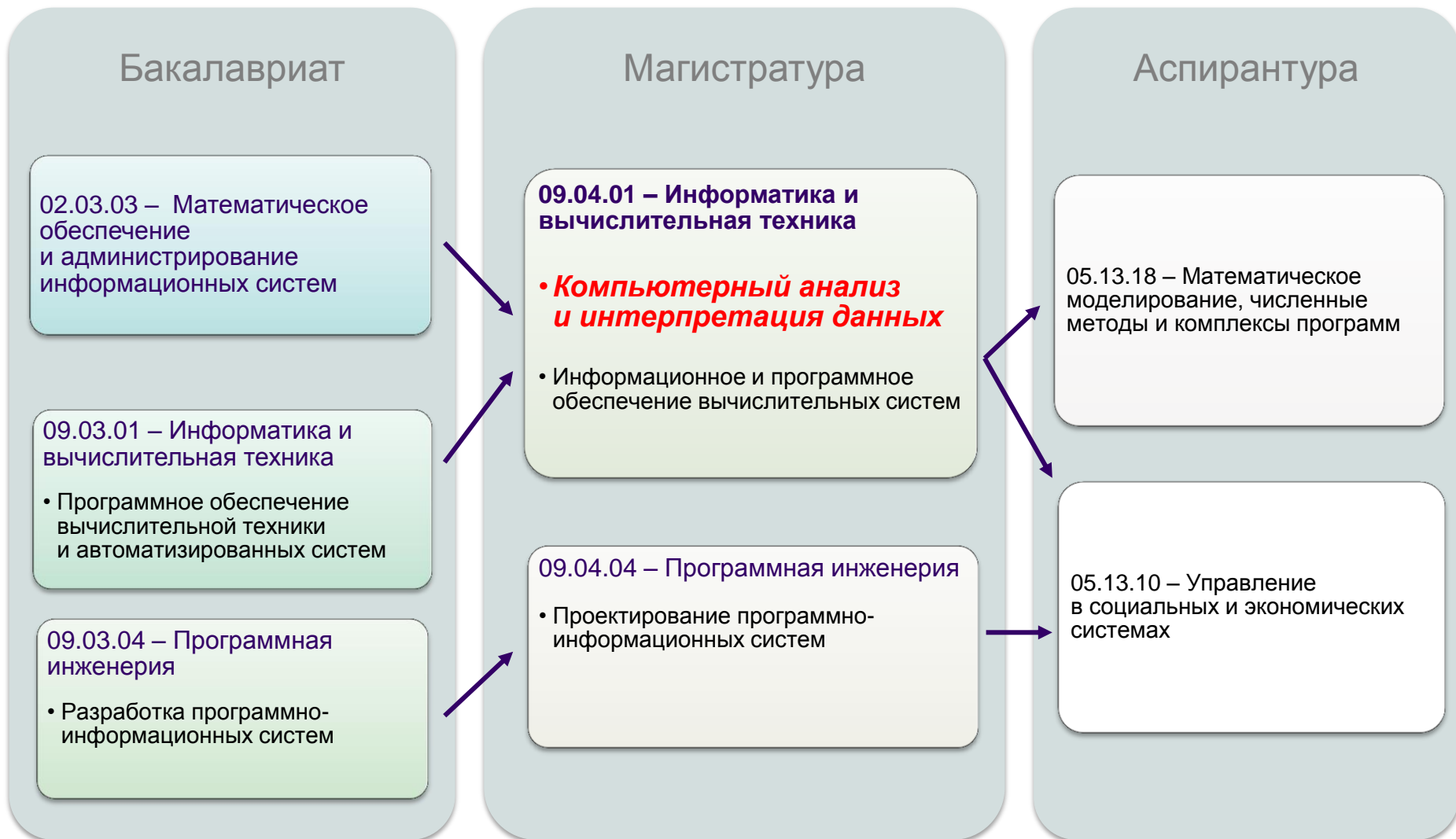
- Введение
- Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - Программная поддержка полного цикла социологического исследования
  - Поиск единомышленников в социальной сети VK
- Итоги работы. Дальнейшие планы и пожелания



# Содержание

- **Введение**
- Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - Программная поддержка полного цикла социологического исследования
  - Поиск единомышленников в социальной сети VK
- Итоги работы. Дальнейшие планы и пожелания

# Кафедра «Информатика и программное обеспечение»: образовательный процесс



# Магистерская программа «Компьютерный анализ и интерпретация данных»



- Реализуется в рамках направления «Информатика и вычислительная техника» с 2009 г.
- К настоящему моменту выпущено 34 магистра
- Основной целью является углубленная (на базе соответствующих направлений бакалавриата) подготовка профессиональных разработчиков программного обеспечения и системных аналитиков со специализацией в следующих областях
  - **обработка и анализ больших объемов данных, методы машинного обучения**
  - **модели и методы поддержки принятия решений**
  - **интеллектуальные системы на основе мягких вычислений**
  - **обработка и анализ изображений, машинное зрение**
  - **цифровая обработка сигналов**

# Ключевые дисциплины учебного плана



## 1. Базовые дисциплины направления

- Методы оптимизации
- Теория принятия решений
- Компьютерное моделирование
- Теория систем и системный анализ
- Технология проектирования, разработки и верификации программного обеспечения

## 2. Математический аппарат, методология и инструменты анализа данных

- Статистический анализ данных
- **Интеллектуальный анализ данных**
- Теория нейронных сетей
- Интеллектуальные системы

## 3. Специальные разделы и приложения методологии анализа данных

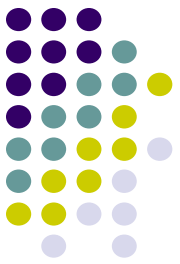
- **Хранилища данных**
- Системы с параллельной обработкой данных
- Обработка и анализ изображений
- Системы машинного зрения
- Цифровая обработка сигналов

# Интеллектуальный анализ данных: общие сведения о курсе



- *Цели курса*
  - изучение современных информационных технологий, предназначенных для интеллектуального анализа данных
  - формирование целостного представления об анализе и интерпретации данных, как о процессе поиска, так и о методологии применения скрытых в них закономерностей для достижения поставленных целей
- *Структура курса*
  - изучается на 1-м курсе во 2-м семестре
  - лекции – 17 часов
  - лабораторные работы – 34 часа
  - курсовой проект, экзамен
- *Теоретическая часть:*
  - структура и возможности аналитических систем
  - основные методы интеллектуального анализа и предобработки данных
  - процесс ETL
  - ансамбли моделей
- *Практическая часть:*
  - использование **Deductor Studio Academic / Professional / Enterprise**

# История использования платформы Deductor на кафедре



- С начала 2000-х – эпизодическое использование Deductor Studio Academic в рамках отдельных лабораторных работ по некоторым дисциплинам
- **2009 – начало систематического использования Deductor Studio Academic в курсе «Интеллектуальный анализ данных» в магистратуре**
- 2015 – начало регулярного взаимодействия с компанией BaseGroup в рамках партнерской программы
- 2016, март – приобретение Deductor Studio Professional
- **2016, июнь – участие с докладом в III межвузовской конференции**
  - <https://basegroup.ru/system/files/events/lagerev.pdf>
- **2016, август – получение лицензии на использование Deductor Studio Enterprise, Deductor Integration Server, Deductor Analytic Server**
- 2016, октябрь – прохождение доцентом Лагеревым Д.Г. тренинга «Разработка скоринговых моделей»
- 2016-2017 – обучение и сертификация преподавателей
  - **доцент Лагерева Д.Г.** – прохождение электронного курса, **сдача экзамена**
  - ассистент Бабурина А.Н. – прохождение электронного курса





# Содержание

- Введение
- **Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise**
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - Программная поддержка полного цикла социологического исследования
  - Поиск единомышленников в социальной сети VK
- Дальнейшие планы



# Наша команда



Бондарева Инна



Козлов Евгений



Герасимчук Иван



Тупикина Екатерина



Журин Владислав



Самородов Илья



# Содержание

- Введение
- **Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise**
  - **Выявление групп риска в рамках мониторинга наркоситуации в Брянской области**
  - Программная поддержка полного цикла социологического исследования
  - Поиск единомышленников в социальной сети VK
- Итоги работы. Дальнейшие планы и пожелания

# Проект «Мониторинг наркоситуации».

## Описание проекта



- Исходные положения
  - Анализ и оценка состояния наркоситуации в Брянской области
  - Проводится ежегодно по заказу УФСКН РФ по Брянской области
  - Исполнитель – кафедра социально-гуманитарных дисциплин Брянского филиала Российской академии народного хозяйства и государственной службы
    - Мы приняли участие благодаря имеющимся научным связям с данной кафедрой
  - Основной метод – социологический опрос населения, в том числе пациентов наркологических клиник
    - Единая форма анкеты (в рамках РФ)
    - Рекомендована статистическая обработка результатов анкетирования
  - Социологический опрос проводился в 2013-2017 гг.
    - 2013 г. – 1875 респондентов
    - 2014 г. – 2211 респондентов
    - 2015 г. – 1915 респондентов
    - 2016 г. – 1208 респондентов
    - 2017 г. – проект планируется продолжить, но данные пока недоступны
- Цель проекта – повышение эффективности обработки данных за счет применения методов интеллектуального анализа

# Проект «Мониторинг наркоситуации».

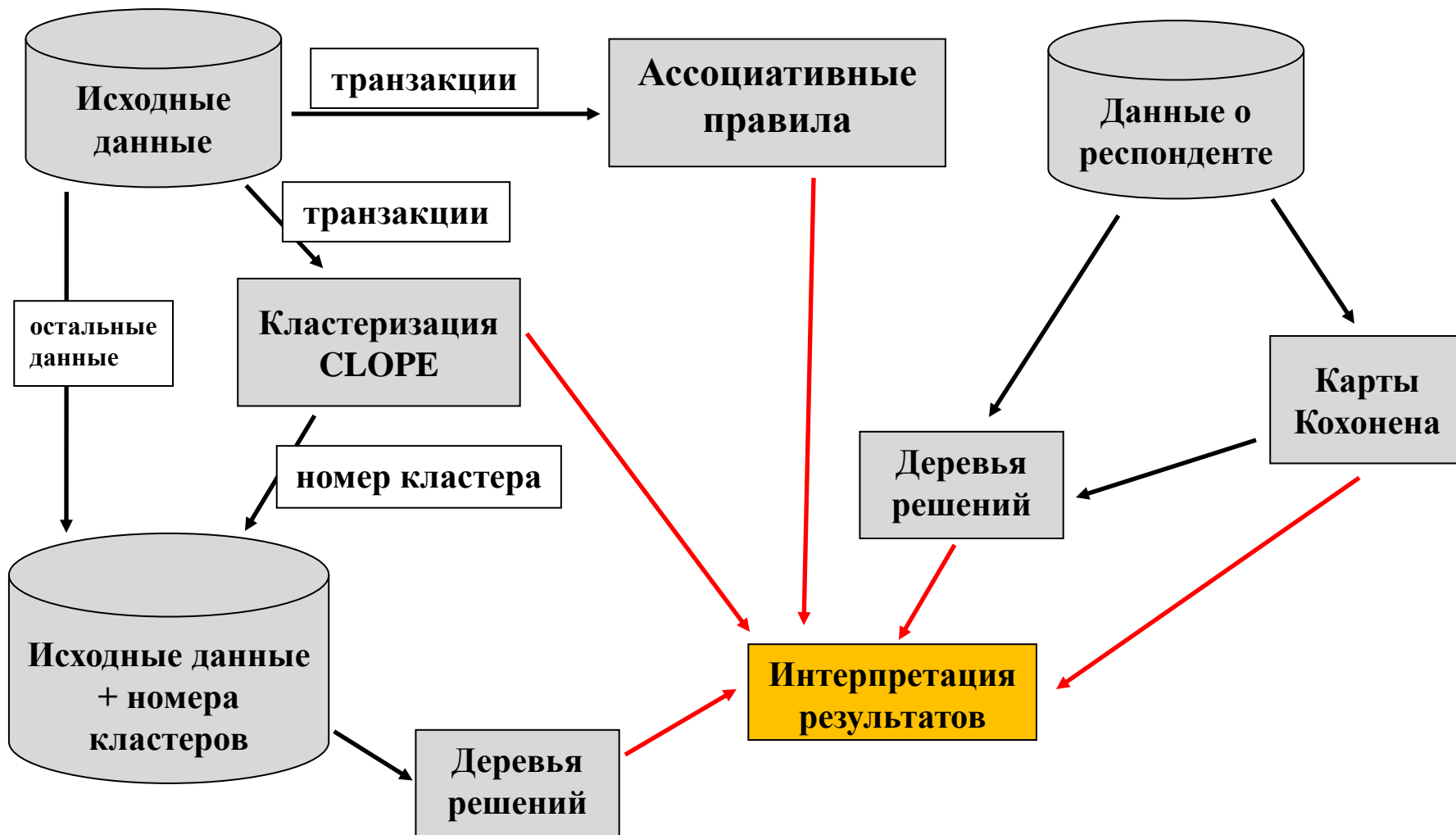
## Структура данных



- Анкета содержит следующие типы вопросов
  - Данные о респонденте (пол, возраст, образование)
  - Жизненные ориентиры респондента (наиболее острые проблемы, жизненные ценности, проведение свободного времени)
  - Вопросы, касающиеся здоровья респондента (оценка здоровья, наличие вредных привычек)
  - Отношение респондента к наркотикам и наркомании
- Структура анкеты в 2013-2016 гг. незначительно менялась. Общее число вопросов – от 37 до 42, среди которых
  - с ответами в категориальной шкале – от 32 до 36 вопросов
  - с ответами в порядковой шкале – от 1 до 2 вопросов
  - с ответами в числовой шкале – от 3 до 4 вопросов

