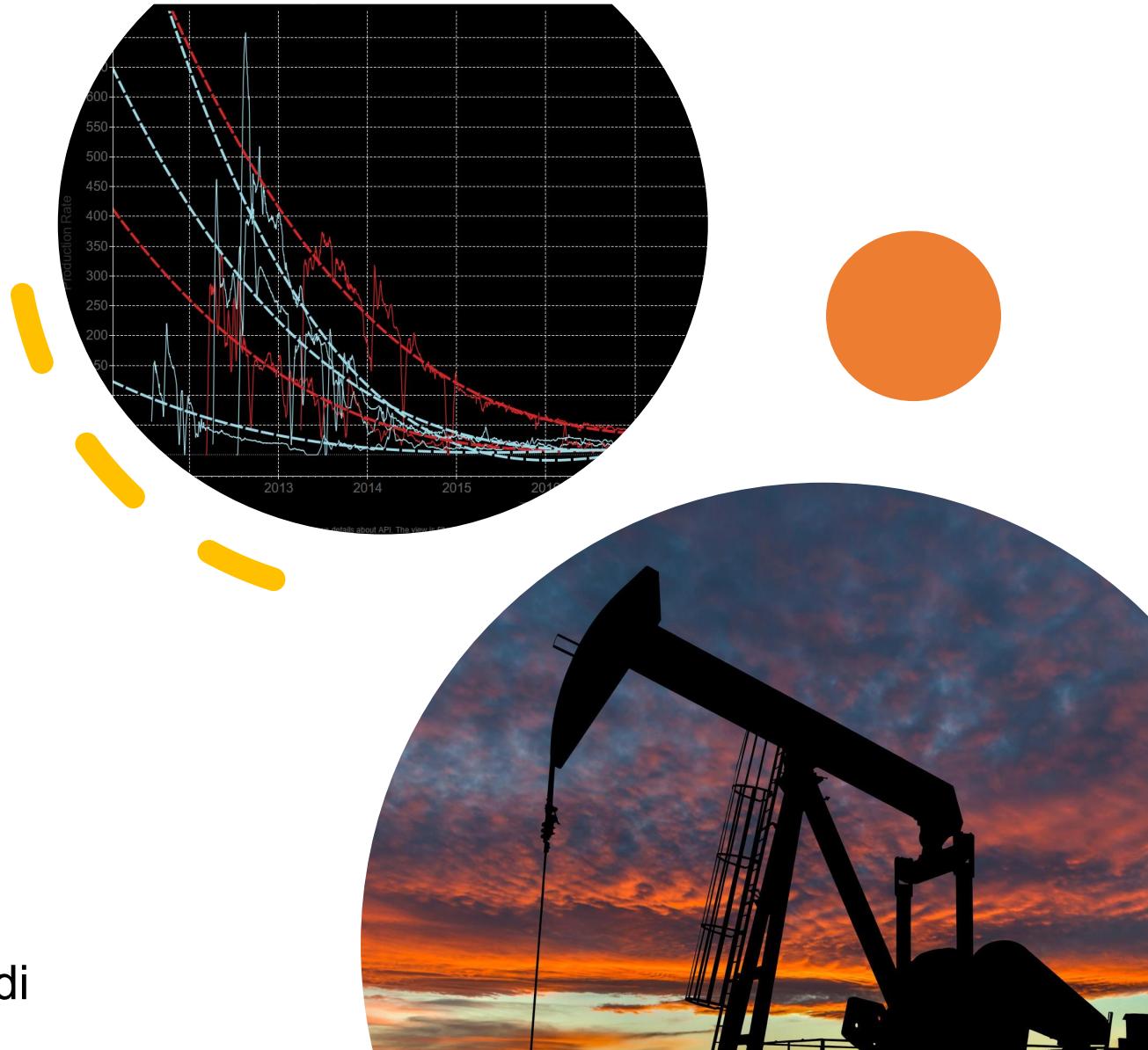


Combining Machine Learning and Empirical Engineering Methods Towards Improving Oil Production Forecasting

Andrew Allen

Thesis Defense
June 19th, 2020

Thesis Committee Chair: Dr. Roy Jafari-Marandi



CAL POLY
College of Engineering

Outline



Introduction and Significance



Background



Literature Review



Case Study



Conclusions



Questions

About Student

- *Academic Background*
 - B.S. and M.S Industrial Engineering
 - Graduate courses focused on data science and analytics
- *Career Interests*
 - Business Intelligence and Data Analytics
 - Energy, sustainability, and continuous improvement
- *Research Motivation*
 - *Learn by doing* approach to solve real-engineering problem
- *Research Goals*
 - Build off existing skills and prepare for entering workforce



Introduction and Significance

An **unconventional resource** is a fossil fuel energy extracted from **unconventional reservoirs**

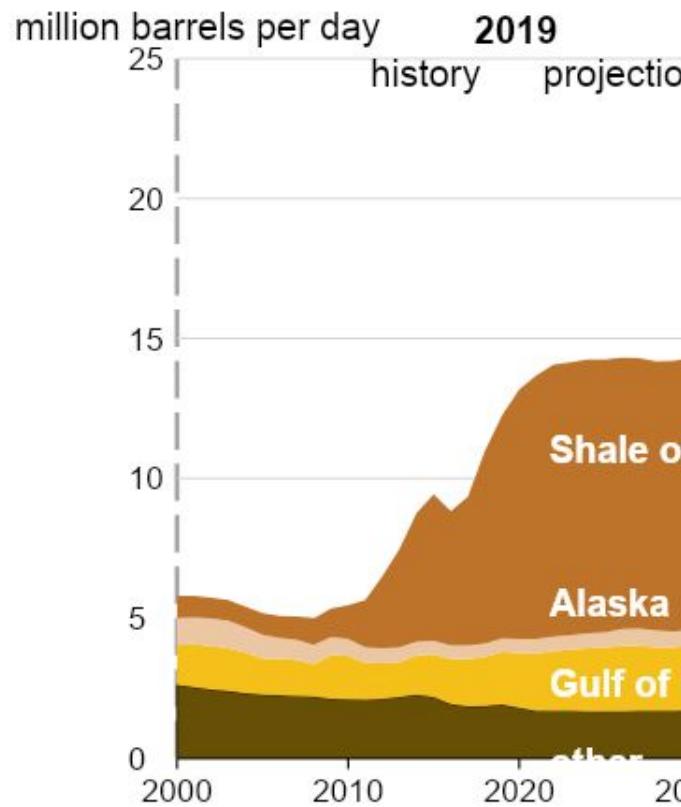
Technology advancements have enabled **extraction** of **shale gas and oil**

Hydraulic fracturing and horizontal drilling has made extracting resources from shale reservoirs possible



