# Instructor: Adam Breindel

LinkedIn: **https://www.linkedin.com/in/adbreind**
Email: adbreind@gmail.com

- 20+ years building systems for startups and large enterprises
- 10+ years teaching front- and back-end technology

- Fun large-scale data projects…
    - Streaming neural net + decision tree fraud scoring
    - Realtime & offline analytics for banking
    - Music synchronization and licensing for networked jukeboxes

- Industries
    - Finance, Insurance
    - Travel, Media / Entertainment
    - Energy, Government
    - Various Others…

# High-Level Plan for Today

- Morning (Part 1)
  - Intro to Spark
  - Spark ML
  - Feature engineering, modeling, evaluating, tuning
- Afternoon (Part 2)
  - Integrating Spark with Python (without sacrificing performance)
  - Productionizing Spark Models (without sacrificing performance)
  - Deep Learning and Future Directions for Integrations (DL, GPU Analytics...)

# Today's Class – Informal Survey

- Today is my first day using Spark

- I've run a few operations in Spark shell

- I've used Spark for 1-2 months in my job

- My job is 50%+ Spark, or I've been using Spark for 6+ months

# Setup with Databricks

Create a Databricks account

- Sign up for **free Community Edition** now at http://tinyurl.com/databricks-ce

- Use **Firefox, Chrome or Safari** (Internet Explorer / MS Edge not fully supported)
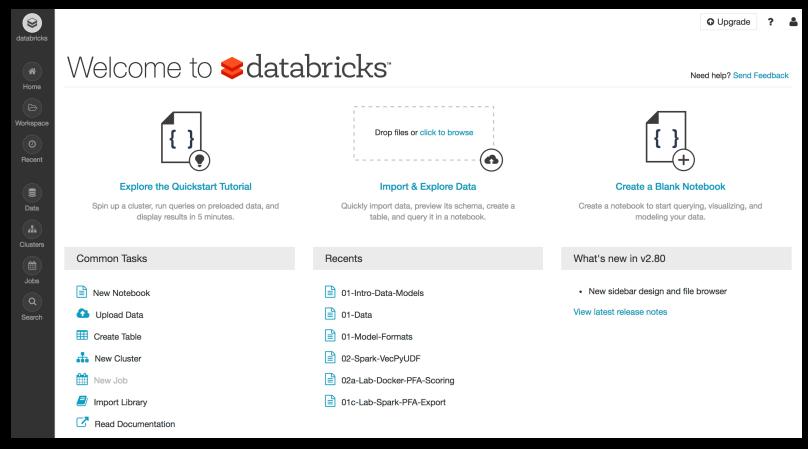
Getting Started

These steps are **illustrated** on subsequent pages; this is the summary:

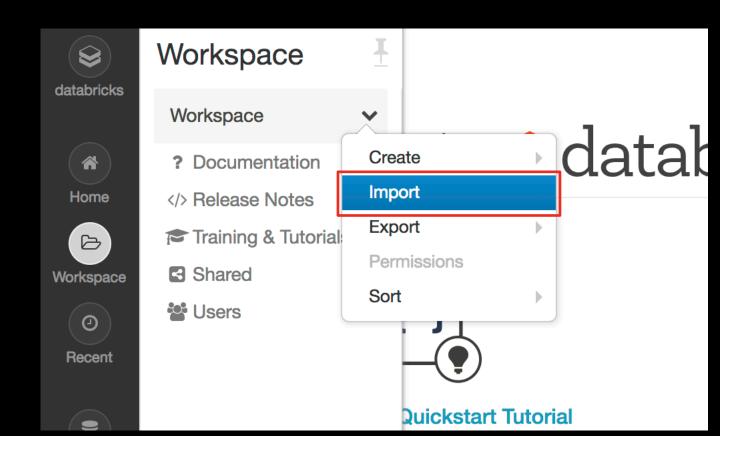1. Copy the courseware link or prepare to type it ☺

   https://materials.s3.amazonaws.com/2019/odsc/east.dbc

2. Import that file into your Databricks account per the instructions on the following slides.

3. Create a cluster: choose `Databricks Runtime 5.2` (also illustrated in the following slides)

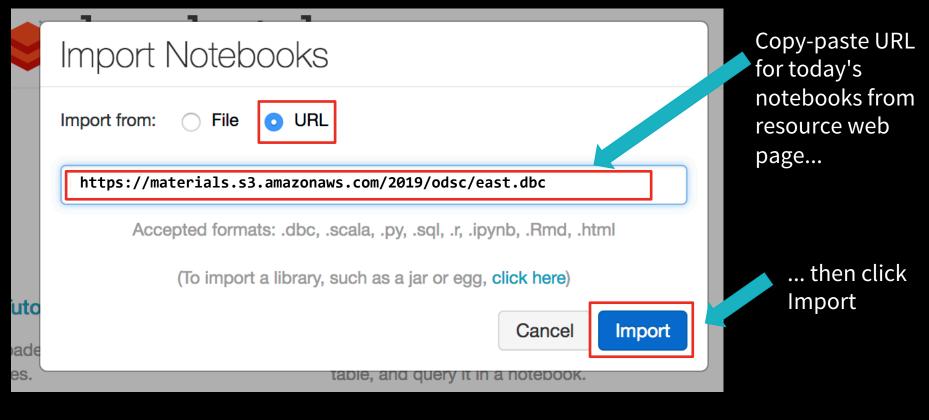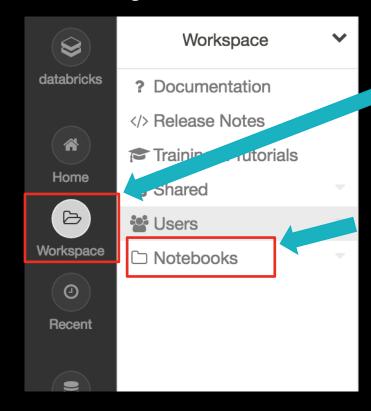# Log in to Databricks

# Import Notebooks…

# Import Notebooks for Today…



Copy-paste URL for today's notebooks from resource web page…

… then click Import

# Find your notebook(s) here...

Click Workspace

Here are your notebook(s).

*The file/folder will have a different name!*

# Create a Cluster

# Create a Cluster

# All set: let's go!