Machine learning for data-driven discovery in solid Earth geoscience

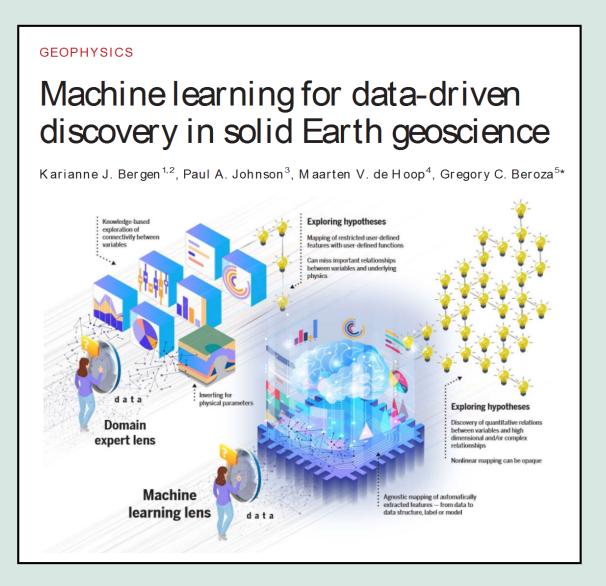
Karianne J. Bergen, Ph.D.

Harvard Data Science Initiative Postdoctoral
Fellow



Harvard John A. Paulson School of Engineering and Applied Sciences

FAST: Earthquake detection by efficient similarity search Single-channel continuous time series data Sparse Similarity **Fingerprint Extraction Efficient Similarity Search** Matrix Waveform Database Waveform Fingerprint **Detection Results**



What is Machine Learning?

Machine learning (ML)

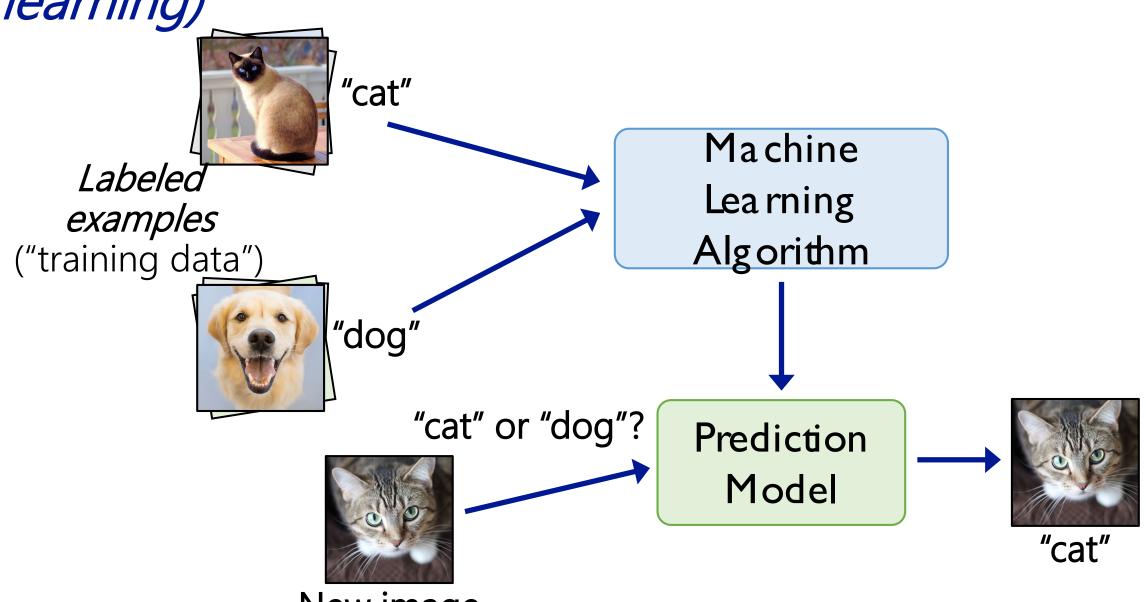
a set of tools for extracting patterns and building predictive models from data.

