# Making the Most of Data: Feature Engineering for Applied Supervised Machine Learning
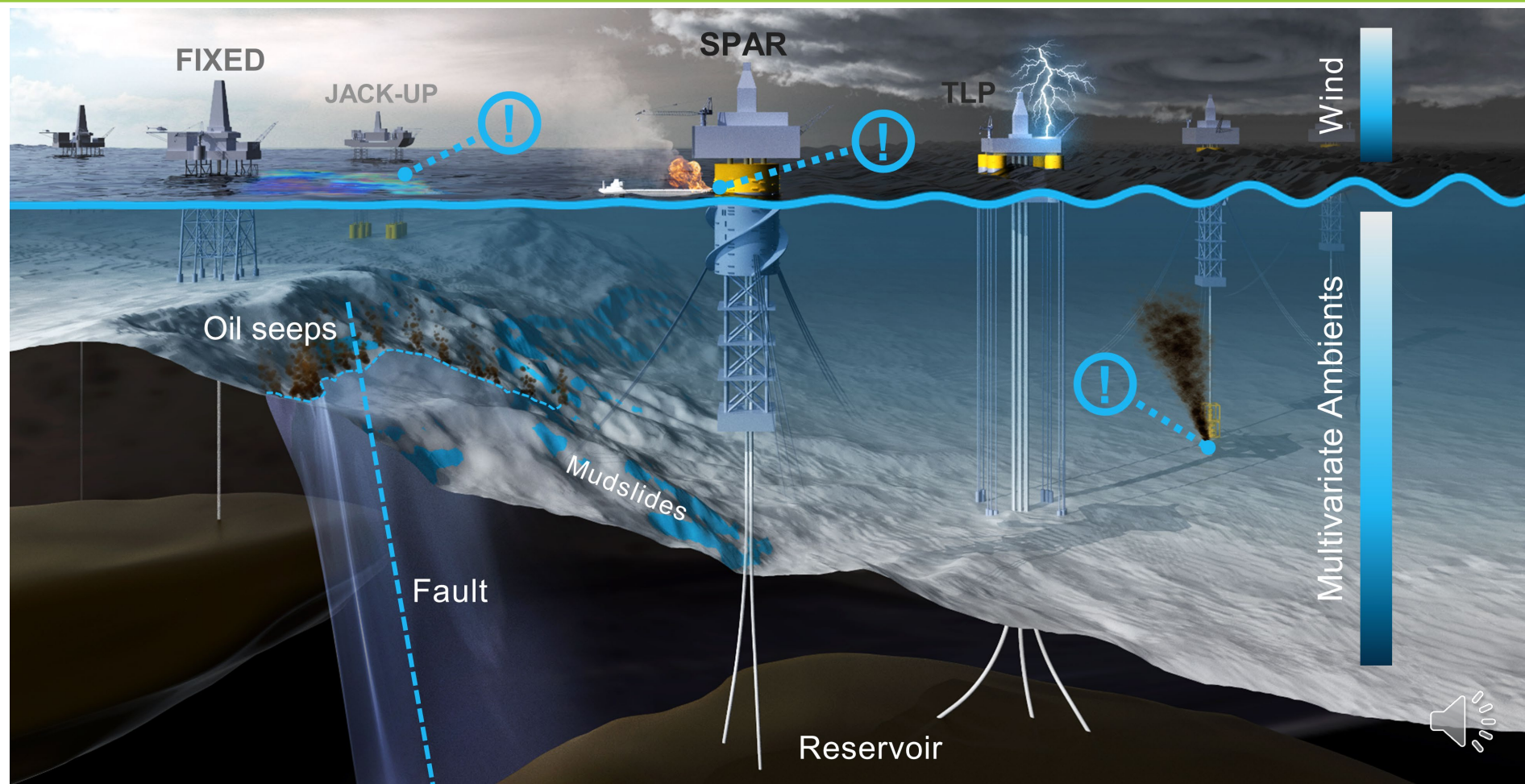
**Thomas Martin**

*Research Scientist*
*Research & Innovation Center*

NATIONAL ENERGY TECHNOLOGY LABORATORY

AIIM

**A**dvanced
**I**nfrastructure
**I**ntegrity
**M**odeling

*Presented to DOE Data Days*
*June 2, 2022 - LLNL*

# Disclaimer

# Authors and Key Personnel

**PIs**

Lucy Romeo[1,2] and Jennifer Bauer[1]

