

# Artiom



## Data Scientist, BI Engineer, Data Engineer

### Key strengths:

- ML and data science programming using Python, Scala, Golang, R.
- Data analysis and machine learning skills. DWH, ETL processes.
- Deep learning skills.
- Knowledge of Telecom, E-commerce, Banking & Finance and Retail domains. Requirements gathering and documenting skills

**7+ years of experience**

## EDUCATION

<b>University</b>	Belarusian State Economic University
<b>Department</b>	Applied mathematics
<b>Degree</b>	Master

## MAIN PROFESSIONAL SKILLS

<b>Programming Languages</b>	Python Golang Matlab R SQL Scala
	Java
<b>Databases, cloud platforms</b>	Mongo BigTable HBase Cassandra Oracle RDBMS PostgreSQL MySQL MariaDB MS SQL Server ClickHouse Vertica Elasticsearch Google BigQuery Google Cloud Platform AWS EC2
<b>Tools</b>	SPSS Eviews Microstrategy BI Oracle BI Tableau Git Jenkins SVN

<b>Technologies</b>	Apache Spark Flask Pandas Numpy Scipy Scikit-learn Keras PyTorch TensorFlow
---------------------	--

## LANGUAGES

Language	Listening	Speaking	Reading	Writing
<b>English</b>	Upper-Intermediate	Upper-Intermediate	Upper-Intermediate	Upper-Intermediate
<b>French</b>	Basic knowledge	Basic knowledge	Basic knowledge	Basic knowledge

## RECENT PROJECTS

### Project #1

<b>Name of the project</b>	<b>Real-time invoice processing data pipeline</b>
<b>Description</b>	Real-time data pipeline built in AWS cloud (step functions + lambda) to parse and process invoice PDF files. Invoice insights (analytics) implemented using Apache Spark. Postgres was used as the main database. The backend API was built using Python FastAPI.
<b>Role</b>	Senior Data Engineer/Backend Engineer
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>• OCR/NLP - correct invoice parsing</li> <li>• Canonical mapping of invoice vendors/line items</li> <li>• Data pipeline with reprocessing and validation functionality and logging to Elasticsearch</li> </ul>
<b>Technologies used</b>	AWS cloud (step functions, lambda-serverless, ECS, Fargate), Docker, Spark, Postgres, Elastic

## Project #2

<b>Name of the project</b>	<b>Recommender system for OSS components based on GitHub data</b>
<b>Description</b>	One of the main goals of the product was to implement a recommendation system which suggests a potential user to follow, compare or analyze similar components with high quality to those already chosen previously. Recommendations were based on both structural similarities between components (NLP of components readmes, title, files & folders structure, then keywords extraction) and usage analytics with defined patterns of behavior on the frontend extracted from yandex metrics raw data and preprocessed in python. ML algorithms and data science stuff were developed on python. Data storage and collection were developed in Mongo for totals and github crawled data and Bigtable for time series.
<b>Role</b>	Chief Data Scientist/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"><li>● Accurate and reliable NLP of readmes and title</li><li>● Specific keywords extraction without general and commonly used words</li><li>● Detection of patterns of behavior in usage analytics raw data, segmentation of users based on these patterns</li></ul>
<b>Technologies used</b>	Python, Mongo, Bigtable, BigQuery, Docker, GitHub

## Project #3

<b>Name of the project</b>	<b>Keywords generation for OSS components based on github data</b>
<b>Description</b>	One of the key goals of the product was development of categories/tags for search engine to make the search of components for users more user-friendly and easier. It was decided to extract these set of tags from the title and readme file of github components. As soon as they were crawled and stored in Mongo, they were processed using several packages like mistune, beautiful soup, nltk and textacy to get only specific tags related to the component without commonly used words and junk words. These keywords were then used in search engine. ML algorithms and data science stuff were developed on python. Data storage and collection were developed in Mongo for totals and github crawled data and Bigtable for time series.
<b>Role</b>	Chief Data Scientist/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"><li>● Accurate and reliable NLP of readmes and title</li><li>● Specific keywords extraction without general and commonly used words</li></ul>
<b>Technologies used</b>	Python, PyTorch Mongo, BigTable, BigQuery, Docker, GitHub

#### Project #4

<b>Name of the project</b>	<b>Road crash prediction model</b>
<b>Description</b>	Classification model to predict an incident. Roads were aggregated to clusters, each cluster had a bunch of metrics, like road sensors, weather information, road characteristics, road constructions, etc. Severe class balancing problem was solved using an imblearn combination of under and oversampling methods.
<b>Role</b>	Chief Data Scientist
<b>Tasks executed</b>	<ul style="list-style-type: none"><li>• Solved class balancing problem</li><li>• Low level of false negatives in prediction</li></ul>
<b>Technologies used</b>	Python, Mongo, Bitbucket

