

# Развертывание ML-моделей в продакшн

**Ирина Степановна Трубчик**

<https://t.me/+PsC-JDrwrvsxNmVi>

Лекция 7

# Цели занятия

- 1 Понять жизненный цикл модели от обучения до продакшена
- 2 Изучить паттерны развертывания (batch, online, streaming)
- 3 Освоить инструменты: Docker, Kubernetes, облачные платформы
- 4 Научиться мониторить и обновлять модели в production



*Сталкивались ли вы с развертыванием моделей?  
Какие проблемы возникали?*



# Что такое deployment ML-модели

Model Deployment — процесс интеграции обученной ML-модели в production-окружение, где она будет делать предсказания на реальных данных.

# КЛЮЧЕВЫЕ КОМПОНЕНТЫ:

- **Model artifacts** — сериализованная модель (pickle, joblib, ONNX)
- **Inference service** — сервис, обрабатывающий запросы
- **Dependencies** — библиотеки и окружение
- **Infrastructure** — серверы, контейнеры, оркестраторы

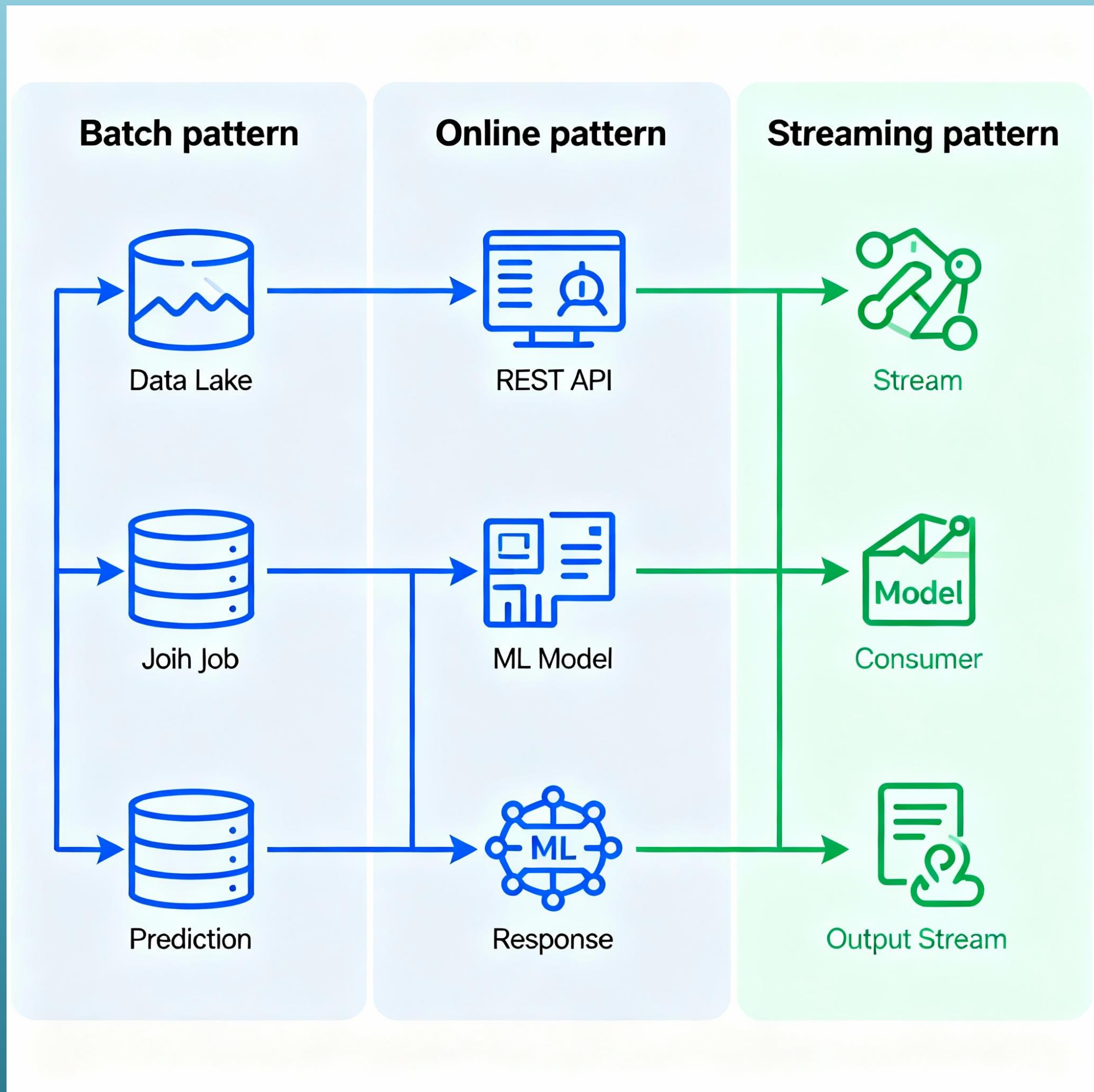
# Схема жизненного цикла



# Паттерны развертывания

Batch Prediction (пакетная обработка)	Online/Real-time Prediction (реалтайм)	Streaming Prediction
<ul style="list-style-type: none"><li>➤ Обработка данных по расписанию (daily, hourly)</li><li>➤ Подходит для некритичных по времени задач</li><li>➤ Пример: рекомендации email-рассылок</li></ul>	<ul style="list-style-type: none"><li>➤ API-сервис отвечает на запросы немедленно</li><li>➤ Низкая latency (&lt;100ms)</li><li>➤ Пример: fraud detection, рекомендации в e-commerce</li></ul>	<ul style="list-style-type: none"><li>➤ Обработка потоковых данных (Kafka, Kinesis)</li><li>➤ Для непрерывных потоков событий</li><li>➤ Пример: мониторинг IoT-датчиков</li></ul>

# Диаграмма паттернов



# Serving архитектуры

Model-in-Service (embedded)	Model-as-Service (dedicated)	Sidecar Pattern