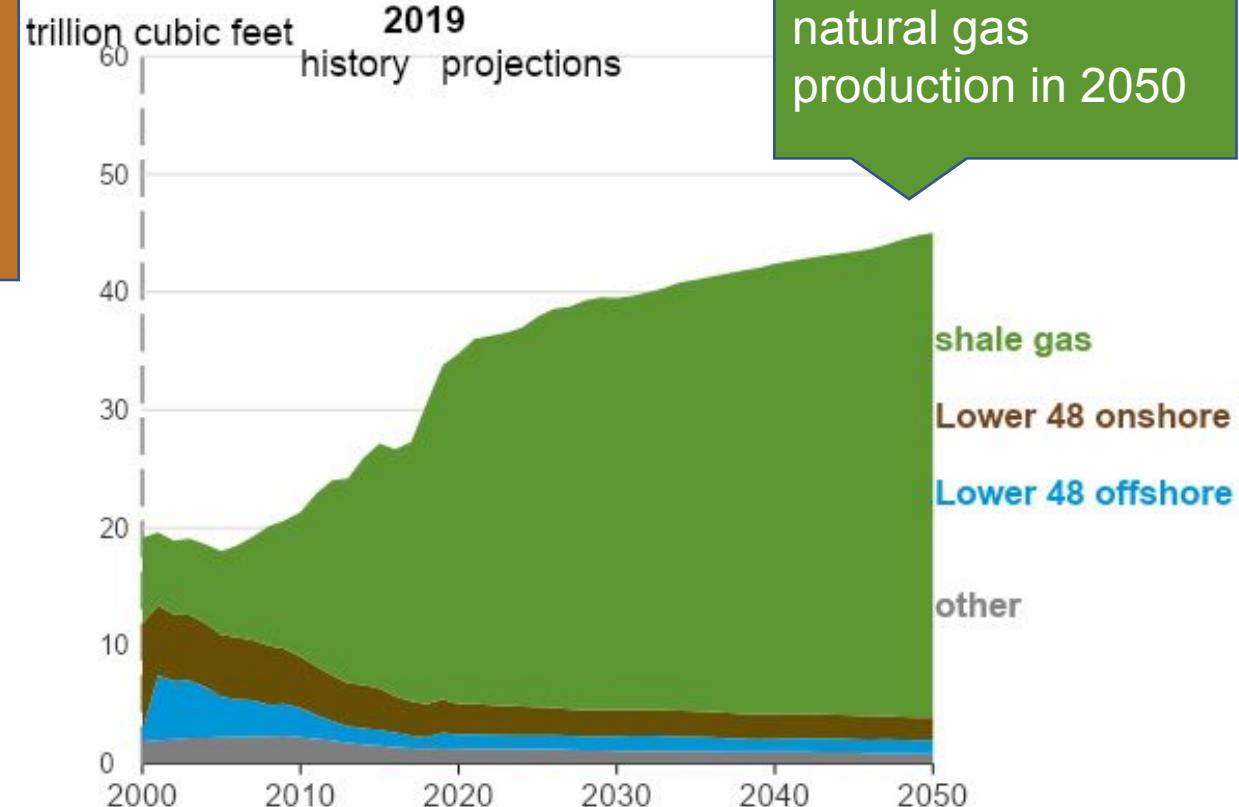
The United States continues to product historically high levels of...

Crude Oil



Shale oil reaches
73% of U.S.
Crude Oil
production in
2050

Natural gas

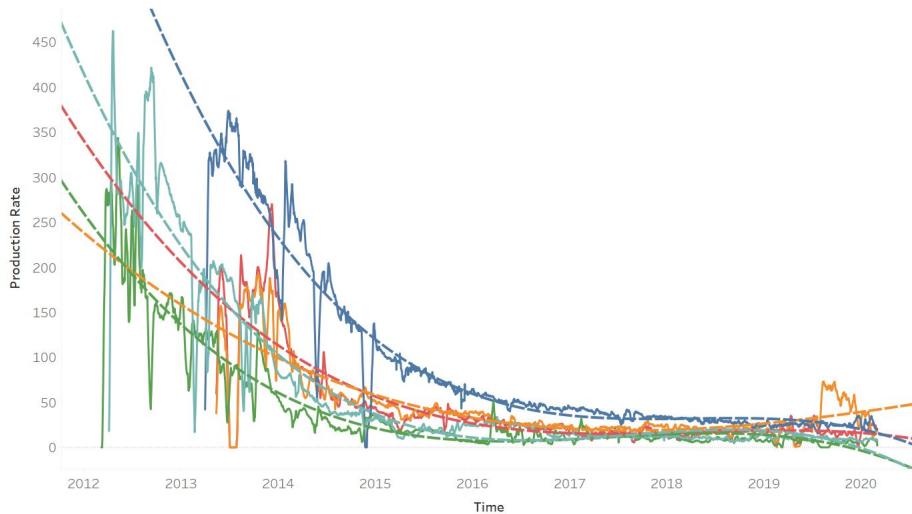


Shale gas reaches
91% of domestic
natural gas
production in 2050

Production Forecasting: Current Industry Methods

(1) Decline Curve Analysis (DCA)

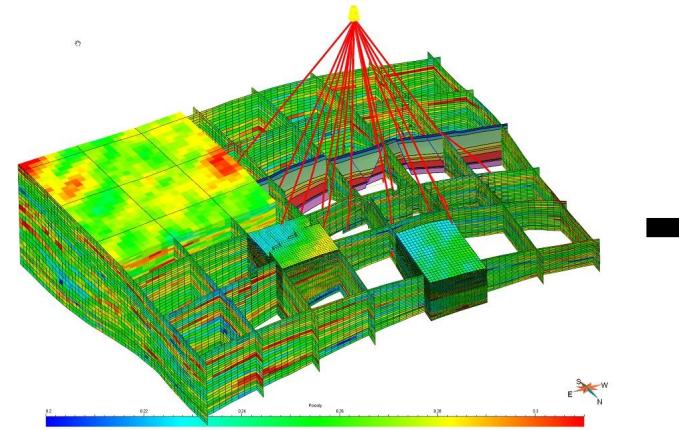
- “Empirical Method”
- Forecasting via graphical curve-fit procedure of historical data



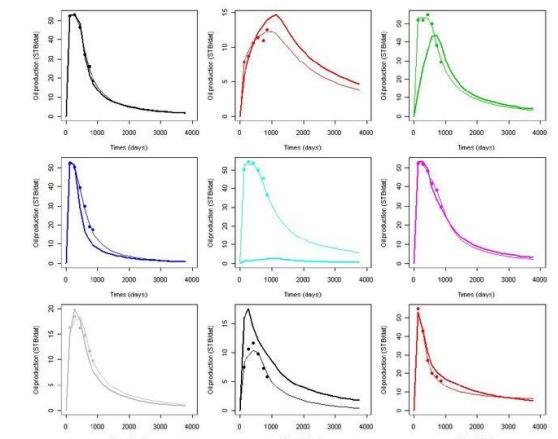
Fit equation through actual field data

(2) Reservoir Simulation

- “Analytical Method”
- Forecasting by simulating production data based on 3D modeling and simulation



3D Reservoir Model



Simulate Production

Problem Overview

Decline Curve Analysis (DCA)

Model Development

Assumptions often misrepresent reality

Forecasting

Tendency to overestimate

Need rich historical production data in given area

Applicable for time-rate data only

Reservoir Simulation

Model Development

Parameter uncertainty now in multiple dimensions

Forecasting

Inconsistent results

Complete measurements and expert interpretation

Computationally expensive

Thesis Problem Statement

Problem

There is need for accurate production forecasting techniques

Solution

Can combining current techniques with machine learning techniques help improve production forecasting?

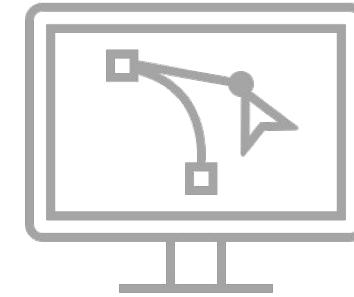
Who benefits?

Enterprises
Systems
People
Planet

What is Machine Learning?

