# First-time Qwiklab Account Setup

- Go to https://h2oai.qwiklab.com/
- Click on "JOIN"
- Create a new account with a valid email address
- You will receive a confirmation email
  - Click on the link in the confirmation email
- Go back to http://h2oai.qwiklab.com and log in
- Go to the Catalog on the left bar
- Choose "H2O and Sparking Water Workshop"
- Wait for instructions

# H2O-3 and Sparkling Water Workshop

Megan Kurka / Tom Kraljevic
H2O.ai

# Agenda

- Introduction to H2O-3

- Hands On
  - Data Cleaning and Supervised Learning

- Introduction to Sparkling Water

- Hands On
  - Feature Engineering for Time Series and AutoML

- Questions and Answers

# H2O Overview

# H2O Products

In-Memory, Distributed Machine Learning Algorithms with H2O Flow GUI

H2O AI Open Source Engine Integration with Spark

Lightning Fast machine learning on GPUs

Automatic feature engineering, machine learning and interpretability
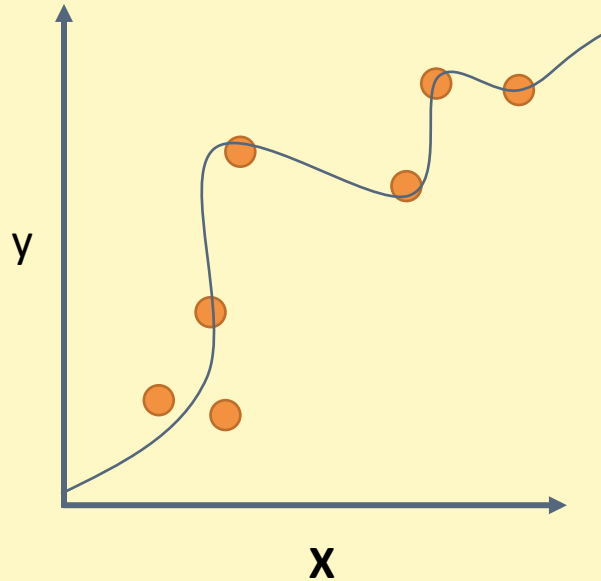
Secure multi-tenant H2O clusters

# What is Machine Learning?

# Supervised Learning

**Regression:**
**How much will a customer spend?**



**H$_2$O algos:**
**Penalized Linear Models**
**Random Forest**
**Gradient Boosting**
**XGBoost**
**Neural Networks**
**Stacked Ensembles**

**Classification:**
**Will a customer churn?**



**H$_2$O algos:**
**Penalized Linear Models**
**Naïve Bayes**
**Random Forest**
**Gradient Boosting**
**XGBoost**
**Neural Networks**
**Stacked Ensembles**

H$_2$O.ai

# Unsupervised Learning

**Clustering:**
Grouping rows – e.g. creating groups of similar customers

$x_j$

$x_i$

**$H_2O$ algos:**
k – means

**Feature extraction:**
Grouping columns – Create a small number of new representative dimensions

$x_j$

$PC_1 = -0.3\ x_i - 0.4\ x_i$

$x_i$

**$H_2O$ algos:**
Principal components
Generalized low rank models
Autoencoders
Word2Vec

**Anomaly detection:**
Detecting outlying rows - Finding high-value, fraudulent, or weird customers

Fraudster

$x_j$

Billionaire

$x_i$

**$H_2O$ algos:**
Principal components
Generalized low rank models
Autoencoders

H$_2$O.ai

# Simplified Typical Machine Learning Pipeline

H2O-3

# What is H2O?

**Math Platform** — Open source in-memory AI engine

- Parallelized and distributed algorithms
- GLM, Random Forest, GBM, Deep Learning, etc.

