# Beyond Inference:
# Bringing ML into Production

DevConf.CZ 2021

Isabel Zimmerman

Data Science Intern

Forward Deployed AI Engineering

Red Hat

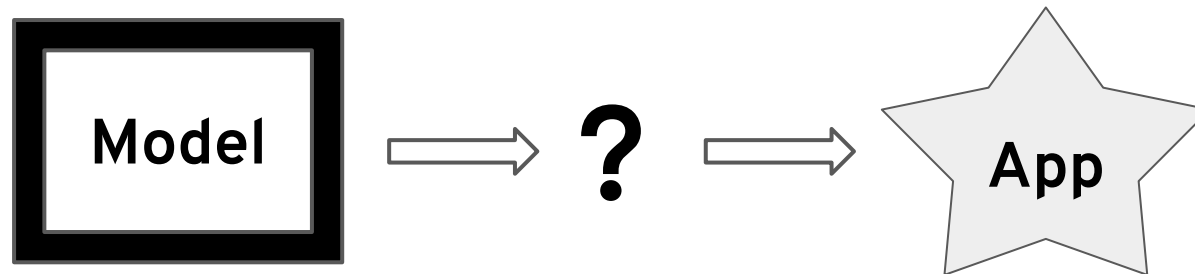# Things to learn

What is model serving?

