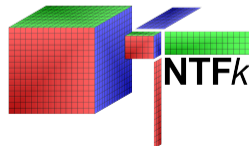


# Predicting Oil and Gas Production from Unconventional Tight-Rock Reservoirs Using Machine Learning

**Velimir V. Vesselinov (monty)** (vuv@lanl.gov)

Earth and Environmental Sciences Division  
Los Alamos National Laboratory, NM, USA

<http://tensors.lanl.gov>



- ▶ **Supervised** ML: learns everything from data
  - ⇒ requires big training datasets
  - ⇒ highly impacted by noise
- ▶ **Physics-informed** ML: learns from data but includes preconceived knowledge about the governing processes
  - ⇒ requires smaller training datasets
  - ⇒ produces better predictability with lower uncertainty
  - ⇒ robust to data noise
- ▶ **Unsupervised** ML: extracts features from data that can be applied for categorization and prediction
  - ⇒ unbiased analyses not impacted by data labeling, subject-matter-expert opinions, and physics assumptions
  - ⇒ however, physics constraints can be added

- ▶ **Supervised** ML: learns everything from data
  - ⇒ requires big training datasets
  - ⇒ highly impacted by noise
- ▶ **Physics-informed** ML: learns from data but includes preconceived knowledge about the governing processes
  - ⇒ requires smaller training datasets
  - ⇒ produces better predictability with lower uncertainty
  - ⇒ robust to data noise
- ▶ **Unsupervised** ML: extracts features from data that can be applied for categorization and prediction
  - ⇒ unbiased analyses not impacted by data labeling, subject-matter-expert opinions, and physics assumptions
  - ⇒ however, physics constraints can be added

- ▶ **Supervised** ML: learns everything from data
  - ⇒ requires big training datasets
  - ⇒ highly impacted by noise
- ▶ **Physics-informed** ML: learns from data but includes preconceived knowledge about the governing processes
  - ⇒ requires smaller training datasets
  - ⇒ produces better predictability with lower uncertainty
  - ⇒ robust to data noise
- ▶ **Unsupervised** ML: extracts features from data that can be applied for categorization and prediction
  - ⇒ unbiased analyses not impacted by data labeling, subject-matter-expert opinions, and physics assumptions
  - ⇒ however, physics constraints can be added

- ▶ **Supervised** ML: requires “labeling” (prior categorization (knowledge) about the processed data)

**Example:** Recognizes images of cats and dogs after extensive training; but cannot recognize horses if not trained

**Cannot discover something that we do not know already**

- ▶ **Unsupervised** ML: extracts hidden features (signals) in the processed data without any prior information (**exploratory analysis; data-driven science**)

**Example:** Identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.)

- ▶ **Supervised** ML: requires “labeling” (prior categorization (knowledge) about the processed data)

**Example:** Recognizes images of cats and dogs after extensive training; but cannot recognize horses if not trained

Cannot discover something that we do not know already

- ▶ **Unsupervised** ML: extracts hidden features (signals) in the processed data without any prior information (**exploratory analysis; data-driven science**)

**Example:** Identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.)

- ▶ **Supervised** ML: requires “labeling” (prior categorization (knowledge) about the processed data)

**Example:** Recognizes images of cats and dogs after extensive training; but cannot recognize horses if not trained

**Cannot discover something that we do not know already**

- ▶ **Unsupervised** ML: extracts hidden features (signals) in the processed data without any prior information (**exploratory analysis; data-driven science**)

**Example:** Identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.)

- ▶ **Supervised** ML: requires “labeling” (prior categorization (knowledge) about the processed data)

**Example:** Recognizes images of cats and dogs after extensive training; but cannot recognize horses if not trained

**Cannot discover something that we do not know already**

- ▶ **Unsupervised** ML: extracts hidden features (signals) in the processed data without any prior information (**exploratory analysis**; **data-driven science**)

**Example:** Identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.)

- ▶ **Supervised** ML: requires “labeling” (prior categorization (knowledge) about the processed data)

**Example:** Recognizes images of cats and dogs after extensive training; but cannot recognize horses if not trained

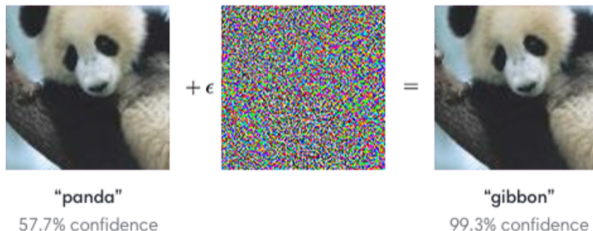
**Cannot discover something that we do not know already**

- ▶ **Unsupervised** ML: extracts hidden features (signals) in the processed data without any prior information (**exploratory analysis**; **data-driven science**)

**Example:** Identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.)

## ▶ Supervised ML

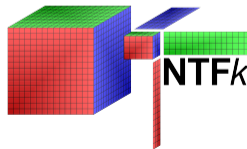
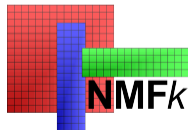
- ▶ introduces subjectivity (through the labeling process)
- ▶ does not provide insights why horses are different from dogs / cats
- ▶ cannot make predictions (that we do not know already)
- ▶ requires huge training (labeled) datasets
- ▶ we do not know why it works
- ▶ is impacted by “adversarial examples”



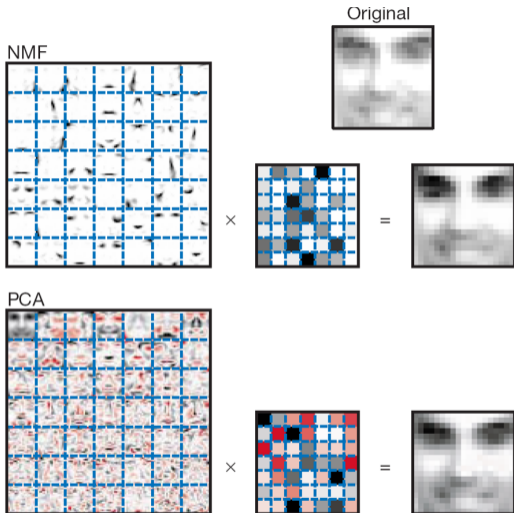
⇒ major limitations of the **supervised** ML methods for **science** applications