# Проект «Мониторинг наркоситуации».

## Ансамбль моделей



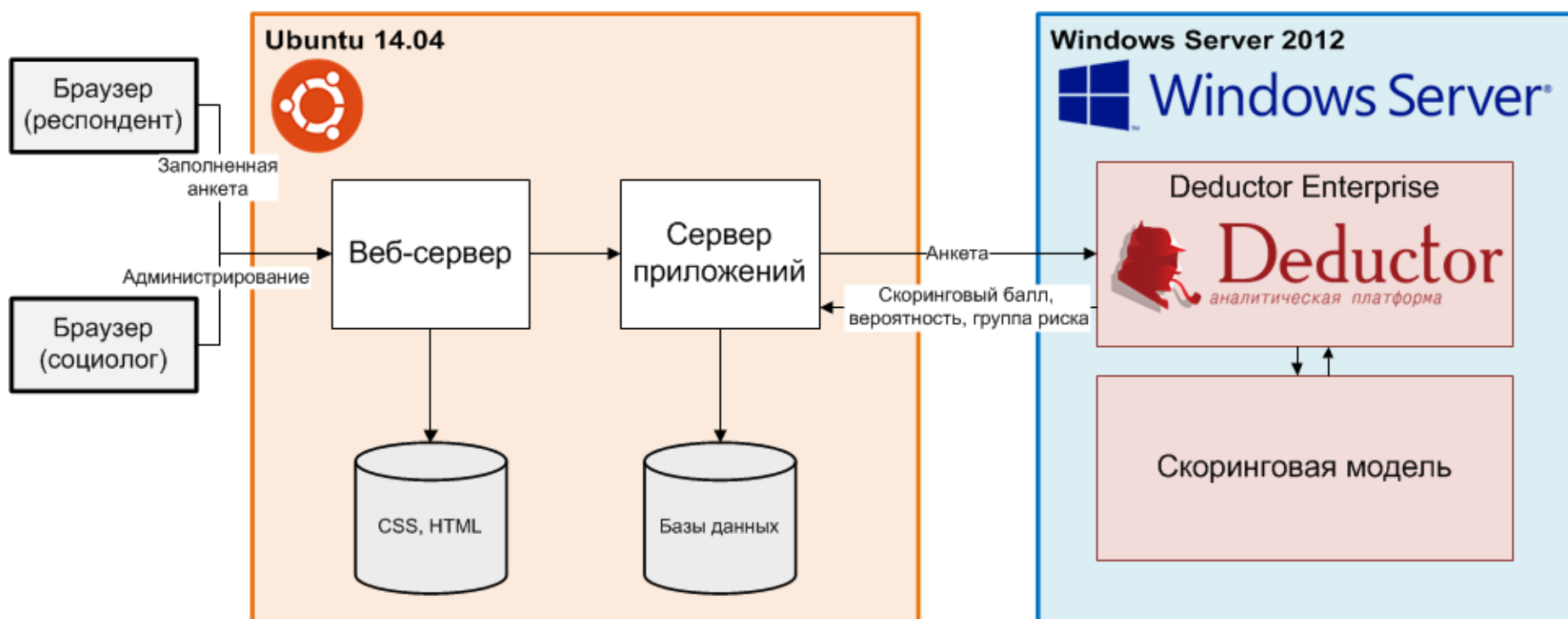


# Постановка задачи

- Проблемы:
  - социально приемлемые ответы респондентов (особенно подростков)
  - ложные ответы на вопросы, касающиеся наркотиков и проблемы наркомании
  - ...
- **Идея:** построить скоринговую модель, которая выделяла бы группы риска относительно наркозависимости на основе анкеты, не содержащей вопросов, напрямую связанных с наркотикам и проблемой наркомании
  - эти вопросы исключены из анкеты
  - добавлены фиктивные вопросы на тему компьютерных игр (для отвлечения внимания), которые не будут подавать в сценарий скоринговой модели



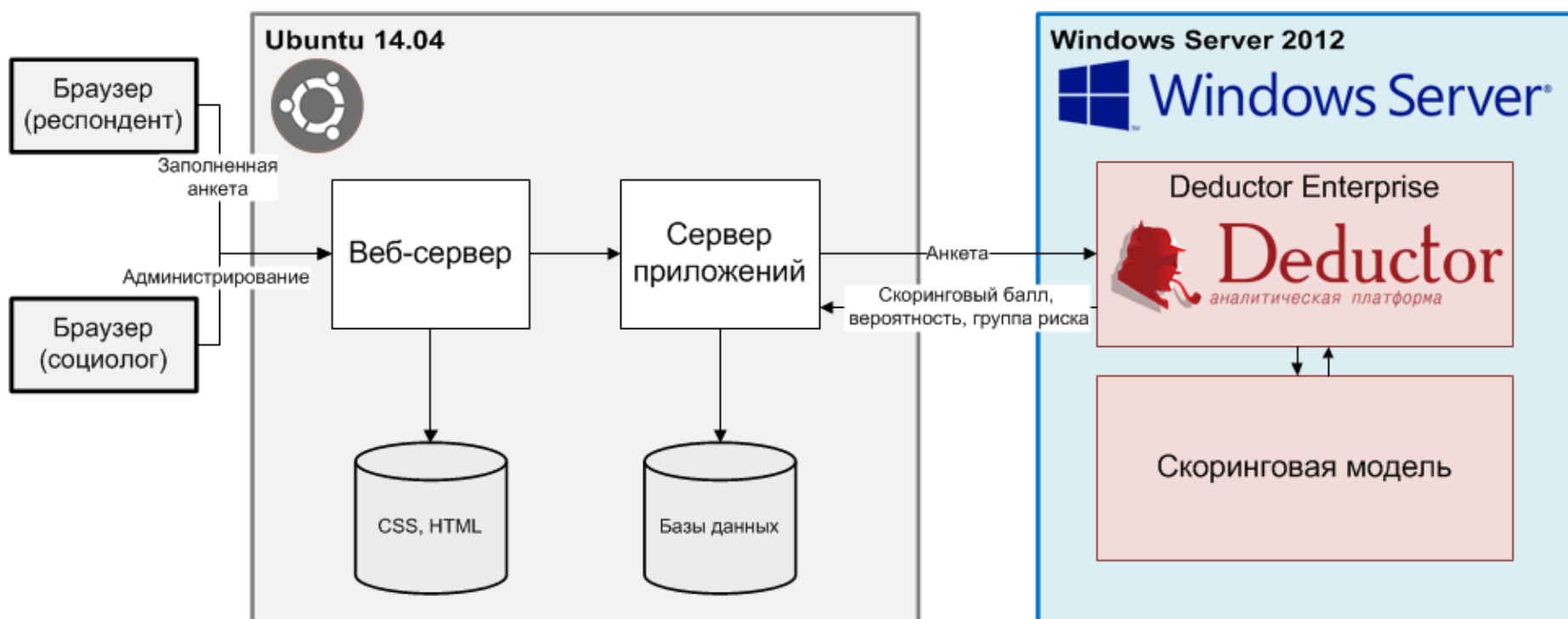
# Архитектура приложения







# Архитектура приложения



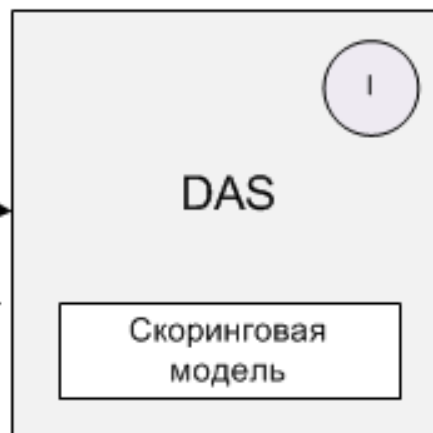
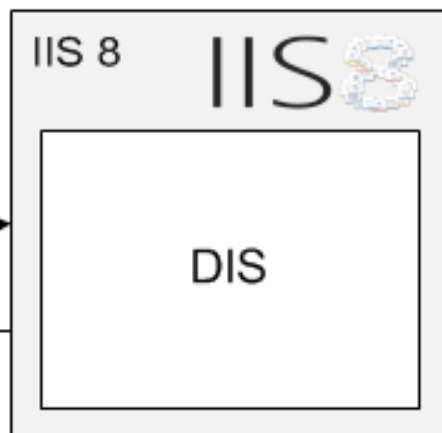
# Архитектура подсистемы взаимодействия с внешними клиентами