<div>python</div> <pre># Модель загружена в памяти сервиса from fastapi import FastAPI model = load_model("model.pkl")  @app.post("/predict") def predict(data):     return model.predict(data)</pre> <div>✓ Простота, низкая задержка (latency) ✗ Сложность обновления</div>	<div>text</div> <pre>[Client] → [Gateway] → [Model Service 1]                         → [Model Service 2]</pre> <div>✓ Независимое масштабирование ✗ Дополнительная latency</div>	<div>text</div> <pre>[App Container] ↔ [Model Container]</pre> <div>✓ Изоляция, легкое обновление ✗ Дополнительные ресурсы</div>



# Контейнеризация с Docker

**Зачем контейнеры для ML?**

- **Изоляция зависимостей**
- **Воспроизводимость окружения**
- **Портабельность между средами**

# Минимальный Dockerfile:

text

```
FROM python:3.10-slim

WORKDIR /app
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

COPY models/ models/
COPY src/ src/

EXPOSE 8080
CMD ["uvicorn", "src.api:app", "--host", "0.0.0.0", "--port", "8080"]
```

## Сборка и запуск:

bash

```
docker build -t ml-model:v1 .
docker run -p 8080:8080 ml-model:v1
```



**Какие еще  
зависимости  
нужно включить  
в образ?**

# Kubernetes для ML

## Почему Kubernetes?

- Автоматическое масштабирование (HРА)
- Self-healing (перезапуск упавших pods)
- Обнаружение служб и балансировка нагрузки
- Управление secrets и configs

# Минимальный Deployment manifest:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: ml-model
spec:
  replicas: 3
  selector:
    matchLabels:
      app: ml-model
  template:
    metadata:
      labels:
        app: ml-model
    spec:
      containers:
      - name: model
        image: ml-model:v1
        ports:
        - containerPort: 8080
      resources:
        requests:
          memory: "512Mi"
          cpu: "500m"
        limits:
          memory: "1Gi"
          cpu: "1000m"
```