**Tech and API** — Easy to use and adopt

- Written in Java – perfect for Java Programmers
- Install is lightweight
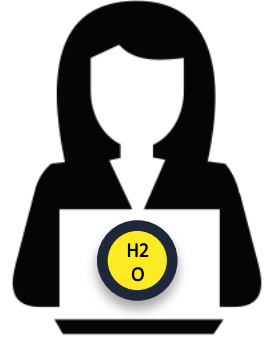- REST API (Java) – run H2O from R, Python, WebUI

**Big Data** — More data? Or better models? BOTH

- Use all of your data – model without sampling
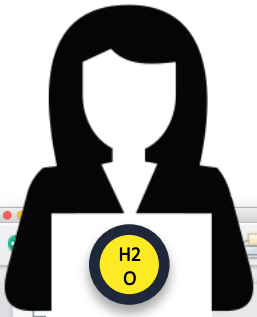- More Data + Better Models = Better Predictions

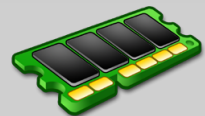# H2O Cluster

H2O Clients

H2O Cluster

H2O

H2O Cluster

REST / JSON
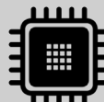
```r
library(h2o)

h2o.init(ip="localhost", port=54321)

h2o.ls()

df_allyears2k <- h2o.getFrame("allyears2k.hex")
deeplearning_model <- h2o.getModel("deeplearning_model")

summary(df_allyears2k)
```