Data mining

tools for extracting unknown patterns or information from large data sets

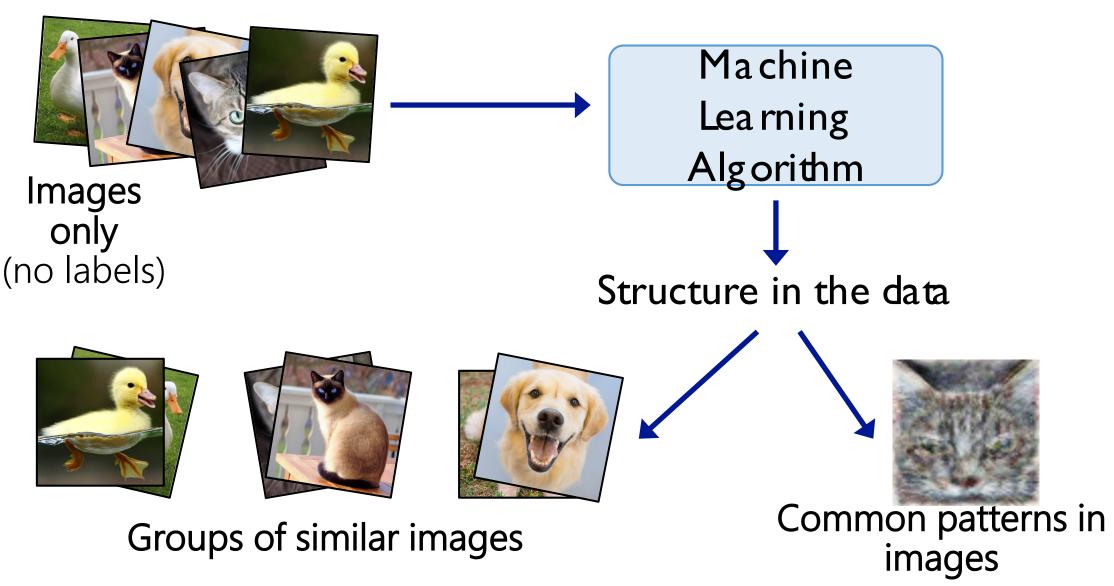


Building models from examples (Supervised learning)



New image K. J. Bergen | COSG Meeting

Finding patterns in data (Unsupervised learning)



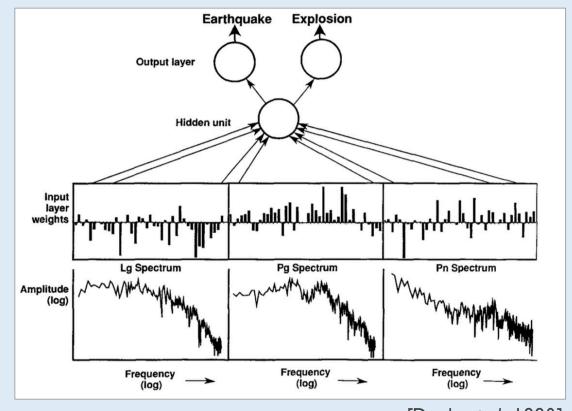
Deep Neural Networks Fast simulations & surrogate models **Recurrent Neural Networks** Inverse problems Deep **Autoencoder Networks** Generative **Convolutional Neural Networks Featurization** Models Dynamic decisions Reinforcement **Dictionary Learning Artif cal Neural Networks** Learning **Feature Learning** Learn joint **Support Vector Machines** probability distribution Prediction Clustering & Random Forests & Ensembles Self-organizing maps Detection & classif cation **Graphical Models** Determine optimal boundary Sparse representation Logistic Regression Semi-Feature representation Domain adaptation Supervised Dimensionality reduction Learning Supervised Learning **Unsurpervised Learning**

Machine learning can help geoscientists extract more knowledge & insights from larger data sets than ever before.

Geoscientists have been using ML for decades.

Artificial Neural Networks

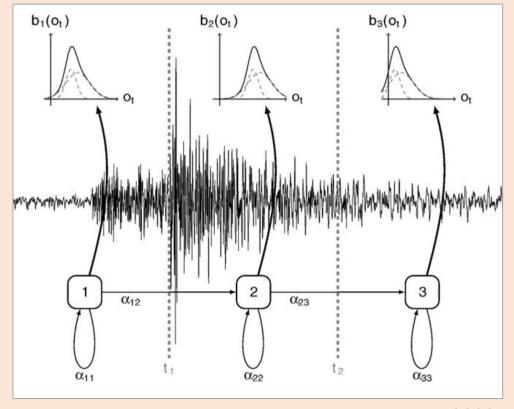
(e.g. Dowla et al., 1990; Dysart & Pulli, 1990)