**Key Personnel**

Rodrigo Duran[3], Alec Dyer[1,2], Isabelle Pfander[1,2], Thomas Martin[1,2], Chukwuemeka Okoli[1,2], Kelly Rose[1], Michael Sabbatino[1,2], Madison Wenzlick[1,2], Patrick Wingo[1,2], Dakota Zaengle[1,2]

1) National Energy Technology Laboratory, 1450 Queen Avenue SW, Albany, OR 97321, USA
2) NETL Support Contractor, 1450 Queen Avenue SW, Albany, OR 97321, USA
3) Theiss Research, 7411 Eads Avenue, La Jolla, CA 92037, USA

Thomas.Martine@netl.doe.gov
Lucy.Romeo@netl.doe.gov
Jennifer.Bauer@netl.doe.gov

# Setting the Stage

**Machine Learning** (ML)

- **Supervised ML** – machine is trained, taught with labeled examples

- **Unsupervised ML** – machine creates its own labels (i.e. clustering)

- **Big Data & Big Data Computing** – Large volumes, variety, variability, velocity of data and the computing engineering & systems to handle them

**Features** – Variables or attributes (ex. continuous or categorical)

**Feature Engineering –** Select, transform, process, and visualize input features of a given dataset

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*

# Offshore Infrastructure Hazards

- Aging infrastructure

- Operational wear-and-tear

- Offshore environment:
  - Extreme weather
  - Climate change
  - Corrosion hazards
  - Geohazards

- **Need:**
  - Identify & prevent hazards
  - Inform safe lifespan extension strategies
  - Environmentally prudent planning in low-carbon economy

Typical platform design life, *20-30 years*
**>60% of platforms >30 years old**

https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/

# How AIIM Operates

Utilizes **big data**, **big data computing**, and **multiple ML models** to forecast infrastructure lifespan and risk.

**Key points:**

- Data **analysis** and **visualization** at every step

- Subject Matter Expert **QAQC**

- A focus on the final product being **explainable**, **logical,** and **defendable**



Do the results make sense?

Yes → **Critical Insights!**

No

Visualize Results

Explore Data

Feature Engineering → Singular Value Decomposition (SVD)

→ Self-Organizing Maps (SOM)

→ Other Methods!

Apply ML

Built a dataset of **>11k platforms** with **>2k features** representing *natural-engineered offshore system*

| | |
|---|---|
| Structures | Metocean |
| Incidents | Biochemical |
| Geohazards | Production |

Nelson et al., 2021

Marine Structures
Volume 83, May 2022, 103152

Applied machine learning model comparison: Predicting offshore platform integrity with gradient boosting algorithms and neural networks

