



SISTEMI ZA OBRADU I ANALIZU VELIKE KOLIČINE PODATAKA - Projekat 3

Filip Trajković 1574

Spisak kontejnera

project3	-	Running (13/1)	⋮
producer-3 6ec4f81b6671	project3-producer-3	Running	2 hours ago
grafana 71d45dc7da40	grafana/grafana	Running	3000:3000
influxdb 99e7fc7f22ae	influxdb:2.1.1	Running	8086:8086
spark-worker-1-3 9862025a776f	bde2020/spark-worker:3.1.2-hadoop3.2	Running	8071:8071
kafka-3 6116e0ba6bfc	wurstmeister/kafka:2.13-2.7.0	Running	9091:9091
spark-worker-2-3 f17857b32418	bde2020/spark-worker:3.1.2-hadoop3.2	Running	8072:8071
datanode-3 fcc091cff35e	bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8	Running	2 hours ago
zookeeper-3 dfffe329403ad	wurstmeister/zookeeper:latest	Running	2181:2181
namenode-3 660c9b46b75a	bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8	Running	9000:9000 9870:9870
nodemanager-3 8b289c347142	bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8	Running	2 hours ago
resourcemanager-3 06cb7d146f32	bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8	Running	1 hour ago
historyserver-3 0f10928ac545	bde2020/hadoop-historyserver:2.0.0-hadoop3.2.1-java8	Running	2 hours ago
spark-master-3 40077f2a7138	bde2020/spark-master:3.1.2-hadoop3.2	Running	7077:7077 8070:8070



Izgled Producer skripte



```
8 producer = KafkaProducer(  
9     bootstrap_servers=[os.environ["KAFKA_HOST"]],  
10    value_serializer=lambda v: json.dumps(v).encode("utf-8"),  
11    api_version=(0, 11),  
12 )  
13  
14 while True:  
15     with open(os.environ["DATA"], "r") as file:  
16         reader = csv.reader(file, delimiter=",")  
17         headers = next(reader)  
18         for row in reader:  
19             value = {headers[i]: row[i] for i in range(len(headers))}  
20             value["lat"] = float(value["lat"])  
21             value["lon"] = float(value["lon"])  
22             value["alt"] = float(value["alt"])  
23             value["user"] = int(value["user"])  
24             value["ts"] = int(time.time())  
25             value["year"] = int(value["year"])  
26             value["month"] = int(value["month"])  
27             value["day"] = int(value["day"])  
28             value["hour"] = int(value["hour"])  
29             value["minute"] = int(value["minute"])  
30             value["second"] = int(value["second"])  
31             print(value)  
32             producer.send(os.environ["KAFKA_TOPIC"], value=value)  
33             time.sleep(float(os.environ["KAFKA_INTERVAL"]))
```

Aplikacija za kreiranje modela za klasifikaciju



Dockerfile

Project 3 > model-creation > Dockerfile > ...

```
1  FROM bde2020/spark-python-template:3.1.2-hadoop3.2
2
3  RUN apk add --no-cache py3-numpy
4
5  ENV CLASSIFICATION_MODEL='hdfs://namenode:9000/model/RfModel'
6  ENV SCALER_LOC="hdfs://namenode:9000/model/Scaler"
7  ENV INDEXER_PIPELINE_LOC="hdfs://namenode:9000/model/IndexerPipeline"
8
9  ENV SPARK_APPLICATION_PYTHON_LOCATION /app/model_creation.py
10 ENV SPARK_APPLICATION_ARGS "hdfs://namenode:9000/user/root/geolife_gps_sorted_without_nodata.csv"
11 ENV SPARK_SUBMIT_ARGS --executor-memory 4G --executor-cores 4
```

Aplikacija za kreiranje modela za klasifikaciju



Shell skripta za pokretanje

```
Project 3 > model-creation > $ start-streaming-app.sh
1  #! /bin/bash
2
3  docker container rm ModelCreationApp
4
5  docker build --rm -t bde/model_creation_app .
6
7  docker run --name ModelCreationApp --net FT_project3_network -p 4042:4042 bde/model_creation_app
8
9
```