1:1    (Top Level)                                    R Script

Console ~/Library/Mobile Documents/com~apple~CloudDocs/0_h2o_docs/demo
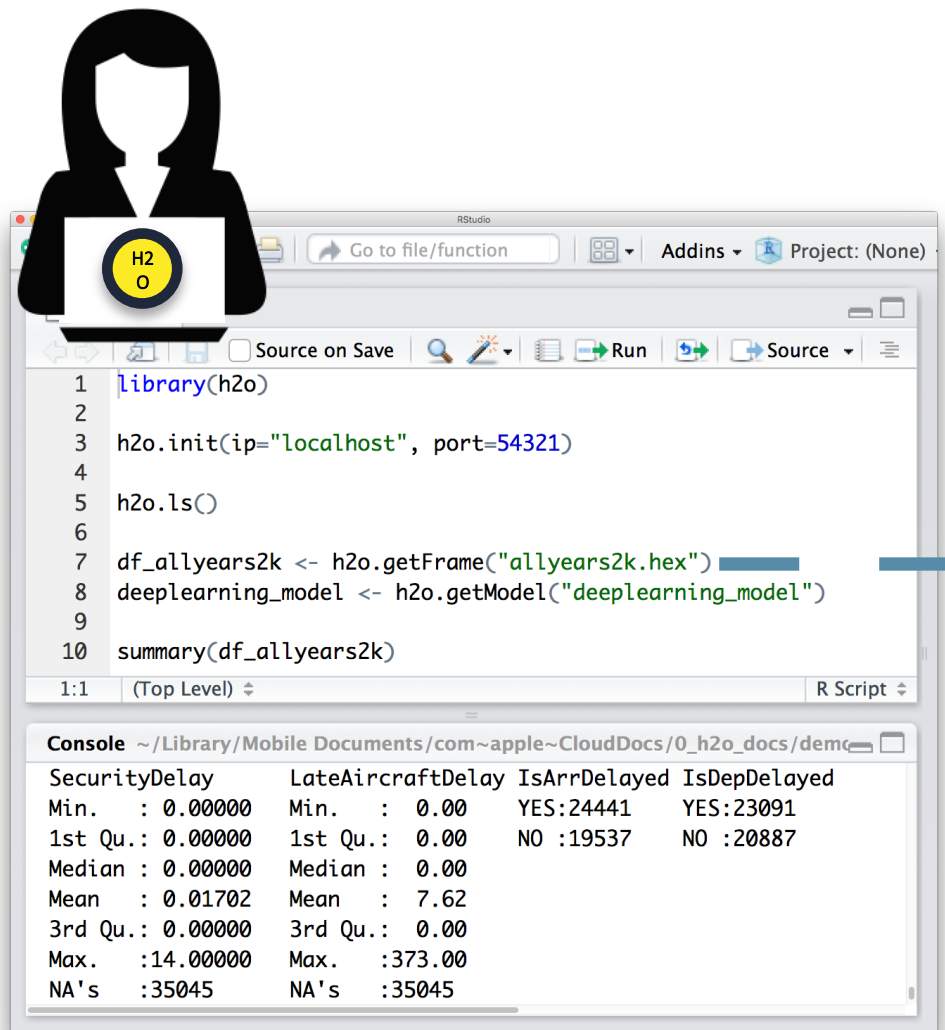
```
SecurityDelay      LateAircraftDelay IsArrDelayed IsDepDelayed
Min.   : 0.00000   Min.   :  0.00    YES:24441    YES:23091
1st Qu.: 0.00000   1st Qu.:  0.00    NO :19537    NO :20887
Median : 0.00000   Median :  0.00
Mean   : 0.01702   Mean   :  7.62
3rd Qu.: 0.00000   3rd Qu.:  0.00
Max.   :14.00000   Max.   :373.00
NA's   :35045      NA's   :35045
```
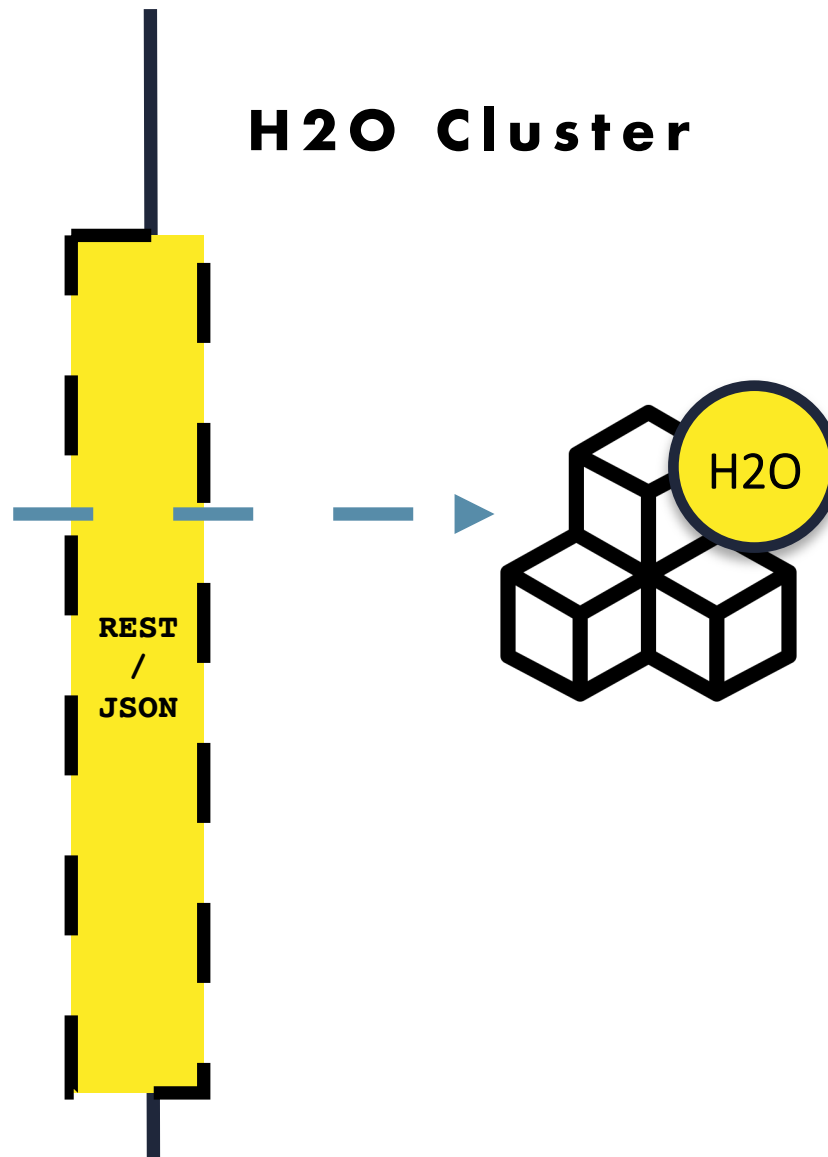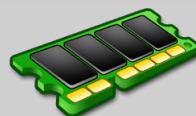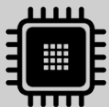
Local Machine