Alec S. Dyer, Dakota Zaengle, Jake R. Nelson, Rodrigo Duran, Madison Wenzlick, Patrick C. Wingo, Jennifer R. Bauer, Kelly Rose, Lucy Romeo
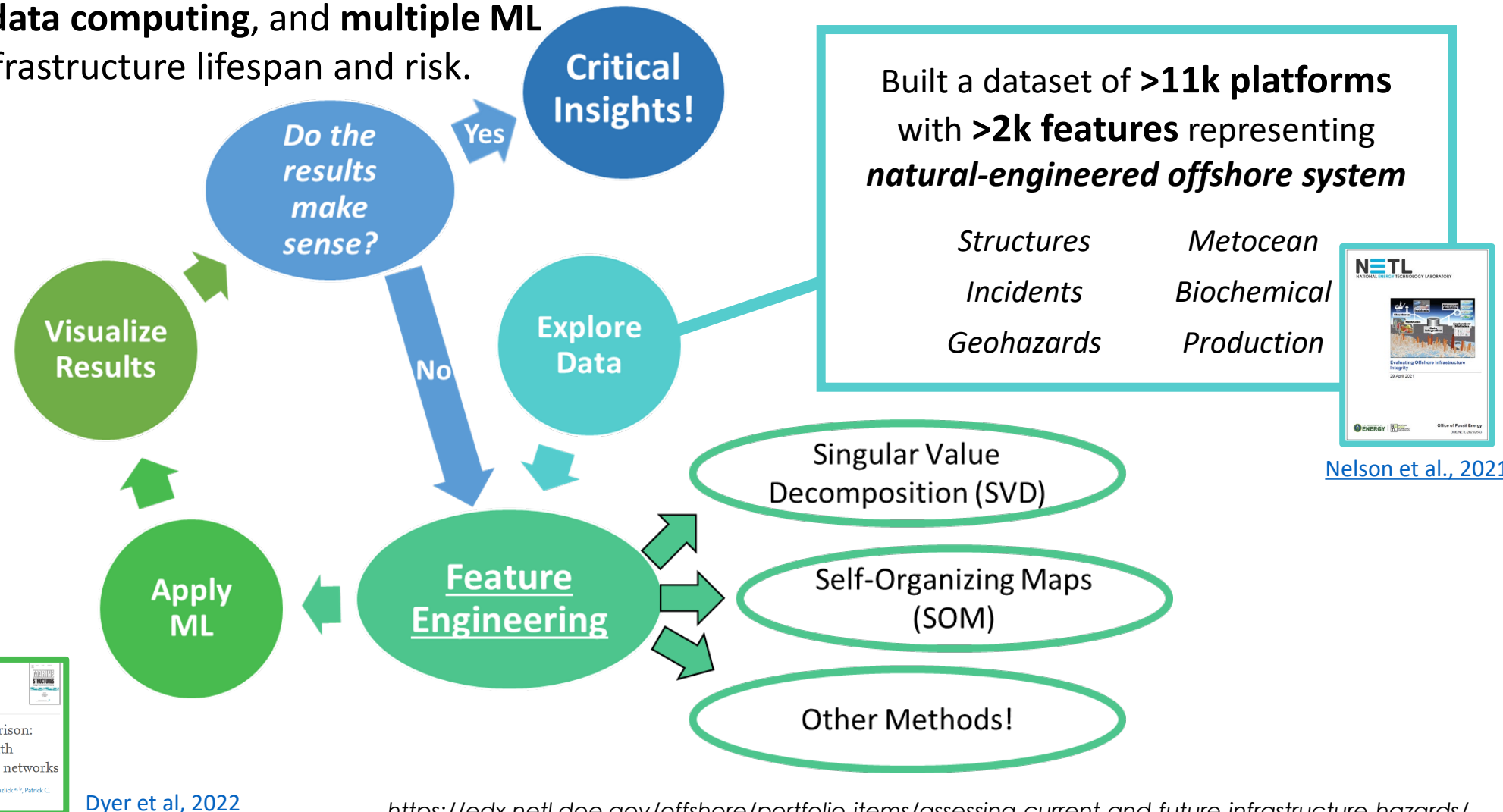
Dyer et al, 2022

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*
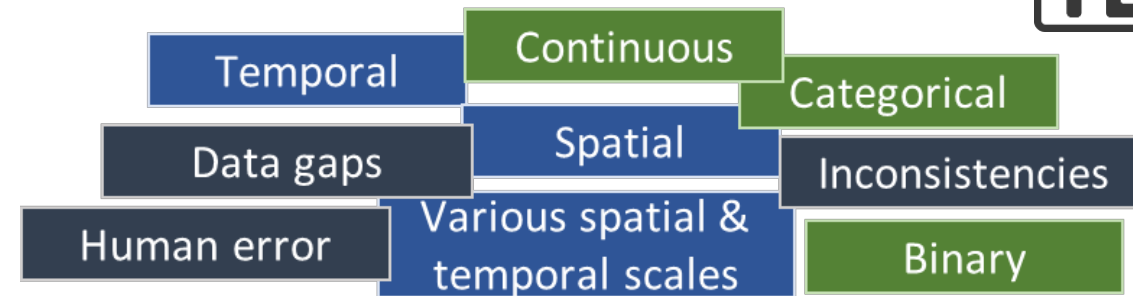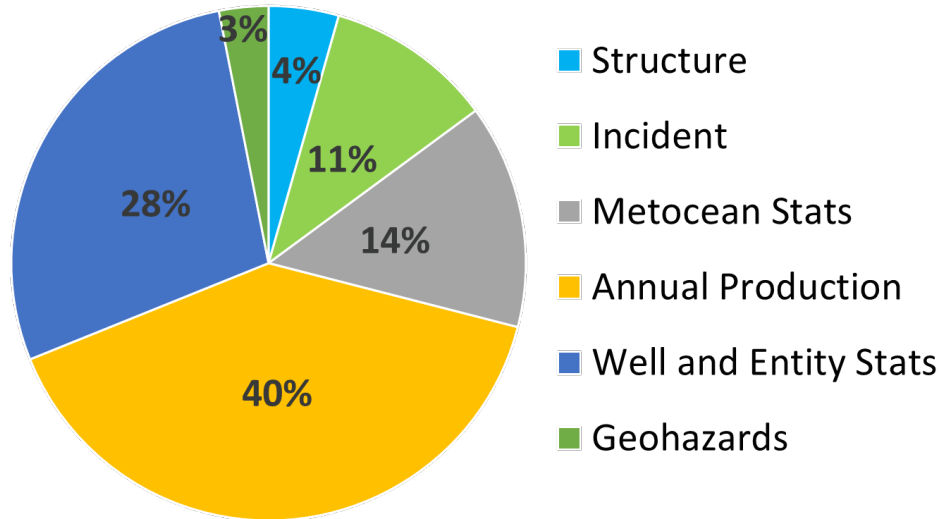
U.S. DEPARTMENT OF ENERGY