Windows Server 2012. deductor.iipo.tu-bryansk.ru

Deductor Server

Deductor Enterprise 5.3



SOAP-запрос

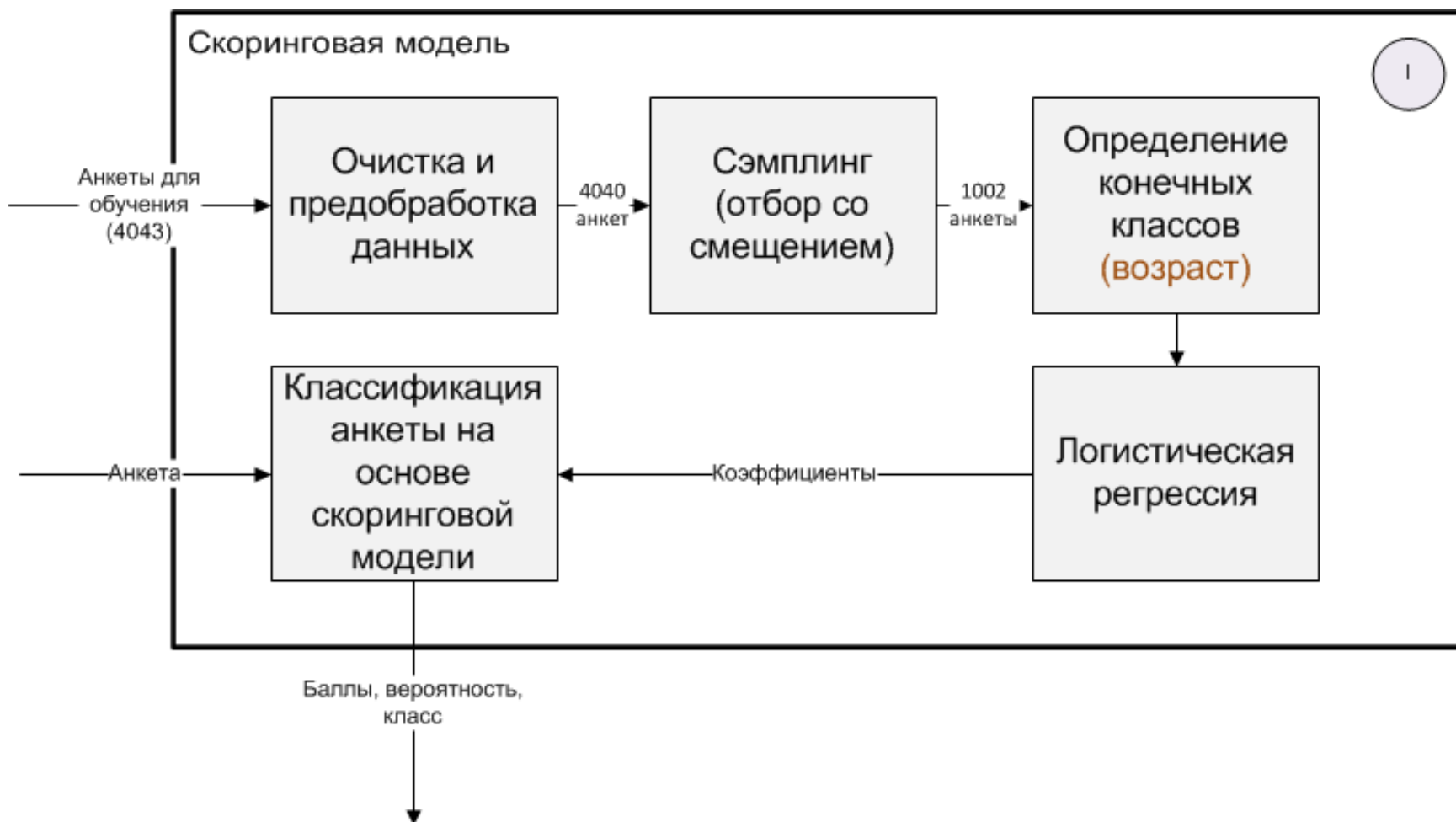
SOAP-ответ

XML

XML



# Структура скоринговой модели



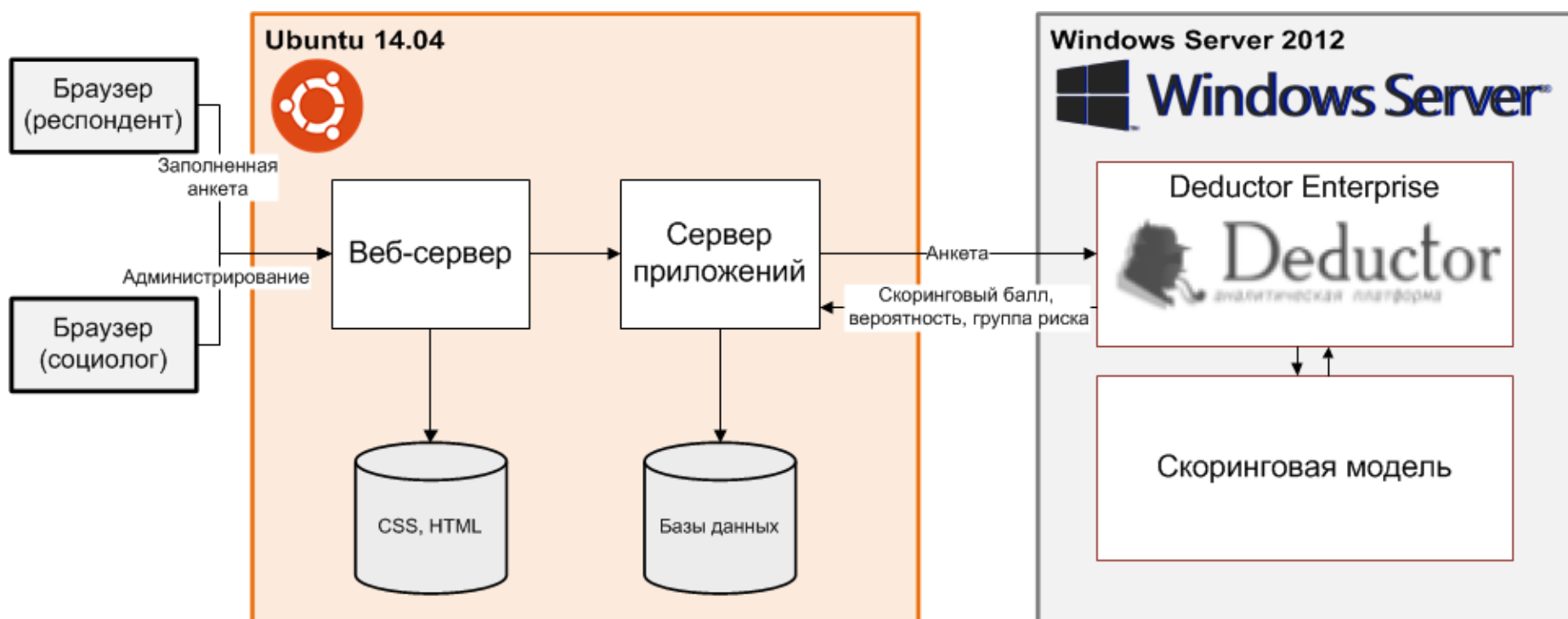


# Скоринговые баллы

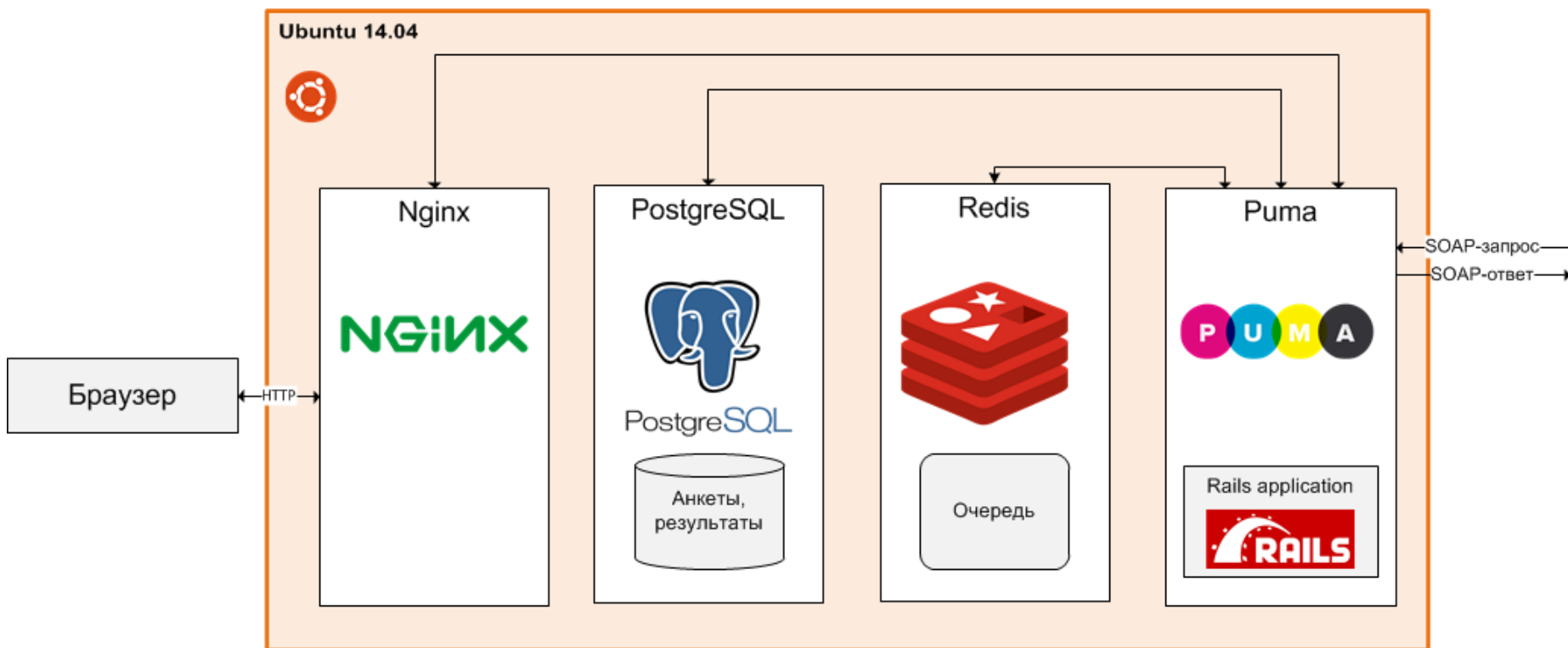
Вариант ответа	Скоринговый балл
Материальное положение – высоко обеспеченная семья	53,69
Материальное положение – не обеспеченная самым необходимым	48,95
Оценка уровня здоровья (настроение, «жизненные силы») – плохое, скорее плохое	48,58
Социальное положение – руководитель отдела, подразделения	43,90
Пол – мужской	39,08
Наличие вредных привычек – имеются	30,68
....	....



# Архитектура приложения



# Архитектура визуальной и административной подсистем





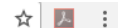
# Схема данных XML

```
<?xml version="1.0" encoding="UTF-8"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
            xmlns:tns="drug_add" targetNamespace="drug_add"
            elementFormDefault="qualified"
            attributeFormDefault="unqualified">
  <xs:complexType name="Input">
    <xs:sequence>
      <xs:element name="q1_1" type="xs:int" />
      <xs:element name="q1_2" type="xs:int" />
      <xs:element name="q1_3" type="xs:int" />
      .....
      <xs:element name="social" type="xs:int" />
      <xs:element name="material" type="xs:int" />
      <xs:element name="Target" type="xs:int" />
    </xs:sequence>
  </xs:complexType>
  <xs:complexType name="Output">
    <xs:sequence>
      <xs:element name="Target_OUT" type="xs:boolean"/>
      <xs:element name="Target_prob" type="xs:float"/>
      <xs:element name="Target_ball" type="xs:float"/>
    </xs:sequence>
  </xs:complexType>
  <xs:element name="input" type="tns:Input"/>
  <xs:element name="output" type="tns:Output"/>
</xs:schema>
```



# Описание WSDL-сервиса

deductor.iipo.tu-bryansk.ru/DIS\_drug/Service.svc?wsdl



This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
<?xml version="1.0" encoding="utf-8" ?>
<wsdl:definitions xmlns:wsdl="http://schemas.xmlsoap.org/wsdl/" xmlns:wsx="http://schemas.xmlsoap.org/ws/2004/09/mex" xmlns:wsu="http://docs.oasis-open.org/wss/2004/01/oasis-200401-wss-wssecurity-utility-1.0.xsd" xmlns:wsa10="http://www.w3.org/2005/08/addressing" xmlns:wsp="http://schemas.xmlsoap.org/ws/2004/09/policy"
xmlns:wsap="http://schemas.xmlsoap.org/ws/2004/08/addressing/policy" xmlns:msec="http://schemas.microsoft.com/ws/2005/12/wsdl/contract"
xmlns:soap12="http://schemas.xmlsoap.org/wsdl/soap12/" xmlns:wsa="http://schemas.xmlsoap.org/ws/2004/08/addressing" xmlns:wsam="http://www.w3.org/2007/05/addressing/metadata"
xmlns:xsd="http://www.w3.org/2001/XMLSchema" xmlns:tns="http://www.basegroup.ru/DeductorIntegrationServer" xmlns:soap="http://schemas.xmlsoap.org/wsdl/soap/"
xmlns:wsaw="http://www.w3.org/2006/05/addressing/wsdl" xmlns:soapenc="http://schemas.xmlsoap.org/soap/encoding/" name="ServiceContract"
targetNamespace="http://www.basegroup.ru/DeductorIntegrationServer">
  <wsdl:types>
    <xsd:schema targetNamespace="http://www.basegroup.ru/DeductorIntegrationServer/Imports">
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/DIS_drug/Service.svc?xsd=xsd0" namespace="drug_add"/>
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/DIS_drug/Service.svc?xsd=xsd1" namespace="http://www.basegroup.ru/warehouses/XsdDbConnection4"/>
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/DIS_drug/Service.svc?xsd=xsd2" namespace="http://www.basegroup.ru/DeductorIntegrationServer"/>
    </xsd:schema>
  </wsdl:types>
  <wsdl:message name="drug_addMessage">
    <wsdl:part name="parameters" element="tns:drug_add"/>
  </wsdl:message>
  <wsdl:message name="drug_addResponseMessage">
    <wsdl:part name="parameters" element="tns:drug_addResponse"/>
  </wsdl:message>
  <wsdl:message name="ServiceClass_drug_add_SOAPExceptionFault_FaultMessage">
    <wsdl:part name="detail" element="tns:SOAPException"/>
  </wsdl:message>
  <wsdl:portType name="ServiceClass">
    <wsdl:operation name="drug_add">
      <wsdl:input wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/drug_add" name="drug_addMessage" message="tns:drug_addMessage"/>
      <wsdl:output wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/drug_addResponse" name="drug_addResponseMessage"
message="tns:drug_addResponseMessage"/>
      <wsdl:fault wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/drug_addSOAPExceptionFault" name="SOAPExceptionFault"
message="tns:ServiceClass_drug_add_SOAPExceptionFault_FaultMessage"/>
    </wsdl:operation>
  </wsdl:portType>
  <wsdl:binding name="BasicHttpBinding_ServiceClass" type="tns:ServiceClass">
    <soap:binding transport="http://schemas.xmlsoap.org/soap/http"/>
    <wsdl:operation name="drug_add">
      <soap:operation soapAction="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/drug_add" style="document"/>
      <wsdl:input name="drug_addMessage">
        <soap:body use="literal"/>
      </wsdl:input>
      <wsdl:output name="drug_addResponseMessage">
        <soap:body use="literal"/>
      </wsdl:output>
      <wsdl:fault name="SOAPExceptionFault">
        <soap:fault name="SOAPExceptionFault" use="literal"/>
      </wsdl:fault>
    </wsdl:operation>
  </wsdl:binding>
  <wsdl:service name="ServiceContract">
    <wsdl:port name="BasicHttpBinding_ServiceClass" binding="tns:BasicHttpBinding_ServiceClass">
      <soap:address location="http://deductor.iipo.tu-bryansk.ru/DIS_drug/Service.svc/ServiceClass"/>
    </wsdl:port>
  </wsdl:service>
</wsdl:definitions>
```





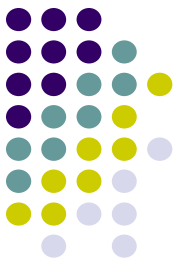
# Тело HTTP-запроса (SOAP)