# H2O Cluster

```r
library(h2o)

h2o.init(ip="localhost", port=54321)

h2o.ls()

df_allyears2k <- h2o.getFrame("allyears2k.hex")
deeplearning_model <- h2o.getModel("deeplearning_model")

summary(df_allyears2k)
```

```
SecurityDelay       LateAircraftDelay   IsArrDelayed    IsDepDelayed
Min.   : 0.00000    Min.   :  0.00      YES:24441       YES:23091
1st Qu.: 0.00000    1st Qu.:  0.00      NO :19537       NO :20887
Median : 0.00000    Median :  0.00
Mean   : 0.01702    Mean   :  7.62
3rd Qu.: 0.00000    3rd Qu.:  0.00
Max.   :14.00000    Max.   :373.00
NA's   :35045       NA's   :35045
```

REST / JSON

Local Machine
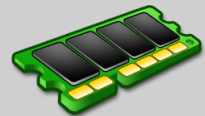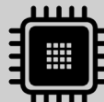
# H2O Cluster

```
1   library(h2o)
2
3   h2o.init(ip="localhost", port=54321)
4
5   h2o.ls()
6
7   df_allyears2k <- h2o.getFrame("allyears2k.hex")
8   deeplearning_model <- h2o.getModel("deeplearning_model")
9
10  summary(df_allyears2k)
```

```
1:1   (Top Level)                                    R Script

Console  ~/Library/Mobile Documents/com~apple~CloudDocs/0_h2o_docs/demo

SecurityDelay      LateAircraftDelay IsArrDelayed IsDepDelayed
Min.   : 0.00000   Min.   :  0.00    YES:24441    YES:23091
1st Qu.: 0.00000   1st Qu.:  0.00    NO :19537    NO :20887
Median : 0.00000   Median :  0.00
Mean   : 0.01702   Mean   :  7.62
3rd Qu.: 0.00000   3rd Qu.:  0.00
Max.   :14.00000   Max.   :373.00
NA's   :35045      NA's   :35045
```
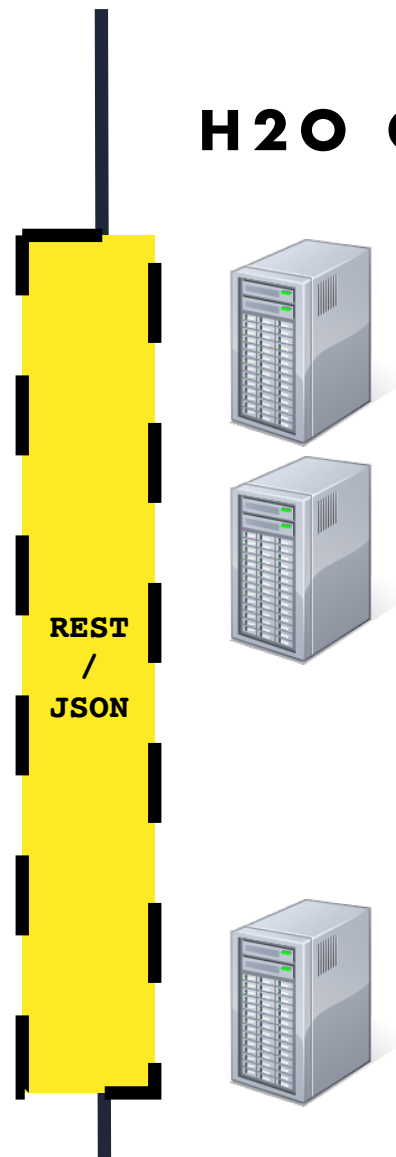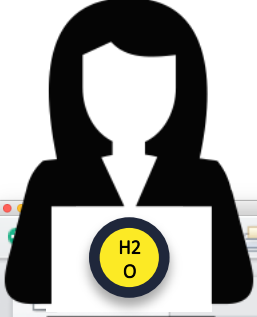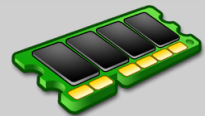
REST / JSON