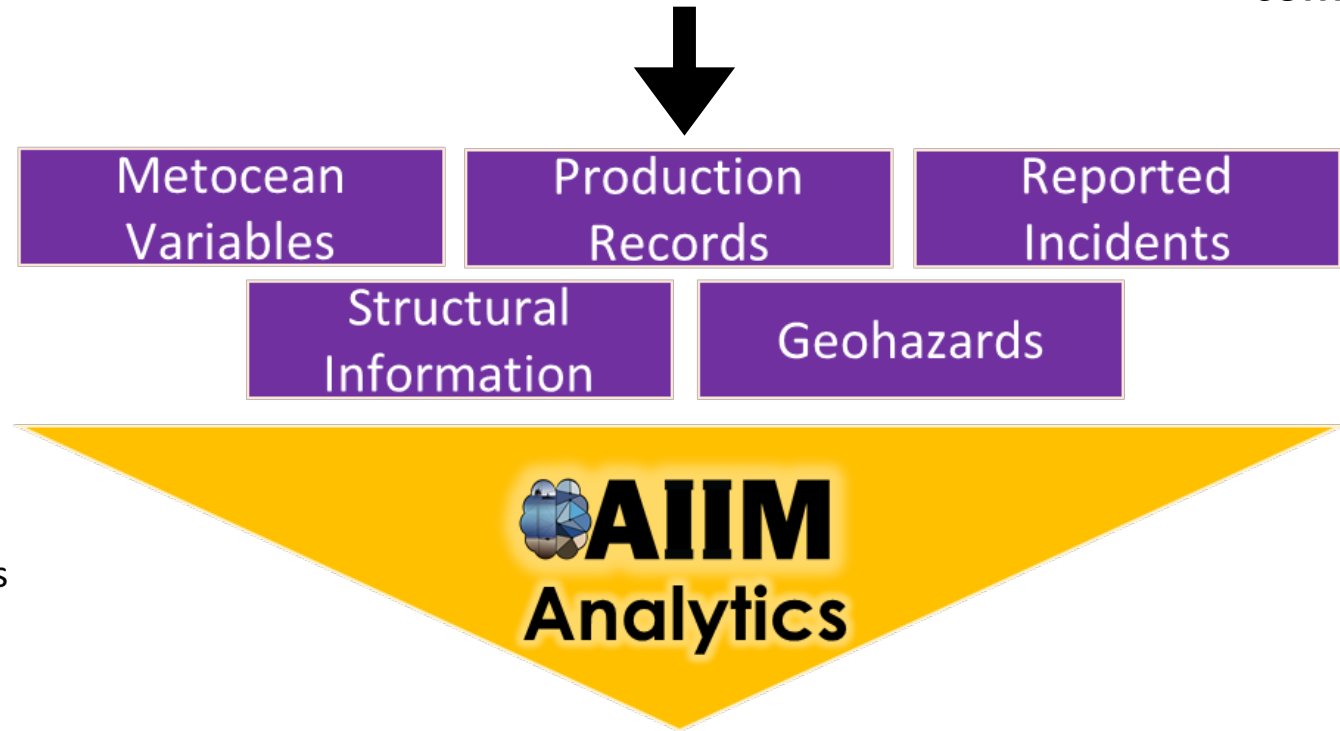# Getting to Know the *Integrated* Data

## Challenges & Opportunities

> 11,000 platform records * > 2,000 features
= **>22,000,000 data values**

*~50% of the dataset has ~90% coverage on a per feature basis*

*Integration has* ***increased data complexity***

| Temporal | Continuous | |
| Data gaps | Spatial | Categorical |
| Human error | Various spatial & temporal scales | Inconsistencies |
| | | Binary |