#### Project #5

<b>Name of the project</b>	<b>Real estate assets price prediction</b>
<b>Description</b>	Prediction of next real estate price using regression model approach. The model was implemented using scikit-learn pipelines.
<b>Role</b>	Chief Data Scientist
<b>Technologies used</b>	Python, PostgreSQL, GitHub

#### Project #6

<b>Name of the project</b>	<b>Prediction of sell intent for real estate assets</b>
<b>Description</b>	Binary classification task regarding prediction for a specific property to be sold or not. Data available: 70K properties with historical information about loans, previous owners, other statistics.
<b>Role</b>	Chief Data Scientist
<b>Technologies used</b>	Python, PostgreSQL, GitHub

#### Project #7

<b>Name of the project</b>	<b>Categorization of OSS components based on GitHub data</b>
----------------------------	--

<b>Description</b>	Github currently doesn't provide a categories feature for components. The task was to label the initial dataset using a set of predefined categories, depending on the content of the repository and then build ML model, predicting the specific class of the repository based on text data available: repository description, readme file, files, and folders hierarchical structure of the repository.
<b>Role</b>	Chief Data Scientist/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>• Neural net model written in python using PyTorch framework, predicting the class of repo with avg accuracy of 75%.</li> </ul>
<b>Technologies used</b>	Python, Mongo, GitHub

### Project #8

<b>Name of the project</b>	<b>Marketing mix modeling for an advertising agency</b>
<b>Description</b>	Identify an optimal distribution of investments between various advertising channels like TV, Radio, Internet. The final model had to work efficiently for different products and brands of various clients of the agency.
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>• Singular spectrum analysis for time series smoothing and future trend prediction</li> <li>• Partial least squares (PLS) regression to get the value of variables impact</li> <li>• Non-linear optimization to get optimal resources allocation</li> <li>• Established reliable functions to explain nonlinear dependencies between variables</li> <li>• Developed a reliable and accurate PLS regression model for preprocessed data</li> </ul>
<b>Technologies used</b>	Python (sklearn, scipy, numpy, pandas)

### Project #9

<b>Name of the project</b>	<b>Churn prediction model for a telecom company</b>
<b>Description</b>	Implement a mathematical model with automatically gathered data from various sources with the ability to give as an output a probability of churn for each client in the upcoming month.

<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>● Big data analysis of different sources: billing, network data, crm, client location info, mnp, call center data</li> <li>● Implementation of a new architecture for data collection and Preprocessing</li> <li>● Developing of ETL process for gathering necessary data for the model</li> <li>● Classification algorithms based on python sklearn: Log Regression, SVM, Decision Trees, Random Forest, Gradient Boosting Algorithms</li> <li>● Clustering based on python sklearn built-in algorithms</li> </ul>
<b>Technologies used</b>	Python, Oracle RDBMS (main data storage), Cloudera, MySQL

### Project #10

<b>Name of the project</b>	<b>Scoring model for a telecom company</b>
<b>Description</b>	The goal of the project was to predict the client's default and not to sell obligatory devices to unreliable clients (future debtors). The output of the model was a probability of being a debtor in a shortcoming perspective. Raw data we collected from Oracle RDBMS. Data analysis was performed using python pandas, numpy, re, matplotlib, seaborn, further data preprocessing and ETL was made on pl/sql. The mathematical model was developed using python imblearn & sklearn. The best results in terms of database integration velocity and accuracy were obtained using Random Forest. The results were also used for segmentation based on the client's financial status and paying discipline.
<b>Role</b>	Chief Data Scientist/Business Analyst/Project Manager
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>● Big data analysis of different sources: billing, network data, crm, mnp, clients' photos, National Bank credit info data</li> <li>● Implementation of a new architecture of data collection and preprocessing</li> <li>● Creation of GUI interface for the scoring machine in the CRM system for sellers</li> <li>● Validation test of several algorithms: log regression, decision trees, SVM, random forest, gradient boosting (chosen the best one)</li> <li>● Developing ETL process for gathering necessary data for the model</li> </ul>
<b>Technologies used</b>	Python, Oracle RDBMS (main data storage), MS Office, Windows

### Project #11

<b>Name of the project</b>	<b>Multiclass classification model for an e-commerce company</b>
----------------------------	--