```
<soap:Envelope
  xmlns:soap="http://schemas.xmlsoap.org/soap/envelope/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema">
  <soap:Body>
    <drug_add xmlns="http://www.basegroup.ru/DeductorIntegrationServer">
      <variables xmlns=""/>
      <data xmlns="">
        <input xmlns="http://www.basegroup.ru/warehouses/XsdDbConnection4">
          <q1_1 xmlns="drug_add">1</q1_1>
          <q1_2 xmlns="drug_add">1</q1_2>
          <q1_3 xmlns="drug_add">0</q1_3>
          <q1_4 xmlns="drug_add">1</q1_4>
          .....
          <sex xmlns="drug_add">2</sex>
          <age xmlns="drug_add">19</age>
          <education xmlns="drug_add">4</education>
          <social xmlns="drug_add">1</social>
          <material xmlns="drug_add">3</material>
          <Target xmlns="drug_add">1</Target>
        </input>
      </data>
    </drug_add>
  </soap:Body>
</soap:Envelope>
```



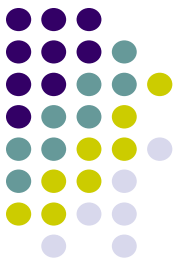
# Тело HTTP-ответа (SOAP)

```
<s:Envelope xmlns:s="http://schemas.xmlsoap.org/soap/envelope/">
  <s:Body xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xmlns:xsd="http://www.w3.org/2001/XMLSchema">
    <drug_addResponse
xmlns="http://www.basegroup.ru/DeductorIntegrationServer">
      <output xmlns="http://www.basegroup.ru/warehouses/XsdDbConnection4"
xmlns:DRGD="drug_add">
        <DRGD:Target_OUT>>false</DRGD:Target_OUT>
        <DRGD:Target_prob>0.0720488804057169</DRGD:Target_prob>
        <DRGD:Target_ball>413.382858296961</DRGD:Target_ball>
      </output>
    </drug_addResponse>
  </s:Body>
</s:Envelope>
```

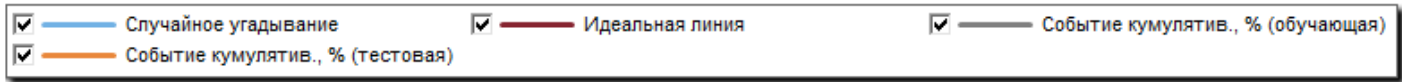
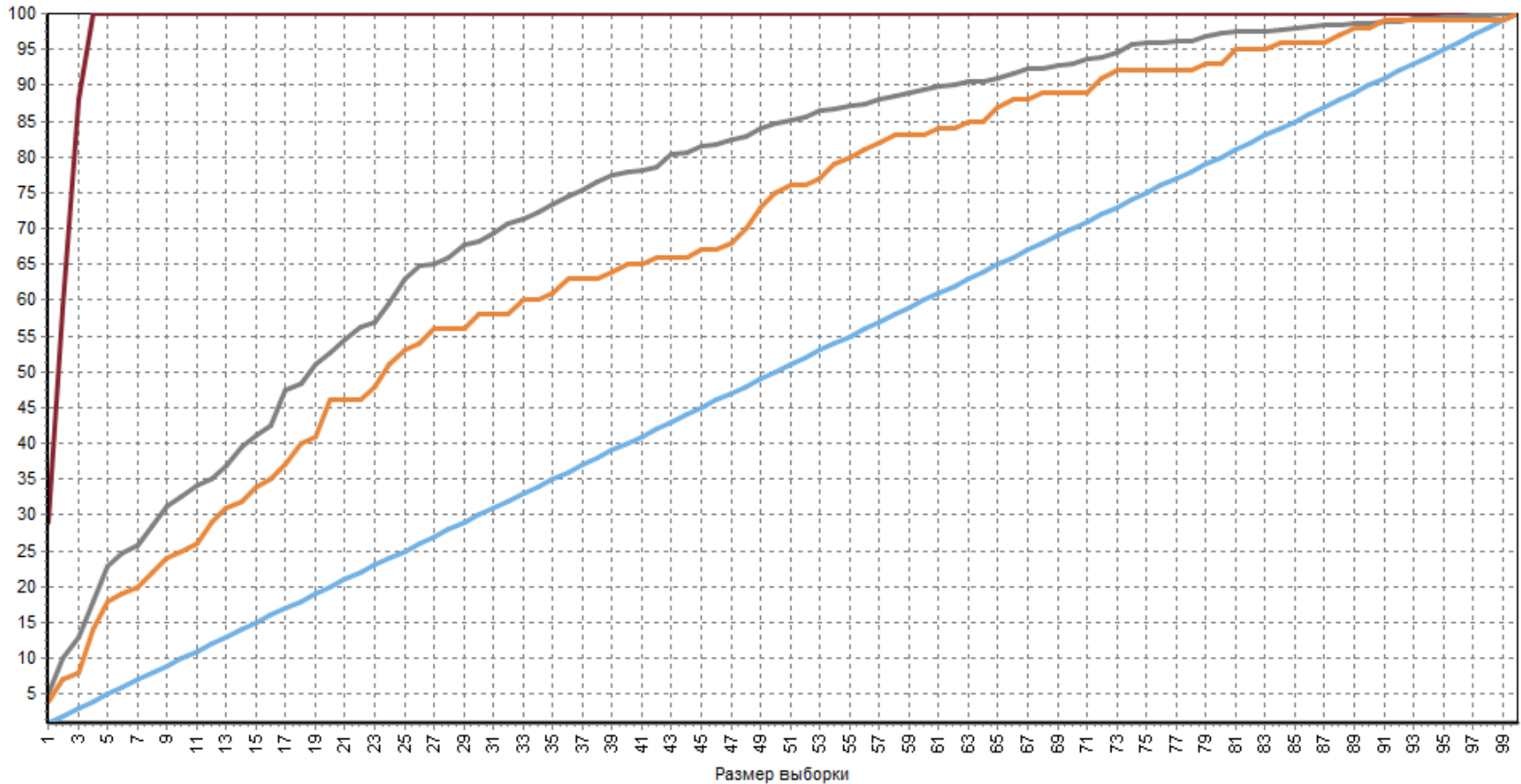


# Оценка качества моделей

- Построено 4 скоринговых модели, которые отличались методом отбора переменных
- Для каждой модели были построены
  - ROC-кривые
  - CAP-кривые
  - диаграммы, отражающие шансы события/не-события
- Выбрана скоринговая модель с лучшими показателями
  - Метод отбора переменных – прямой отбор
  - Использована в сценарии Deductor Enterprise в качестве рабочей модели

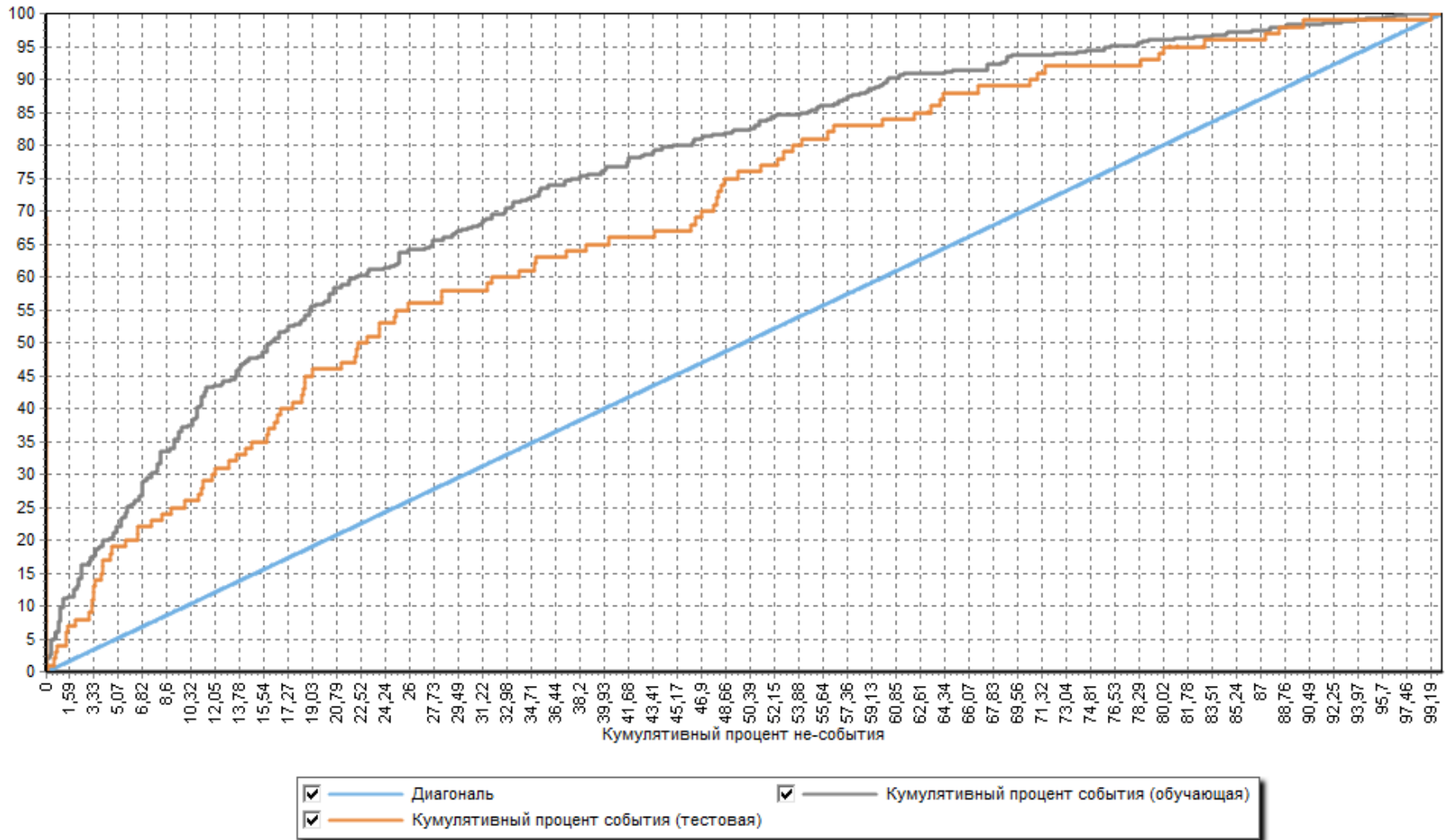


# САР-кривые



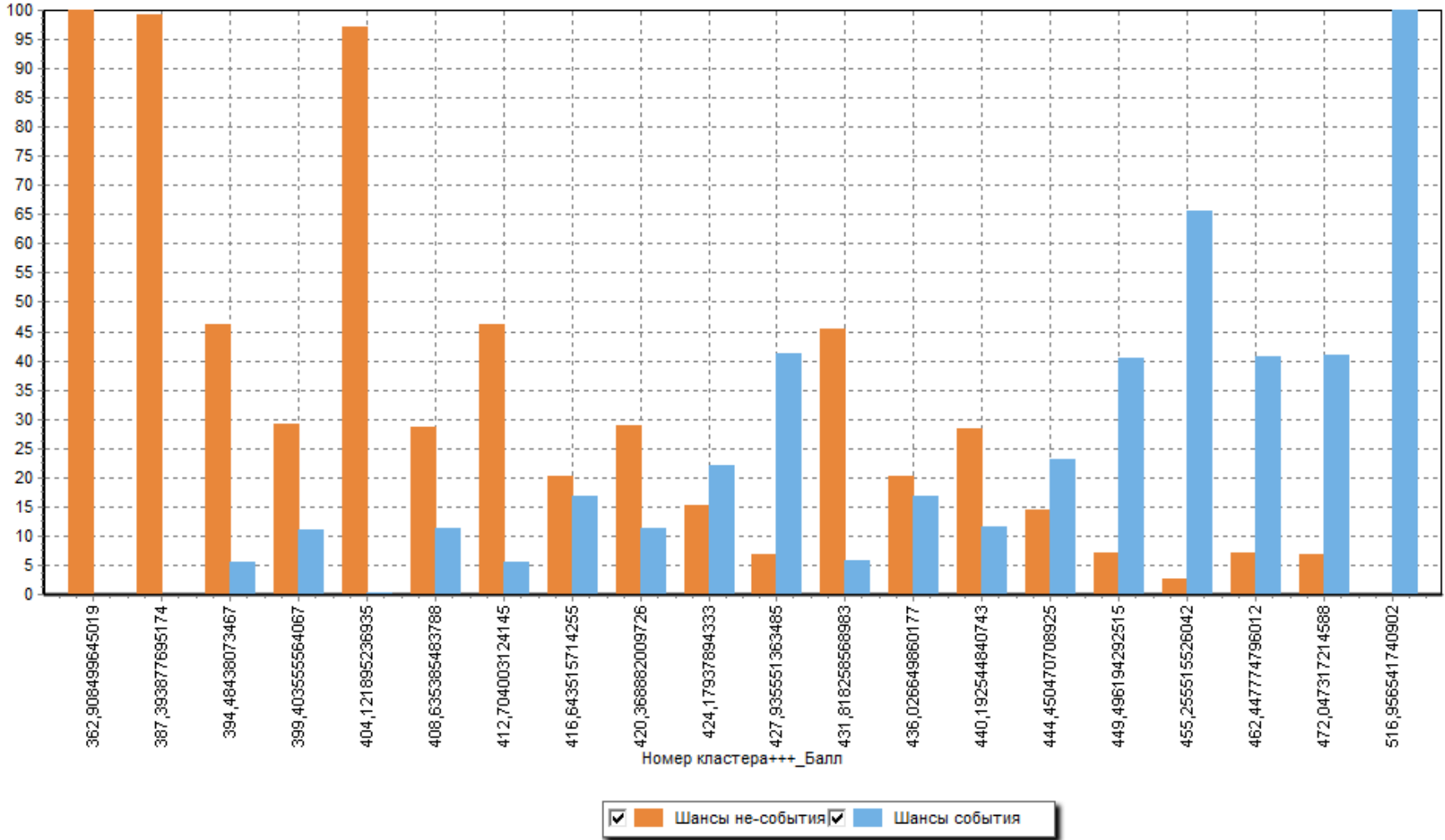


# ROC-кривые





# Шансы события/не-события





# Панель администрирования

Questionnaires | Deducto x Анкетирование x

social.iipo.tu-bryansk.ru/admin/questionnaires

Deductor Панель управления Admin Users Comments **Questionnaires**

ADMIN /

## Questionnaires

Групповые операции ▾

<input type="checkbox"/>	Uid	Is Analyze	Status	Risk	Percent	Points	Created At	Updated At				
<input type="checkbox"/>	c66bed3dc7d7ee2dbcee2646410d9e	НЕТ	ЗАПОЛНЯЕТСЯ...	НЕИЗВЕСТНО	0	0	24 июня 2017, 20:20	24 июня 2017, 20:20	<a href="#">Открыть</a>	<a href="#">Изменить</a>	<a href="#">Удалить</a>	<a href="#">Прямая ссылка</a>
<input type="checkbox"/>	71dd030bbe2535890b30feb0611b6d	ДА	УСПЕШНО ЗАВЕРШЕНА	ВСЕ ХОРОШО	25	455	24 июня 2017, 13:59	24 июня 2017, 14:00	<a href="#">Открыть</a>	<a href="#">Изменить</a>	<a href="#">Удалить</a>	<a href="#">Прямая ссылка</a>
<input type="checkbox"/>	7c7b53e39cdb550830e74ff02bb485	ДА	УСПЕШНО ЗАВЕРШЕНА	ВСЕ ХОРОШО	6	412	23 июня 2017, 15:55	23 июня 2017, 15:56	<a href="#">Открыть</a>	<a href="#">Изменить</a>	<a href="#">Удалить</a>	<a href="#">Прямая ссылка</a>
<input type="checkbox"/>	da120618559f77c9552512123dcc0e	ДА	УСПЕШНО ЗАВЕРШЕНА	В ЗОНЕ РИСКА	68	509	22 июня 2017, 23:07	23 июня 2017, 15:47	<a href="#">Открыть</a>	<a href="#">Изменить</a>	<a href="#">Удалить</a>	<a href="#">Прямая ссылка</a>

Загрузка: [CSV](#) [XML](#) [JSON](#)

Результат: 4 Questionnaires



# Фрагмент анкеты

Questionnaires | Deduct... x Анкетирование x

social.iipo.tu-bryansk.ru/q/c66bed3dc7d7ee2dbcee2646410d9e

## Здравствуйте!

Мы проводим всероссийское исследование, посвященное изучению привычек и убеждений граждан России. В исследовании участвуют жители всех субъектов Российской Федерации, проживающие в городских и сельских населенных пунктах. Мы просим Вас высказать свое мнение по ряду вопросов.

**Анкета анонимная, Вам не нужно указывать фамилию. Все данные будут использованы только в обобщенном виде для научных целей.**

После каждого вопроса написано, сколько ответов нужно дать. Обведите в кружок цифру около варианта ответа, который выражает Ваше мнение, или напишите свой вариант.

**Пожалуйста, ответьте на ВСЕ вопросы. Ваше мнение очень важно для нас!**

### Для начала просим Вас ответить на вопросы, характеризующие Ваши жизненные ориентиры

**№1. Укажите, пожалуйста, ПЯТЬ наиболее острых проблем, требующих решения в первую очередь в Вашем населенном пункте (возможно несколько вариантов ответа):**

- Нехватка жилья
- Качество дорог
- Алкоголизм
- Безработица
- Состояние жилищно-коммунальной сферы
- Наркомания
- Качество медицинского обслуживания
- Преступность

**№2. Выберите, пожалуйста, из ниже перечисленного списка не более ПЯТИ наиболее значимых для Вас ценностей:**

- Активная, деятельная жизнь
- Жизненная мудрость
- Здоровье
- Красота природы и искусства



# Информация о заполненной анкете



Questionnaire #26 | Deductor x Анкетирование x

social.iipo.tu-bryansk.ru/admin/questionnaires/26

Deductor Панель управления Admin Users Comments Questionnaires ababurin@bk.ru Выйти

ADMIN / QUESTIONNAIRES /

## Questionnaire #26

Изменить Questionnaire Удалить Questionnaire

Questionnaire подробнее

UID	da120618559f77c9552512123dccc0e
IS ANALYZE	ДА
STATUS	УСПЕШНО ЗАВЕРШЕНА