Aplikacija za kreiranje modela za klasifikaciju

Uključene biblioteke,
šema podataka i
environment
promenljive



```
Project 3 > model-creation > model_creation.py > ...
1  import os
2  from pyspark.sql import SparkSession
3  from pyspark.sql.types import StructType, StructField, StringType, FloatType, IntegerType, DoubleType
4  from pyspark.sql.functions import col, from_json, window, count, mean, max, min
5  from pyspark import SparkConf
6  import sys
7  from pyspark.ml.feature import StringIndexer, VectorAssembler, MinMaxScaler, StringIndexerModel
8  from pyspark.ml.classification import RandomForestClassifier
9  from pyspark.ml.evaluation import MulticlassClassificationEvaluator
10 from pyspark.ml import PipelineModel
11 from pyspark.ml import Pipeline
12
13 schema = StructType(
14     [
15         StructField("lat", FloatType()),
16         StructField("lon", FloatType()),
17         StructField("alt", FloatType()),
18         StructField("label", StringType()),
19         StructField("user", StringType()),
20         StructField("year", IntegerType()),
21         StructField("month", IntegerType()),
22         StructField("day", IntegerType()),
23         StructField("hour", IntegerType()),
24         StructField("minute", IntegerType()),
25         StructField("second", IntegerType())
26     ]
27 )
28
29 if __name__ == '__main__':
30     if len(sys.argv) < 2:
31         print("Usage: main.py <input folder> ")
32         exit(-1)
33
34 SCALER_LOCATION = os.getenv('SCALER_LOC')
35 INDEXER_PIPELINE_LOC = os.getenv('INDEXER_PIPELINE_LOC')
36 CLASSIFICATION_MODEL = os.getenv('CLASSIFICATION_MODEL')
37
38
```

→ Aplikacija za kreiranje modela za klasifikaciju

Preprocesiranje podataka

```
62     label_indexer = StringIndexer(  
63         inputCol="label", outputCol="label_index", handleInvalid="keep")  
64     user_indexer = StringIndexer(  
65         inputCol="user", outputCol="user_index")  
66     # pipeline = Pipeline(stages=[label_indexer])  
67  
68     pipeline = Pipeline(stages=[label_indexer, user_indexer])  
69  
70     fitted_indexer = pipeline.fit(df)  
71  
72     fitted_indexer.write().overwrite().save(INDEXER_PIPELINE_LOC) #("hdfs://namenode:9000/model/IndexerPipeline")  
73  
74     indexer_model = PipelineModel.load(INDEXER_PIPELINE_LOC) #("hdfs://namenode:9000/model/IndexerPipeline")  
75  
76     indexedDf = indexer_model.transform(df)  
77  
78     # indexedDf = fitted_label_indexer.transform(df)  
79  
80     indexedDf = indexedDf.drop("label")  
81     indexedDf = indexedDf.drop("user")  
82  
83     features = ["lat", "lon", "alt", "user_index",  
84                  "year", "month", "day", "hour", "minute", "second"]  
85  
86     featureAssembler = VectorAssembler(  
87         inputCols=features, outputCol="features")  
88  
89     assembleredDf = featureAssembler.transform(indexedDf)  
90  
91     scaler = MinMaxScaler(inputCol="features", outputCol="scaledFeatures")  
92  
93     fitScaler = scaler.fit(assembleredDf)  
94  
95     fitScaler.write().overwrite().save(SCALER_LOCATION)  
96  
97  
98     scaledData = fitScaler.transform(assembleredDf)  
99  
100    scaledData.show()  
101  
102    finalDf = scaledData.select("scaledFeatures", "label_index")  
103
```

Aplikacija za kreiranje modela za klasifikaciju

Kreiranje i evaluacija modela

```
112
113     train_data, test_data = finalDf.randomSplit([0.8, 0.2], seed=42)
114
115     rf = RandomForestClassifier(
116         featuresCol='scaledFeatures', labelCol='label_index', maxDepth=10)
117
118     rfModel = rf.fit(train_data)
119
120     predictions = rfModel.transform(test_data)
121
122
123     # metricName="f1"
124     # metricName="precisionByLabel"
125     # metricName="recallByLabel"
126
127     evaluatorAccuracy = MulticlassClassificationEvaluator(
128         labelCol="label_index", predictionCol="prediction", metricName="accuracy")
129     evaluatorF1 = MulticlassClassificationEvaluator(
130         labelCol="label_index", predictionCol="prediction", metricName="f1")
131     evaluatorByLabel = MulticlassClassificationEvaluator(
132         labelCol="label_index", predictionCol="prediction", metricName="precisionByLabel")
133     evaluatorRecallByLabel = MulticlassClassificationEvaluator(
134         labelCol="label_index", predictionCol="prediction", metricName="recallByLabel")
135
136     print(evaluatorAccuracy.evaluate(predictions))
137     print(evaluatorAccuracy.evaluate(predictions))
138     print(evaluatorByLabel.evaluate(predictions))
139     print(evaluatorRecallByLabel.evaluate(predictions))
140
141     rfModel.write().overwrite().save(CLASSIFICATION_MODEL) #("hdfs://namenode:9000/model/RfModel")
142
```