[Dowla et al., 1990]

Hidden Markov Models

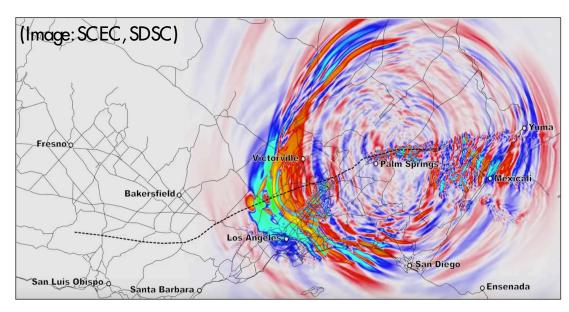
(e.g. Ohrnberger, 2001; Beyreuther et al., 2008)



Recent developments have created new opportunities for scientific discovery with ML in solid Earth geoscience.

- 1) Massive geoscience data sets
- 2) New ML algorithms and models
- 3) Improvements in computing technology & tools

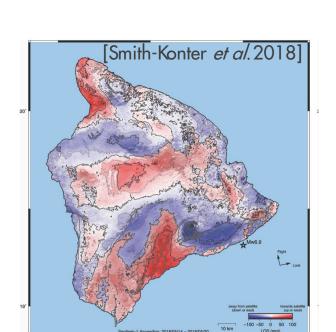
Massive geoscience data sets



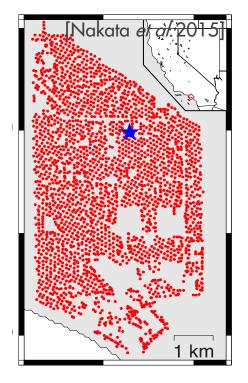
Large-scale simulations







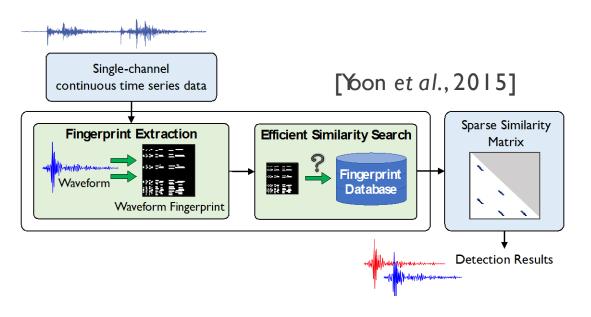
Spatiotemporal & remote sensing data



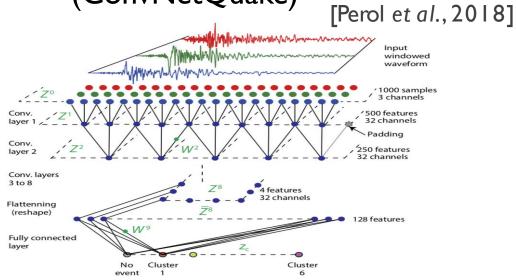
Dense sensor arrays

New ML algorithms and models

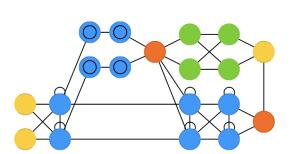
Large-scale Data Mining (FAST)



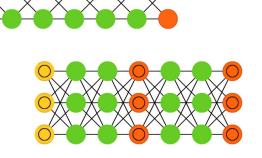
Convolutional Neural Network (ConvNetQuake)