- ▶ Feature extraction (**FE**)
- ▶ Blind source separation (**BSS**)
- ▶ Detection of disruptions / anomalies
- ▶ Image recognition
- ▶ Separate physics processes
- ▶ Discover unknown dependencies and phenomena
- ▶ Develop reduced-order/surrogate models
- ▶ Identify dependencies between model inputs and outputs
- ▶ Guide development of physics models representing the data
- ▶ Make predictions
- ▶ Optimize data acquisition
- ▶ “Label” datasets for supervised ML analyses

- ▶ Novel LANL-patented, open-source, unsupervised Machine Learning (ML) methods and computational techniques
- ▶ Based in matrix/tensor factorization coupled with custom  $k$ -means clustering and nonnegativity/sparsity constraints:
  - NMF $k$ : Nonnegative **Matrix** Factorization
  - NTF $k$ : Nonnegative **Tensor** Factorization
  - <https://github.com/TensorDecompositions>
- ▶ Capable to efficiently process large datasets (TB's) utilizing GPU's, TPU's & FPGA's  
⇒ **Julia**, Flux.jl, AutoOffLoad.jl, TensorFlow, PyTorch, MXNet

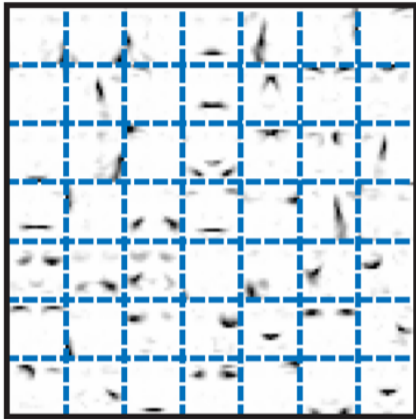


- ▶ NMF vs PCA (Lee & Seung, 1999)
- ▶ NMF: Nonnegative Matrix Factorization
- ▶ PCA: Principal Component Analysis

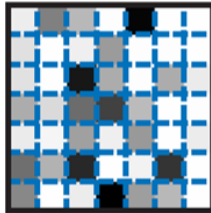


**Nonnegativity constraints provide meaningful and interpretable results (+sparsity)**

- ▶ **Tensors** (multi-dimensional/multi-modal/multi-way datasets) are everywhere:
  - ▶ observational data are typically a 5-D tensor (x, y, z, t, attributes)
  - ▶ model outputs are typically a 5-D tensor (x, y, z, t, attributes)
  - ▶ data dependency to  $N$  parameters will form a  $(N + 5)$ -D tensor

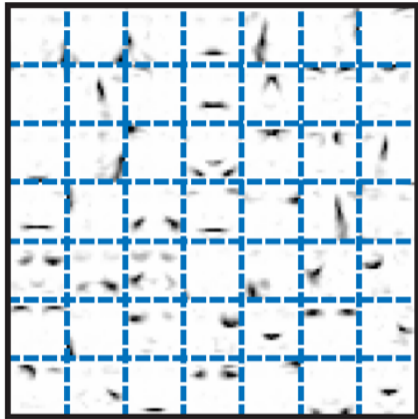


$\times$

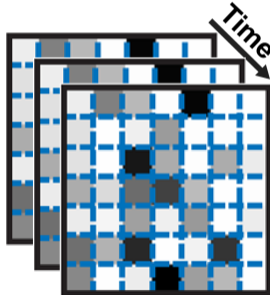


$=$

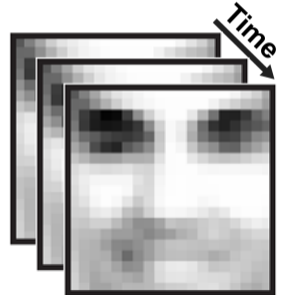


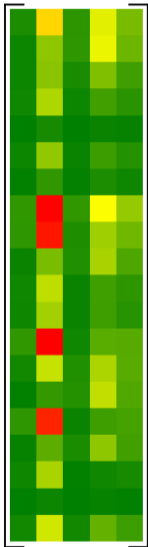


$\otimes$



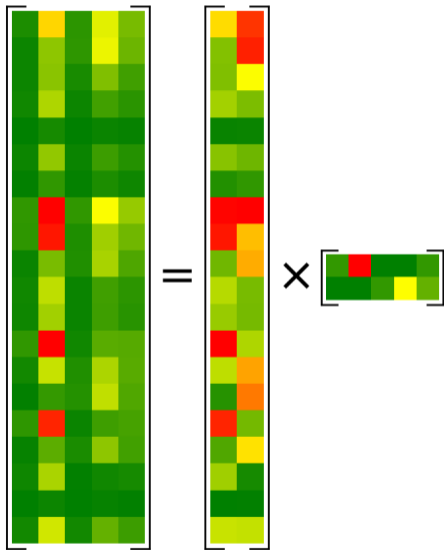
$=$





$X$   
[20 × 5]

$X$  – data matrix  
[attributes × observations]



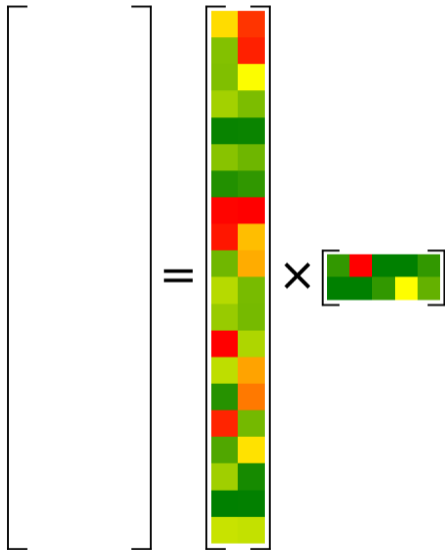
$$X = W \times H$$

$$[20 \times 5] = [20 \times 2] \times [2 \times 5]$$

$X$  – **data** matrix  
 [attributes  $\times$  observations]