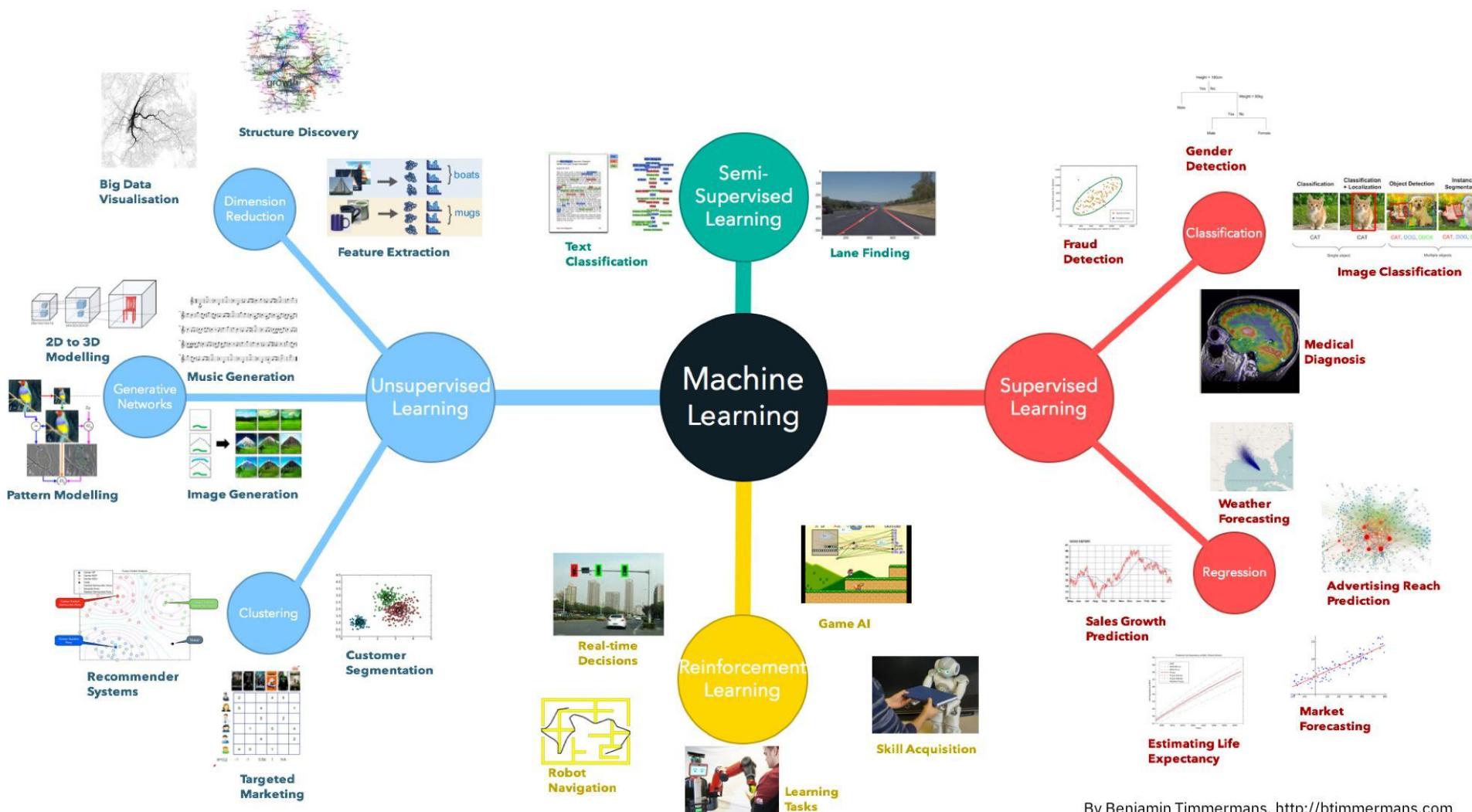


Human Learning: Improving at doing a task with experience



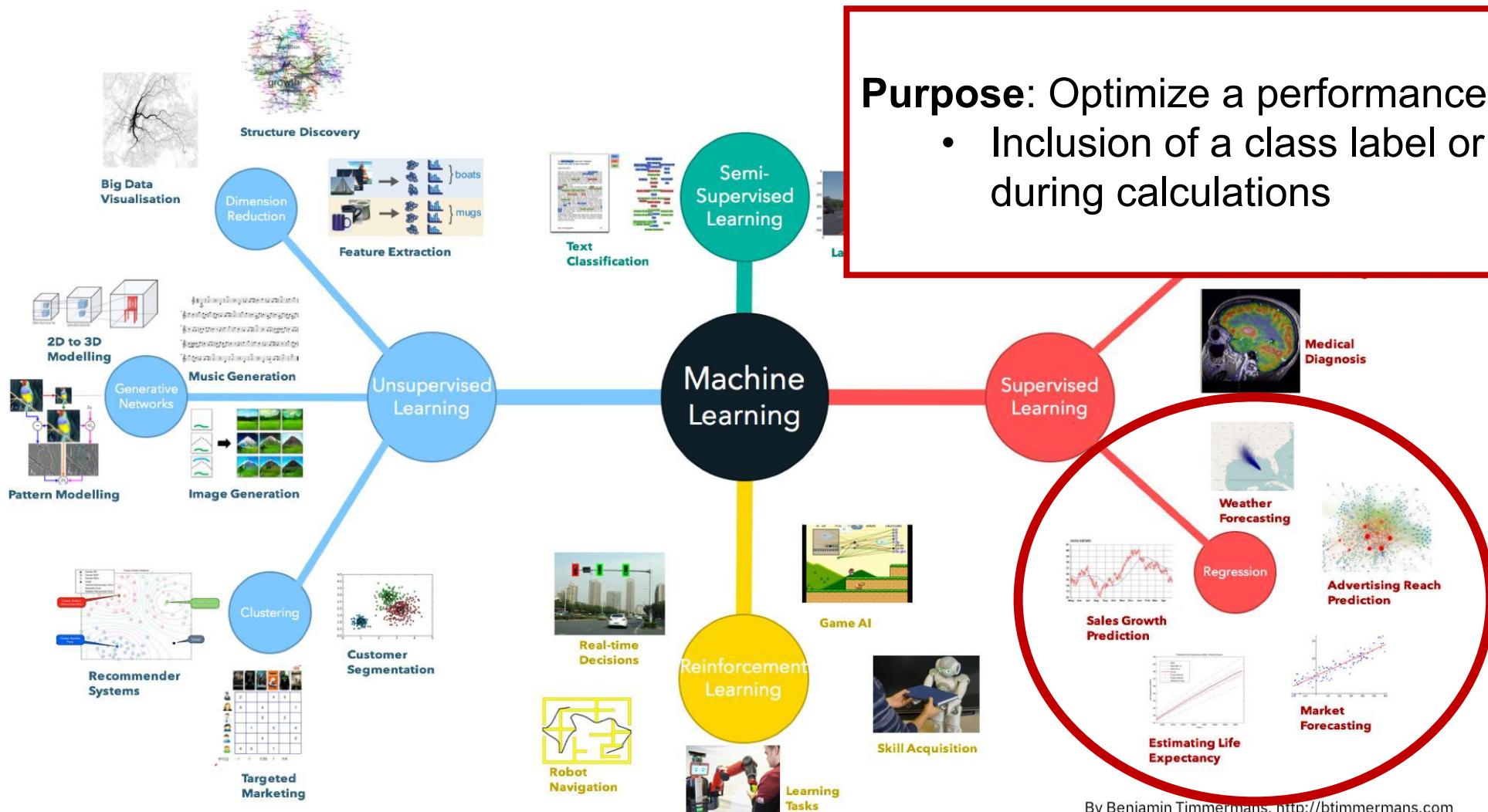
Machine Learning: Estimating parameters of a model so it best fits the dataset

What is Machine Learning?

