# Hit a Home Run

*Discover and operationalize data-analytics for winning*

*Artificial Intelligence and Machine Learning*

**IBM** + **aginity** + **H₂O**.ai

# The Moneyball All-Stars

**David Kearns**
Data Science

@IBMAnalytics

**Ari Kaplan**
Mr. Moneyball

@Aginity

**Chris Coad**
Story Teller

@Aginity

**Erin LeDell**
Chief ML Scientist

@H20.ai

IBM + aginity + H₂O.ai

# What is Moneyball?

*Moneyball is the practice of using data, analytics, artificial intelligence and machine learning to find the best baseball players and build the optimal team to win championships and increase organizational revenue*

**IBM** + ○ aginity + **H₂O**.ai

# Organizational AI Challenges

Data-related challenges are hindering 96% of organizations from taking full advantage of AI

IBM + aginity + H₂O.ai

# Organizational AI Challenges

# 80% of organizations face collaboration challenges due to silos

# Organizational AI Challenges

Data systems don't "do AI" and AI technologies don't "do data." Organizations use 7 disparate tools

IBM + aginity + H₂O.ai

# How Moneyball Relates to Business?

- Diverse datasets living in different systems and formats
- Collaboration across consistent data and analytics to power predictive insights and capabilities
- AI & ML is built on engineered features and "analytics-ready" data
- Technology drives more data  and competition drives need for innovation and fast and clear ways to consume these insights

# Moneyball Project Goals

- Create a predictive player performance model

- Work with a distributed set of data sources in different data repositories

- Build consistent data and derived attributes (analytics) across that can be shared across the world

- Share domain expertise

- Operationalize ML/AI

IBM + aginity + H₂O.ai

# Combine Public and Proprietary Databases

## Lahman Database

- Public database from 1871 to 2017

- Aggregate pitching and batting statistics

| Attribute | Description |
|-----------|-------------|
| playerID | Player ID code |
| AB | At Bats |
| R | Runs |
| H | Hits |
| SO | Strike Outs |
| + More | … |

http://www.seanlahman.com/baseball-archive/statistics/

## Ari Database

- Private database from 2012-2017

- Pitch-by-pitch play for each MLB game

| Attribute | Description |
|-----------|-------------|
| Pitch_Type | Two - character code of type of pitch |
| Spin_rate | Spin of the pitch in rotations per minute. |
| Start_speed | The velocity of the pitch in MPH |
| End_speed | The velocity of the pitch when it arrives at the plate in MPH |
| Spray_des | Classification of type of hit |
| + More | … |

IBM + aginity + H$_2$O.ai

# The Manual Way...



| Phase | Challenge |
|---|---|
| Model | Limited traceability and full picture of how it was built |
| Model Building | Manual model tests & specific expertise |
| Modeling Table | |
| Data Quality & Transformation | Manual data prep, re-coding & no model reuse |
| Data Integration | Disparate data sources & small data sets |

IBM + aginity + H₂O.ai

# To the Enterprise Way: IBM + Aginity + H2O

**Challenge**

**Solve**

Limited traceability and full picture of how it was built

Who, what, when and how the analytic was created

Manual model tests & specific expertise

AutoML

Manual data prep, re-coding & no model reuse

Re-usable analytics operationalized and universally shareable

Disparate data sources & small data sets

Rich data sets in from multiple enterprise data solutions

**Application Layer**

H2O.ai    IBM Immersive Insights

**Analytics Layer**

Aginity Amp

**Data Layer**

IBM    DB2    Hortonworks
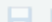
IBM + aginity + H2O.ai

# *Demos*

File    Edit    Code    View    Plots    Session    Build    Debug    Profile    Tools    Help          Disk usage:    24%

Project: (None)

ui.R    H2oml.R

Source on Save          Run    Source

```r
82
83      # H2O AutoML with Lahman only
84      automl_lahman = h2o.automl(x = features,
85                                 y = targets[n_target],
86                                 training_frame = h_train,
87                                 validation_frame = h_valid,
88                                 max_models = 10, # increase this to allow more models
89                                 max_runtime_secs = 120, # increase this to allow more time
90                                 stopping_metric = "RMSE",
91                                 stopping_rounds = 3,
92                                 seed = n_seed,
93                                 exclude_algos = c("DeepLearning"), # you can exclude any algo
94                                 project_name = paste0("AutoML_Lahman", targets[n_target]))
95
96
97      # Extract model
98      model_best_lahman = automl_lahman@leader
99      # print(automl_lahman@leaderboard)
100
101
102     # Make predictions for all records in one go
103     tmp_yhat_lahman = as.data.frame(h2o.predict(model_best_lahman, h_all))
104
105
106     # Store Results
107     colnames(tmp_yhat_lahman) = paste0("pred_lahman_", targets[n_target])
108     d_all_with_pred = cbind(d_all_with_pred, tmp_yhat_lahman)
109
110
111     }
```