Applikacija za predikciju novih podataka



Shell skripta za pokretanje

```
Project 3 > spark-streaming-prediction > $ start-streaming-app.sh
1  #! /bin/bash
2
3  docker container rm StreamingPredictionApp
4
5  docker build --rm -t bde/streaming_prediction_app .
6
7  docker run --name StreamingPredictionApp --net FT_project3_network -p 4043:4043 bde/streaming_prediction_app
8
9
```

Aplikacija za predikciju novih podataka

Dockerfile

```
Project 3 > spark-streaming-prediction > Dockerfile > ...
1  FROM bde2020/spark-python-template:3.1.2-hadoop3.2
2
3  RUN apk add --no-cache py3-numpy
4  RUN cd /app pip install -r requirements.txt
5
6  ENV INFLUXDB_HOST "influxdb"
7  ENV INFLUXDB_PORT "8086"
8  ENV INFLUXDB_USERNAME "admin"
9  ENV INFLUXDB_PASSWORD "admin"
10 ENV INFLUXDB_TOKEN "n6zAU61C3paPpmfY3dvbQLoVORcqcd0STCCw811VB7MUxEKbppao5CSUg8sMaRzSlWiNI090H6iHMzAwfU_1yw=="
11 ENV INFLUXDB_ORG "bigdata_projekat"
12 ENV INFLUXDB_BUCKET "locations"
13
14 ENV CLASSIFICATION_MODEL='hdfs://namenode:9000/model/RfModel'
15 ENV SCALER_LOC="hdfs://namenode:9000/model/Scaler"
16 ENV INDEXER_PIPELINE_LOC="hdfs://namenode:9000/model/IndexerPipeline"
17
18 ENV KAFKA_HOST=kafka-3:9092
19 ENV KAFKA_TOPIC=locations
20 ENV KAFKA_CONSUMER_GROUP=Spark-Group
21
22 ENV SPARK_APPLICATION_PYTHON_LOCATION /app/streaming_prediction.py
23 ENV SPARK_APPLICATION_ARGS "hdfs://namenode:9000/user/root/geolife_gps_sorted_without_nodata.csv"
24 ENV SPARK_SUBMIT_ARGS --packages org.apache.spark:spark-streaming-kafka-0-10_2.12:3.1.2,org.apache.spark:spark-sql-kafka-0-10_2.12:3.1.2 --executor-memory 4G --executor-cores 4
```

Aplikacija za predikciju

novih podataka

Uključene biblioteke,

šema podataka i

environment

promenljive

```
1 import os
2 from pyspark.sql import SparkSession
3 from pyspark.sql.types import StructType, StructField, StringType, FloatType, IntegerType, DoubleType
4 from pyspark.sql.functions import col, from_json, window, count, mean, max, min
5 from pyspark import SparkConf
6 import sys
7 from pyspark.ml.feature import StringIndexer, VectorAssembler, MinMaxScalerModel
8 from pyspark.ml.classification import RandomForestClassifier, RandomForestClassificationModel
9 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
10 from pyspark.ml import PipelineModel
11 from pyspark.ml import Pipeline
12
13 # InfluxDB libraries
14 from influxdb_client import InfluxDBClient, Point, WritePrecision
15 from influxdb_client.client.write_api import SYNCHRONOUS
16 from datetime import datetime
17
18 schema = StructType(
19     [
20         StructField("lat", FloatType()),
21         StructField("lon", FloatType()),
22         StructField("alt", FloatType()),
23         StructField("label", StringType()),
24         StructField("user", StringType()),
25         StructField("year", IntegerType()),
26         StructField("month", IntegerType()),
27         StructField("day", IntegerType()),
28         StructField("hour", IntegerType()),
29         StructField("minute", IntegerType()),
30         StructField("second", IntegerType())
31     ]
32 )
33
34 INFLUXDB_ORG = os.getenv('INFLUXDB_ORG')
35 INFLUXDB_BUCKET = os.getenv('INFLUXDB_BUCKET')
36 INFLUXDB_HOST = os.getenv('INFLUXDB_HOST')
37 INFLUXDB_PORT = os.getenv('INFLUXDB_PORT')
38 INFLUXDB_TOKEN = os.getenv('INFLUXDB_TOKEN')
39
```