**Feature Breakdown**

- 4% Structure
- 11% Incident
- 14% Metocean Stats
- 40% Annual Production
- 28% Well and Entity Stats
- 3% Geohazards

| Metocean Variables | Production Records | Reported Incidents |
| Structural Information | | Geohazards |

**AIIM Analytics**

https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/

# Singular Value Decomposition (SVD)

SVD efficiently **identifies** and **summarizes** **important information** in a **correlation** or **covariance** matrix
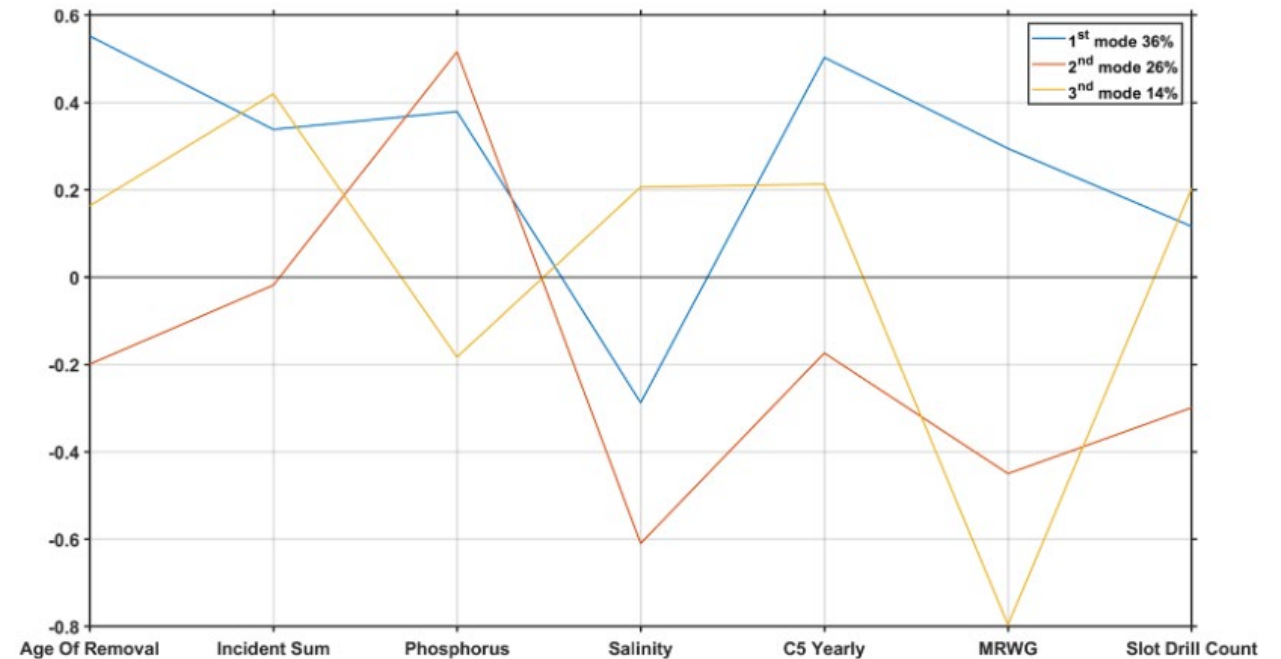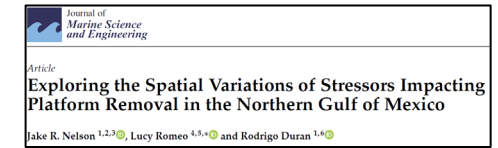
Nelson et al, 2021

Journal of *Marine Science and Engineering*

*Article*
**Exploring the Spatial Variations of Stressors Impacting Platform Removal in the Northern Gulf of Mexico**

Jake R. Nelson [1,2,3], Lucy Romeo [4,5,*] and Rodrigo Duran [1,6]

## Pros

- Interpretable
- Appropriate for time series and continuous spatial data
- Most efficient way to summarize data in a matrix ( Eckart-Young Theorem)

## Cons

- Does not work with categorical features
- Incomplete data requires pre-processing
- Expert opinion needed to select features



**First three right singular vectors** of a data correlation matrix, showing relations between **input variables** and the **target variable** "Age of Removal". 76% percent of features explained by 6 features.