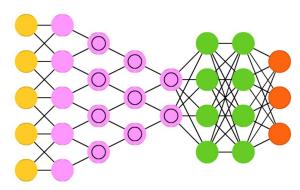






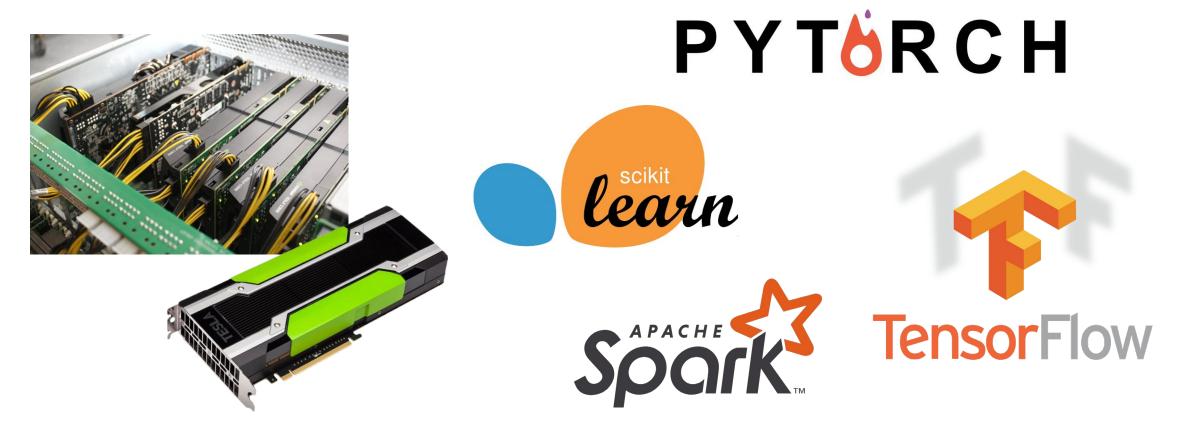






[Van Veen & Leijnen, 2019]

Improvements in computing technology & tools





Where has ML been particularly successful?

[Treml et al. 2016] [DeepMind]

