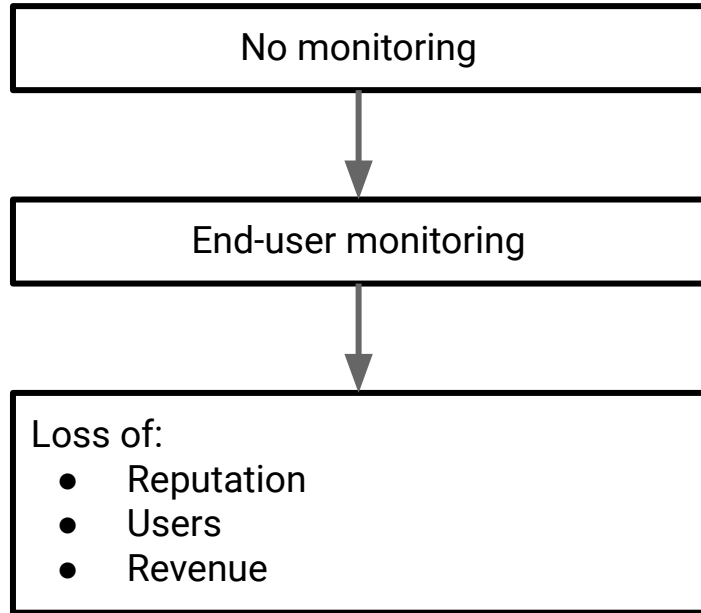# Data Drift Monitoring for ML Projects

Alex Kim @alex000kim

# Introduction

- Why monitor?
- Application monitoring vs ML monitoring
- Common causes of model and data drift
- What to monitor?

# Why monitor?

No monitoring

⬇

End-user monitoring

⬇

Loss of:
- Reputation
- Users
- Revenue

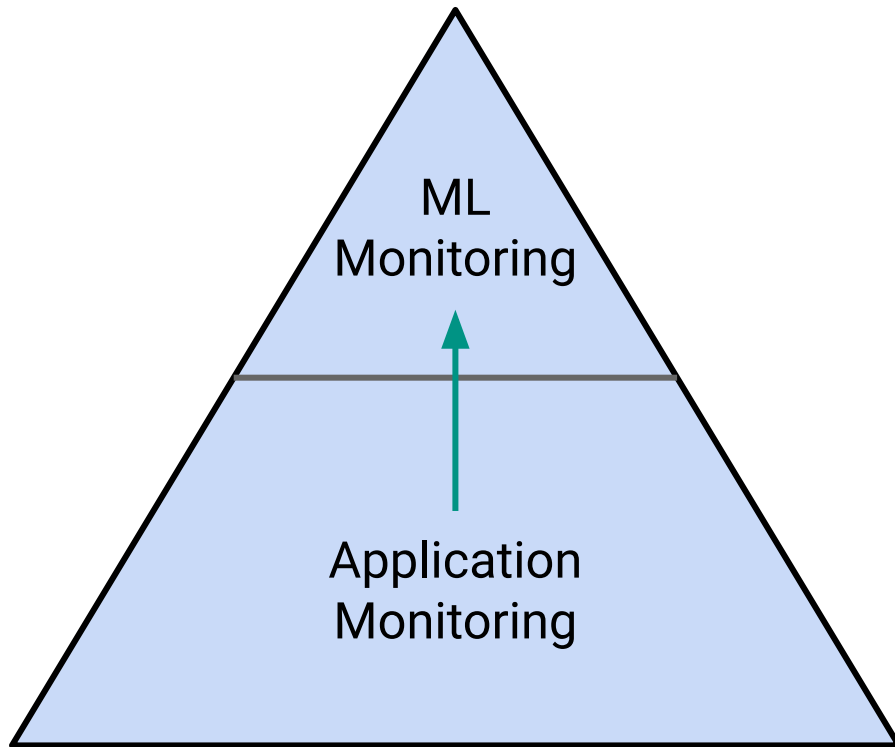# Application monitoring vs ML monitoring