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*
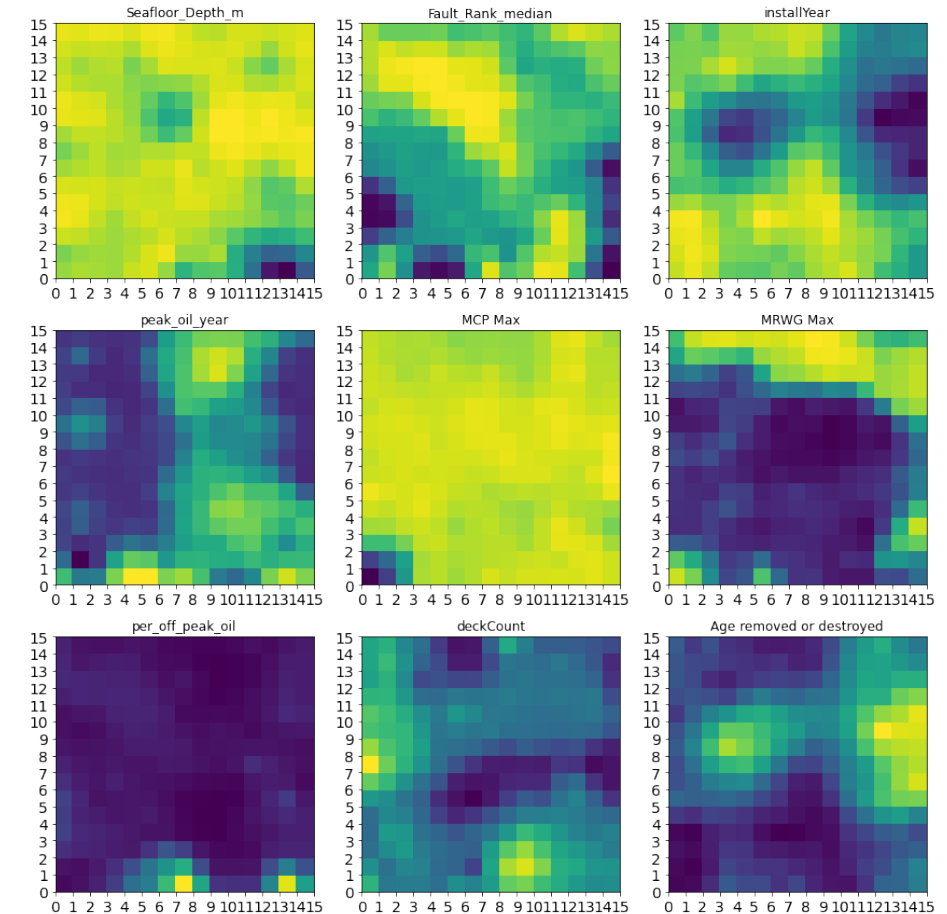
# Self Organizing Maps (SOM)

SOM is an *unsupervised* ML technique that is a specific type of neural network. SOMs identify non-linear feature relationships.

**Pros**

- Can be used with nonlinear features
- Relatively fast
- Threshold for different features is user-selected
- Can be used to create composite features

**Cons**

- Can't be used with categorical data
- Expert opinion needed to select features
- Like all neural networks, complete and pre-processed data helps with convergence and speed



Self Organizing Map Weights compared to target variable (lower right)