$W$  – **feature (signal)** matrix  
 [attributes  $\times$  features]

$H$  – **mixing** matrix  
 [features  $\times$  observations]



$$X = W \times H$$

$$[20 \times 5] = [20 \times 2] \times [2 \times 5]$$

$X$  – **data** matrix

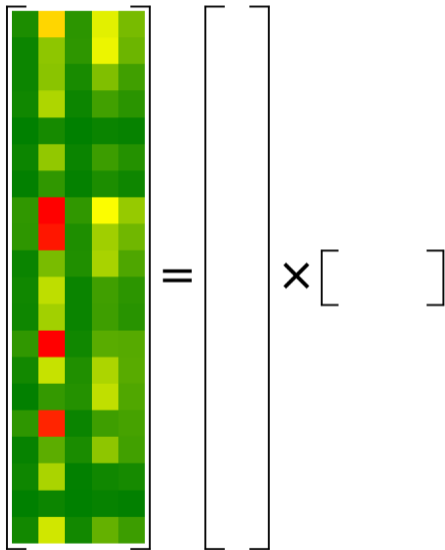
[attributes × observations]

$W$  – **feature (signal)** matrix

[attributes × features]

$H$  – **mixing** matrix

[features × observations]



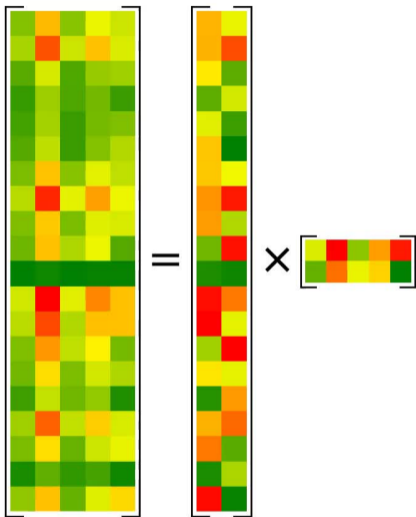
$$X = W \times H$$

$$[20 \times 5] = [20 \times ?] \times [? \times 5]$$

⇒ 100 **knowns**

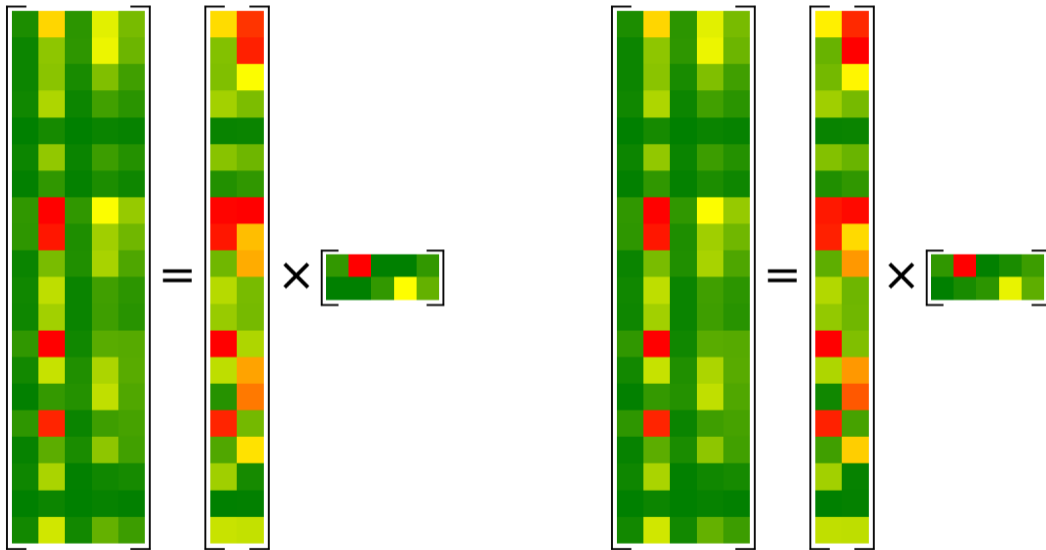
⇒ **unknown** number of features  
(2 or more)

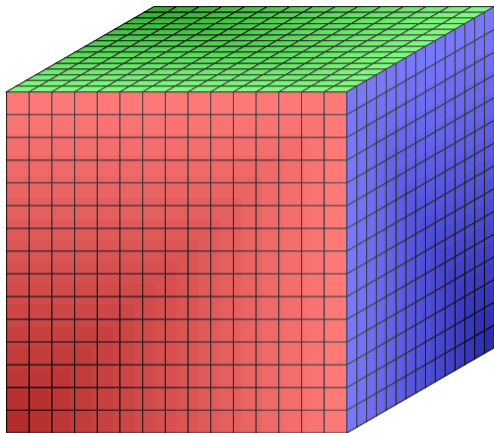
⇒ **unknown** matrix elements of  $W$  and  $H$   
(50 or more)



0001

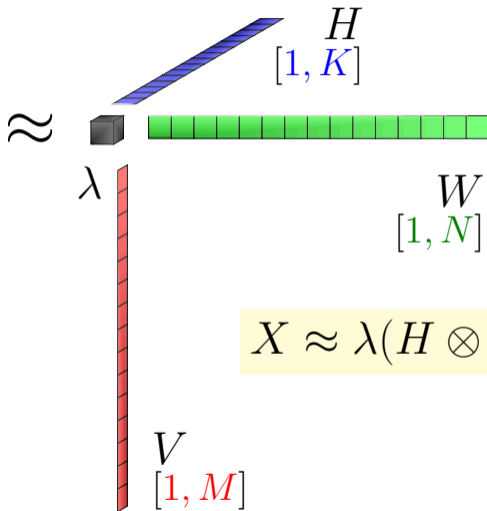
# NMF<sub>k</sub>: true vs. estimated matrix factorization