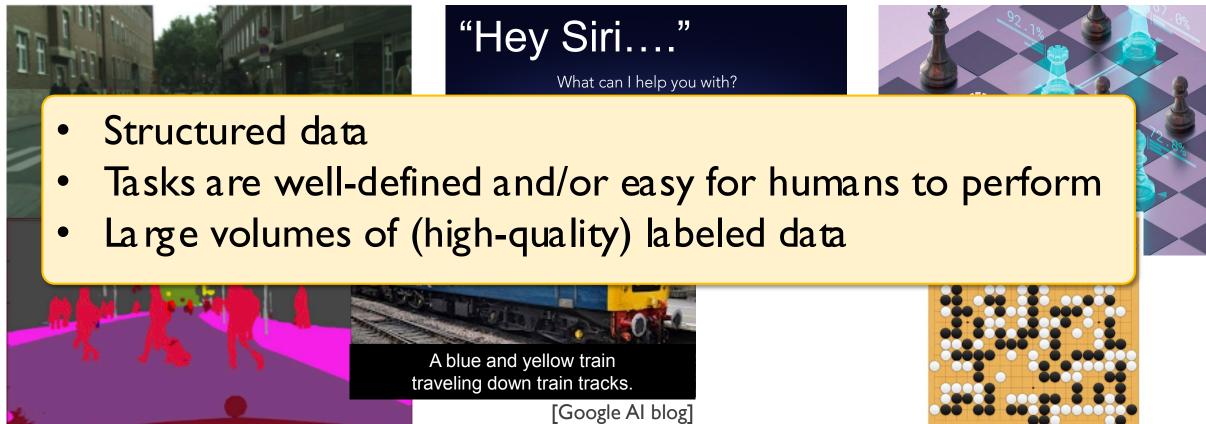


Image Processing data on a grid: images

La nguage Processing sequential data: text, audio

Game Play

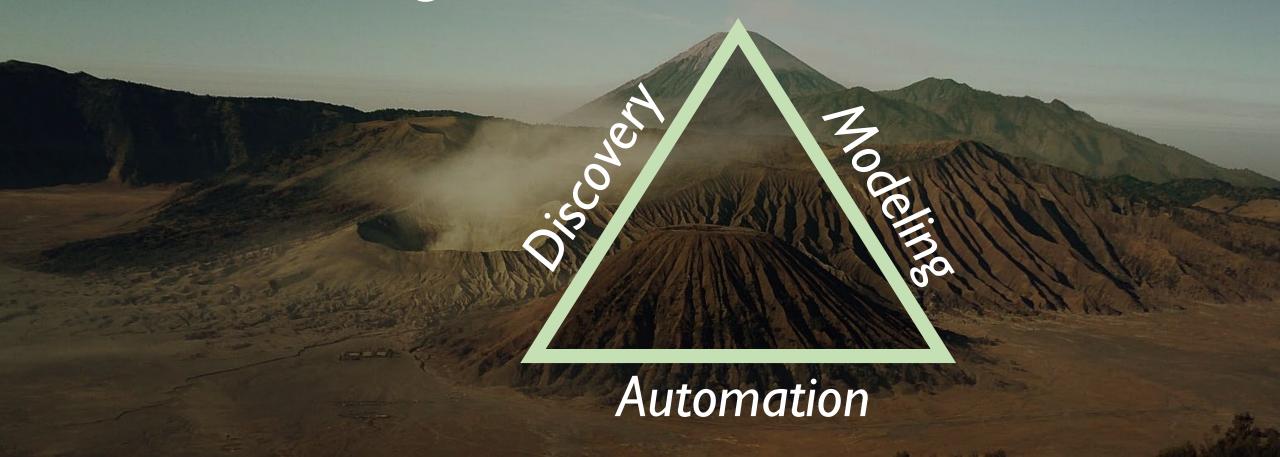
Data sets in solid Earth geoscience

- Low signal-to-noise data
- Heterogenous noise
- Limited or low-quality labels
- Ground truth often unavailable
- Massive data sets require scalable methods
- Modeling across multiple scales (spatial, temporal)



How is machine learning being used in solid Earth geoscience today?

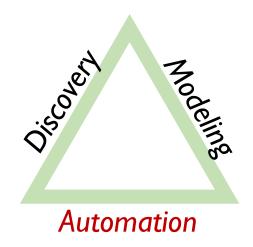
How might it be used in the near future?



Automated prediction & analysis

Goal: Perform a complex or repetitive task

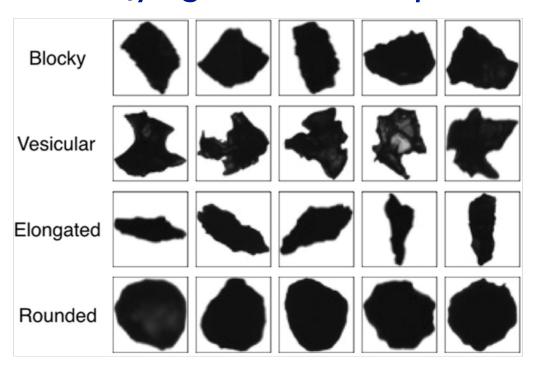
- a) challenging, tedious or infeasible for an analyst to perform, OR
- b) difficult to express as a set of explicit commands
- ML as a tool for high accuracy predictions





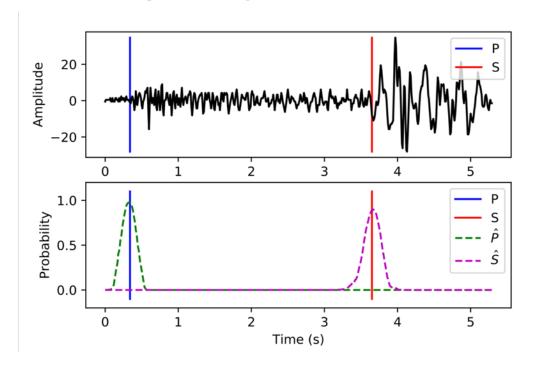
Automated data analysis

Classifying volcanic ash particles



Tedious task for an analyst, well-suited to automation with ML

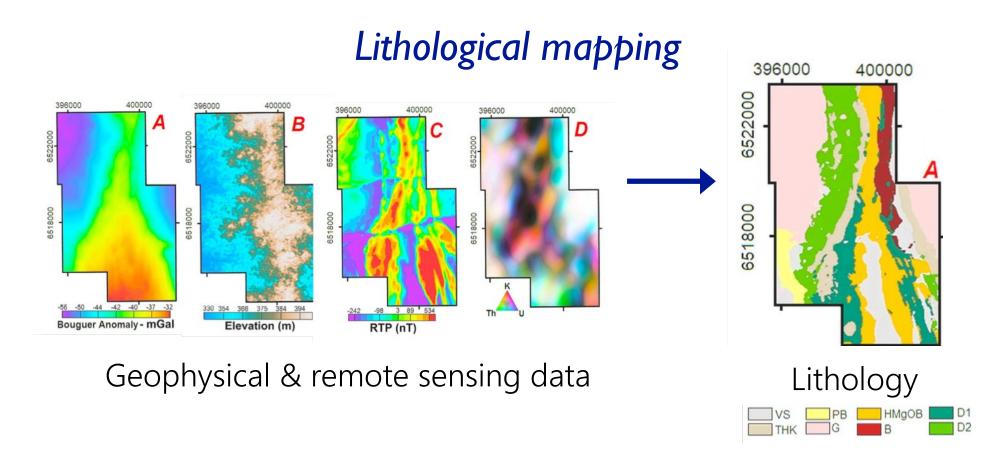
Phase-picking in seismic data



Automating, improving accuracy of existing analysis pipeline

[Zhu & Beroza (2018), Geophys J.