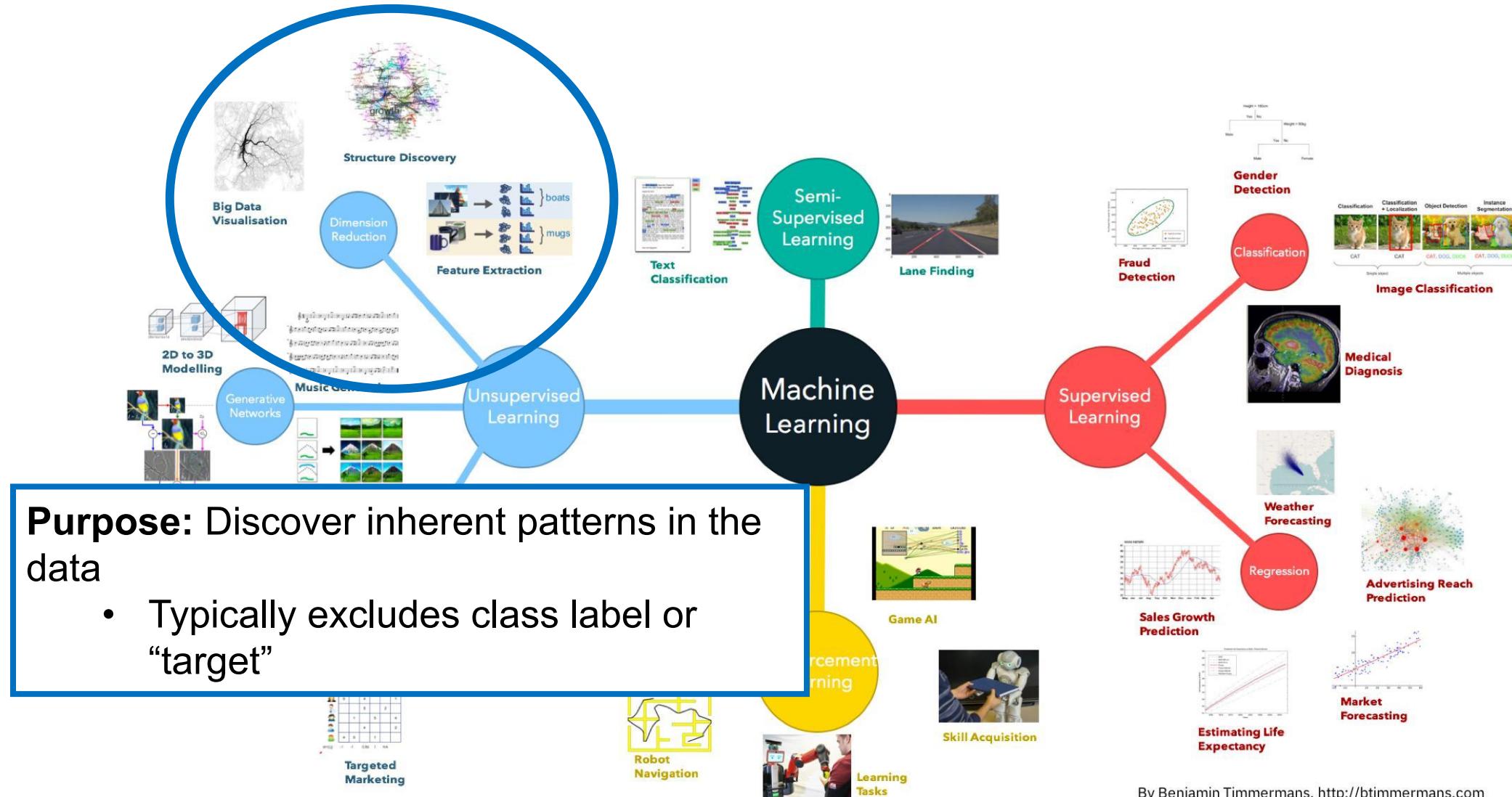


By Benjamin Timmermans. <http://btimmermans.com>

What is Machine Learning?



What is Machine Learning?



Supervised Learning: Linear Regression

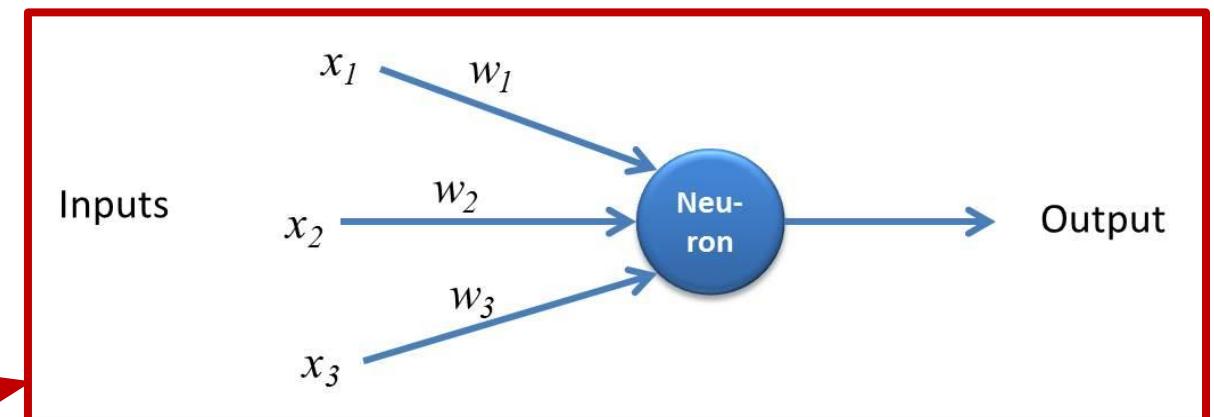
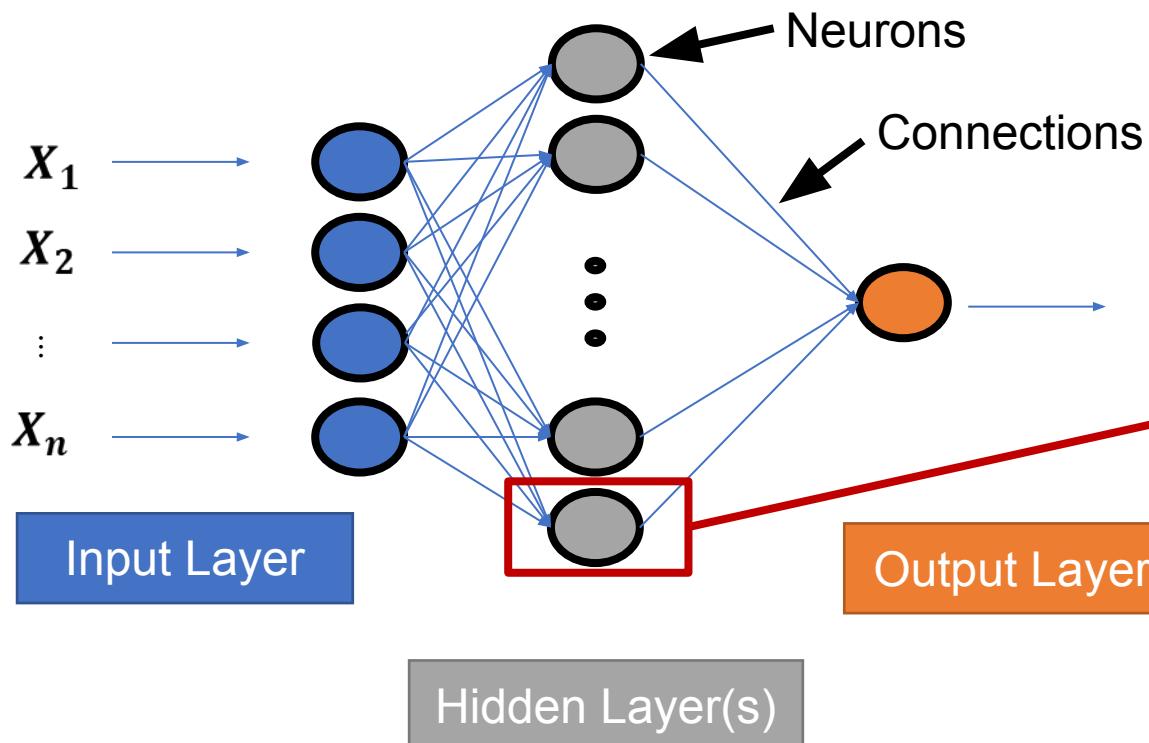
Linear Regression (LR)