https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/
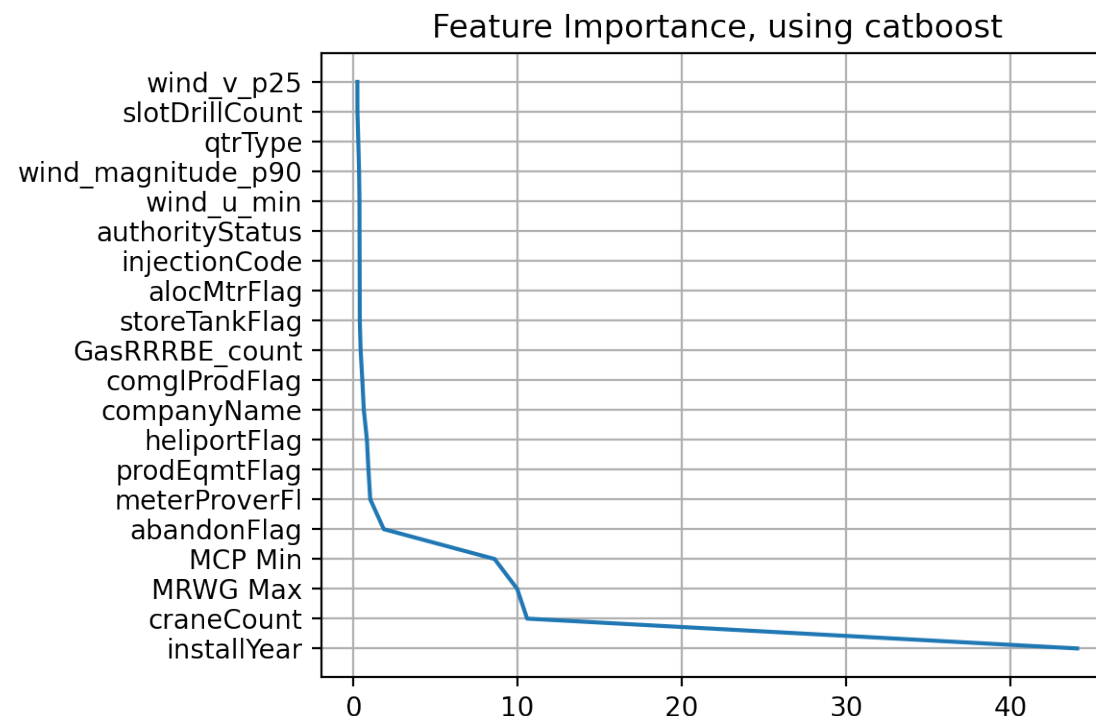
# Letting the Model Decide – Feature Importance

**Gradient Boosted Decision Trees (GBDTs)** are a common and well-used ML algorithm. This is one method to assess every features importance.

| Pros | Cons |
|---|---|
| • Handles all data types<br>• Easily interpretable<br>• Great ML model to be used for prediction as well | • Shared feature importance (potentially collinearity w/other input variables).<br>• Scores are presented quantitively, easy to overinterpret.<br>• Model accuracy has limited impact on feature importance. |



Feature Importance, using catboost

Feature Importance from a GBDT (using CatBoost). This specific model removed many features while still retaining similar accuracy. Low importance features could be further removed.

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*

# Comparison of Methods

## Overcoming incomplete, complex, multivariate data

| SVD | SOM | GBDT | *Using any method alone will give an incomplete picture* |
|---|---|---|---|
| • **Best** for **numeric dat**a (time series, spatial) | • **Best** for deciding between **closely related non-linear features** | • Best for **categorical, incomplete data** | |
| Identified variables containing duplicate information.<br><br>Highlighted storm-related features as important. | Confirmed age variables are key.<br><br>Confirmed findings from SVD testing. | Using top 5–10 variables does not degrade model performance.<br><br>Continued interpretation of environmental loadings and age variables is key. | |