Automated prediction



Learning to predict lithology from sparse ground-truth measurements

[Kuhn et al., (2018), Geophysics.]

Machine learning challenges in solid Earth geoscience

• Data set shift, cova ria te shift

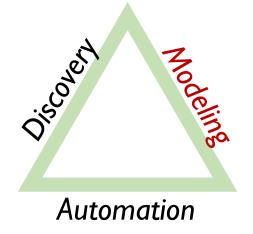
Automation

- Biases in data collection, labeling
- Evaluating performance

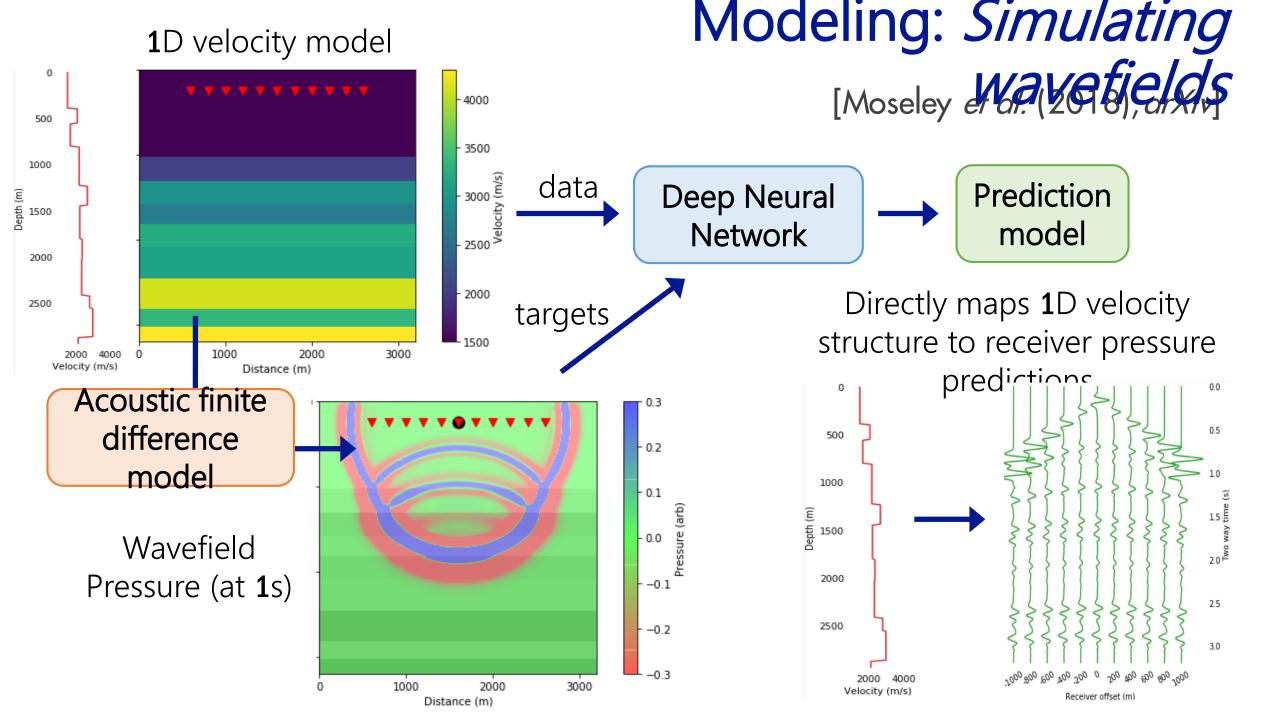
Modeling & inversion

Goal: Create a representation that captures relationships or structure in a data set.