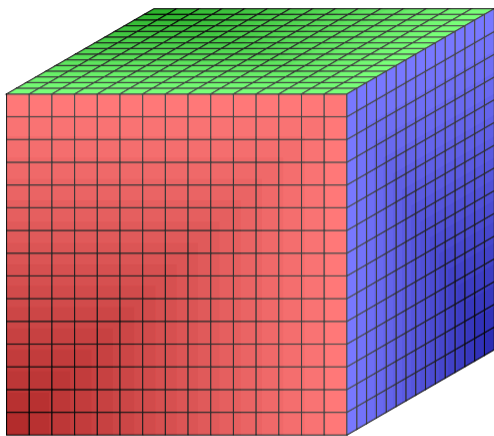


$$X$$

$$[K, M, N]$$

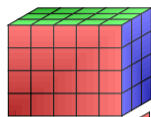


$$X \approx \lambda(H \otimes W \otimes V)$$

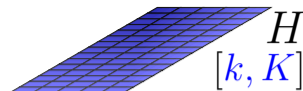


$X$   
 $[K, M, N]$

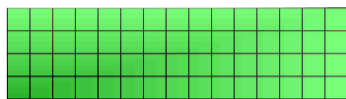
$\approx$



$G$   
 $[k, m, n]$



$H$   
 $[k, K]$

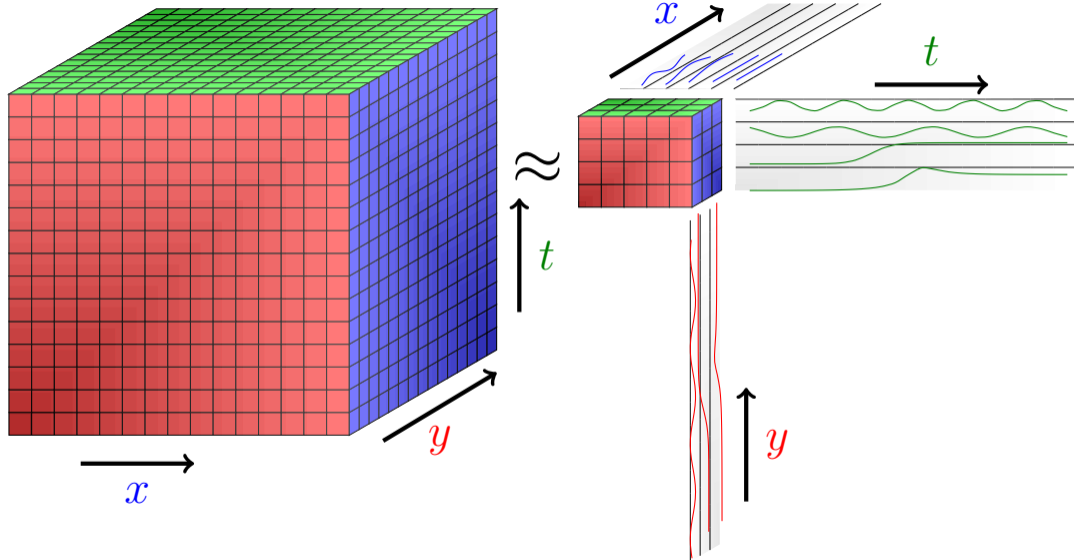


$W$   
 $[n, N]$



$V$   
 $[m, M]$

$$X \approx G \otimes H \otimes W \otimes V$$



▶ **Field Data:**

- ▶ Contamination
- ▶ Climate
- ▶ Geothermal
- ▶ Seismic
- ▶ Oil/gas production
- ▶ CO<sub>2</sub> sequestration
- ▶ **Wildfires** (California 2020)
- ▶ **COVID-19**

▶ **Lab Data:**

- ▶ X-ray Spectroscopy
- ▶ UV Fluorescence Spectroscopy
- ▶ Microbial population analyses
- ▶ Isotope fractionation

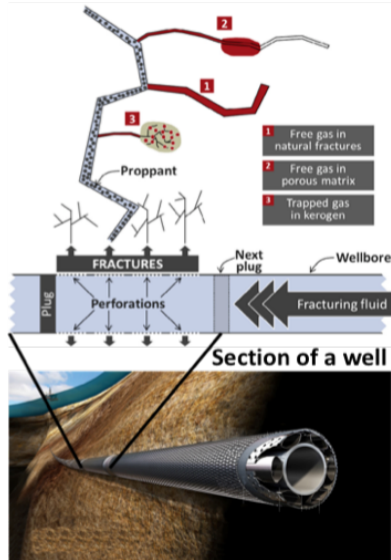
▶ **Operational Data:**

- ▶ LANSCE: Los Alamos Neutron Accelerator
- ▶ Oil/gas production
- ▶ CO<sub>2</sub> sequestration

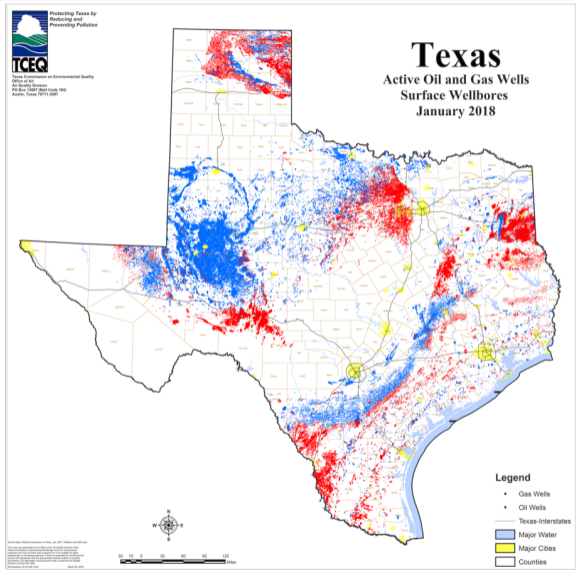
▶ **Model Outputs:**

- ▶ Reactive mixing  $A + B \rightarrow C$
- ▶ Phase separation of co-polymers
- ▶ Molecular Dynamics of proteins
- ▶ Climate
- ▶ CO<sub>2</sub> sequestration
- ▶ Wildfires