- Fit set of predictor variables \mathbf{X} to a target variable \mathbf{Y} using linear predictor functions
- Estimate β s using the data
- Multivariate approach to production forecasting

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n$$

Artificial Neural Networks (ANN)

- Experimental branch of science to simulate human learning
 - Multilayer Perceptron (MLP) most famous and widely used ANN
 - Nonlinear approach to production forecasting



$$\mathbf{y} = \left(\sum_{i=1}^n w_i x_i + b \right) = \varphi(\mathbf{w}^T \mathbf{x} + b)$$

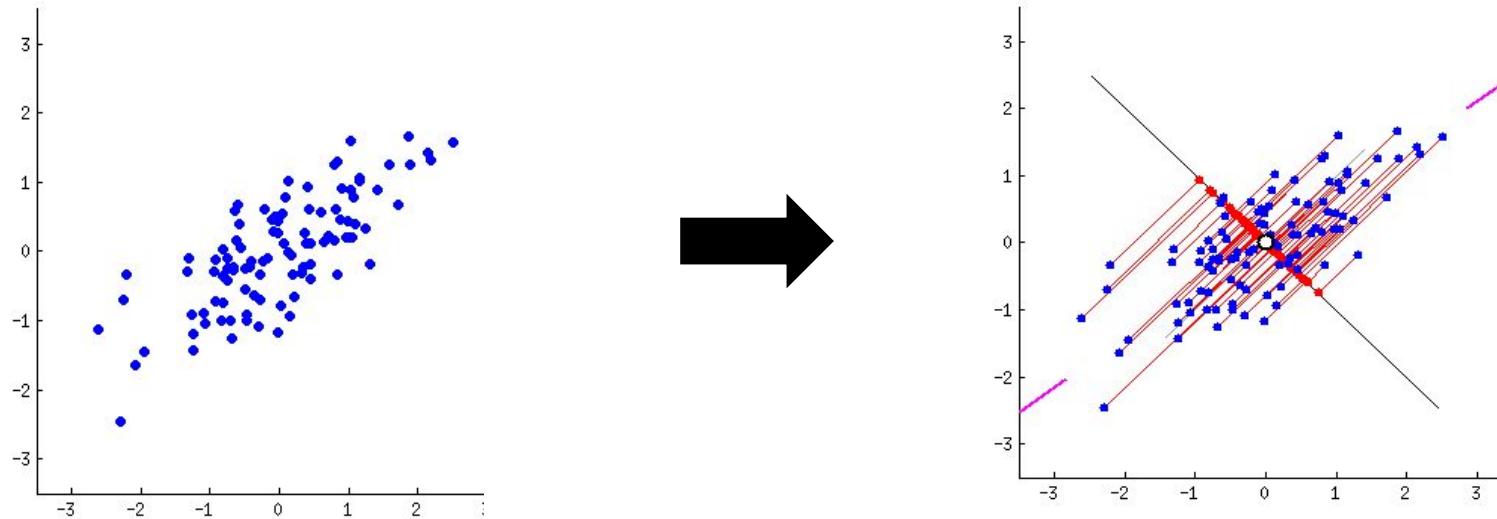
MLP Activation Function

Unsupervised Learning via Dimension Reduction

- Represent the data with fewer dimensions
 - Retain the most statistical variance possible
 - More tractable subset of original data, easier for learning
 - Discover intrinsic patterns in the data
- *Feature selection* = creating a **subset** of the original feature set
- *Feature extraction* = creating a **new set** of features from the original feature set.

Feature Extraction via Principal Components Analysis (PCA)

- PCA is dimension reduction technique that maps original data onto a *linear subspace*, such that the retained statistical variance of the projected data is maximized

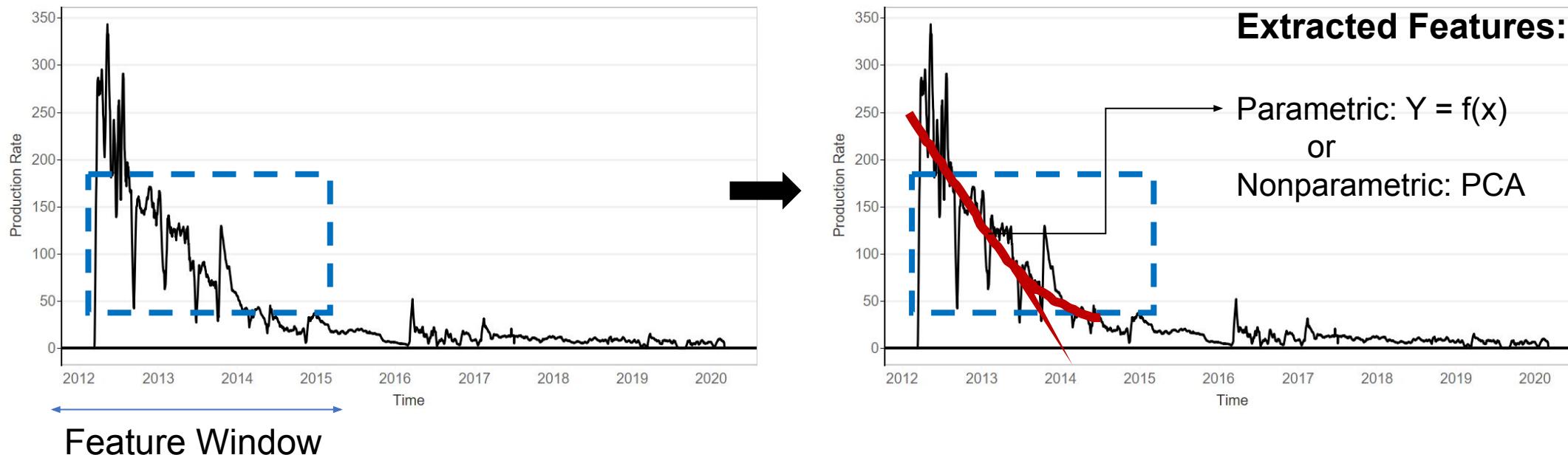


- The principal components (PCs) represent a **new subset** of the data

Feature Extraction of Time-Series

Decompose time-series into components that represent some pattern of the series

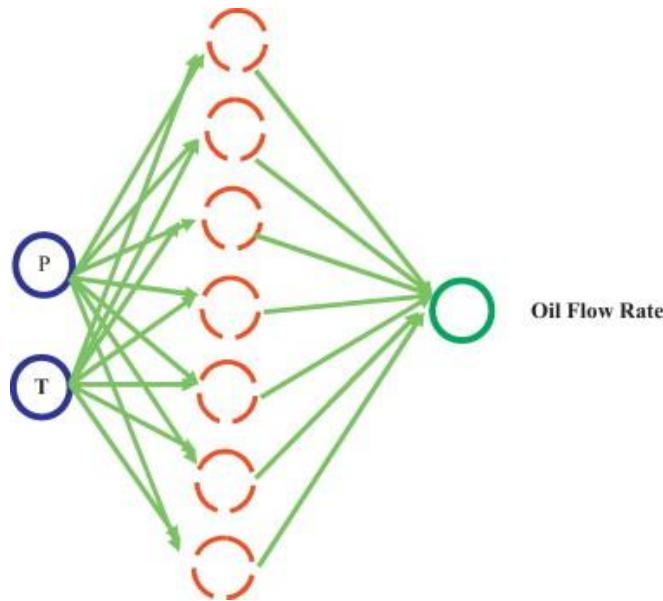
- (1) understand time series patterns (2) forecast new values
- Parametrically via curve-fitting
- Nonparametrically via PCA