ANSWER	<pre>&lt;?xml version="1.0" encoding="UTF-8"?&gt; &lt;hash&gt;   &lt;q1&gt;     &lt;q1-1&gt;1&lt;/q1-1&gt;     &lt;q1-2 type="integer"&gt;0&lt;/q1-2&gt;     &lt;q1-3&gt;1&lt;/q1-3&gt;     &lt;q1-4&gt;1&lt;/q1-4&gt;     &lt;q1-5 type="integer"&gt;0&lt;/q1-5&gt;     &lt;q1-6&gt;1&lt;/q1-6&gt;     &lt;q1-7 type="integer"&gt;0&lt;/q1-7&gt;     &lt;q1-8&gt;1&lt;/q1-8&gt;   &lt;/q1&gt;   &lt;q2&gt;     &lt;q2-1 type="integer"&gt;0&lt;/q2-1&gt;     &lt;q2-2 type="integer"&gt;0&lt;/q2-2&gt;     &lt;q2-3 type="integer"&gt;0&lt;/q2-3&gt;     &lt;q2-4 type="integer"&gt;0&lt;/q2-4&gt;     &lt;q2-5 type="integer"&gt;0&lt;/q2-5&gt;     &lt;q2-6 type="integer"&gt;0&lt;/q2-6&gt;     &lt;q2-7 type="integer"&gt;0&lt;/q2-7&gt;     &lt;q2-8 type="integer"&gt;0&lt;/q2-8&gt;     &lt;q2-9&gt;1&lt;/q2-9&gt;     &lt;q2-10 type="integer"&gt;0&lt;/q2-10&gt;     &lt;q2-11 type="integer"&gt;0&lt;/q2-11&gt;     &lt;q2-12 type="integer"&gt;0&lt;/q2-12&gt;     &lt;q2-13 type="integer"&gt;0&lt;/q2-13&gt;     &lt;q2-14 type="integer"&gt;0&lt;/q2-14&gt;     &lt;q2-15 type="integer"&gt;0&lt;/q2-15&gt;   &lt;/q2&gt; &lt;/hash&gt; </pre>
RISK	В ЗОНЕ РИСКА
PERCENT	68
POINTS	509



# Итоги и перспективы

- Построена скоринговая модель для выявления группы риска относительно наркозависимости
- Разработаны соответствующие веб-приложения, которые используют серверные компоненты аналитической платформы Deductor Enterprise:
  - веб-приложение для анкетирования респондентов
  - веб-приложение для социолога (просмотр результатов опроса и классификация респондентов)
- В ближайшее время планируется
  - представить проект заказчику
  - выполнить экспериментальную проверку на реальном социологическом опросе
  - подготовить пресс-релиз для сайта компании BaseGroup Labs
  - опубликовать полученные результаты



# Содержание

- Введение
- **Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise**
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - **Программная поддержка полного цикла социологического исследования**
  - Поиск единомышленников в социальной сети VK
- Итоги работы. Дальнейшие планы и пожелания



# О проекте

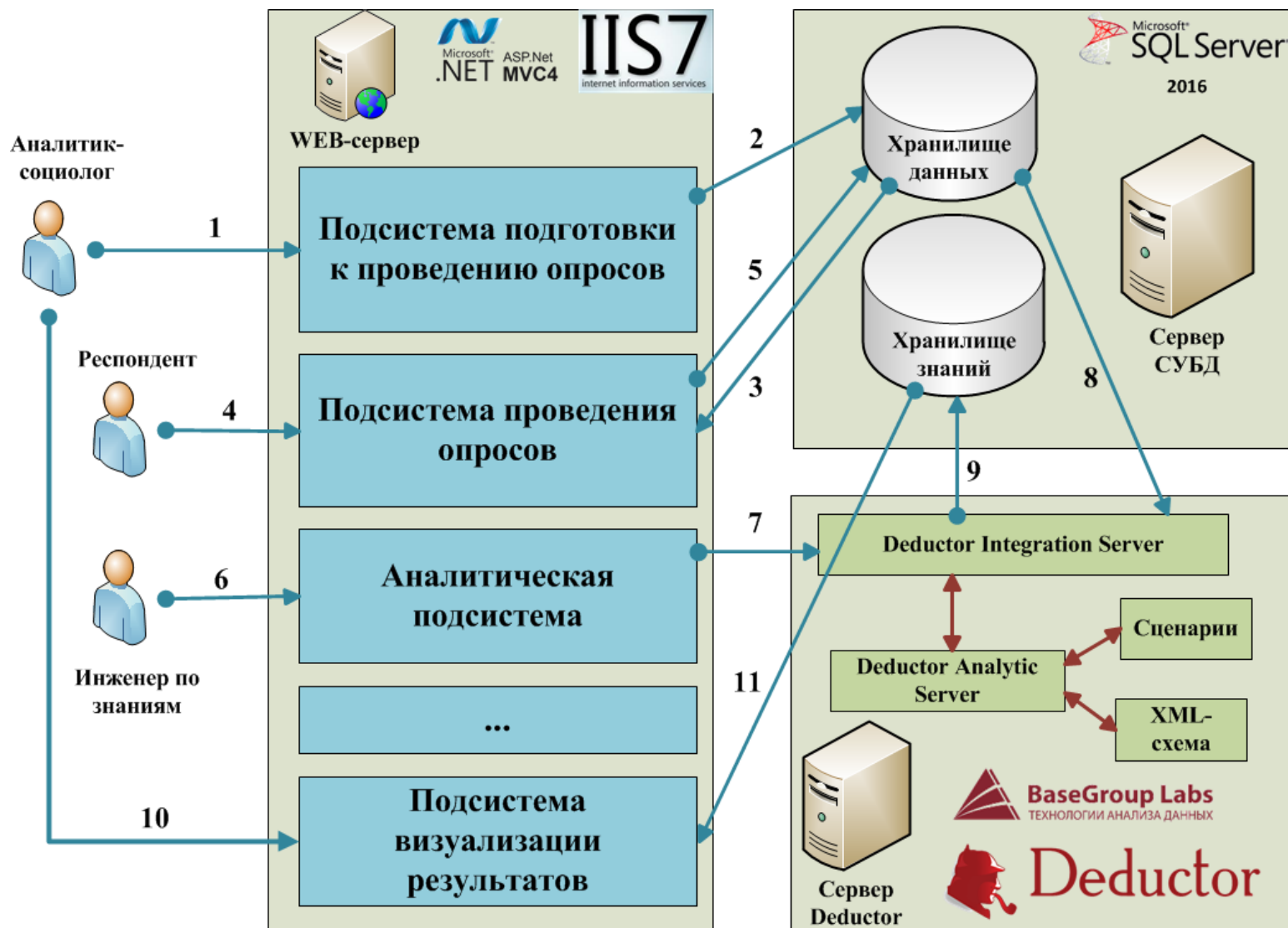
- Выполняется в рамках кандидатской диссертации Бабурина А.Н.
- Рабочее название «Разработка методов и программных средств комплексной автоматизации социологических исследований с использованием ансамблей моделей интеллектуального анализа данных»
  - Специальность 05.13.10 «Управление в социальных и экономических системах»
- **Основная идея:** комплексная автоматизация полного цикла социологического исследования, начиная от составления анкеты и заканчивая выдачей результатов интеллектуального анализа заказчику.

# Диаграмма вариантов использования

## ИСПОЛЬЗОВАНИЯ



# Архитектура программного комплекса

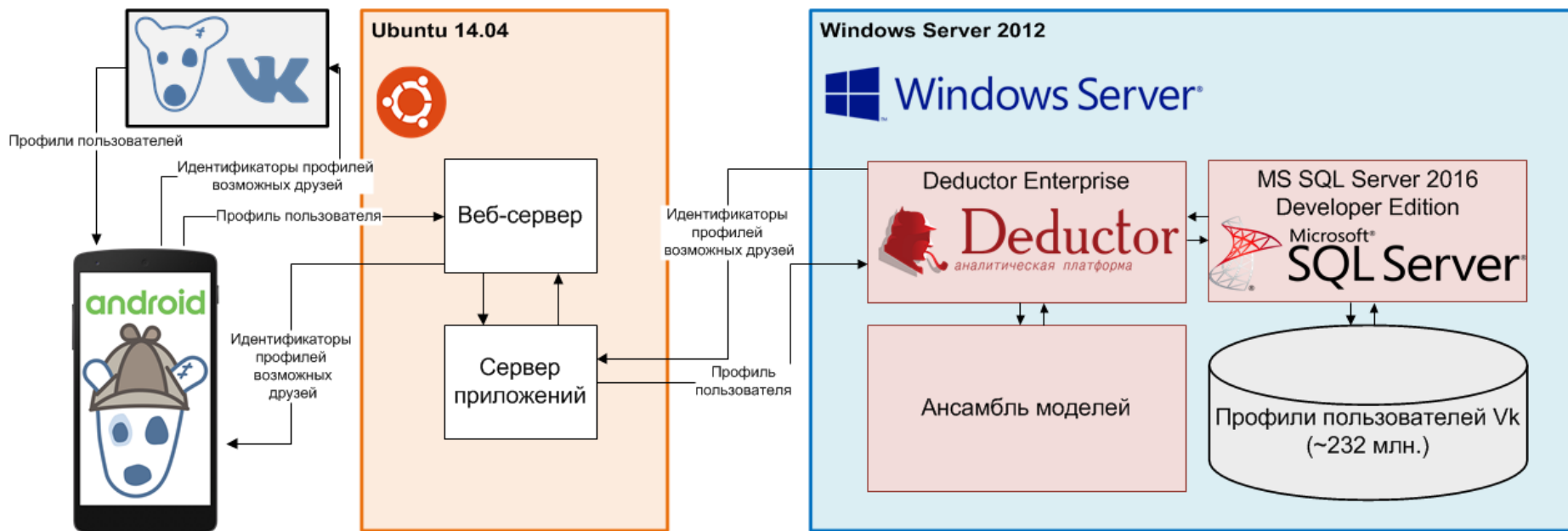




# Содержание

- Введение
- **Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise**
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - Программная поддержка полного цикла социологического исследования
  - **Поиск единомышленников в социальной сети VK**
- Итоги работы. Дальнейшие планы и пожелания

# Архитектура программного комплекса



\* Выражаем благодарность компании «АйТи Про» ([www.itprocomp.ru](http://www.itprocomp.ru)) и лично **Бондареву Борису Игоревичу** (главному архитектору) за помощь в получении данных из социальной сети VK



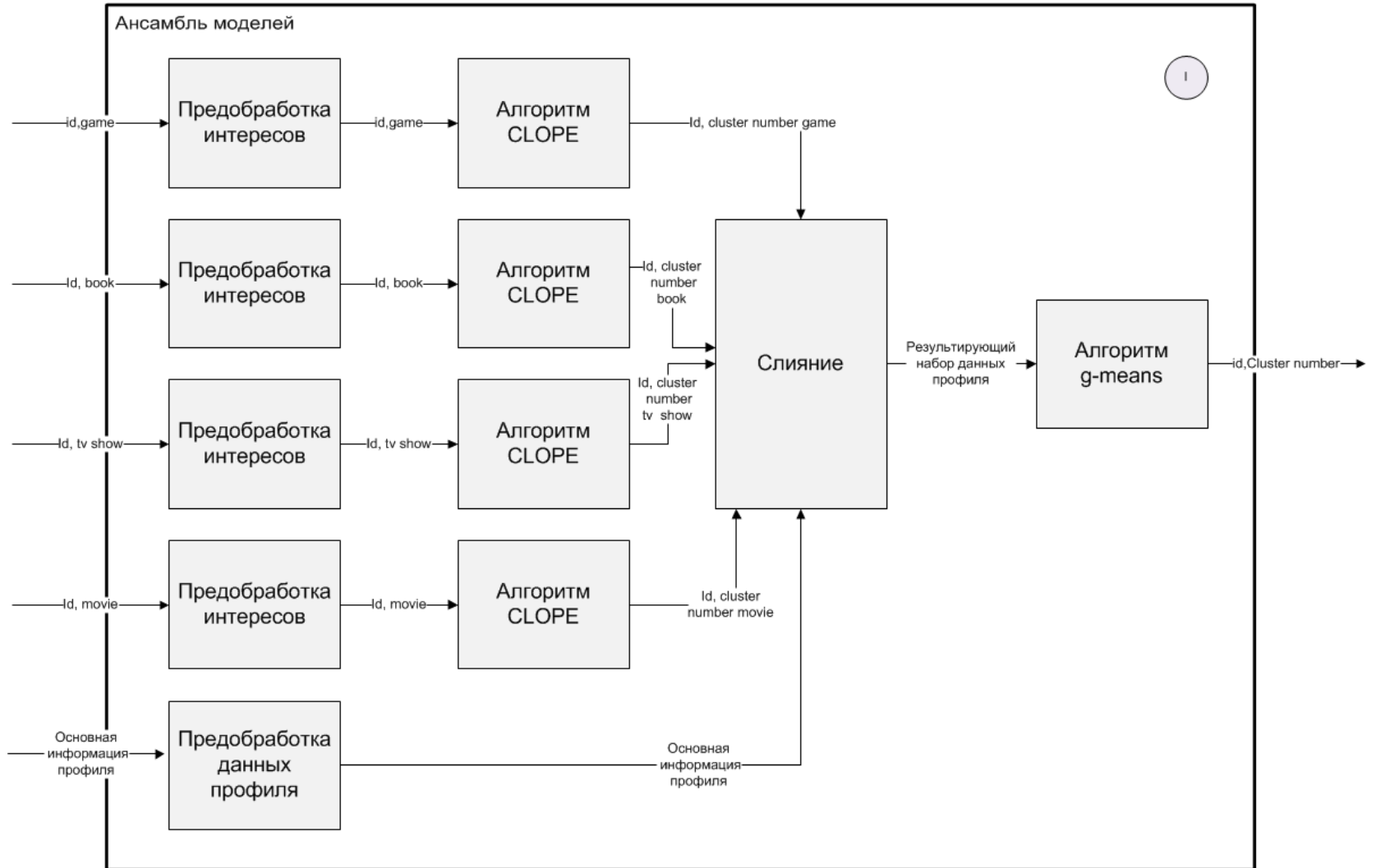


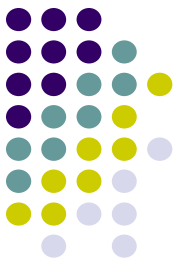
# Конфигурация сервера

- Процессор: Intel Xeon E5-2620 v2 2,1 ГГц
- ОЗУ: 16 Гб
- Объем жесткого диска: 1 Тб
- ОС: MS Windows Server 2016 R2
- СУБД: MS SQL Server 2016 Developer Edition
- Аналитическая платформа: Deductor Enterprise



# Ансамбль моделей





# Предобработка интересов

- Все операции выполнялись средствами Deductor, кроме преобразования строк в формате JSON (в VK так хранятся данные о интересах пользователя) в табличный вид, пригодный для импорта в Deductor
- Операция осложнялась вольным написанием названий интересов: (WoW, World of Warcraft, worldofwarcraft, World of Warcraft – BurningCrusade, Wor croft, Wor craft, Wor of Kraft, wor workrawt, worald of warkraft, ...)
- Очистка и предобработка данных:
  - Строки приводятся к нижнему регистру.
  - Замена незаполненных интересов словом «нет»
  - Удаляются все символы кроме букв, так же удаляются буквы «i» и «v», так как пользователи часто используют их в качестве цифр
  - Удаляются пустые строки
  - Замена по словарю наиболее частых ошибок. Например Wor craft → warcraft
  - Отбрасываются редкие варианты (менее 1000 вхождений)

# Примеры результатов работы алгоритма CLOPE



- Кластер 13:
  - Футбол
  - Твистер
  - Теннис
  - Бильярд
  - Волейбол
  - Шахматы
- Кластер 148:
  - Ролевые
  - Взрослые
  - Сексуальные
- Кластер 280:
  - Зомбиферма
  - Need for Speed Carbon
  - Паук
  - Warcraft
- Кластер 309:
  - Баскетбол
  - Гарри Поттер
  - RPG
  - Теннис
  - Приключения
  - Компьютерные игры
  - Метро
  - Футбол
  - Секс
  - Игры
  - Бильярд
  - Рыбалка
  - Sims
  - Спорт
  - Жизнь
  - Маджонг



# Предобработка для G-Means

- Данная таблица состоит из следующих столбцов:
  - Идентификатор пользователя
  - Дата рождения
  - Пол
  - Город проживания
  - Страна проживания
- Алгоритм предобработки:
  - На основе даты рождения определяются:
    - возраст
    - знак зодиака
    - знак животного по восточному календарю
  - Выполняется слияние с таблицами, полученными по результатам работы алгоритмов CLOPE

# Создание веб-сервиса для обработки набора записей

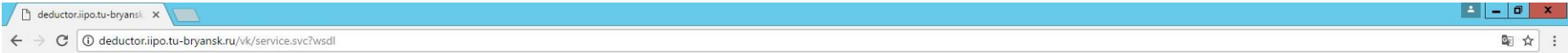