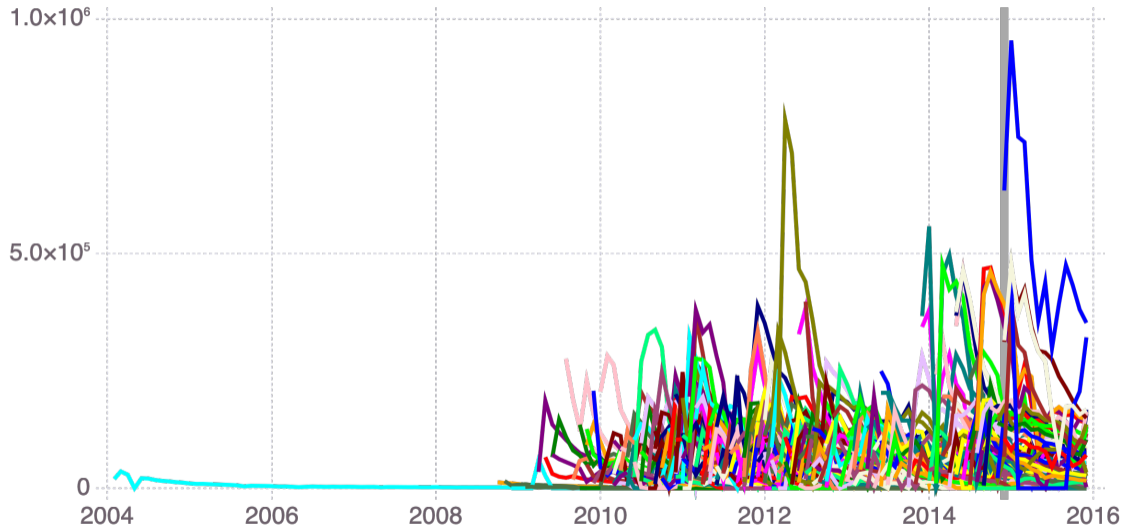
- ▶ Oil/Gas production from unconventional reservoirs extracts a small portion of the available resources (<10%)
- ▶ Oil/Gas production is challenging to predict and optimize
- ▶ Physics processes during well development (including hydrofracking) and extraction are poorly understood and challenging to simulate
- ▶ Alternative is to learn to predict system behavior based on the observed oil/gas production at existing wells



- ▶ Large public datasets are available representing unconventional oil and gas production (U.S. and world wide)
- ▶ Data represent monthly production rates (oil, gas, water) + many other well attributes
- ▶ ~ 2,000,000 wells in U.S.
- ▶ > 300,000 wells in Texas
- ▶ > 20,000 wells in Eagle Ford Shale Play
- ▶ 327 gas wells in Eagle Ford Shale Play selected for preliminary analyses