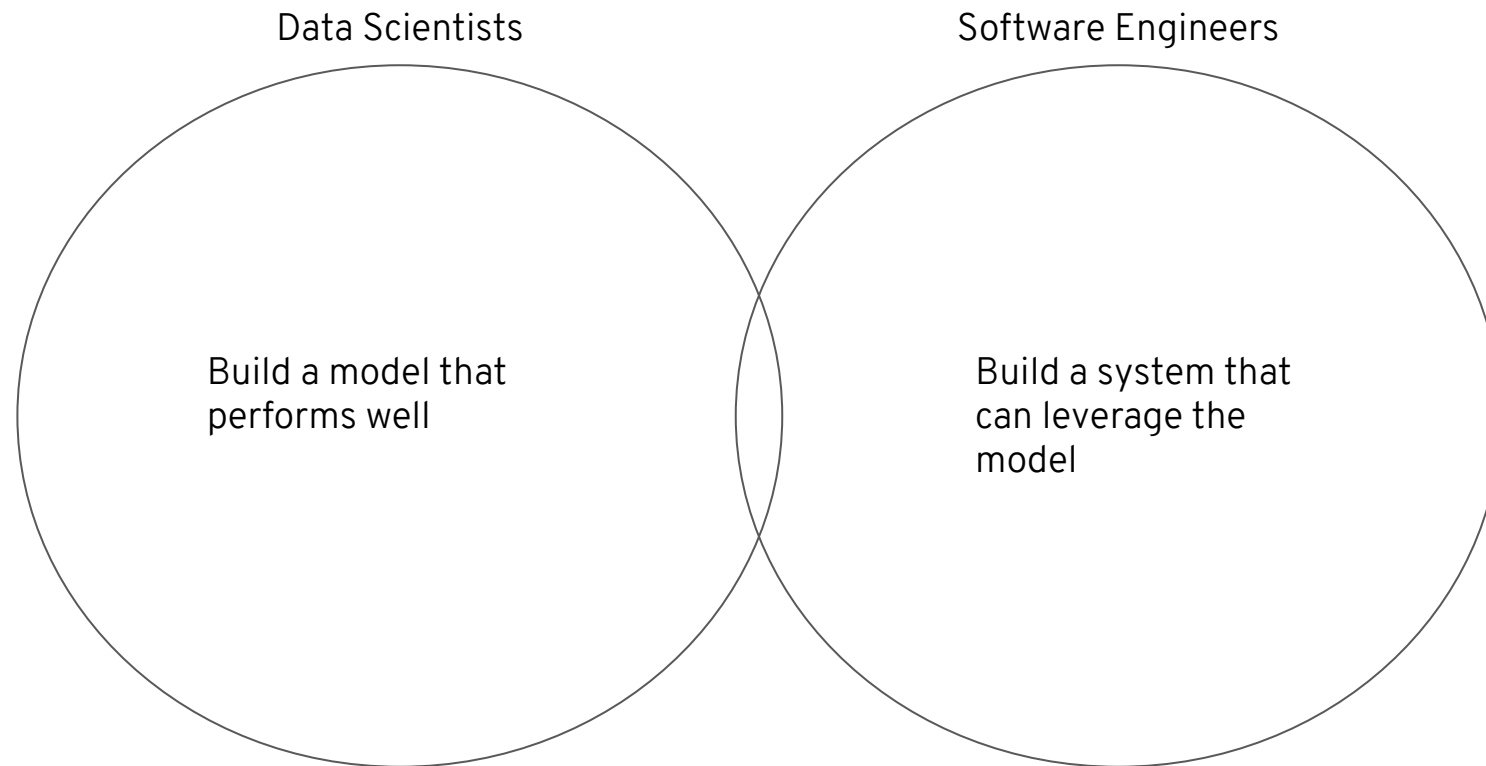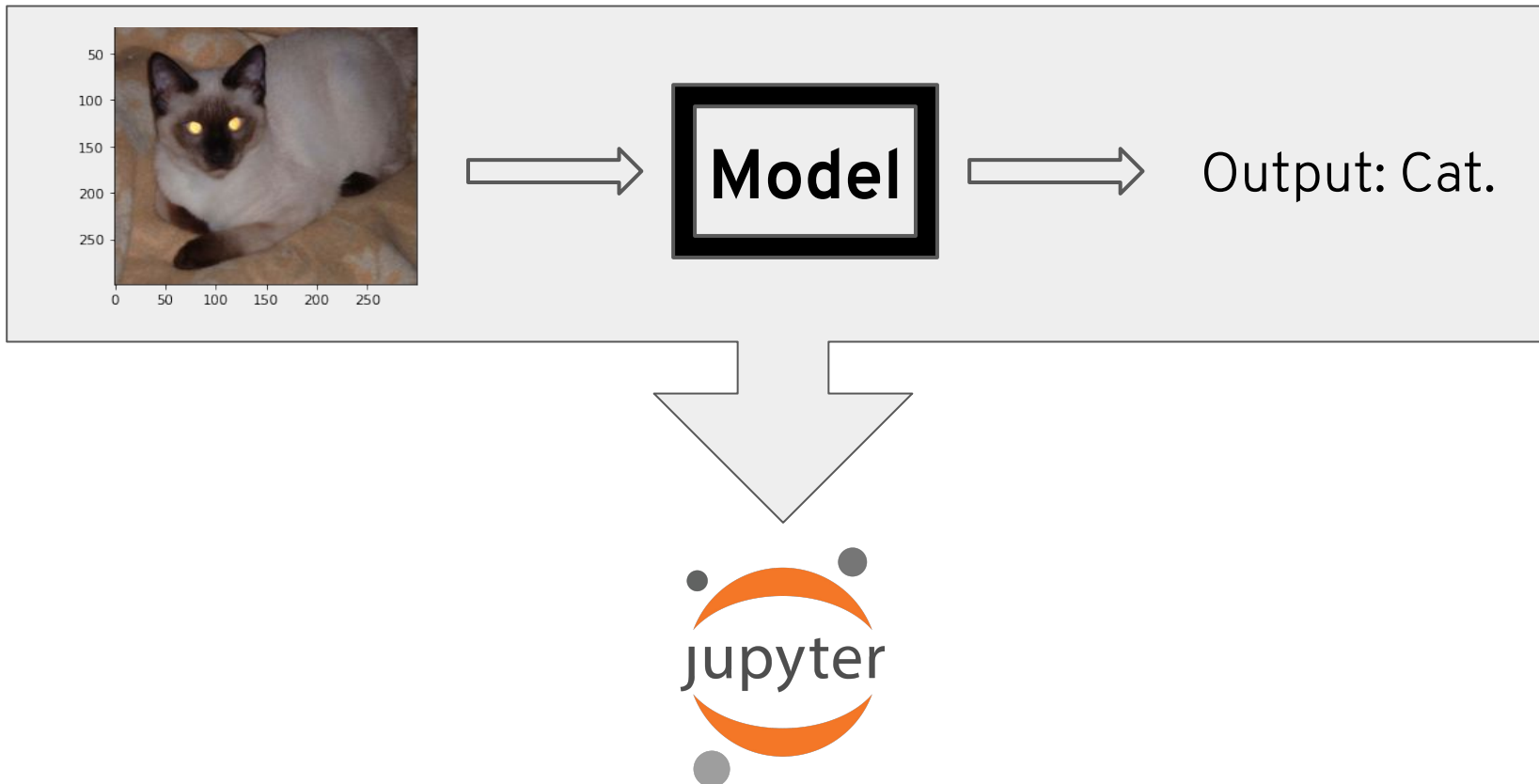What features of model serving can I utilize?