Applikacija za predikciju



novih podataka

Klasa za upis u
InfluxDB
bazu podataka

```
39
40     class InfluxDBWriter:
41         def __init__(self):
42             self._org = INFLUXDB_ORG
43             self._token = INFLUXDB_TOKEN
44             self.client = InfluxDBClient(url=f"http:///{INFLUXDB_HOST}:{INFLUXDB_PORT}", token=self._token, org=self._org)
45             self.write_api = self.client.write_api(write_options=SYNCHRONOUS)
46
47         def open(self, partition_id, epoch_id):
48             print(f"Opened {partition_id}, {epoch_id}")
49             return True
50
51         def process(self, row):
52             self.write_api.write(bucket=INFLUXDB_BUCKET, record=self._row_to_point(row))
53
54         def close(self, error):
55             self.write_api.__del__()
56             self.client.__del__()
57             print(f"Closed with error: {error}")
58
59         def _row_to_point(self, row):
60             print(row)
61
62             # timestamp = datetime.utcnow().strftime('%Y-%m-%dT%H:%M:%S')
63
64             point = (
65                 Point.measurement("locations_measurement2")
66                 .tag("measure", "locations_measurement2")
67                 .time(datetime.utcnow(), WritePrecision.NS)
68                 .field("Latitude", float(row['lat']))
69                 .field("Longitude", float(row['lon']))
70                 .field("Altitude", float(row['alt']))
71                 .field("Label_index", int(row['label_index'])))
72                 .field("User_index", int(row['user_index'])))
73                 .field("Label_prediction", int(row['prediction'])))
74             )
75             return point
76
```

Aplikacija za predikciju



novih podataka

Postavljanje environment
promenljivih i učitavanje
podataka sa Kafka topic-a

```
76
77 if __name__ == '__main__':
78     if len(sys.argv) < 2:
79         print("Usage: main.py <input folder> ")
80         exit(-1)
81
82 INDEXER_PIPELINE_LOC = os.getenv('INDEXER_PIPELINE_LOC')
83 SCALER_LOCATION = os.getenv('SCALER_LOC')
84 CLASSIFICATION_MODEL = os.getenv('CLASSIFICATION_MODEL')
85
86 INFLUXDB_HOST = os.getenv('INFLUXDB_HOST')
87 INFLUXDB_PORT = os.getenv('INFLUXDB_PORT')
88 INFLUXDB_USERNAME = os.getenv('INFLUXDB_USERNAME')
89 INFLUXDB_PASSWORD = os.getenv('INFLUXDB_PASSWORD')
90 INFLUXDB_DATABASE = os.getenv('INFLUXDB_DATABASE')
91
92 KAFKA_HOST = os.getenv('KAFKA_HOST')
93 KAFKA_TOPIC = os.getenv('KAFKA_TOPIC')
94 KAFKA_CONSUMER_GROUP = os.getenv('KAFKA_CONSUMER_GROUP')
95
96
97 conf = SparkConf()
98 conf.setMaster("spark://spark-master-3:7077")
99 #conf.setMaster("local")
100 # conf.set("spark.driver.memory", "4g")
101 # conf.set("spark.cassandra.connection.host", "cassandra")
102 # conf.set("spark.cassandra.connection.port", 9042)
103
104 spark = SparkSession.builder.config(
105     conf=conf).appName("StreamingPredictionApp").getOrCreate()
106
107 spark.sparkContext.setLogLevel("ERROR")
108
109 df = (
110     spark.readStream.format("kafka")
111     .option("kafka.bootstrap.servers", KAFKA_HOST)
112     .option("subscribe", KAFKA_TOPIC)
113     .option("startingOffsets", "latest")
114     .option("groupIdPrefix", KAFKA_CONSUMER_GROUP)
115     .load()
116 )
117
```