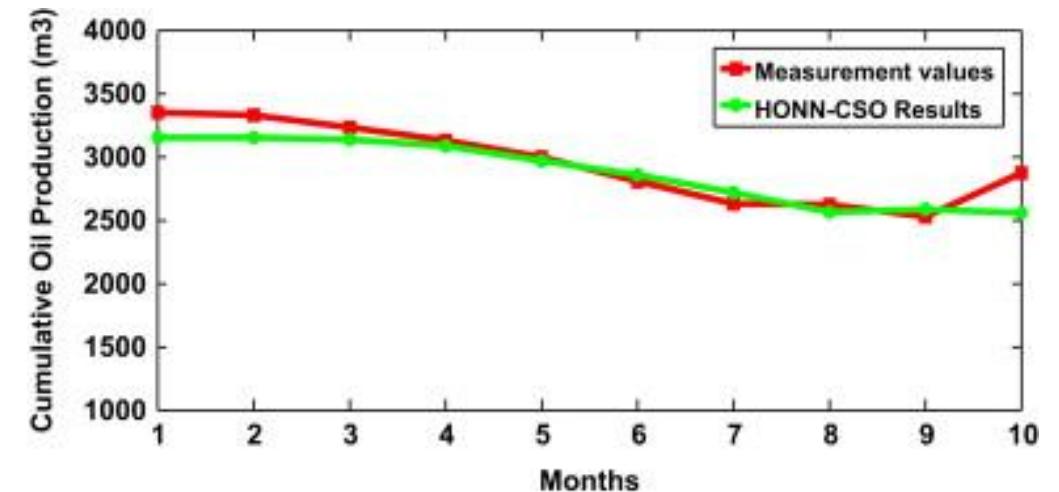
Literature Review

1

Literature Review: ANN in Oil/Gas Industry



Ahmadi and Ebadi (2017) used pressure (P) and temperature (T) data as inputs to MLP to predict oil flow rates



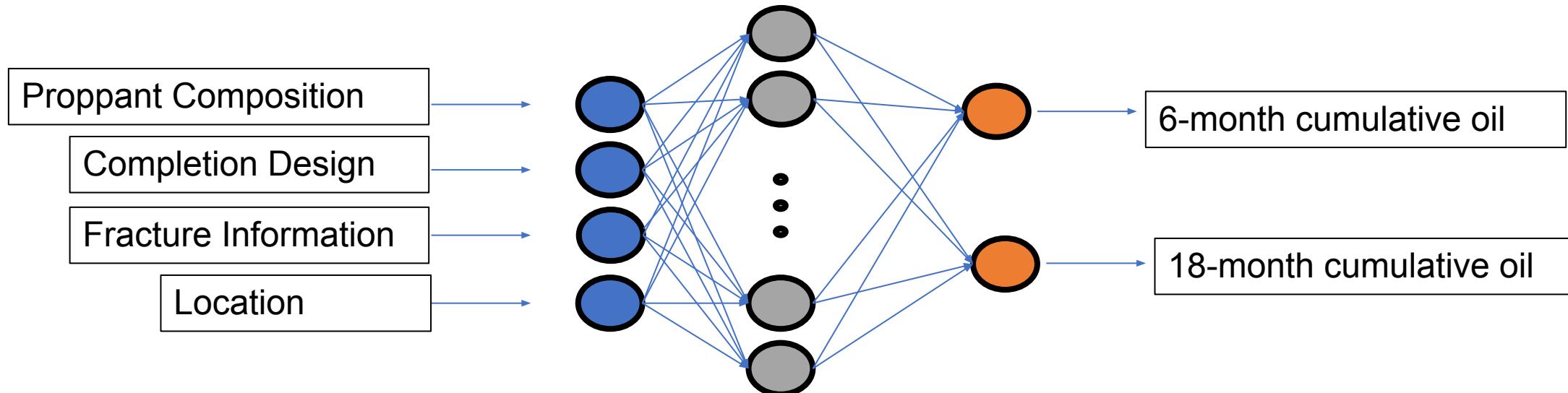
Chakra and Song (2013) applied an ANN (HONN-CSO) using correlation matrix of oil, gas, and water data as inputs to predict 10-month cumulative oil production

- Production yielded < 4.0% error
- Did not incorporate static well parameters

Literature Review: ANN in Oil/Gas Industry

Wang and Chen (2019)

- Applied MLP to predict cumulative oil production in Bakken Shales
- Test set R^2 71% and 72% for 6- and 18-month cumulative oil

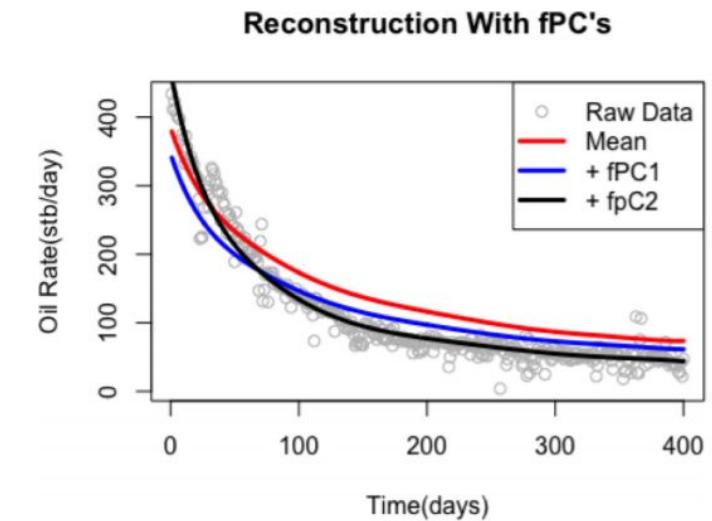
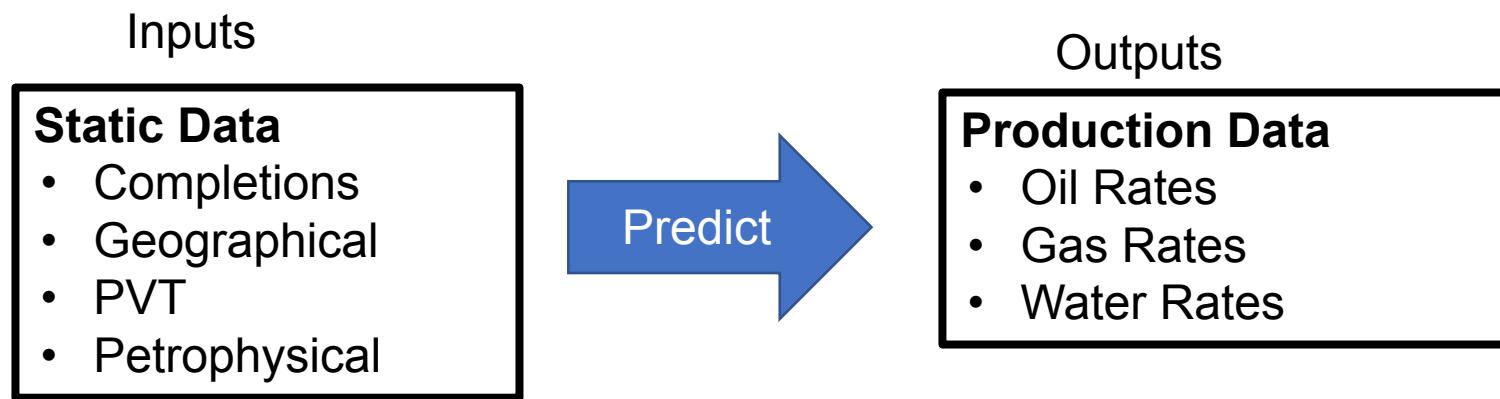


Luo and Tian (2019)

- Compared LR and MLP to predict first 6-months cumulative oil
- Measured depth, TVD, fluid, and proppant to be most important variables