How can I deploy a model?

Red Hat

# Why model serving?

Red Hat

# Why do we need model serving?

# Do we need model serving?

Data Scientists
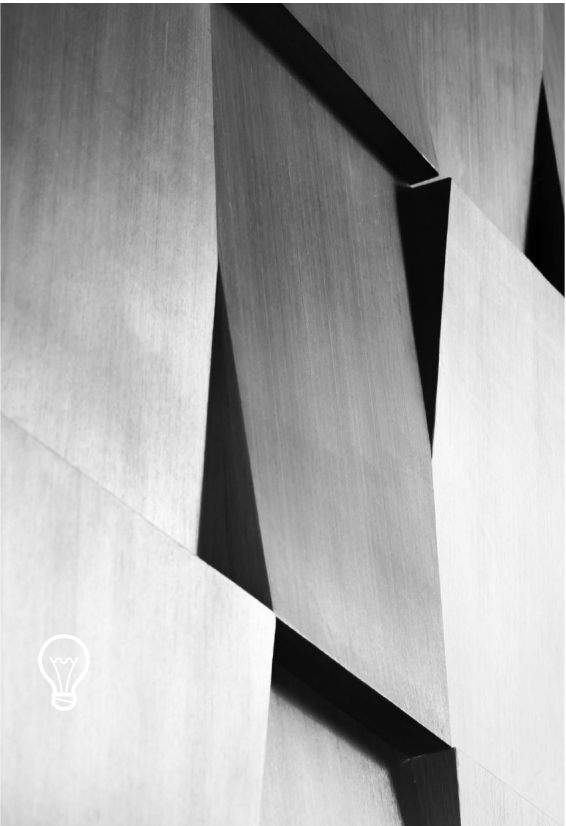
Software Engineers

Build a model that performs well

Build a system that can leverage the model

Red Hat

# How do models act?

# What is model serving?

Data Scientists

Software Engineers

**Model** → **?** ← **App**

Red Hat

# What is model serving?
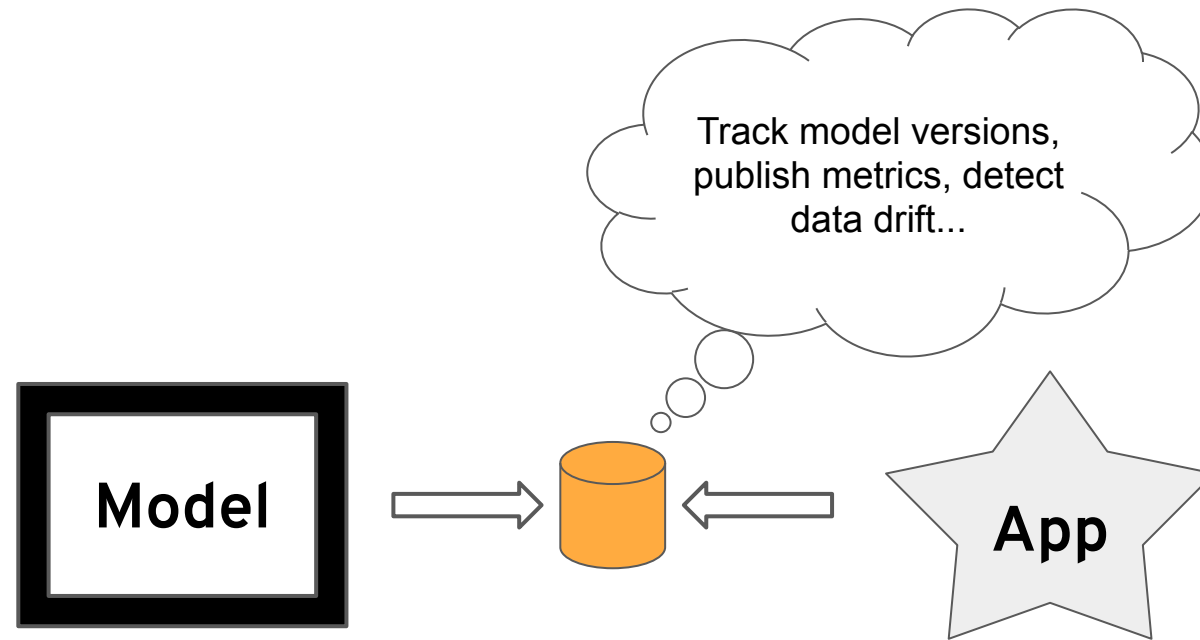
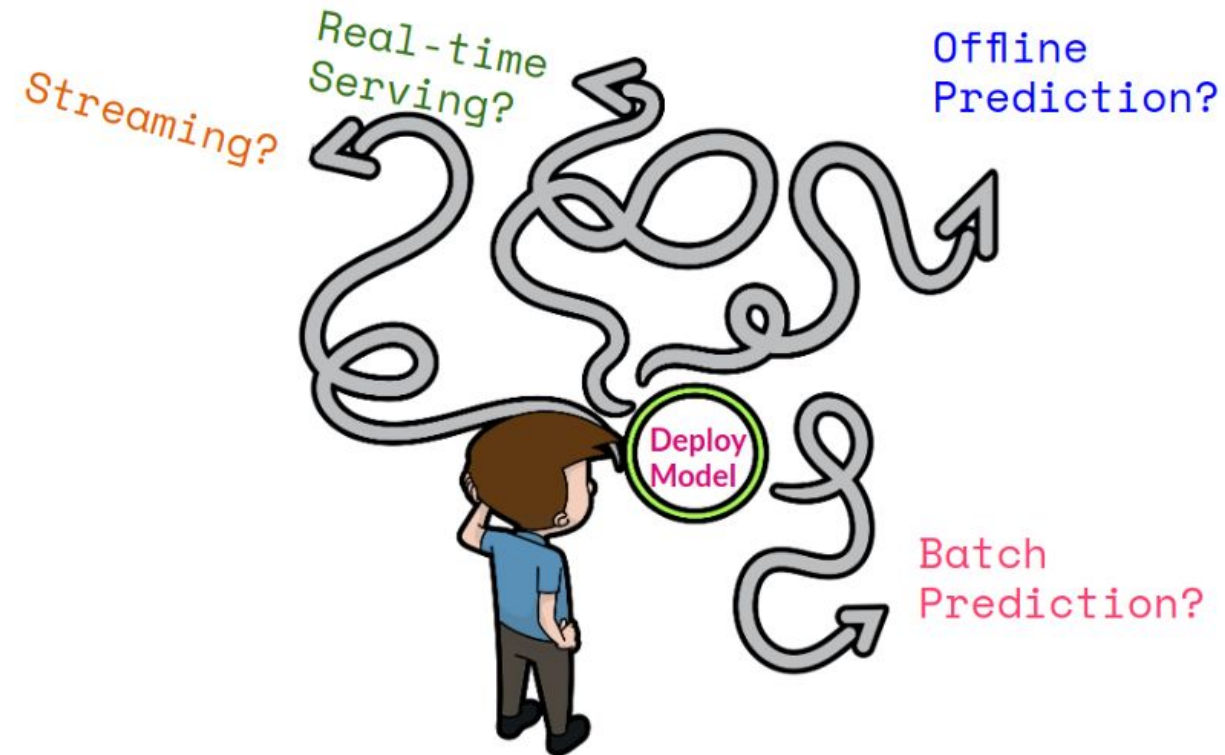Data Scientists

Software Engineers

Model serving

**Model**

**App**

What features should we look for?

# How can we make model servers more useful?

Track model versions, publish metrics, detect data drift...