Applikacija za predikciju

novih podataka

Primena

preprocesiranja

nad novim podacima,

izvršenje predikcije

i upis

u InfluxDB

bazu podataka

```
117
118     loaded_values = df.select(
119         from_json(col("value").cast(
120             "string"), schema).alias("data")
121     )
122
123     parsed_values = loaded_values.selectExpr("data.*")
124     parsed_values.printSchema()
125
126     indexer_model = PipelineModel.load(INDEXER_PIPELINE_LOC) #("hdfs://namenode:9000/model/IndexerPipeline")
127     scaler_model = MinMaxScalerModel.load(SCALER_LOCATION)
128     classification_model = RandomForestClassificationModel.load(CLASSIFICATION_MODEL)
129
130
131     indexedDf = indexer_model.transform(parsed_values)
132
133     indexedDf = indexedDf.drop("label")
134     indexedDf = indexedDf.drop("user")
135
136     features = ["lat", "lon", "alt", "user_index",
137                 "year", "month", "day", "hour", "minute", "second"]
138
139     featureAssembler = VectorAssembler(
140         inputCols=features, outputCol="features")
141
142     assembleredDf = featureAssembler.transform(indexedDf)
143
144     finalData = scaler_model.transform(assembleredDf)
145
146     prediction = classification_model.transform(finalData)
147
148     prediction.printSchema()
149
150     query = prediction.writeStream \
151         .foreach(InfluxDBWriter()) \
152         .start()
153     query.awaitTermination()
```