**Application monitoring:**

- Latency
- Response error rate
- CPU
- RAM
- Disk space

**ML monitoring:**
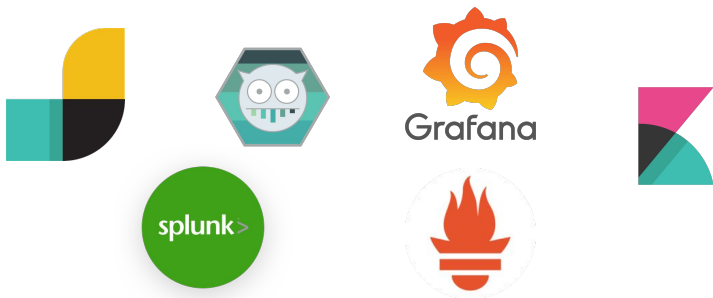
- Model performance metrics
- Data drift
- Concept drift

# Application Monitoring Tools

- **Instrumentation & metrics**: statsd, prometheus, etc.
- **Event logging & tracing**: logstash, splunk, etc.
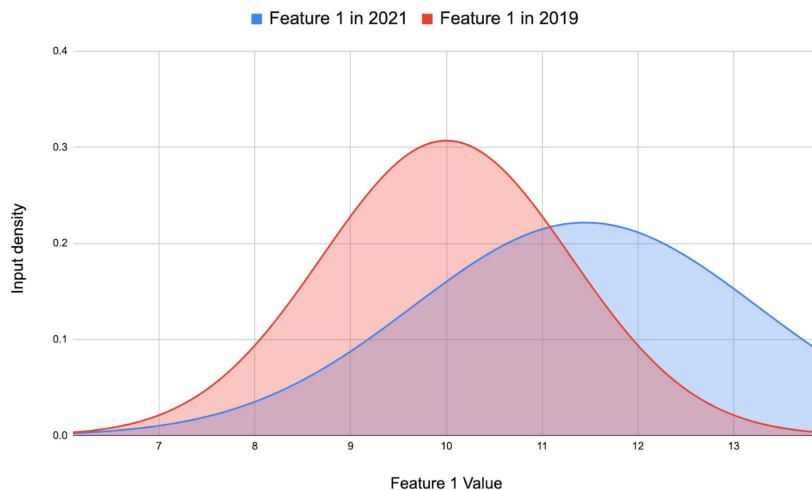- **Dashboards**: grafana, kibana, graphite, etc.

# ML Monitoring Tools

- Alibi Detect
- Arize AI
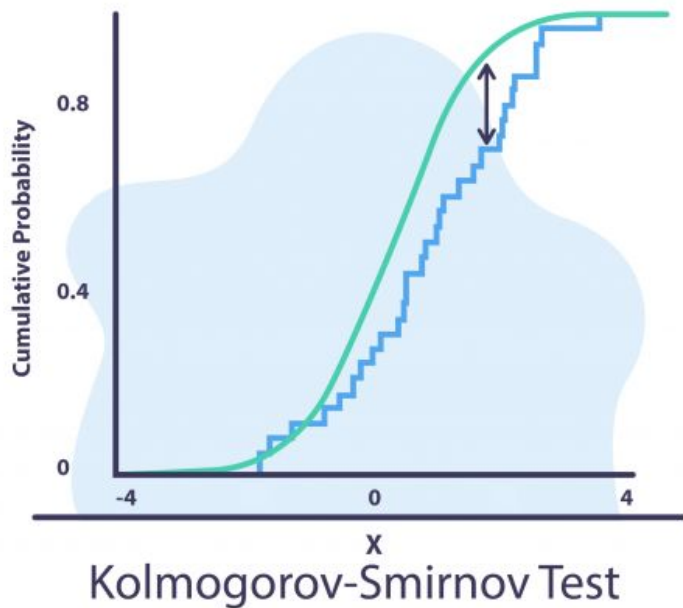- Evidently AI
- WhyLabs
- Fiddler

# Causes of data drift

- changes data source (e.g. ℃ vs ℉, broken sensor)
- data preprocessing pipeline (e.g. variable scaling, data imputation)
- market conditions lead to changes in user behaviour (e.g. change in disposable income, consumer preferences)
- regulations change user behaviour (e.g. GDPR)
- upstream system or company policy (e.g. change in UI, opt-out vs opt-in data collection)

# What to monitor?

- Model performance on new data
- Input data summary stats (% missing, min/max, etc)
- Data distribution
- Statistical distances between training data and new data
  (e.g. Chi-Squared, Kolmogorov-Smirnov)



Kolmogorov-Smirnov Test

# Practice time!