Literature Review: PCA in Oil/Gas Industry

- Reconstructing production profiles with FPCA (Bhattacharya and Nikolaou, 2013)
 - Sum products of PCs and associated scores and add each sum to the mean function computed on entire ensemble



Literature Review: PCA in Oil/Gas Industry

- Makinde and Lee (2019)
 - Feature extraction .5 to 3 years of production data for forecasting 30 years time-horizon

Profile Reconstruction using .5 to 3 years of data
(Bhattacharya and Nikolaou, 2013)



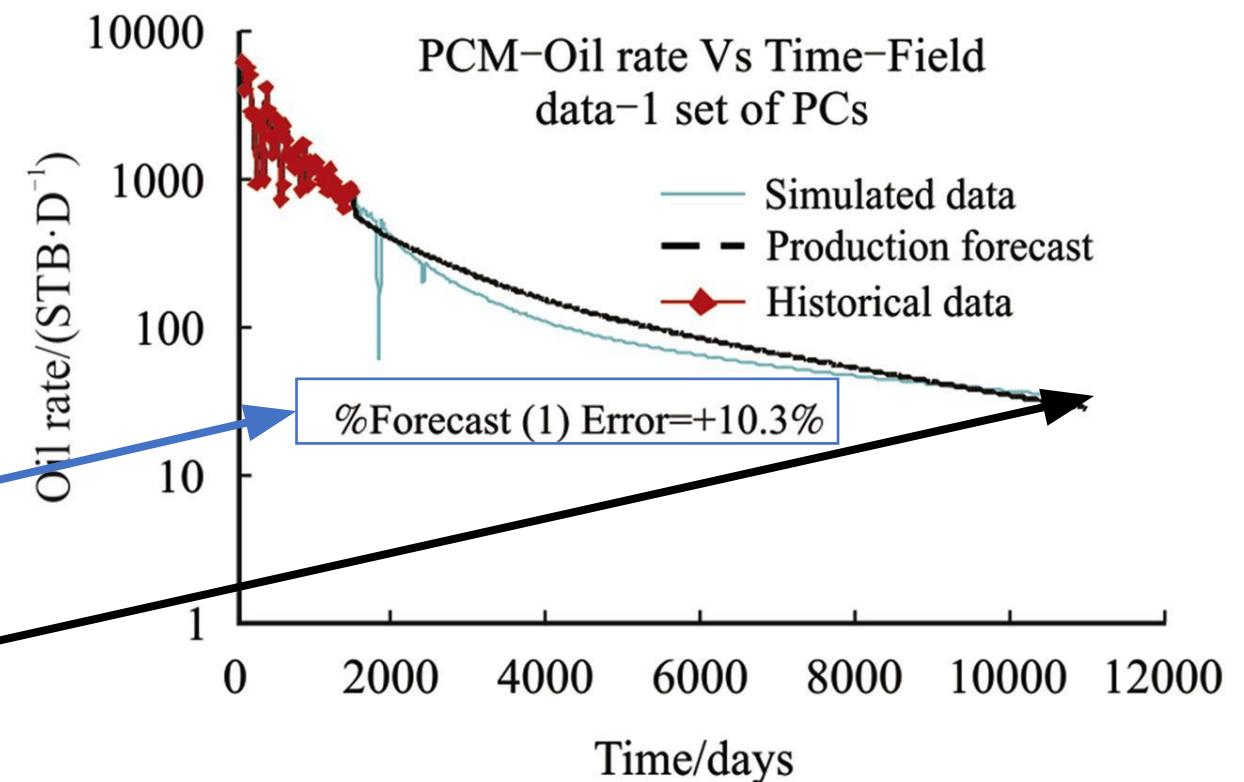
Evaluate PCA reconstruction



Forecast 30 years



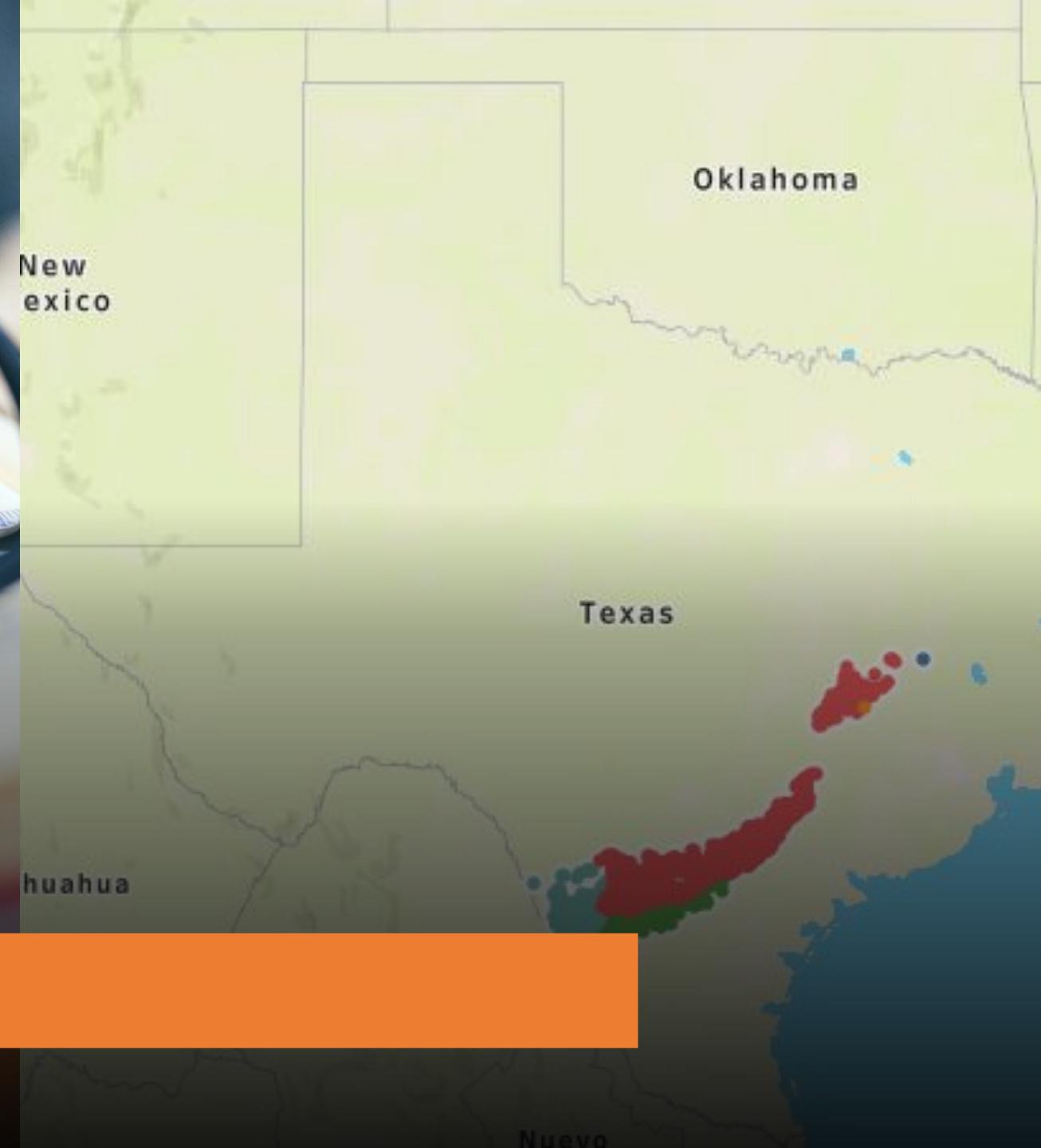
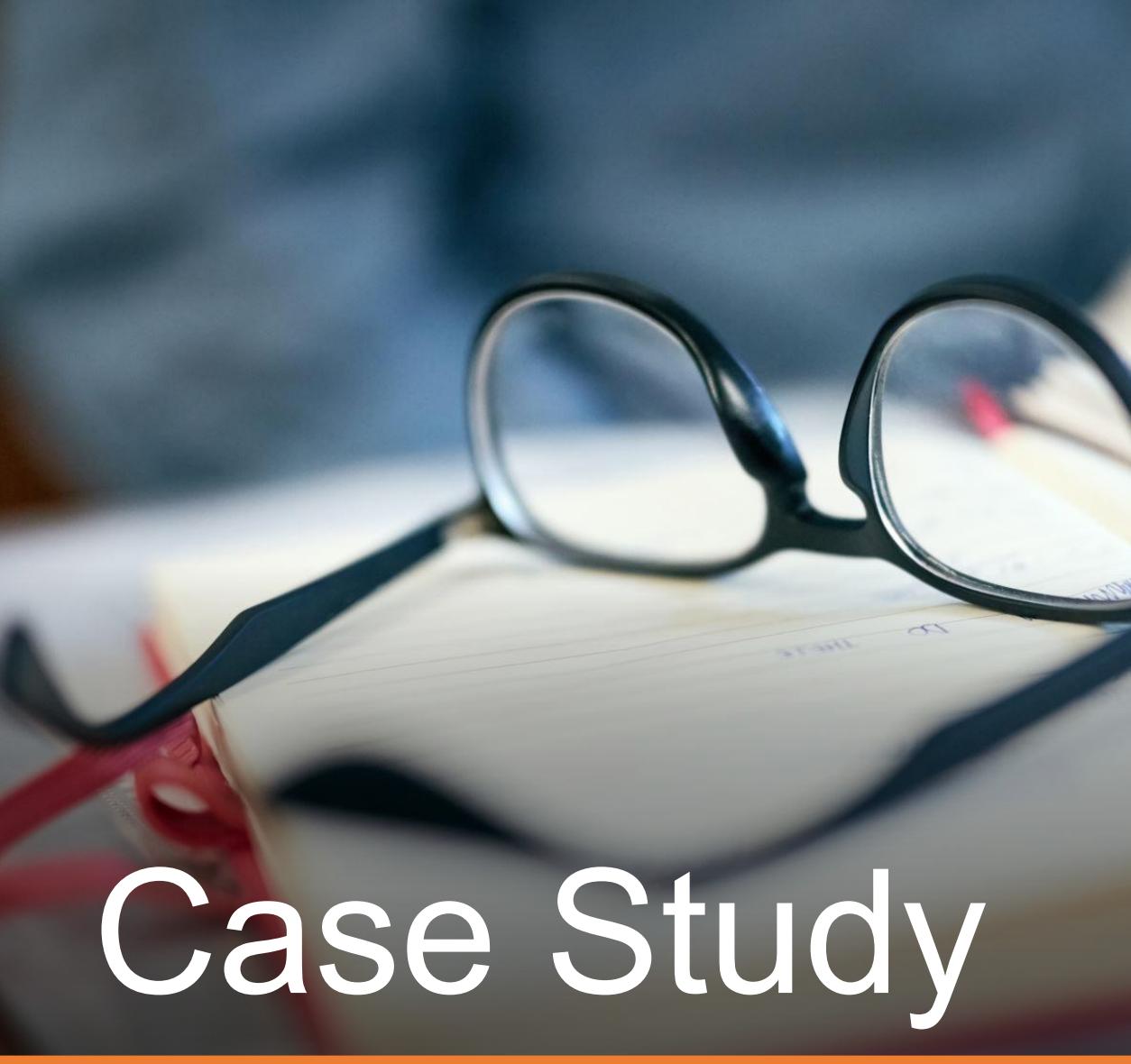
Compare to simulated data