Aplikacija za predikciju novih podataka

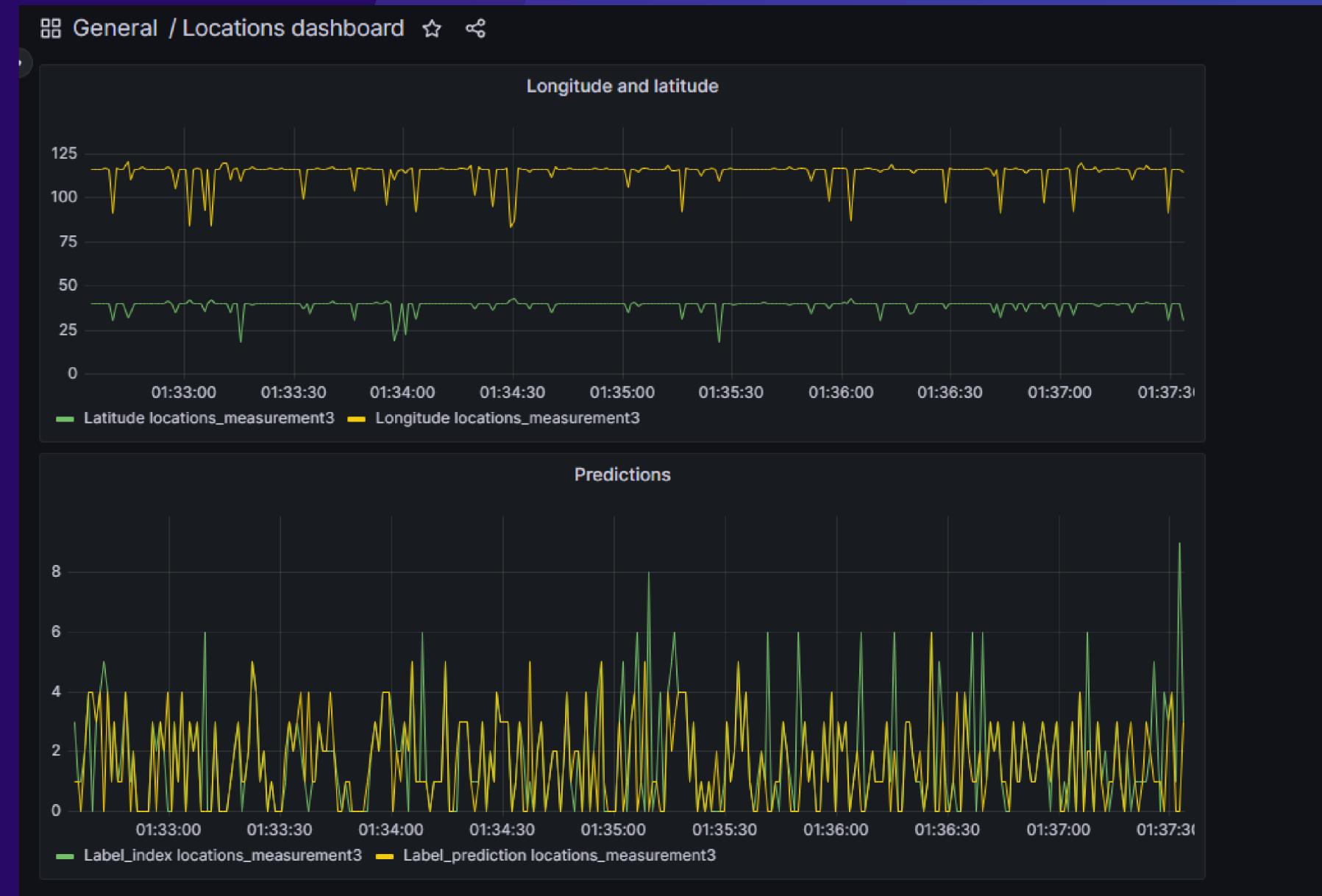
Prikaz podataka u InluxDB

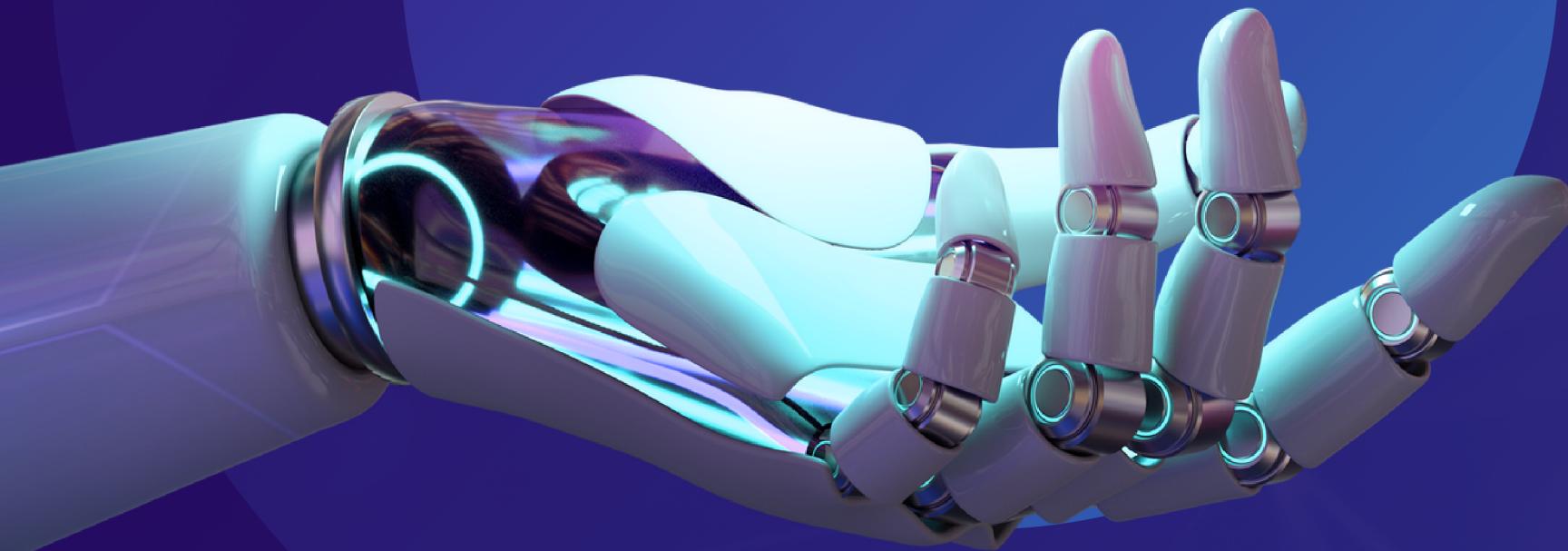


Aplikacija za predikciju

novih podataka

Prikaz podataka preko Grafane





HVALA
NA
PAŽNJI



ficax.ai@elfak.rs