JVM 1

JVM 2

JVM N

Local Machine

# H2O Cluster

```
1   library(h2o)
2
3   h2o.init(ip="localhost", port=54321)
4
5   h2o.ls()
6
7   df_allyears2k <- h2o.getFrame("allyears2k.hex")
8   deeplearning_model <- h2o.getModel("deeplearning_model")
9
10  summary(df_allyears2k)
```

```
SecurityDelay       LateAircraftDelay  IsArrDelayed   IsDepDelayed
Min.   : 0.00000    Min.   :  0.00     YES:24441      YES:23091
1st Qu.: 0.00000    1st Qu.:  0.00     NO :19537      NO :20887
Median : 0.00000    Median :  0.00
Mean   : 0.01702    Mean   :  7.62
3rd Qu.: 0.00000    3rd Qu.:  0.00
Max.   :14.00000    Max.   :373.00
NA's   :35045       NA's   :35045
```
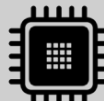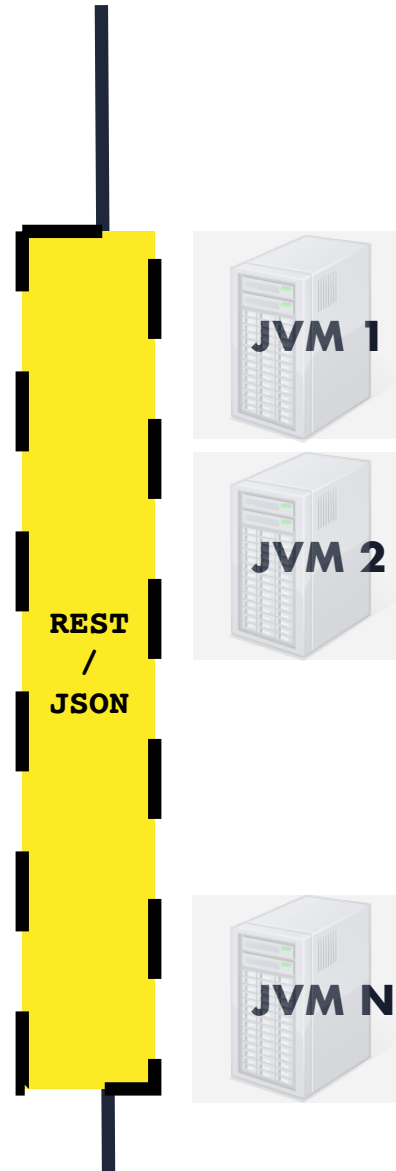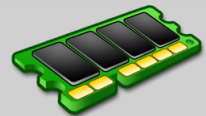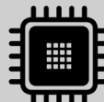
Local Machine

REST / JSON

JVM 1

JVM 2

JVM N

| $X_1$ | $X_2$ | $X_3$ | ... | $X_p$ | $y$ |
|-------|-------|-------|-----|-------|-----|
|       |       |       |     |       |     |