<b>Description</b>	The goal of the project was to predict the probability of purchase for each client of an e-commerce company. Raw data was collected using Logs API Yandex Metrics and analyzed in python using numpy, pandas. For storage and aggregation purposes we developed Clickhouse DB in cloud on Microsoft Azure under Ubuntu 16.04. Using python clickhouse module we preprocessed all the data and developed ETL process for data processing on a daily basis. The mathematical model was developed using python imblearn & sklearn. Algorithms we used for classification on a validation stage: Log Regression, SVM, Decision Trees, Random Forest, Gradient Boosting Algorithms. We divided a set of clients with predicted probability into 3 groups (excellent, good, normal) and sent them as segments to the advertisement platform using API for to retarget them and to attract new users using look-alike ML technology to increase conversion rate and decrease marketing campaign costs.
<b>Role</b>	Data Scientist/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>● Automatic data exchange with Logs API</li> <li>● Implementation of data collection and preprocessing using Clickhouse in the cloud under Ubuntu</li> <li>● Math model developing using python with algorithms validation checking procedure</li> </ul>
<b>Technologies used</b>	Python, Clickhouse DB (main data storage), Ubuntu

### Project #12

<b>Name of the project</b>	<b>DWH transformation and modernization for a telecom company</b>
<b>Description</b>	The goal of the project was an improvement of the ETL process, increasing the quality, accuracy, and velocity of the data collection, creating a new data model and data marts for BI tool and adding data from new sources with the enhancement of current attributes and metrics for the reporting purposes.
<b>Role</b>	Data Analyst/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>● Implementation of a new architecture for data collection</li> <li>● Creation of new tables and views in Oracle RDBMS and new data marts for Microstrategy to meet the requirements of the business</li> <li>● Increasing the velocity of data collection during ETL process</li> <li>● Creation of new reports and dashboards in Microstrategy based on DWH data</li> <li>● Checking system for possible DWH bugs and defects</li> <li>● Update of Microstrategy BI to a newer version</li> </ul>
<b>Technologies used</b>	Python, Microstrategy, Oracle RDBMS

### Project #13

<b>Name of the project</b>	<b>Descriptive analysis for a telecom company</b>
----------------------------	---



<b>Description</b>	The project goal was to identify a typical client of the company and the patterns of the clients' behavior. The analysis was divided into several main parts: overall analysis with different dimensions, zero ARPU clients' analysis, data traffic usage analysis.
<b>Role</b>	Data Scientist/Business Analyst/Architect
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>• The overall analysis of the data in different dimensions: which data we have, what is the level of their accuracy, where are they stored</li> <li>• Creation of a dictionary of features that explained clients behavior with links to corresponding DBs</li> <li>• Zero ARPU analysis using python pandas, numpy, matplotlib, seaborn with recommendations regarding how to decrease the number of such clients</li> <li>• Data usage analysis using python pandas, numpy, matplotlib, seaborn</li> <li>• Analysis of necessary database architecture improvement</li> </ul>
<b>Technologies used</b>	Python, Oracle RDBMS, MySQL, MS Office, Windows

#### Project #14

<b>Name of the project</b>	<b>Call center performance analytical platform for a telecom company</b>
<b>Description</b>	The project goal was to develop a set of KPI for the call center and then create a system of online updated dashboards using the data in several database sources (Oracle, SQL Server, MySQL). Dashboards were implemented using Microstrategy 10 with creation of attributes and metrics and overall warehouse catalog improvement.
<b>Role</b>	Business Analyst/Microstrategy Architect
<b>Tasks executed</b>	<ul style="list-style-type: none"> <li>• Analysis of several data sources (Oracle, SQL Server, MySQL)</li> <li>• Creation of connections from Microstrategy to necessary DBs</li> <li>• Creation of attributes, metrics and OLAP cubes in Microstrategy for implementation of dashboards</li> <li>• Implementation of dashboards using Microstrategy 10 in the customer's native language</li> </ul>
<b>Technologies used</b>	Microstrategy, Oracle RDBMS, SQL Server, MySQL, Windows

#### Project #15

<b>Name of the project</b>	<b>Financial assets forecasting</b>
<b>Description</b>	The goal of the project was to predict financial assets of the bank for accurate forecasting of the business plan KPI. Data analysis was developed in python pandas, numpy, matplotlib and seaborn, math model was developed using matlab. Due to high quasi-periodicity of data and their chaotic behavior, it was decided to use for prediction an algorithm singular spectrum analysis with principal component analysis to decrease the dimensionality of data series and smooth the data.

<b>Role</b>	Data Scientist/Business Analyst
<b>Tasks executed</b>	<ul style="list-style-type: none"><li>• Developing a mathematical model for assets prediction</li><li>• Implementation of a new business plan template with an integrated model in it</li></ul>
<b>Technologies used</b>	Python, Matlab