116:22    (Untitled)          R Script

Environment    History    Connectio

Import Dataset

Global Environment

Files    Plots    Packages    Help

Console    Terminal

~/

Let's talk

File    Edit    Code    View    Plots    Session    Build    Debug    Profile    Tools    Help    Disk usage: 24%

Go to file/function    Addins    Project: (None)

**ui.R**    **H2oml.R**

Reload App    R Script

```
1   source("func_lib.R")
2
3   jscode <- "shinyjs.set_caption_pic = function(ele
4   elem = document.getElementById(element_name);
5   if (elem) {
6   if (elem.classList.contains('downArrow')) {
7   elem.classList.add('upArrow');
8   elem.classList.remove('downArrow');
9   } else {
10  elem.classList.add('downArrow');
11  elem.classList.remove('upArrow');
12  }}
13  }"
14
15  ui <- fluidPage (
16    tags$head(includeCSS("style.css")),
17
18    shinyjs::useShinyjs(),
19    shinyjs::extendShinyjs(text = jscode, functions
20
21    tags$img(src="Aginity_logo_small.jpg"),
22
```

11:34    (Top Level)    R Script

**Environment**    **History**    **Connections**

Import Dataset    List

Global Environment

**Files**    **Plots**    **Packages**    **Help**    **Viewer**
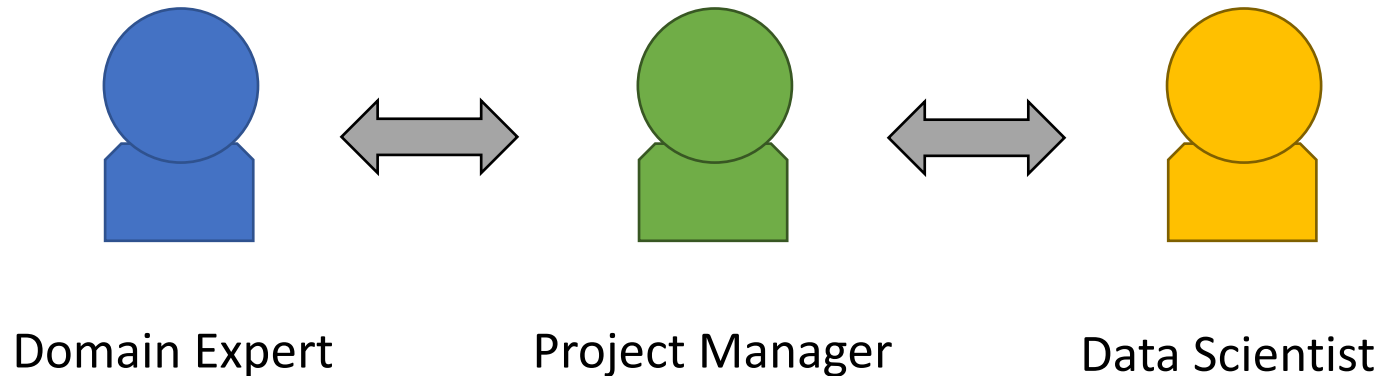
Publish

aginity

**Console**    **Terminal**

~/

```
########"
[1] "Remove: Response_Types"
[1] "refTables after remove: "
NULL
[1] "colNames after remove:"
NULL
[1] "asset_id ####"
[1] "b4f51a23-de4f-48a2-b0db-db48e00c1955"
```

Click feature names to add to modeling set:

∨ Base asset: AriDB2012

**Choose Workspace**

Main Workspace

**Choose Project**

Money Ball

**Choose Analytic Frame**

AriDB2012

**Choose Cluster**

amp-demo

**Choose Target Platform**

AWS Redshift

**Enter Target (e.g. 'DB.Tablename')**

| | Variable/Feature name: | Description |
|---|---|---|
| 1 | mlbid | Player Identifier |
| 2 | gameid | Game Identifier |
| 3 | gameday | Date of game played |
| 4 | batter_team | 3 letter abbreviation for team |
| 5 | pitcher_team | 3 letter abbreviation for team |
| 6 | h_a | Player home or away team |
| 7 | inning_num | Inning of the pitch |
| 8 | atbat_num | nth batter to appear in game |
| 9 | ball | Pitch count of balls |
| 10 | strike | Pitch count of strikes |
| 11 | outs | Number of outs at start of at bat |
| 12 | batter | mlbid of the batter |
| 13 | pitcher | mlb id of the pitcher |