cluster returns frame key

# H2O Cluster

```r
library(h2o)

h2o.init(ip="localhost", port=54321)

h2o.ls()

df_allyears2k <- h2o.getFrame("allyears2k.hex")
deeplearning_model <- h2o.getModel("deeplearning_model")

summary(df_allyears2k)
```
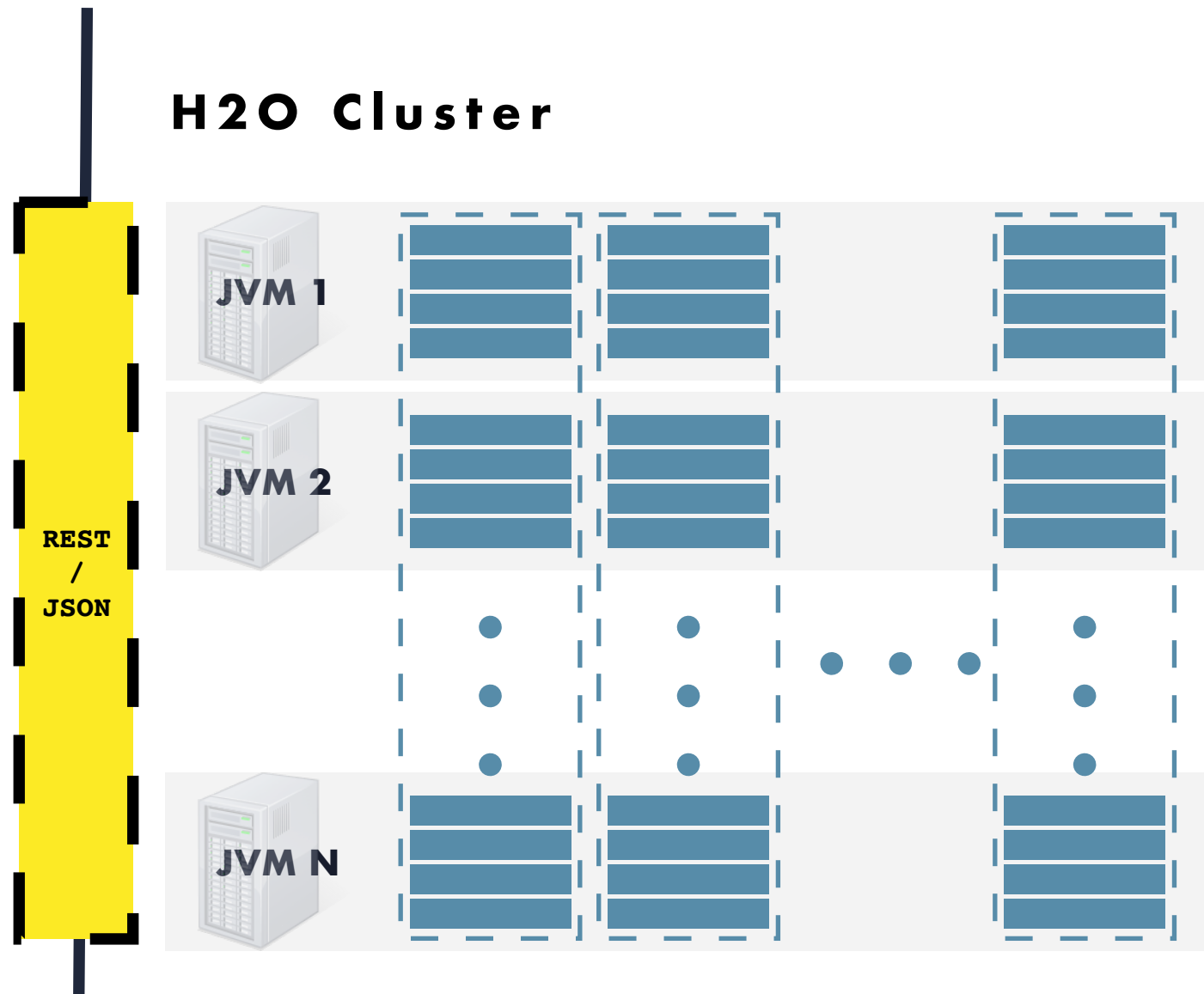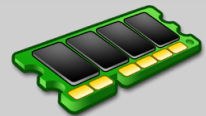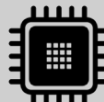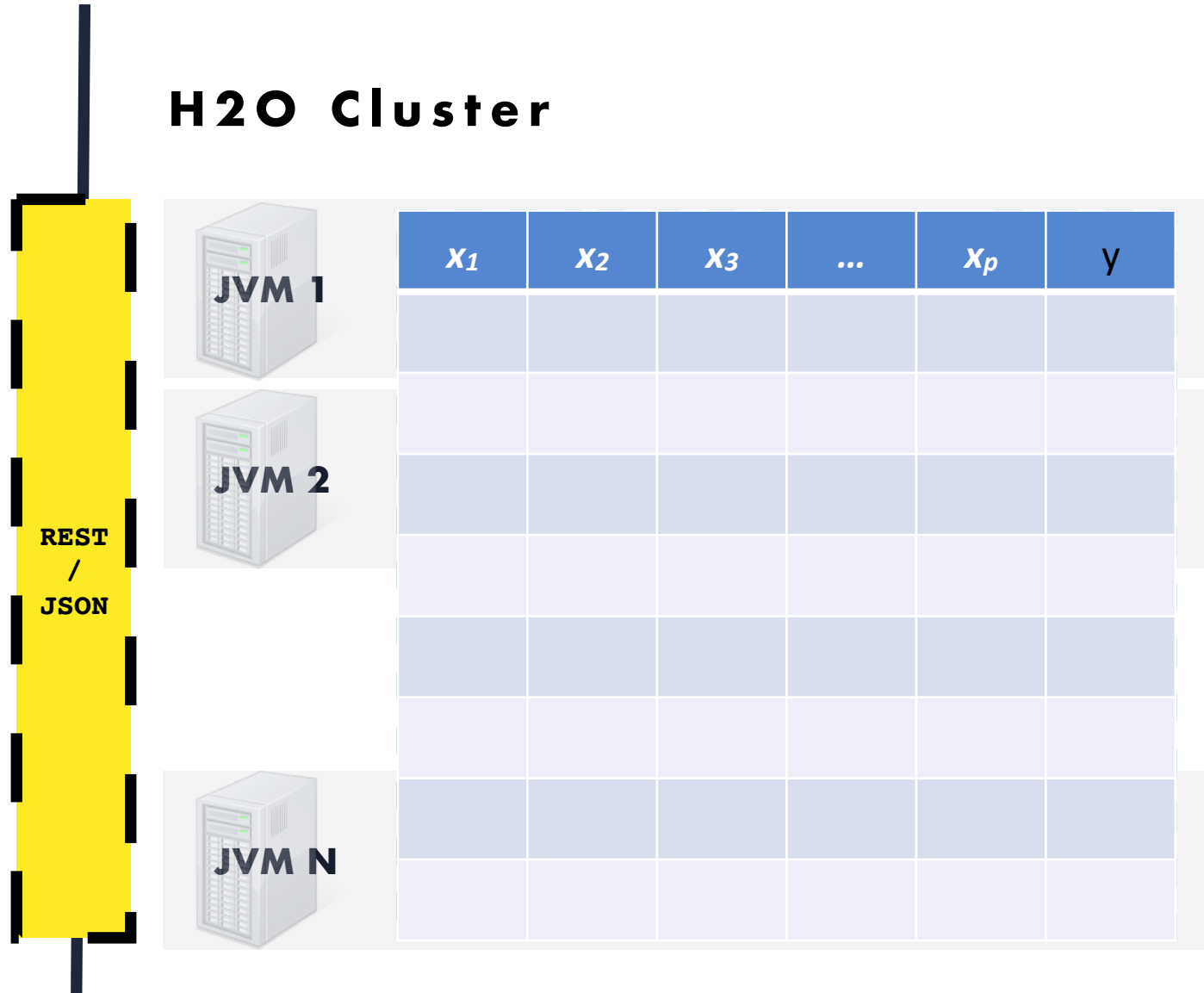
```
SecurityDelay      LateAircraftDelay IsArrDelayed  IsDepDelayed
Min.   : 0.00000   Min.   :  0.00    YES:24441    YES:23091
1st Qu.: 0.00000   1st Qu.:  0.00    NO :19537    NO :20887
Median : 0.00000   Median :  0.00
Mean   : 0.01702   Mean   :  7.62
3rd Qu.: 0.00000   3rd Qu.:  0.00
Max.   :14.00000   Max.   :373.00
NA's   :35045      NA's   :35045
```

JVM 1

JVM 2

REST / JSON

JVM N

Local Machine

| $X_1$ | $X_2$ | $X_3$ | ... | $X_p$ | y |
|-------|-------|-------|-----|-------|---|
|       |       |       |     |       |   |
|       |       |       |     |       |   |
|       |       |       |     |       |   |
|       |       |       |     |       |   |
|       |       |       |     |       |   |
|       |       |       |     |       |   |
|       |       |       |     |       |   |

# H2O-3 Hands On

# The Data

- Lending Club peer-to-peer loan default dataset
- Covers loans issued in years 2007-2011
- 42,538 loans

| Type of information | Column names |
| --- | --- |
| Demographic Information | annual_inc, home_ownership, emp_length |
| Loan Information | purpose, term, desc, int_rate |
| **Response Column** | **bad_loan** |