

Deploy in production

Or, pipeline to save the world

SKB Kontur

- B2B
- 40+ products
- 7700+ employee
- 1000+ devs
- 1200+ tech support

Products

- Tax reporting
- Accounting, salary, personnel
- Electronic document management
- Trade
- Checking counterparties
- Electronic bidding and procurement
- Interaction with GIS

ML Tasks

Добавление товара



[▶ Как заполнить карточку товара](#)

Наименование

Макароны увелка 450г

☐

Весовой товар

Группа

Бакалея > Макаронные изделия

[✎ Изменить группу](#)



Документы

+ Новый документ

Входящие документы



Контрагенты



Сообщения



Справочная

Возможно, это ваши поставщики. Пригласите их и получайте документы быстрее



ИП Юзефпольский А. М.
ОАО «Ростелеком»
ООО «Снег с дождем»

[Перейти к приглашению](#)

асование

Переместить

Аннулирование

Удал

Статус

91-A от 20.02.15

● Требуется подпись

от 20.02.15

Подписан

000,00 ₽



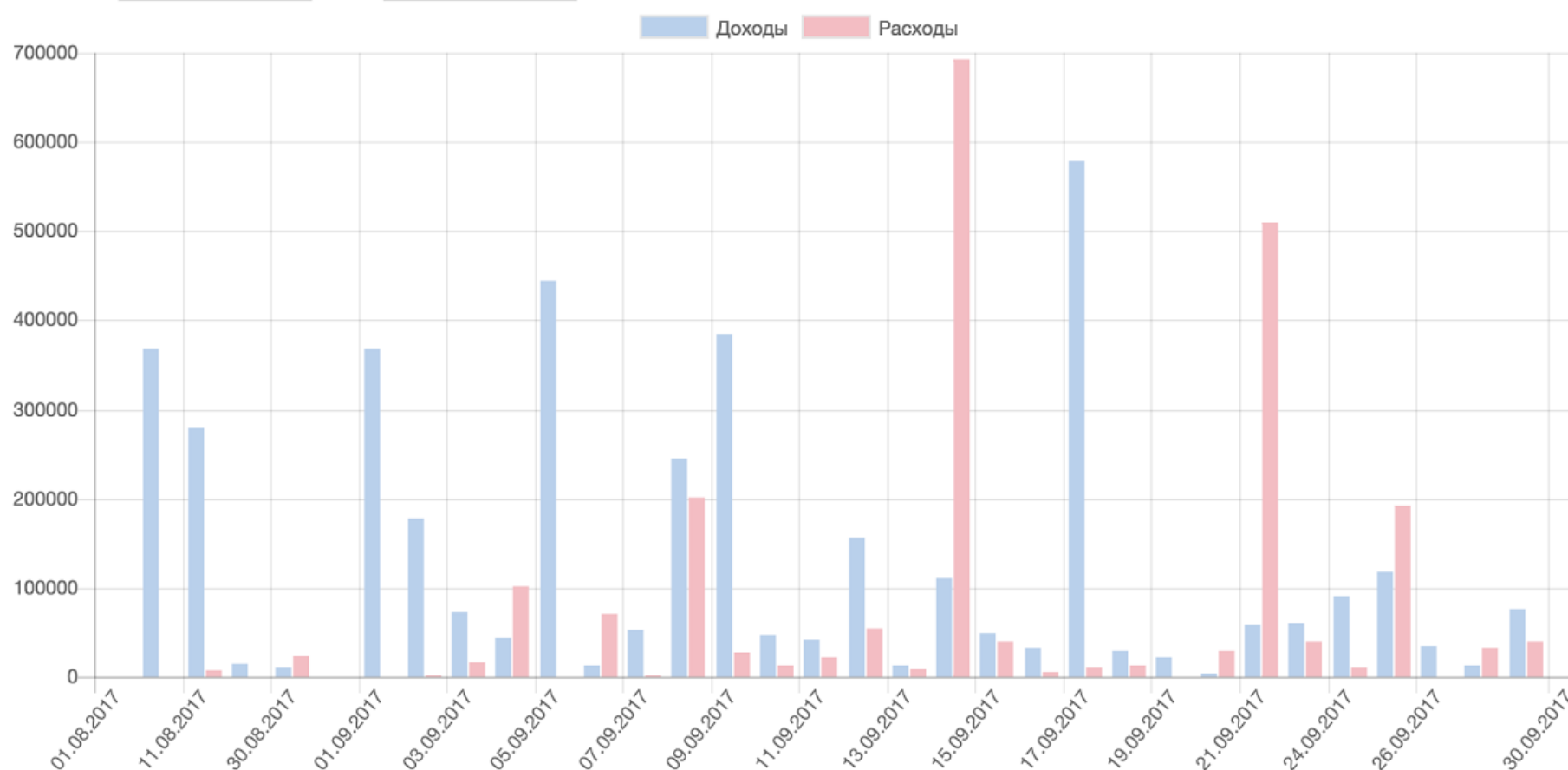
Аналитика

➤ Если у вас есть пожелания по развитию сервиса — напишите нам об этом

Загрузить банковскую выписку

По всем счетам

Период 01.08.2017 — 01.10.2017



Остаток
1 952 200,71 ₽

Доходы
4 006 178,76 ₽

Расходы
2 178 187,08 ₽

Доходы



Расходы



Tech Stack

- Langs
 - Reality: .Net
 - Data Science: + Python, Java/Scala, C++
- Databases
 - Reality: MongoDB, Cassandra, MsSql
 - Data Science: + Hadoop/Spark
- Deploy
 - Reality: Windows, Octopus, TeamCity, Houston, Vostok
 - Data Science: + Linux, Docke

Principle Scheme

- Data Science as a service
 - Task
 - Hypothys
 - Research
 - MVP
 - Evaluation
 - Integration
 - Support and extension

Problem



Research model

!=

Production model

Solution

Scikit-Learn Pipeline

Concept

- Every preprocessing step is sk-learn compatible transformer
- Preprocessing + decision model as Pipeline

Concept

- All params are accessible from outside (easy grid search)
- Model is serialisable

Research



Deploy



Practice

- Create class for your transformation
- Inherit sklearn's BaseEstimator and TransformerMixin
- Implement fit, transform
- Do not compute **ANYTHING** in `__init__`

Practice

```
class MyTransformer(BaseEstimator, TransformerMixin):
    def __init__(self, value):
        # no logic here
        self.value = value

    def fit(self, X: np.array, y=None):
        return self

    def transform(self, X: np.array):
        # logic here
        if self.value < 0:
            self.value = 0

        return X * value
```

Practice

- Use FeatureUnion for parallel transformation
- Use memory arg in Pipeline to cache transformation (for grid search)
- Use sklearn-pandas Mapper for DataFrame

Practice

```
model = Pipeline(  
    steps=[  
        ('fill_na', FillNa()),  
        ('union', FeatureUnion([  
            ('text_vect', CountVectoriser()),  
            ('have_name', TextMatcher())  
        ])),  
        ('clf', LGBMClassifier)  
    ]  
)
```

Practice

- Create separate library and store all transformation in it
- Share this lib across projects

Practice

- Use joblib to serialise model
joblib from sklearn **not compatible** with separate joblib, choose one
- Use database to store model
- Use docker to deploy your service
- Try Pipenv and Pipefile
- Try Kafka for data streams

You good