```
D:\!Analys\VK\scheme.xsd - Notepad++
File Edit Search View Encoding Language Settings Macro Run Plugins Window ?
scheme.xsd
11     <xs:element name="sex" type="xs:int" />
12   </xs:sequence>
13 </xs:complexType>
14 <xs:complexType name="Output">
15   <xs:sequence>
16     <xs:element name="id" type="xs:int"/>
17     <xs:element name="id_animal" type="xs:int"/>
18     <xs:element name="id_horoscope" type="xs:int"/>
19     <xs:element name="cluster_number_game" type="xs:int"/>
20     <xs:element name="cluster_number_book" type="xs:int"/>
21     <xs:element name="cluster_number_movie" type="xs:int"/>
22     <xs:element name="cluster_number_tv" type="xs:int"/>
23     <xs:element name="cluster_number" type="xs:int"/>
24     <xs:element name="distance" type="xs:float"/>
25   </xs:sequence>
26 </xs:complexType>
27 <xs:element name="input" type="tns:Input"/>
28 <xs:element name="output" type="tns:Output"/>
29
30 <xs:complexType name="InputRowsType">
31   <xs:sequence>
32     <xs:element name="RowItem" type="tns:Input" maxOccurs="unbounded"/>
33   </xs:sequence>
34 </xs:complexType>
35 <xs:element name="InputRows" type="tns:InputRowsType"/>
36 <xs:complexType name="OutputRowsType">
37   <xs:sequence>
38     <xs:element name="RowItem" type="tns:Output" maxOccurs="unbounded"/>
39   </xs:sequence>
40 </xs:complexType>
41 <xs:element name="OutputRows" type="tns:OutputRowsType"/>
42
43 </xs:schema>
```

exTensible Markup Language file      length : 1 755   lines : 43      Ln : 42   Col : 1   Sel : 0 | 0      Windows (CR LF)   UTF-8   INS      20:30   24.06.2017



# Описание WSDL-сервиса



This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
<?xml:definitions xmlns:wsdl="http://schemas.xmlsoap.org/wsdl/" xmlns:wsx="http://schemas.xmlsoap.org/ws/2004/09/mex" xmlns:wssu="http://docs.oasis-open.org/wss/2004/01/oasis-200401-wss-wssecurity-utility-1.0.xsd" xmlns:wsa10="http://www.w3.org/2005/08/addressing"
xmlns:wsp="http://schemas.xmlsoap.org/ws/2004/09/policy" xmlns:wsap="http://schemas.xmlsoap.org/ws/2004/08/addressing/policy" xmlns:msec="http://schemas.microsoft.com/ws/2005/12/wsdl/contract" xmlns:soap12="http://schemas.xmlsoap.org/wsdl/soap12/"
xmlns:wsa="http://schemas.xmlsoap.org/ws/2004/08/addressing" xmlns:wsam="http://www.w3.org/2007/05/addressing/metadata" xmlns:xsd="http://www.w3.org/2001/XMLSchema" xmlns:tns="http://www.basegroup.ru/DeductorIntegrationServer"
xmlns:soap="http://schemas.xmlsoap.org/wsdl/soap/" xmlns:wsaw="http://www.w3.org/2006/05/addressing/wsdl" xmlns:soapenc="http://schemas.xmlsoap.org/soap/encoding/" name="ServiceContract" targetNamespace="http://www.basegroup.ru/DeductorIntegrationServer">
  <wsdl:types>
    <xsd:schema targetNamespace="http://www.basegroup.ru/DeductorIntegrationServer/Imports">
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/vk/Service.svc?xsd=xsd0" namespace="DIS_vk"/>
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/vk/Service.svc?xsd=xsd1" namespace="http://www.basegroup.ru/warehouses/vk"/>
      <xsd:import schemaLocation="http://deductor.iipo.tu-bryansk.ru/vk/Service.svc?xsd=xsd2" namespace="http://www.basegroup.ru/DeductorIntegrationServer"/>
    </xsd:schema>
  </wsdl:types>
  <wsdl:message name="DIS_vkMessage">
    <wsdl:part name="parameters" element="tns:DIS_vk"/>
  </wsdl:message>
  <wsdl:message name="DIS_vkResponseMessage">
    <wsdl:part name="parameters" element="tns:DIS_vkResponse"/>
  </wsdl:message>
  <wsdl:message name="ServiceClass_DIS_vk_SOAPExceptionFault_FaultMessage">
    <wsdl:part name="detail" element="tns:SOAPException"/>
  </wsdl:message>
  <wsdl:portType name="ServiceClass">
    <wsdl:operation name="DIS_vk">
      <wsdl:input wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/DIS_vk" name="DIS_vkMessage" message="tns:DIS_vkMessage"/>
      <wsdl:output wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/DIS_vkResponse" name="DIS_vkResponseMessage" message="tns:DIS_vkResponseMessage"/>
      <wsdl:fault wsaw:Action="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/DIS_vkSOAPExceptionFault" name="SOAPExceptionFault" message="tns:ServiceClass_DIS_vk_SOAPExceptionFault_FaultMessage"/>
    </wsdl:operation>
  </wsdl:portType>
  <wsdl:binding name="BasicHttpBinding_ServiceClass" type="tns:ServiceClass">
    <soap:binding transport="http://schemas.xmlsoap.org/soap/http"/>
  </wsdl:binding>
  <wsdl:operation name="DIS_vk">
    <soap:operation soapAction="http://www.basegroup.ru/DeductorIntegrationServer/ServiceClass/DIS_vk" style="document"/>
    <wsdl:input name="DIS_vkMessage">
      <soap:body use="literal"/>
    </wsdl:input>
    <wsdl:output name="DIS_vkResponseMessage">
      <soap:body use="literal"/>
    </wsdl:output>
    <wsdl:fault name="SOAPExceptionFault">
      <soap:fault name="SOAPExceptionFault" use="literal"/>
    </wsdl:fault>
  </wsdl:operation>
</wsdl:definitions>
```