- Learning surrogate models
- Model reduction & coarse-graining
- Intersection of ML & numerical simulations



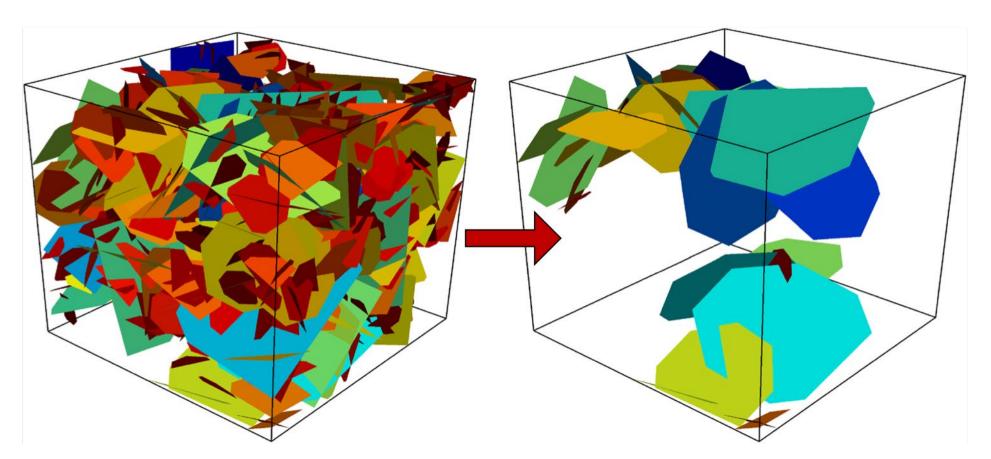




Modeling: *Model reduction for flow in porous media*

[Valera et al. (2018), Comput.

Geosci.]



Full network model
Discrete fracture network

"Backbone" (reduced network model)
Subnetwork capturing flow pattern of full network

K. J. Bergen | COSG Meeting

Machine learning challenges in solid Earth geoscience

• Data set shift, cova ria te shift

Automation

- Biases in data collection, labeling
- Evaluating performance
- Quantifying model uncertainty

Modeling

- Physical constraints, domain knowledge
- Expense of collecting training data (from simulations)

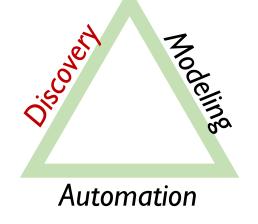
Discovering patterns & insights

Goal: Extract new information (complex patterns, structure, or relationships) from scientific data sets,

especially patterns *not easily* revealed

by conventional analysis techniques.

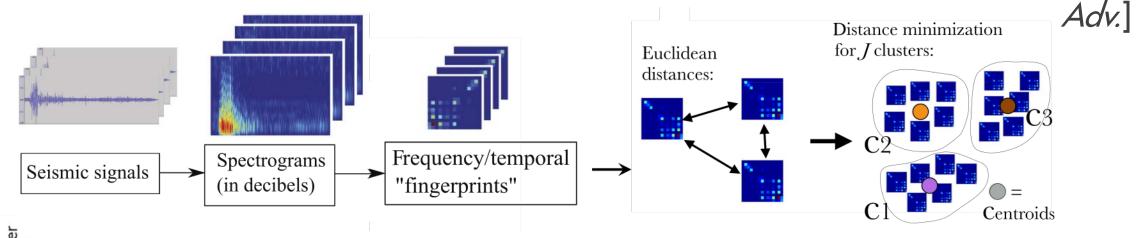
- Unsupervised learning
- Generative models

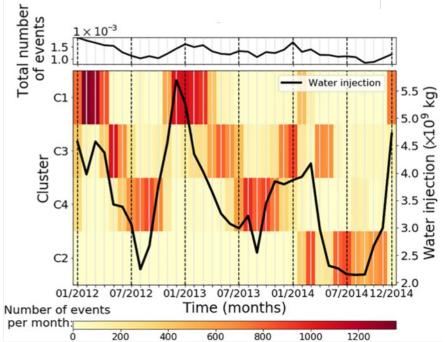




Discovery: finding temporal patterns among seismic signals

[Holtzman et al. (2018), Sci.



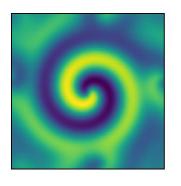


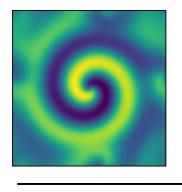
ML identifies spectra-temporal patterns, too subtle for standard methods, that reflect changes in faulting process.