```
---
apiVersion: v1
kind: Service
metadata:
  name: ml-model-service
spec:
  selector:
    app: ml-model
  ports:
  - port: 80
    targetPort: 8080
  type: LoadBalancer
```

# Стратегии развертывания

Blue-Green Deployment	Canary Deployment	A/B тестирование
<div>text</div> <div>[Blue (v1)] ← 100% traffic [Green (v2)] ← ready</div> <div>→ Switch traffic instantly</div> <div>text</div> <div>[Blue (v1)] ← idle [Green (v2)] ← 100% traffic</div> <div>✓ Мгновенный rollback</div> <div>✗ Требуется 2x ресурсов</div>	<div>text</div> <div>[v1] ← 90% traffic [v2] ← 10% traffic</div> <div>→ Gradually increase</div> <div>text</div> <div>[v1] ← 50% traffic [v2] ← 50% traffic</div> <div>✓ Снижение рисков</div> <div>✗ Сложность мониторинга</div>	<div>text</div> <div>Group A → [Model v1] Group B → [Model v2]</div> <div>✓ Бизнес-метрики</div> <div>✗ Требуется аналитика</div>

# Облачные платформы

 **SageMaker**

 **Google Vertex AI**

 **Microsoft  
Azure ML**

 **Yandex Cloud**



python

```
from sagemaker.sklearn import SKLearnModel

model = SKLearnModel(
    model_data="s3://bucket/model.tar.gz",
    role=role,
    entry_point="inference.py"
)

predictor = model.deploy(
    instance_type="ml.m5.large",
    initial_instance_count=2
)
```



python

```
from google.cloud import aiplatform

endpoint = aiplatform.Endpoint.create(display_name="ml-model")
endpoint.deploy(
    model=model,
    machine_type="n1-standard-4",
    min_replica_count=1,
    max_replica_count=10
)
```



python

```
from azureml.core.webservice import AciWebservice

deployment_config = AciWebservice.deploy_configuration(
    cpu_cores=2,
    memory_gb=4
)

service = Model.deploy(
    workspace=ws,
    name="ml-model",
    models=[model],
    deployment_config=deployment_config
)
```





# Российские облачные платформы

Yandex.Cloud	VK Cloud (бывший MCS)	Ростелеком/Ростелеком Cloud
<ul style="list-style-type: none"><li>➤ Container Registry — хранилище Docker образов</li><li>➤ Managed Kubernetes — управляемый K8s сервис</li><li>➤ Полная интеграция с Russian Federation, соответствие законам о локализации данных</li></ul>	<ul style="list-style-type: none"><li>➤ Kubernetes as a Service</li><li>➤ Хранение данных в России (Frankfurt, Moskva regions)</li><li>➤ Интеграция с VK экосистемой</li></ul>	<ul style="list-style-type: none"><li>➤ На базе OpenStack + Kubernetes</li><li>➤ Для госзаказов особенно популярно</li></ul>

```
bash
```

```
# Пример: создание K8S кластера в Yandex.Cloud
yc k8s cluster create \
  --name my-cluster \
  --region ru-central1 \
  --version 1.28
```

# сравнительная таблица

Решение	Тип	Российское	Функциональность	Сложность	Best for
Kubernetes	Open-source	Нет (Google)	★ ★ ★ ★ ★	Высокая	Production, большие системы
Yandex.Cloud K8s	Managed	✓	★ ★ ★ ★ ★	Средняя	Облачные системы в России
Docker Swarm	Open-source	Нет	★ ★ ★	Низкая	Небольшие проекты, быстрый старт
Apache Mesos	Open-source	Нет	★ ★ ★ ★	Высокая	Дата-центры, высоконагруженные
VK Cloud K8s	Managed	✓	★ ★ ★ ★ ★	Средняя	VK экосистема
Ростелеком Cloud	Managed	✓	★ ★ ★ ★	Средняя	Госзаказы, on-premise

# Мониторинг в production

Operational Metrics	Model Quality Metrics	Business Metrics	Инструменты:
<ul style="list-style-type: none"><li>➤ Latency (p50, p95, p99)</li><li>➤ Throughput (requests/sec)</li><li>➤ Error rate</li><li>➤ Resource utilization (CPU, RAM)</li></ul>	<ul style="list-style-type: none"><li>➤ Prediction distribution drift</li><li>➤ Feature drift</li><li>➤ Accuracy degradation</li><li>➤ Outlier detection</li></ul>	<ul style="list-style-type: none"><li>➤ Conversion rate</li><li>➤ Revenue impact</li><li>➤ User satisfaction</li></ul>	<ul style="list-style-type: none"><li>➤ Prometheus + Grafana</li><li>➤ Evidently AI</li><li>➤ WhyLabs</li><li>➤ Seldon Alibi</li></ul>

# Пример Prometheus метрик:

```
from prometheus_client import Counter, Histogram

predictions_total = Counter('predictions_total', 'Total predictions')
prediction_latency = Histogram('prediction_latency_seconds', 'Prediction latency')

@app.post("/predict")
def predict(data):
    with prediction_latency.time():
        result = model.predict(data)
        predictions_total.inc()
    return result
```



# Model Registry и версионирование

```
import mlflow

# Регистрация модели
mlflow.register_model(
    model_uri="runs:/abc123/model",
    name="flight-delay-predictor"
)

# Продвижение в Production
client = mlflow.tracking.MlflowClient()
client.transition_model_version_stage(
    name="flight-delay-predictor",
    version=3,
    stage="Production"
)

# Загрузка для deployment
model = mlflow.pyfunc.load_model(
    model_uri="models:/flight-delay-predictor/Production"
)
```

## Stages:

- ✓ None → только что зарегистрирована
- ✓ Staging → тестируется
- ✓ Production → активно используется
- ✓ Archived → устаревшая



# развертывание модели (demo)

## Подготовка модели

python

```
# train.py
import mlflow
import joblib
from sklearn.ensemble import RandomForestClassifier

with mlflow.start_run():
    model = RandomForestClassifier()
    model.fit(X_train, y_train)

    mlflow.sklearn.log_model(model, "model")
    mlflow.register_model("runs:/.../model", "demo-model")
```

# развертывание модели (demo)

Создание API:

python

```
# api.py
from fastapi import FastAPI
import mlflow.pyfunc

app = FastAPI()
model = mlflow.pyfunc.load_model("models:/demo-model/Production")

@app.post("/predict")
def predict(features: dict):
    prediction = model.predict([list(features.values())])
    return {"prediction": prediction[0]}
```

# развертывание модели (demo)

## Развертывание