# Пример результатов работы в мобильном приложении (1)



12:30

← Ваш профиль

**Владимир Ефремов**

Основная информация

Возраст	Город
29 лет	Москва
Знак зодиака	Китайский календарь
Рак	Кролик

Личная информация

Любимые книги	Любимые телешоу
Журналы	Нет таких
Любимые игры	Любимые фильмы
Футбол	Русские боевики

12:30

← Найденные результаты

- ЕФ Евгений Шарифьянов
- АП Алёша Петров
- СБ Святослав Быстров
- ТН Trevor Hansen
- АВ Aaron Bennett
- АС Abbey Christensen
- АС Ali Connors
- АН Alex Nelson
- АС Anthony Stevens

12:30

←

**Евгений Шарифьянов**

Основная информация

Возраст	Город
29 лет	Мурманск
Знак зодиака	Китайский календарь
Рак	Кролик

Личная информация

Любимые книги	Любимые телешоу
Нет	Камеди клуб, наша раша
Любимые игры	Любимые фильмы
Футбол	Комедии



# Пример результатов работы в мобильном приложении (2)



12:30

← Ваш профиль

**Владимир Ефремов**

Основная информация

Возраст	Город
29 лет	Москва
Знак зодиака	Китайский календарь
Рак	Кролик

Личная информация

Любимые книги	Любимые телешоу
Журналы	Нет таких
Любимые игры	Любимые фильмы
Футбол	Русские боевики

12:30

←

**Алёша Петров**

Основная информация

Возраст	Город
29 лет	Москва
Знак зодиака	Китайский календарь
Козерог	Дракон

Личная информация

Любимые книги	Любимые телешоу
Нет	Камеди
Любимые игры	Любимые фильмы
Танки	Шаг вперёд

12:30

←

**Святослав Быстров**

Основная информация

Возраст	Город
29 лет	Москва
Знак зодиака	Китайский календарь
Телец	Дракон

Личная информация

Любимые книги	Любимые телешоу
Нет времени читать	Очень много
Любимые игры	Любимые фильмы
Не играю	Очень много



# Перспективы развития проекта

- Поддержание в актуальном состоянии базы данных профилей пользователей социальной сети ВКонтакте
- Улучшение качества предобработки
  - Применение более совершенных словарей
  - Применение функции вычисления значения Дамерау-Левенштейна для определения схожести строк с целью идентификации опечаток
- Поиск оптимального алгоритма кластеризации профилей
  - Опыты с другими алгоритмами (EM, карты Кохонена, ...)
  - Экспериментальная проверка с целью выявления наилучшего
- Ранжирование выходного набора данных для повышения релевантности результатов подбора людей со схожими интересами
  - По расстоянию между местами проживания в километрах
  - По длине цепочки друзей (количество промежуточных «звеньев»)



# 8% за 23 часа...

Мастер обработки - Кластеризация (5 из 7)

**Кластеризация набора данных**  
Запуск процесса кластеризации

Обучающее множество

- Максимальная ошибка
- Средняя ошибка
- Распознано (%)

Тестовое множество

- Максимальная ошибка
- Средняя ошибка
- Распознано (%)

Название текущего процесса  
Поиск кластеров

Процент выполнения текущего процесса обучения  
8%

Время обучения 23:10:00

Инициализировать перед обучением

▶ Пуск    || Пауза    ■ Стоп

< Назад    Далее >    Отмена

Windows taskbar: 22:44 12.06.2017



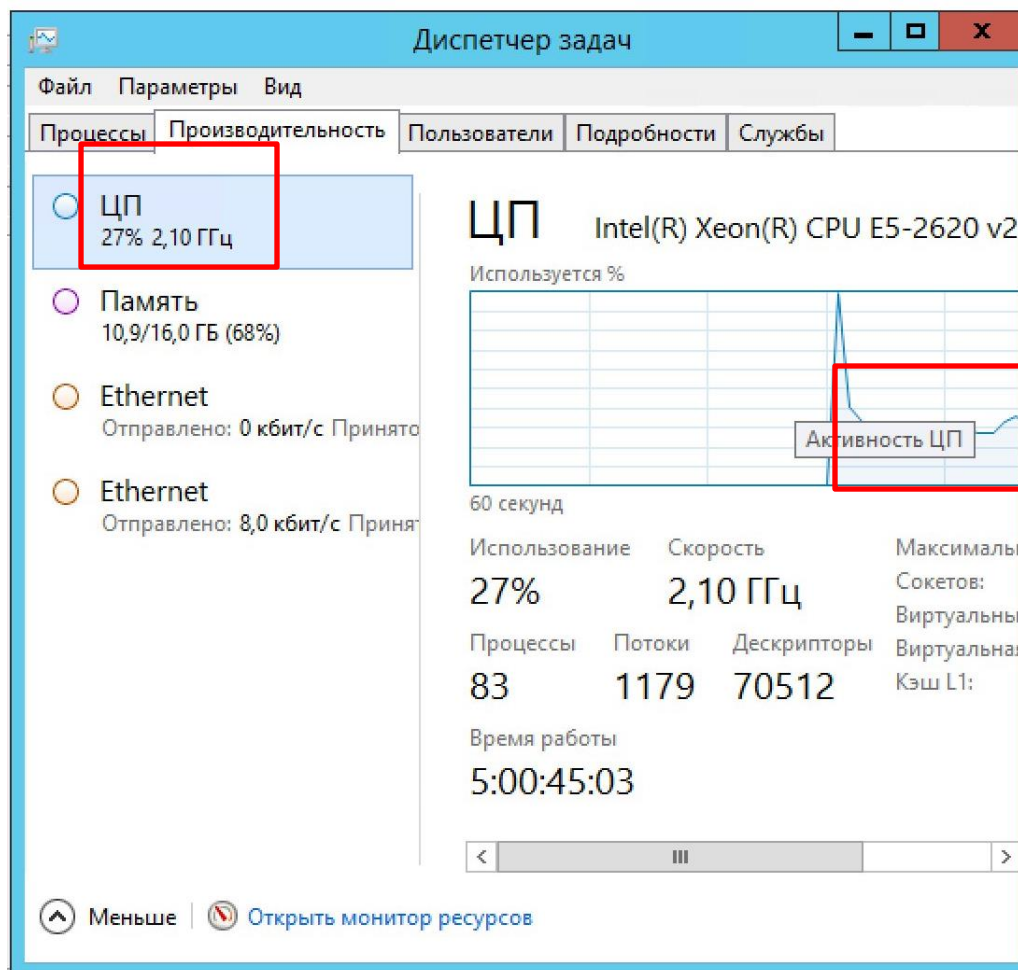
# Затраты времени

- Кластеризация транзакций методом CLOPE занимала **от 24 до 72 часов**
- Кластеризация g-means: 11 столбцов данных, настройки не менялись:
  - 5 000 000 строк, ограничение 1000 кластеров, длительность ~ **54 часа**.
  - 300 000 строк, итоговое количество кластеров 2213, длительность ~ **25 часов**.
- После выбора алгоритма имела место долгая предобработка
  - Занимала несколько часов
  - На это время Deductor «зависал» и не реагировал на действия пользователя
  - При этом не отражался никакой визуальный индикатор