Let's talk

# *Key Findings*

# 1. Diverse Team with Specific Skills

- Important have a team with a wide variety of skills so the full context of what you're trying to accomplish is understood

- Provide the context first

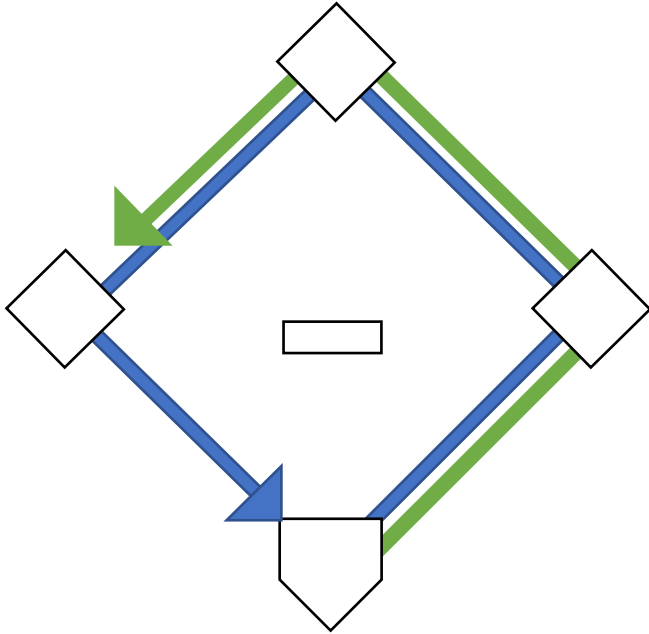Domain Expert          Project Manager          Data Scientist

# 2. Encapsulate Domain Expert Knowledge

- Who do you constantly rely for domain expertise?

- How can you encapsulate their knowledge...

- ...Make it easily accessible to the rest of the team, department or enterprise?

IBM + aginity + H₂O.ai

# 3. Similar Analytics Across Use Cases



**80%** Of enterprise analytics **are highly similar.**

**100%** Of enterprise analytics **are recreated.**

# 4. Enterprise-width Features Enables AI/ML

*AutoML*

**Active** Analytics & Data Catalog

**3rd Party analytics**

**Center of Excellence**

Single-Family Analytics

Financial Engineering

Risk

Other

*Domain Analysts "Contribute" Their Subject Matter Expertise*

IBM + aginity + H₂O.ai

# Final Project Result

## $20M

trade 2 weeks prior to the season beginning



IBM + aginity + H₂O.ai

# Try it Yourself!

GitHub: **https://github.com/woobe/moneyball**

# Get in Touch with Us

**David Kearns**
Data Science

@DaithiOCiaran

**Ari Kaplan**
Mr. Moneyball

@arikaplan

**Chris Coad**
Story Teller

@chriscoad

**Erin LeDell**
Chief ML Scientist

@ledell

IBM + aginity + H$_2$O.ai

# Sports Analytics Podcasts

- Selfish plug (in case you dozed off)
  - http://www.ibmbigdatahub.com/podcast/making-data-simple-hit-home-run-using-ai-machine-learning

- New York Yankees on injury prevention
  - https://itunes.apple.com/gb/podcast/4-optimising-player-training-treatment-strategies-in/id1327803354?i=1000411908873&mt=2
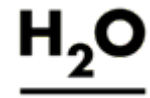
- ML Sports-focused Series
  - https://twimlai.com/aiinsports2018/

# Find Out More About Our Products



- ICP for Data
  - https://www.ibm.com/analytics/cloud-private-for-data
- Watson Studio
  - https://www.ibm.com/cloud/watson-studio



- Aginity Amp
  - https://www.aginity.com/main/products/
- Aginity Workbench
  - https://www.aginity.com/main/workbench/



- AutoML
  - https://www.h2o.ai/products/h2o/
- Driverless AI
  - https://www.h2o.ai/products/h2o-driverless-ai/