**Model**

**App**

Red Hat

# How are you deploying?

# A/B Testing

90%

10%

# A/B Testing

80%

20%

# Ensemble



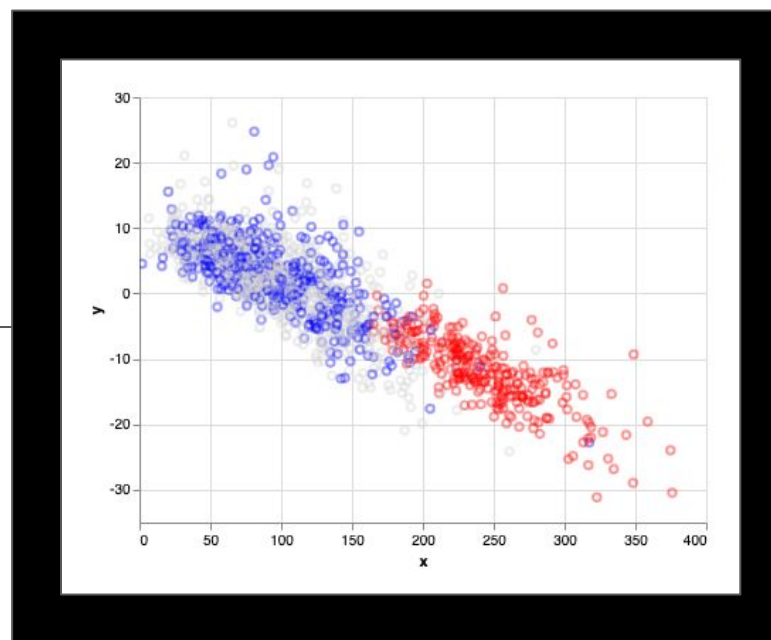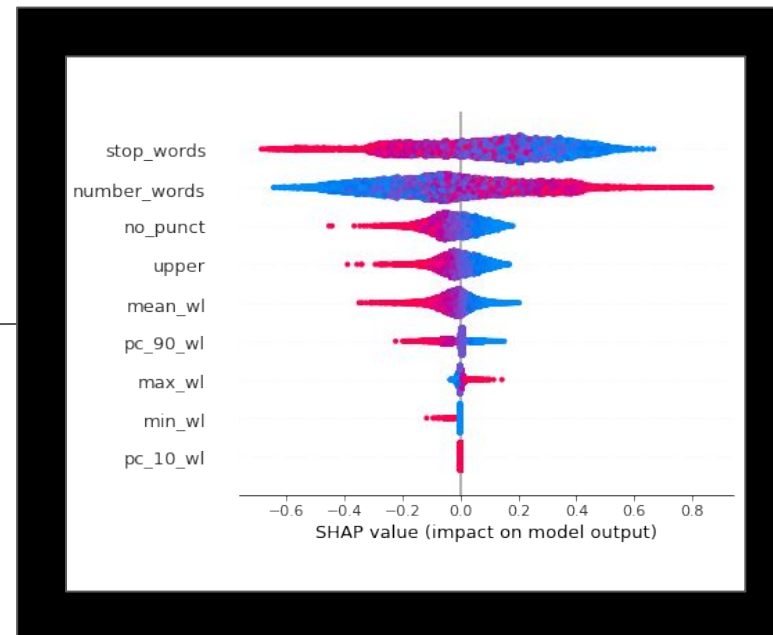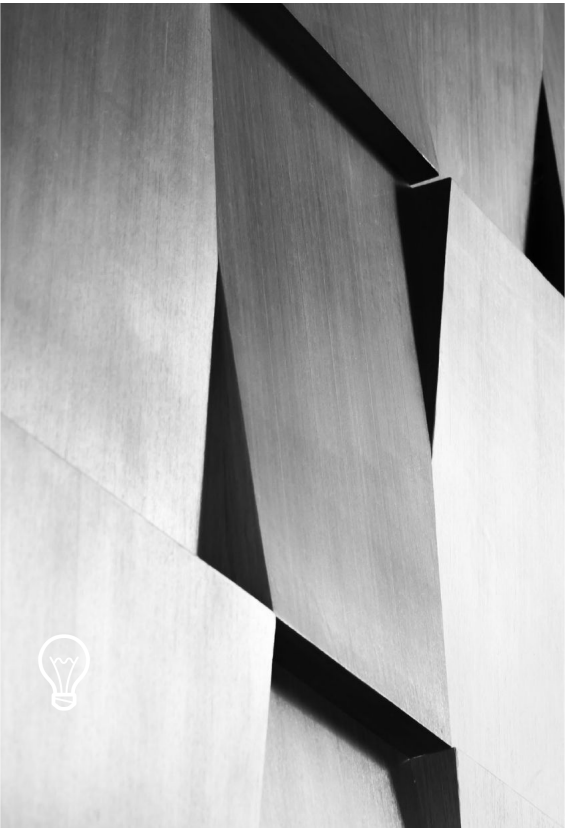Output: Cat.