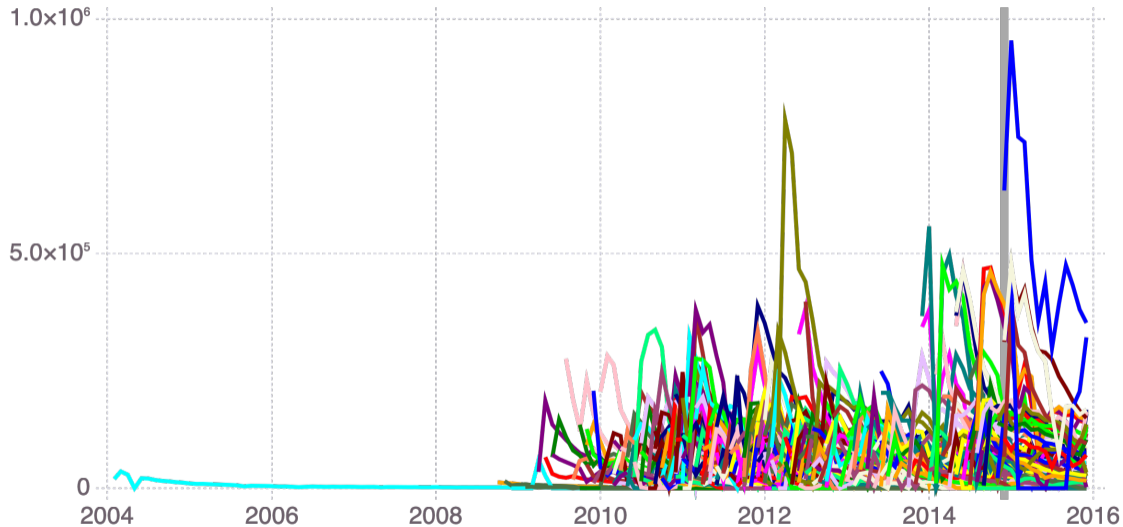


# Eagle Ford Shale Play: Monthly production volumes [MCF] of 327 gas wells

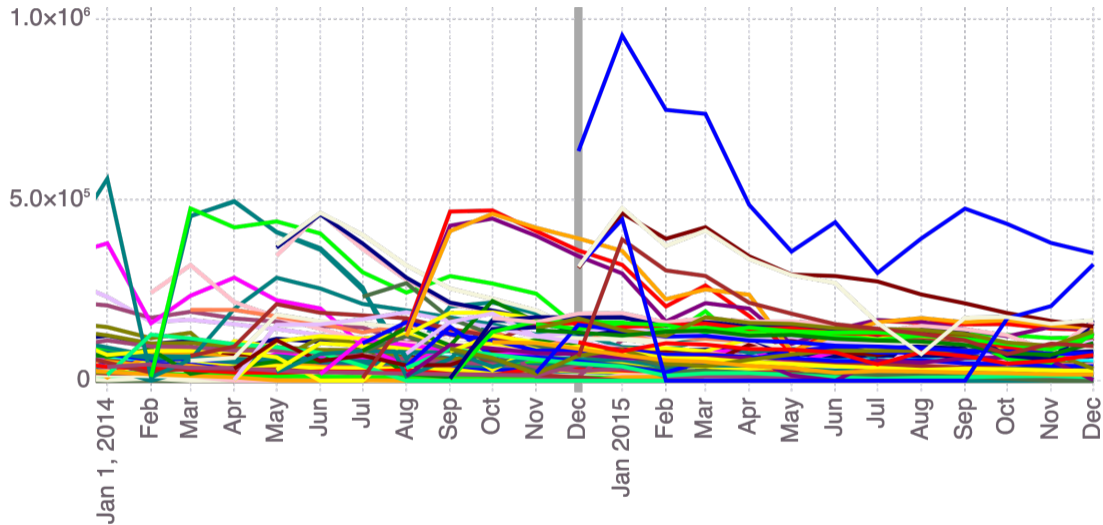


- ▶ Use all the data up to a given cutoff date (e.g. 2015)
- ▶ Apply ML to learn behavior of the “known“ well transients
  - Identify and group wells which behave similarly (having similar production transients)
  - Discover the optimal number of **master decline curves** required to reproduce the observed transients
  - **master decline curves** represent production **features/signatures**
- ▶ Apply ML to predict **blindly** the unknown production transients beyond the cutoff
- ▶ Prediction is obtained by discovering to which type (group) the wells producing beyond the cutoff belong
- ▶ i.e., discovering what combinations of the **master decline curves** can represent the wells producing beyond the cutoff
- ▶ Unsupervised ML analyses performed using **SmartTensors (NMFk/NTFk)**

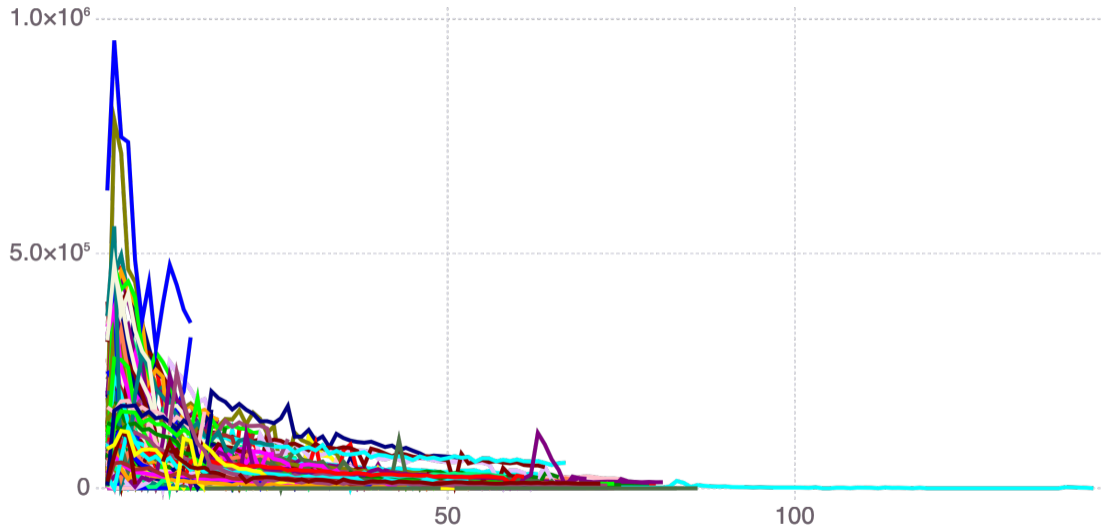
# Eagle Ford Shale Play: Monthly production volumes [MCF] of 327 gas wells



# Eagle Ford Shale Play: Monthly production volumes [MCF] of 327 gas wells

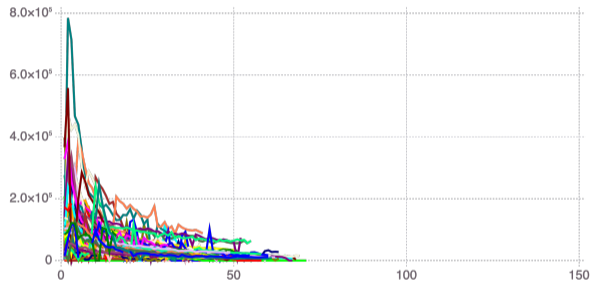
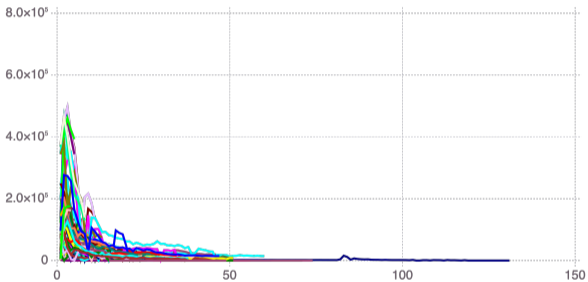