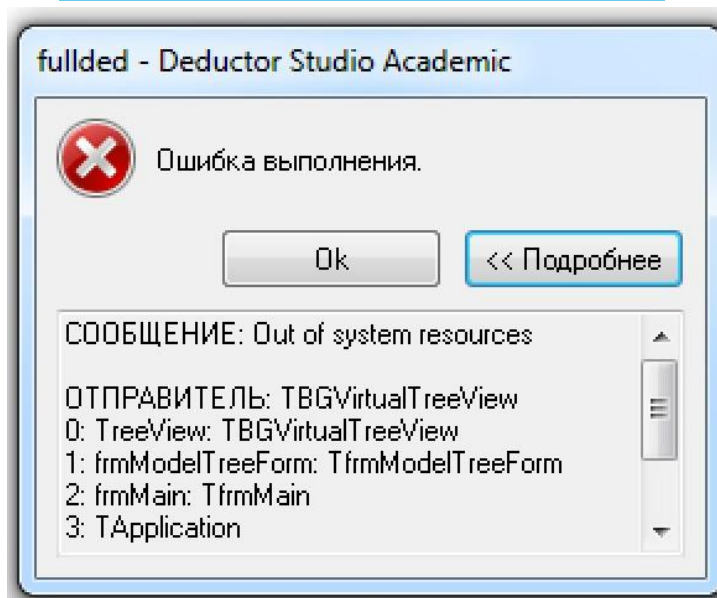
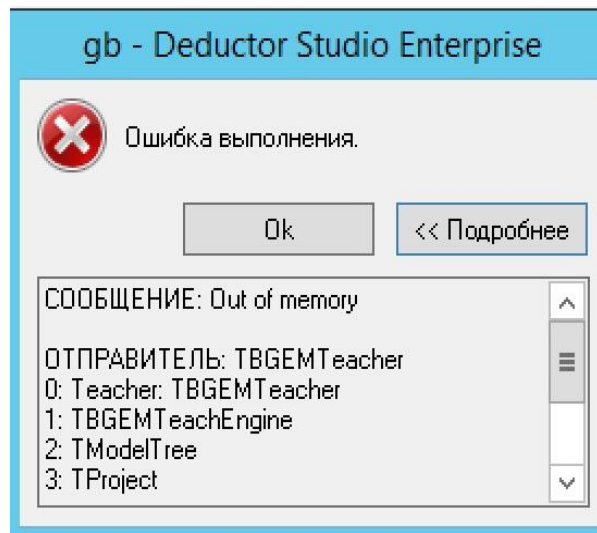
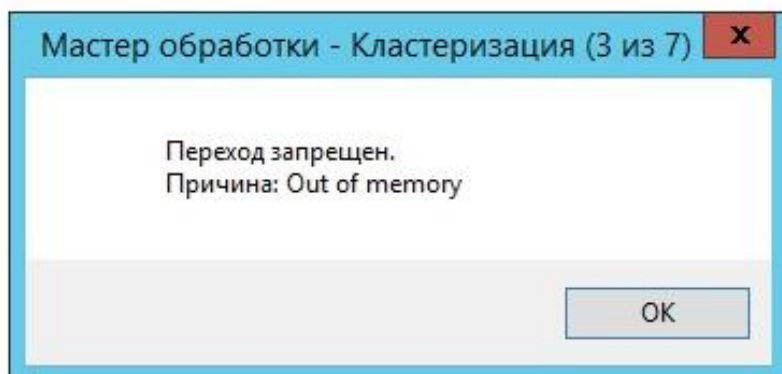
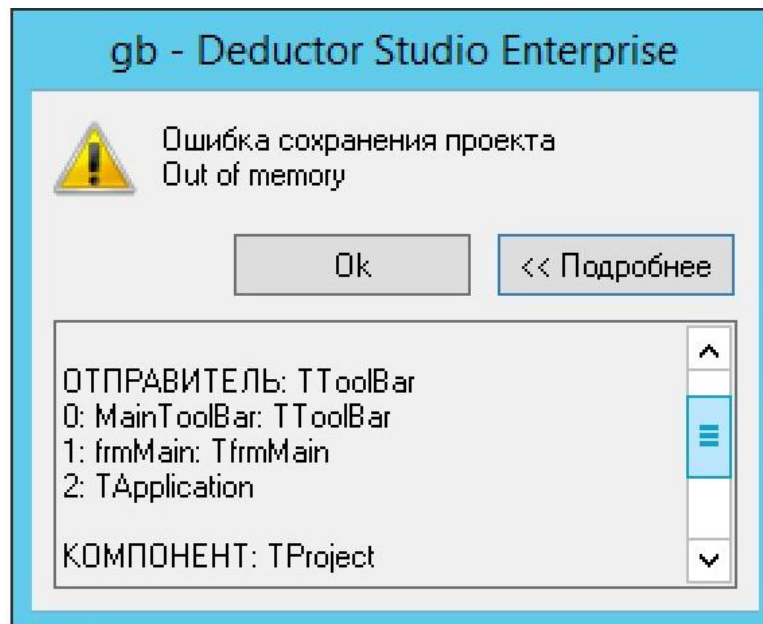
# Выполнение кластеризации

Deductor использует только одно ядро, остальные бездействуют





# Out of memory...





# Настройки файла подкачки...

Настройка файла подкачки

Параметры кэширования    Информация об используемой памяти

Файл подкачки открыт без ошибок

Использование оперативной памяти

Выделено в кэше	<input type="text" value="14%"/>	290 979 840	Байт ▾
Выделено всего	<input type="text" value="38%"/>	1 632 051 200	Байт ▾
Зарезервировано всего	<input type="text" value="81%"/>	3 496 607 744	Байт ▾

Использование файла подкачки

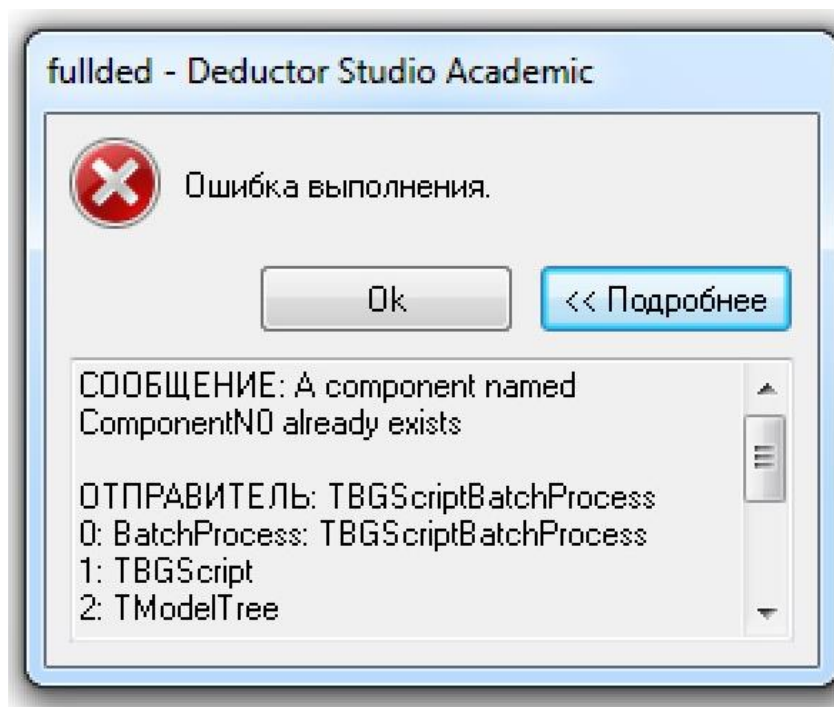
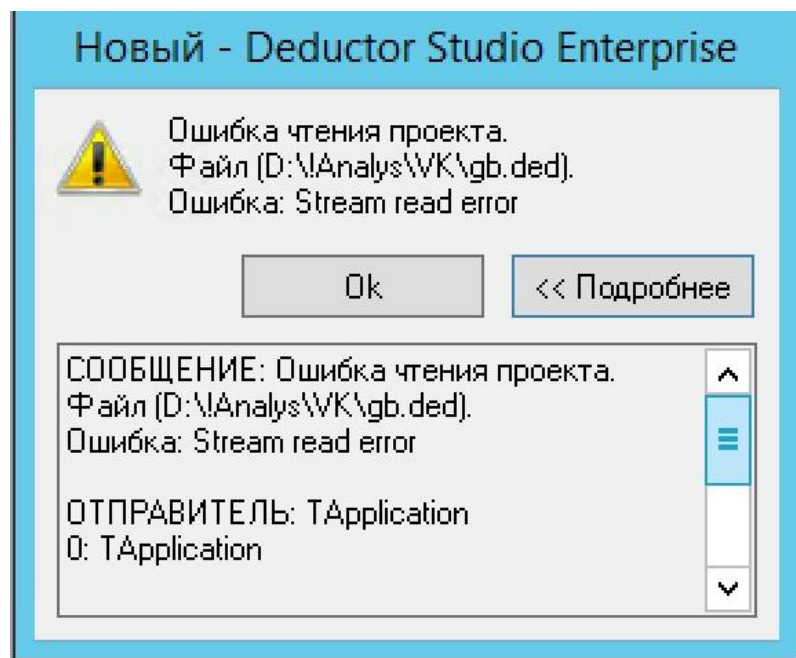
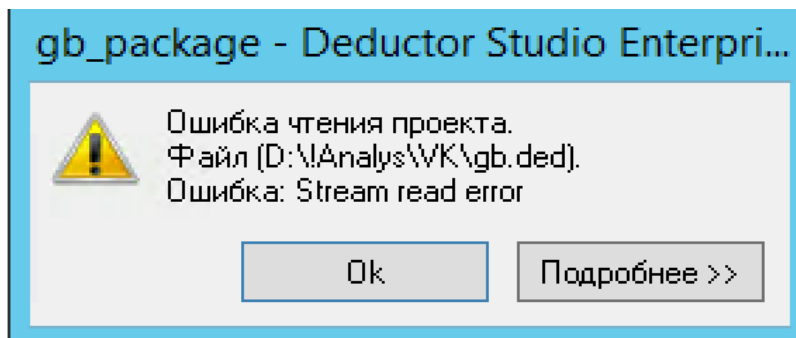
Текущим приложением	<input type="text" value="1%"/>	290 979 840	Байт ▾
Всего	<input type="text" value="33%"/>	11 371 216 896	Байт ▾
Количество подключенных к файлу подкачки приложений	N/A		

Обновить

Ok    Отмена



# Прочие проблемы...





# Проблемы DIS при работе совместно с кластеризацией

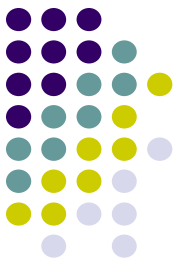


```
DIS_log_2017-06-21_14-20-43 {a9de46aa-877a-4fa4-9487-20e45583b40a} — Блокнот
Файл  Правка  Формат  Вид  Справка
Текст исключения: Не удалось установить подключение к Deductor Server
Класс исключения: System.Exception
Трасса
    в DeductorIntegrationServer.Core.ServiceClass.ConnectDeductorServer(String address, Int32 port, LogWriter
logWriter, Int32& activeApplicationCount)
    в DeductorIntegrationServer.Core.ServiceClass.ExecuteProject(XDocument inputData, LogWriter logWriter,
XDocument& outputData)
    в DeductorIntegrationServer.Core.ServiceClass.InternalExecute(XDocument inputData, LogWriter logWriter,
XDocument& outputData, Exception& exception)
    в DeductorIntegrationServer.Core.ServiceClass.Execute(InputData inputData)
-----
Текст исключения: Подключение не установлено, т.к. конечный компьютер отверг запрос на подключение
82.179.88.29:4386
Класс исключения: System.Net.Sockets.SocketException
Трасса
    в System.Net.Sockets.Socket.Connect(IPAddress[] addresses, Int32 port)
    в System.Net.Sockets.Socket.Connect(String host, Int32 port)
    в DeductorClient.DeductorConnector.Connect()
    в DeductorIntegrationServer.Core.ServiceClass.ConnectDeductorServer(String address, Int32 port, LogWriter
logWriter, Int32& activeApplicationCount)
```



# Содержание

- Введение
- Проекты, выполненные с использованием серверных компонентов платформы Deductor Enterprise
  - Выявление групп риска в рамках мониторинга наркоситуации в Брянской области
  - Программная поддержка полного цикла социологического исследования
  - Поиск единомышленников в социальной сети VK
- **Итоги работы. Дальнейшие планы и пожелания**



# Планы и пожелания прошлого года

(были представлены на конференции в июне 2016)

- **Обучение и сертификация**

- Пройти обучение двум преподавателям дисциплины «Интеллектуальный анализ данных»
- Пройти сертификацию этим преподавателям
- Рассмотреть возможность сертификации студентов в дистанционном режиме

- **Расширение количества лицензий Deductor Studio Professional**

- В настоящее время имеется одна лицензия, и для комфортного выполнения курсовых проектов этого недостаточно

- **Учет упоминавшихся ранее особенностей обучения магистрантов со специализацией в области программирования**

- Курсовые проекты и магистерские диссертации должны быть направлены не просто на анализ данных, а на *разработку приложений анализа данных*, что возможно только при наличии серверных компонентов Deductor
- При получении такого опыта обязуемся рассказать об этом на следующей конференции



# Итоги работы за прошедший год

- **Обучение и сертификация**
  - доцент Лагерев Д.Г. – 3 сертификата BaseGroup
  - ассистент Бабурин А.Н. – 1 сертификат BaseGroup
  - сертификация 5-и магистрантов
- **Расширение количества лицензий Deductor Studio Professional**
  - Получена лицензия на Deductor Studio Enterprise, Deductor Integration Server, Deductor Analytic Server
- **Учет упоминавшихся ранее особенностей обучения магистрантов *со специализацией в области программирования***
  - Представлены 3 проекта, ориентированные на *разработку приложений анализа данных*, с использованием серверных компонентов Deductor
  - Один из них можно отнести к категории Big Data



# Дальнейшие планы и пожелания

- **Продолжение обучения и сертификации преподавателей**
- **Использование Deductor Enterprise и/или Loginom студентами и аспирантами**
  - для выполнения магистерских и кандидатских диссертаций
  - для выполнения курсовых проектов (продолжить практику использования)
- **Получение дополнительных лицензий для Deductor Enterprise**
  - В настоящее время имеется одна лицензия, и для эффективной работы с большими данными этого недостаточно
  - При наличии только одной лицензии невозможно выполнять лабораторные работы по разработке сервисов DIS
- **Переход на платформу Loginom**
  - Обучение работе с платформой (в рамках мастер-класса 28 июня)
  - Получение доступа к платформе
  - Развертывание платформы Loginom на кафедральном сервере (для работы с большими данными)

IV межвузовская конференция  
*Бизнес-аналитика. Использование аналитической  
платформы Logiном (Deductor) в учебном процессе вуза*



**СПАСИБО ЗА ВНИМАНИЕ !**

*Подвесовский Александр Георгиевич  
заведующий кафедрой, к.т.н., доцент  
apodv@tu-bryansk.ru*

*Лагерев Дмитрий Григорьевич  
к.т.н., доцент  
lagerevdg@mail.ru*

*Бабурин Артем Николаевич  
ассистент  
ababurin@bk.ru*

г. Москва, 27 июня 2017 г.