# Data Drift

# Explainability

How do I deploy a model?

```yaml
apiVersion: machinelearning.seldon.io/v1
kind: SeldonDeployment
metadata:
  name: image-default
  namespace: seldon
spec:
  name: income
  predictors:
  - graph:
      children: []
      implementation: SKLEARN_SERVER
      modelUri: gs://seldon-models/sklearn/image/model-0.23.2
      name: classifier
    explainer:
      type: AnchorTabular
      modelUri: gs://seldon-models/sklearn/image/explainer-py36-0.5.2
    name: default
     replicas: 1
```
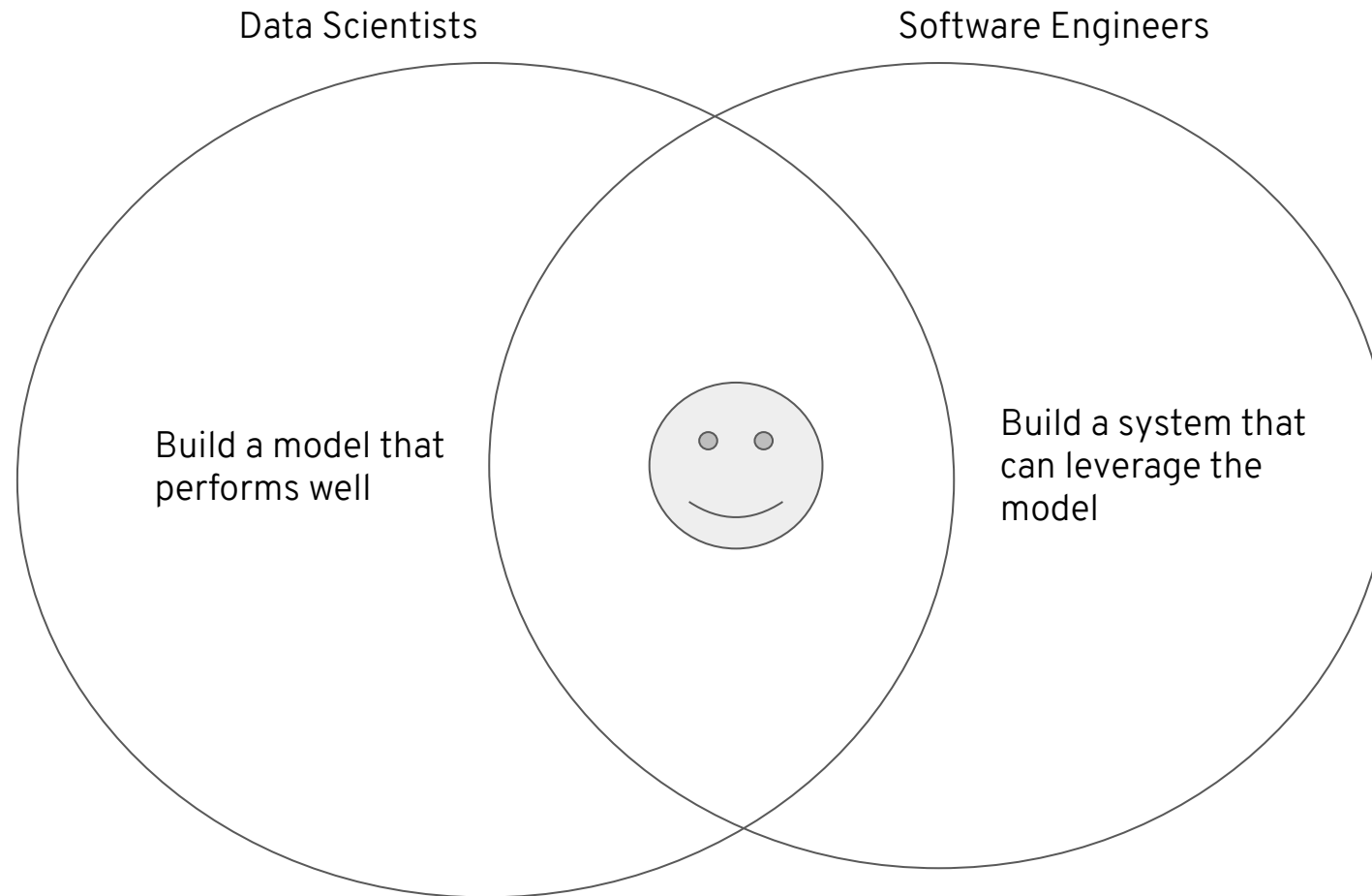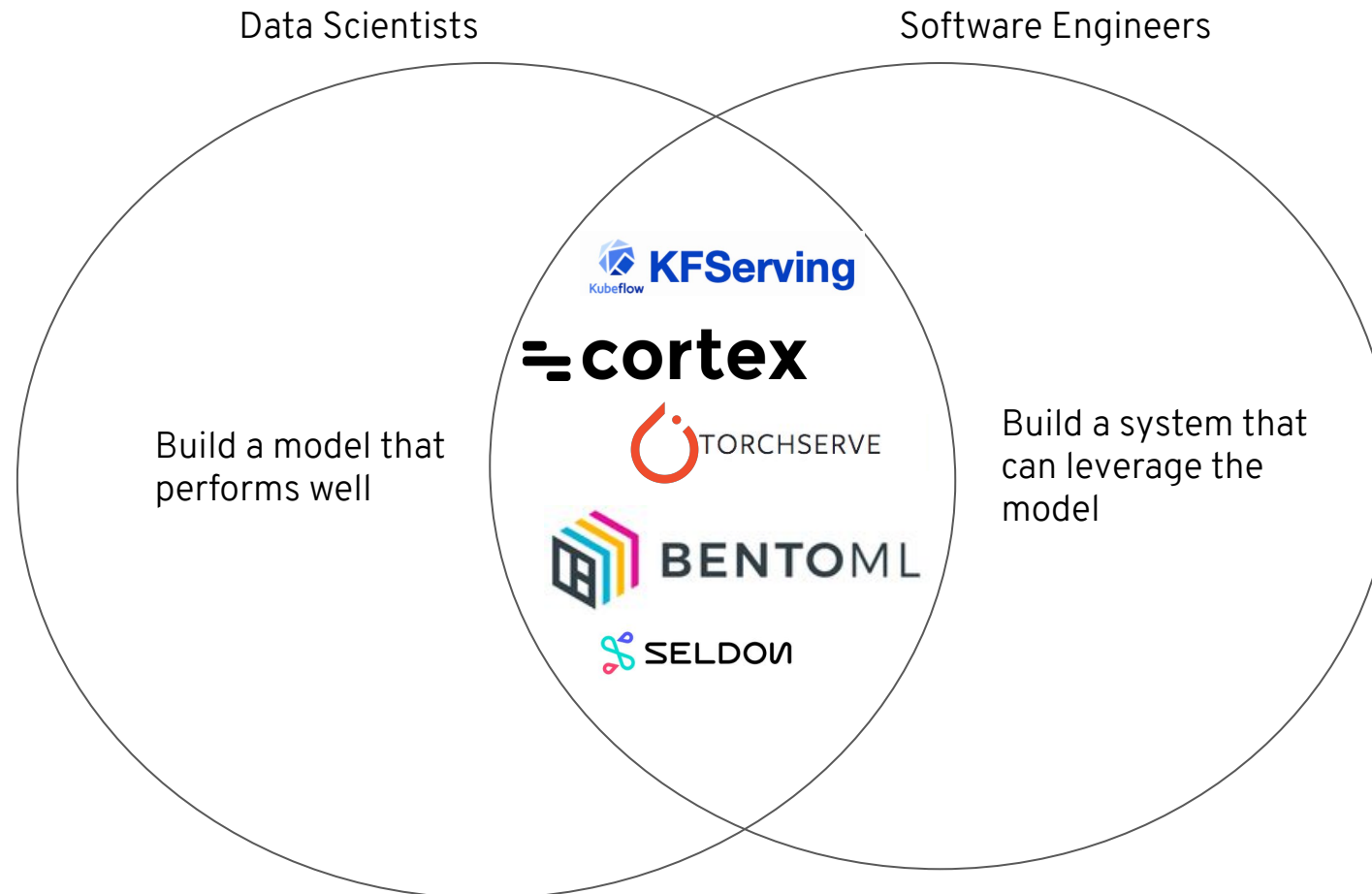
SELDON → jupyterhub → Prometheus → Grafana

Red Hat

# Demo

# Do we need model serving?

Data Scientists

Software Engineers

Build a model that performs well

Build a system that can leverage the model

# Do we need model serving?

Data Scientists

Software Engineers

**KFServing**
Kubeflow

**cortex**

TORCHSERVE

**BENTOML**

SELDON

Build a model that performs well

Build a system that can leverage the model

https://github.com/EthicalML/awesome-production-machine-learning#model-serving-and-monitoring

# Stay Connected.

**in** linkedin.com/isabel-zimmerman

**🐦** twitter.com/isabelizimm