# Eagle Ford Shale Play: Monthly production volumes [MCF] of 327 gas wells

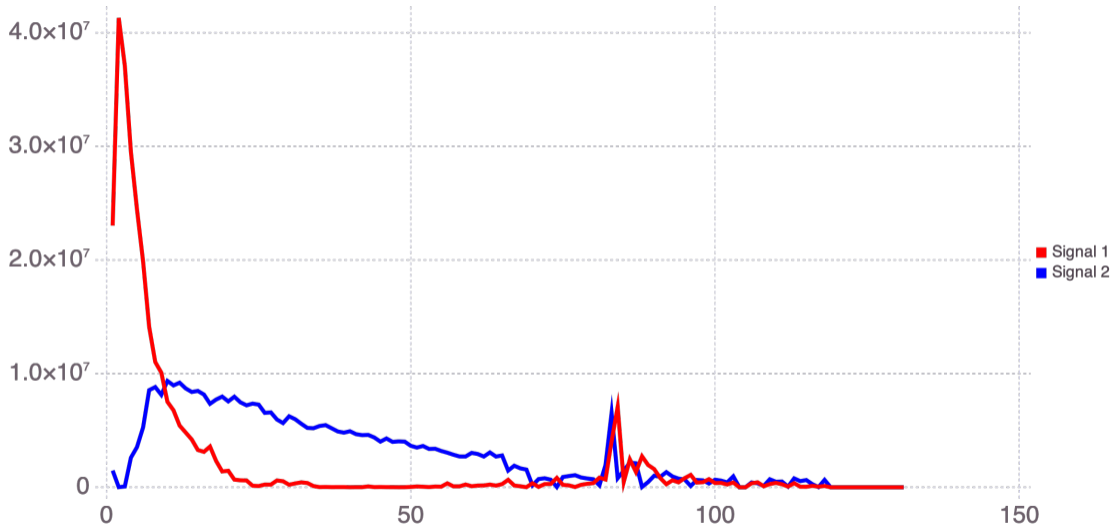


## 'Fast' declining (135)

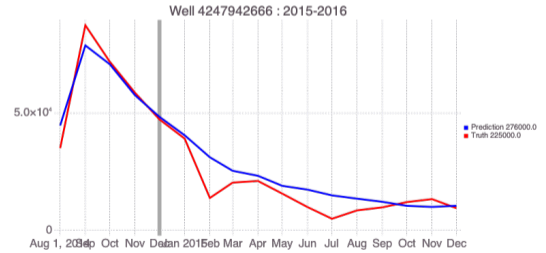
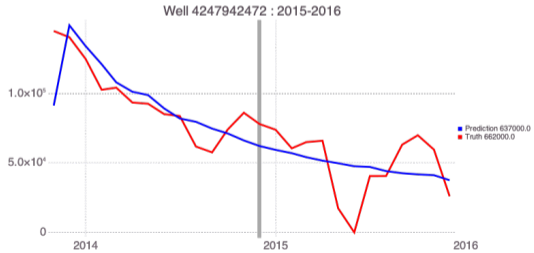
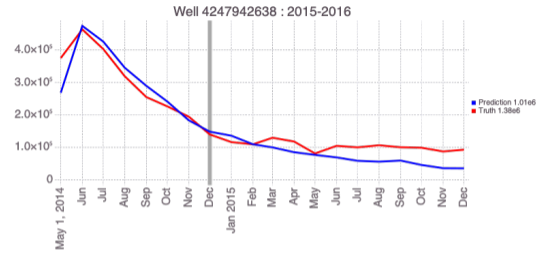
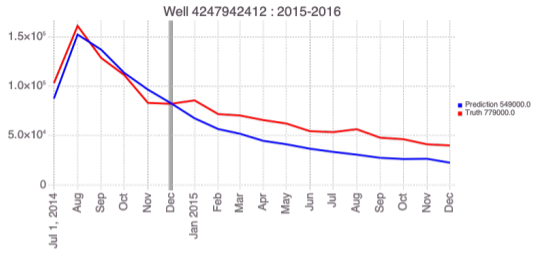
## 'Slow' declining (192)