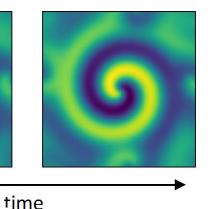
Discovery: learning governing equations from data

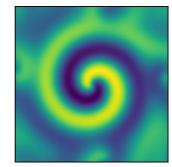
[Champion *et al.* (2019), *PNAS*]

Reaction-diffusion system







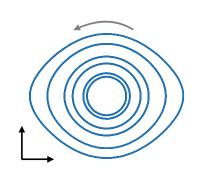


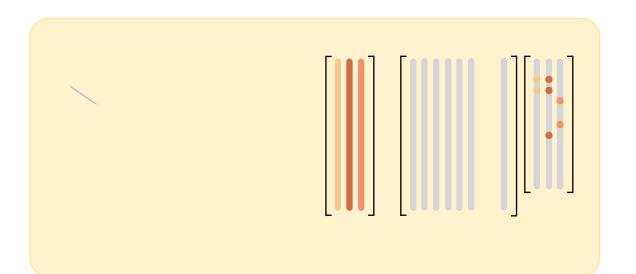
$$\dot{z}_1 = -0.85z_2$$

 $\dot{z}_2 = 0.97z_1$

$$\dot{z}_2 = 0.97z_1$$

Nonlinear pendulum





Machine learning challenges in solid Earth geoscience

• Data set shift, cova ria te shift

Automation

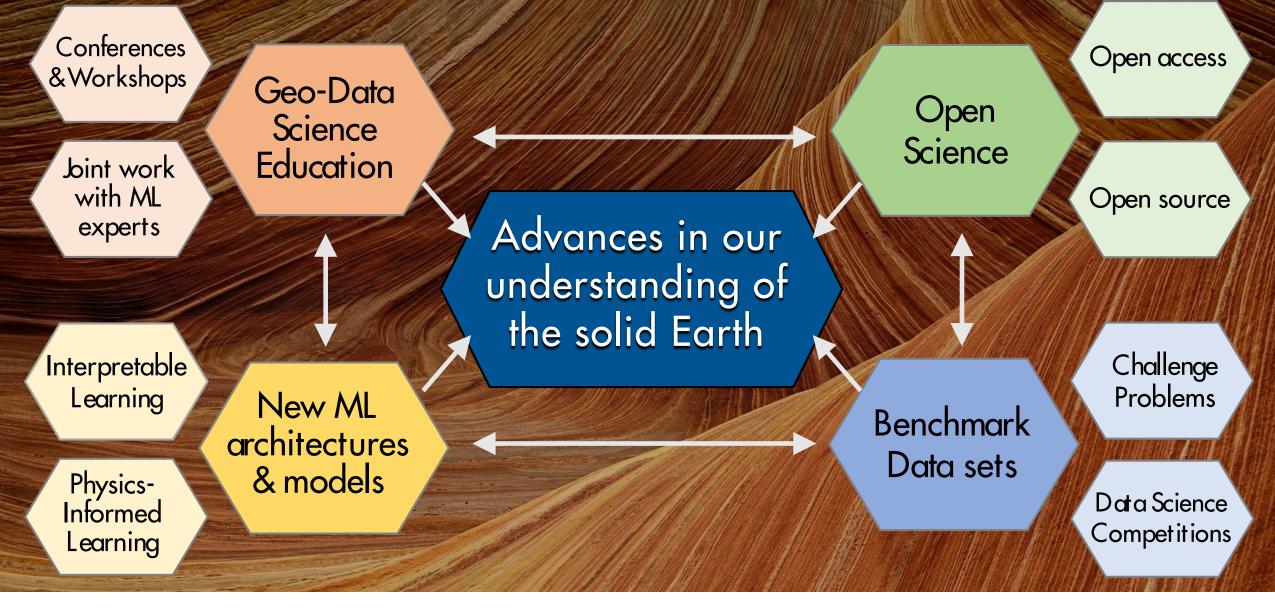
- Biases in data collection, labeling
- Evaluating performance
- Quantifying model uncertainty

Modeling

- Physical constraints, domain knowledge
- Expense of collecting training data (from simulations)
- Interested in outliers, infrequent events, unexpected patterns

Discovery

- Interpolation vs. extrapolation
- Discovery often requires interpreting "black box" models





karianne_bergen@fas.harvard.edu



Modeling: Representing seamount bathymetry

[Valentine *et al.* (2013)]