```
bash
```

```
docker build -t demo-model .
```

```
docker run -p 8080:8080 demo-model
```

```
# Тест
```

```
curl -X POST http://localhost:8080/predict \  
  -H "Content-Type: application/json" \  
  -d '{"feature1": 0.5, "feature2": 0.3}'
```



# Чек-лист best practices

Pre-deployment:	Deployment:	Post-deployment:
<ul style="list-style-type: none"><li><input type="checkbox"/> Модель зарегистрирована в Model Registry</li><li><input type="checkbox"/> Все тесты пройдены (unit, integration, model quality)</li><li><input type="checkbox"/> Документация API обновлена</li><li><input type="checkbox"/> Определены SLA и resource limits</li><li><input type="checkbox"/> Настроен мониторинг и алерт</li></ul>	<ul style="list-style-type: none"><li><input type="checkbox"/> Используется версионирование образов</li><li><input type="checkbox"/> Применяется canary/blue-green стратегия</li><li><input type="checkbox"/> Настроен health check endpoint</li><li><input type="checkbox"/> Логируются все предсказания для аудита</li><li><input type="checkbox"/> Есть rollback plan</li></ul>	<ul style="list-style-type: none"><li><input type="checkbox"/> Мониторинг latency и throughput</li><li><input type="checkbox"/> Отслеживание data drift</li><li><input type="checkbox"/> Сравнение business metrics с baseline</li><li><input type="checkbox"/> Регулярный review логов</li><li><input type="checkbox"/> План обновления модели</li></ul>



# Ключевые вопросы для самопроверки:

1. Что такое model deployment и какие этапы он включает?
2. В чем разница между batch и online inference?
3. Зачем нужен Model Registry?
4. Что такое canary deployment?
5. Какие метрики важно мониторить в production?

# Материалы и ссылки

## Документация и tutorиалы:

1. [MLflow Model Registry](#)
2. [Kubernetes ML Patterns](#)
3. [AWS SageMaker Deployment Guide](#)
4. [Seldon Core — open-source ML deployment](#)

## Книги:

1. ["Building Machine Learning Powered Applications" by Emmanuel Ameisen](#)
2. ["Machine Learning Design Patterns" by Lakshmanan, Robinson, Munn](#)

## Курсы:

1. [DeepLearning.AI: "MLOps Specialization"](#)
2. [DataLens: анализ и визуализация данных: https://yandex.cloud/ru/training/datalens](https://yandex.cloud/ru/training/datalens)



# Материалы и ссылки

## Инструменты:

1. BentoML — framework для packaging и serving
2. TorchServe — deployment для PyTorch
3. TensorFlow Serving
4. KServe (ex-KFServing)
5. <https://datalens.ru/>

## Домашнее чтение:

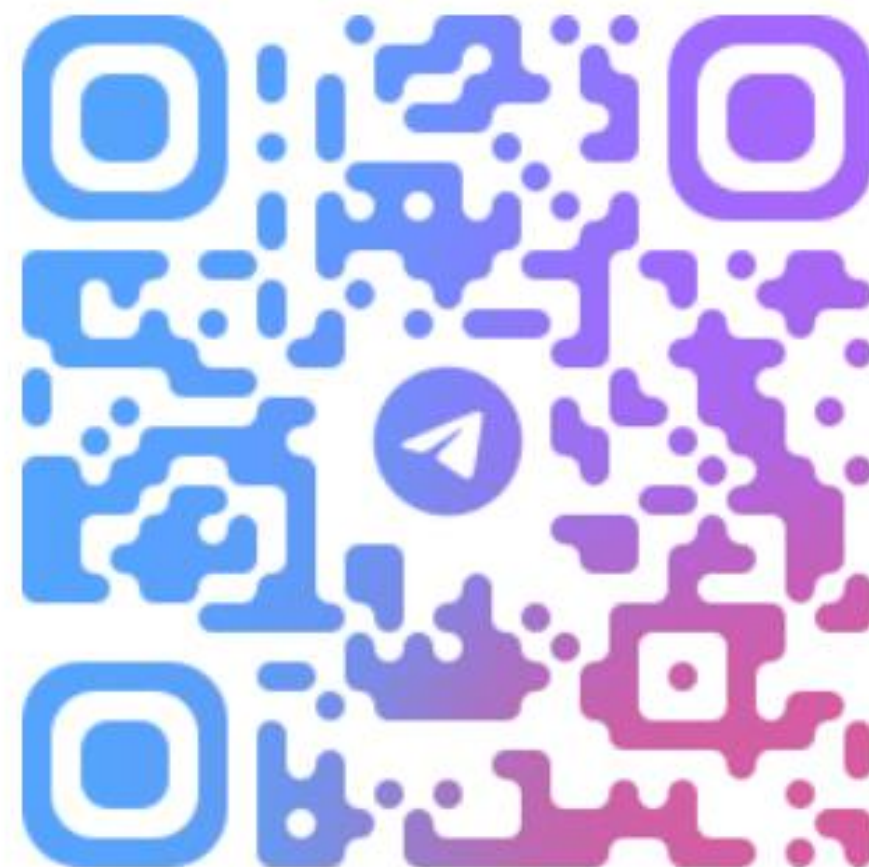
Изучить документацию одной из облачных платформ для ML deployment и написать краткое сравнение.



# Вопросы



Телеграм <https://t.me/+PsC-JDrwrvsxNmVi>



**СКИФ**

**(<https://do.skif.donstu.ru/course/view.php?id=7508>)**