# Eagle Ford Shale Play: Master Decline Curves [MCF over months]

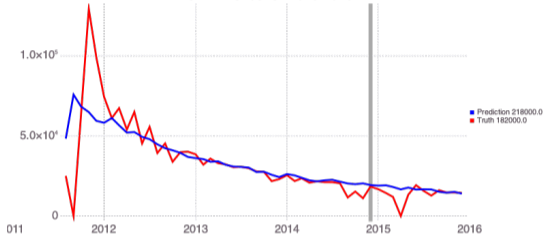


# Eagle Ford Shale Play: Blind predictions beyond 2015

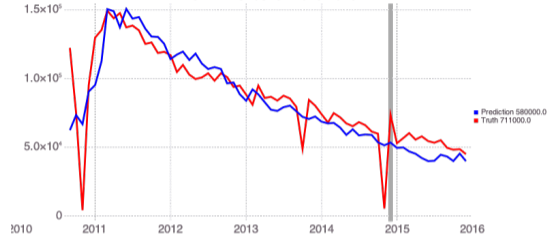


# Eagle Ford Shale Play: Blind predictions beyond 2015

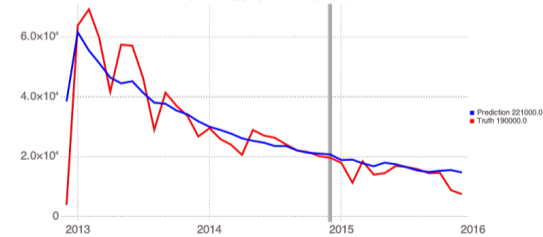
Well 4247940978 : 2015-2016



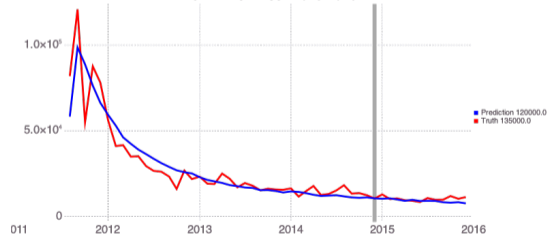
Well 4247940815 : 2015-2016



Well 4212332547 : 2015-2016

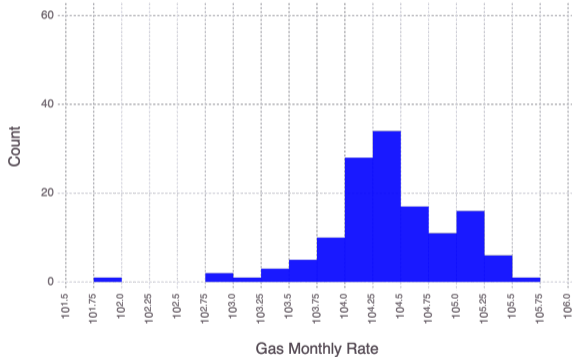


Well 4247941283 : 2015-2016

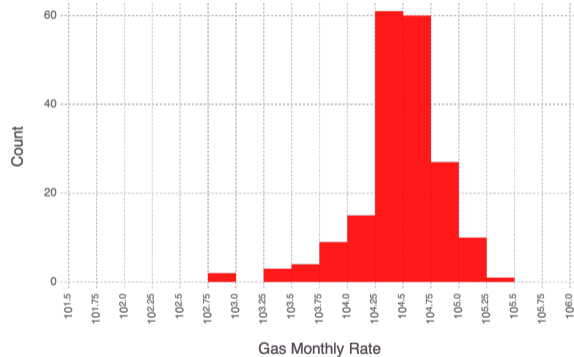


## Monthly rate histograms

### 'Fast' declining (135)

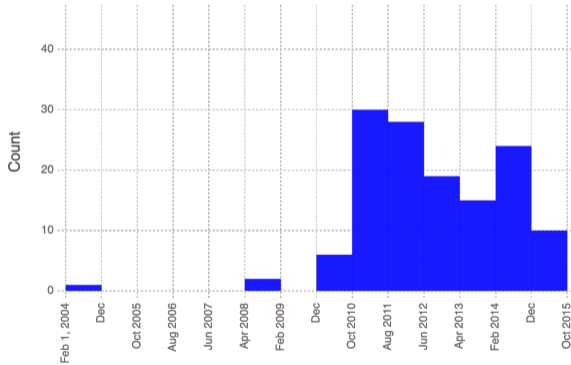


### 'Slow' declining (192)

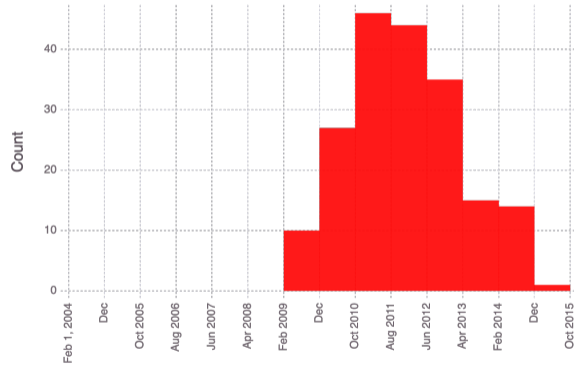


## Drilling date histograms

### 'Fast' declining (135)

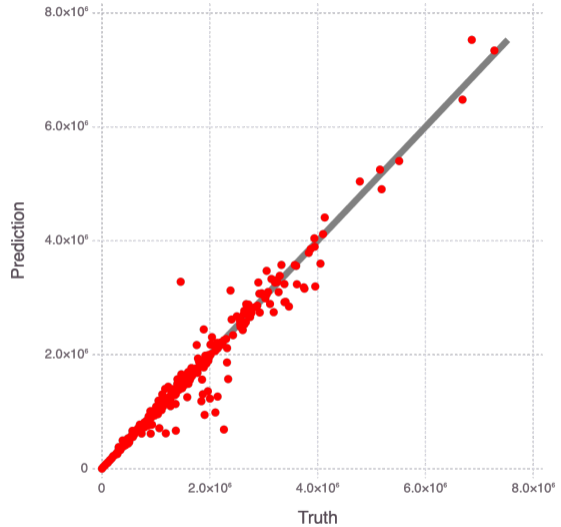


### 'Slow' declining (192)

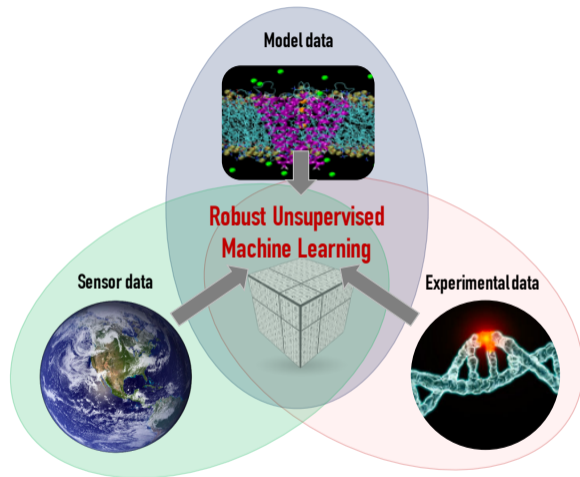


- ▶ Other well attributes also differ between the 2 groups
- ▶ For example:
  - Operators
  - Proppant mass
  - Injected fluid volumes
  - ...

- ▶ 300 wells continue producing beyond 2015
- ▶  $r^2 = 0.96$

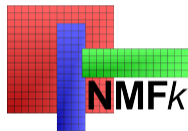


- ▶ Developed **novel** unsupervised and physics-informed ML methods and computational tools
- ▶ Some of our tools have been recently patented
- ▶ Our ML methods have been used to solve various real-world problems (brought breakthrough discoveries related to human cancer research)
- ▶ Several ongoing projects (DOE, ARPA-E, ...)



## ► Codes:

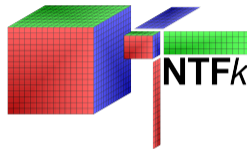
NMF<sub>k</sub>



MADS



NTF<sub>k</sub>



## ► Examples:

[http://madsjulia.github.io/Mads.jl/Examples/blind\\_source\\_separation](http://madsjulia.github.io/Mads.jl/Examples/blind_source_separation)

<http://tensors.lanl.gov>

<http://tensordecompositions.github.io>

<https://github.com/TensorDecompositions>

<https://hub.docker.com/u/montyvesselinov>



- ▶ Vesselinov, Munuduru, Karra, O'Malley, Alexandrov, Unsupervised Machine Learning Based on Non-Negative Tensor Factorization for Analyzing Reactive-Mixing, **Journal of Computational Physics**, Special issue: Machine Learning, 2019.
- ▶ Stanev, Vesselinov, Kusne, Antoszewski, Takeuchi, Alexandrov, Unsupervised Phase Mapping of X-ray Diffraction Data by Nonnegative Matrix Factorization Integrated with Custom Clustering, **Nature Computational Materials**, 2018.
- ▶ Vesselinov, O'Malley, Alexandrov, Nonnegative Tensor Factorization for Contaminant Source Identification, **Journal of Contaminant Hydrology**, 2018.
- ▶ O'Malley, Vesselinov, Alexandrov, Alexandrov, Nonnegative/binary matrix factorization with a D-Wave quantum annealer, **PLOS ONE**, 2018.
- ▶ Vesselinov, O'Malley, Alexandrov, Contaminant source identification using semi-supervised machine learning, **Journal of Contaminant Hydrology**, 2017.
- ▶ Alexandrov, Vesselinov, Blind source separation for groundwater level analysis based on nonnegative matrix factorization, **WRR**, 2014.