https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/

# Key Findings

**Thomas Martin**
**thomas.martin@netl.doe.gov**
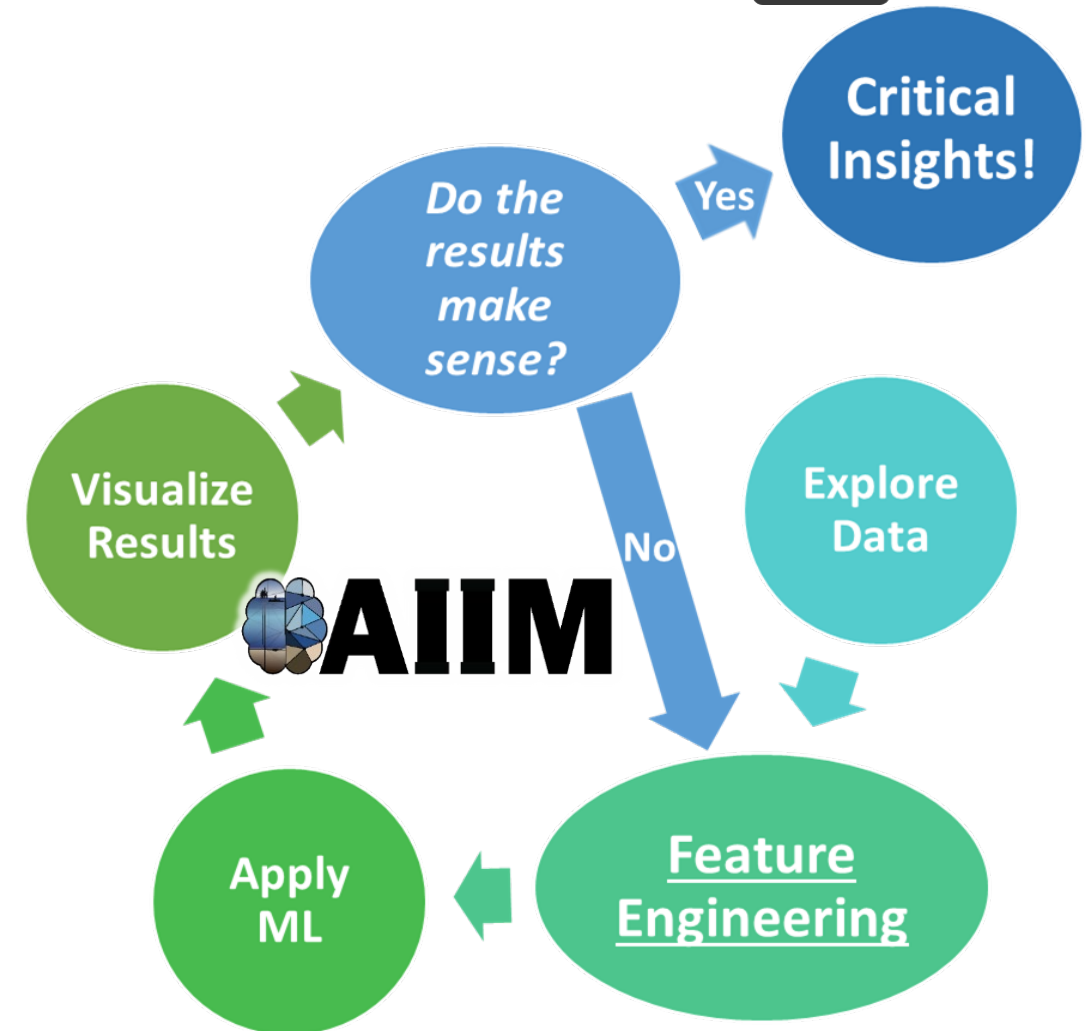
NATIONAL ENERGY TECHNOLOGY LABORATORY

## *Feature Engineering Matters!*

- Reduced input features from >2,000 to >20

- *Minimizing* complexity & error, *maintaining* accuracy

- Insights to inform **safe infrastructure reuse & removal**

- Identify hazards to support **environmental** and **operational risk prevention**

## Next Steps:

- Finalize feature engineering

- Expand to evaluate pipelines & wellbores

- Develop & compare additional models

- Build an interactive AIIM modeling and visualization tool

Critical Insights!

Do the results make sense?  Yes

Explore Data

No

Visualize Results

AIIM

Apply ML

Feature Engineering

EDX Energy Data eXchange

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*

U.S. DEPARTMENT OF ENERGY

# NETL
# Resources

VISIT US AT: www.NETL.DOE.gov

🐦 @NETL_DOE

📷 @NETL_DOE

f @NationalEnergyTechnologyLaboratory

👤 Thomas Martin
   thomas.martin@netl.doe.gov

https://edx.netl.doe.gov/offshore

U.S. DEPARTMENT OF
ENERGY

# References

- Dyer, A.S., Zaengle, D., Nelson, J.R., Duran, R., Wenzlick, M., Wingo, P.C., Bauer, J.R., Rose, K., and Romeo, L. (2022). Applied machine learning model comparison: Predicting offshore platform integrity with gradient boosting algorithms and neural networks, Marine Structures, Volume 83, 103152. https://doi.org/10.1016/j.marstruc.2021.103152.

- Nelson, J. R., Romeo, L., & Duran, R. (2021). Exploring the Spatial Variations of Stressors Impacting Platform Removal in the Northern Gulf of Mexico. *Journal of Marine Science and Engineering*, *9*(11), 1223.

- Nelson, J., Dyer, A., Romeo, L., Wenzlick, M., Zaengle, D., Duran, R., Sabbatino, M., Wingo, P., Barkhurst, A., Rose, K., Bauer, J. Evaluating Offshore Infrastructure Integrity; DOE/NETL-2021/2643; NETL Technical Report Series; U.S. Department of Energy, National Energy Technology Laboratory: Albany, OR, 2020; p 70. doi.org/10.2172/1780656

*https://edx.netl.doe.gov/offshore/portfolio-items/assessing-current-and-future-infrastructure-hazards/*