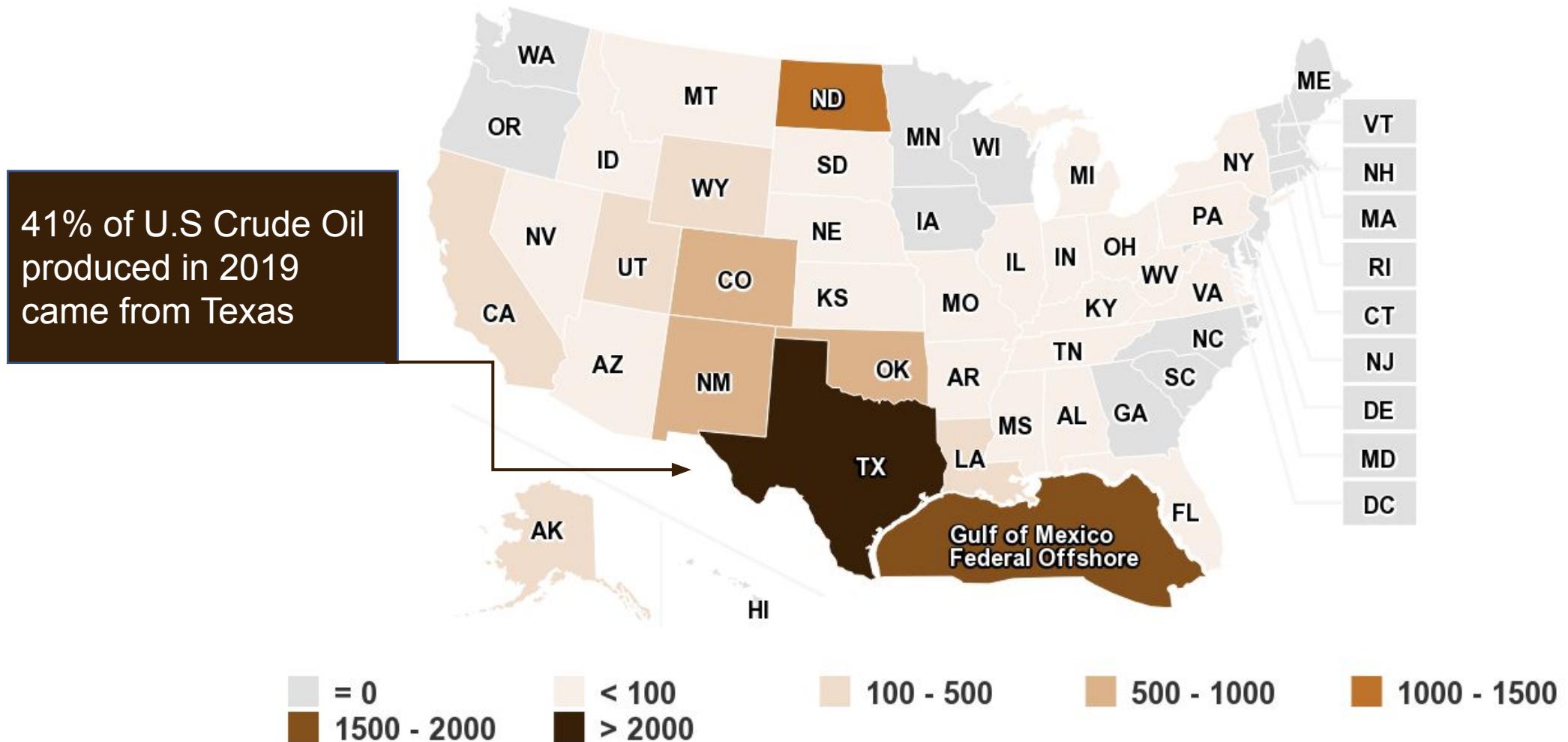
Study Goal and Objectives

- Study Goal
 - *Create multivariate well-proxy model that incorporates both static and dynamic variables to predict cumulative oil production*
- Main Purposes
 - 1) Determine if regression-based predictions using **parametric** and **nonparametric** feature extraction could predict cumulative oil
 - 2) Determine efficacy of including static well variables into modeling
 - 3) Determine if model complexity aids prediction results

Case Study



Texas Oil

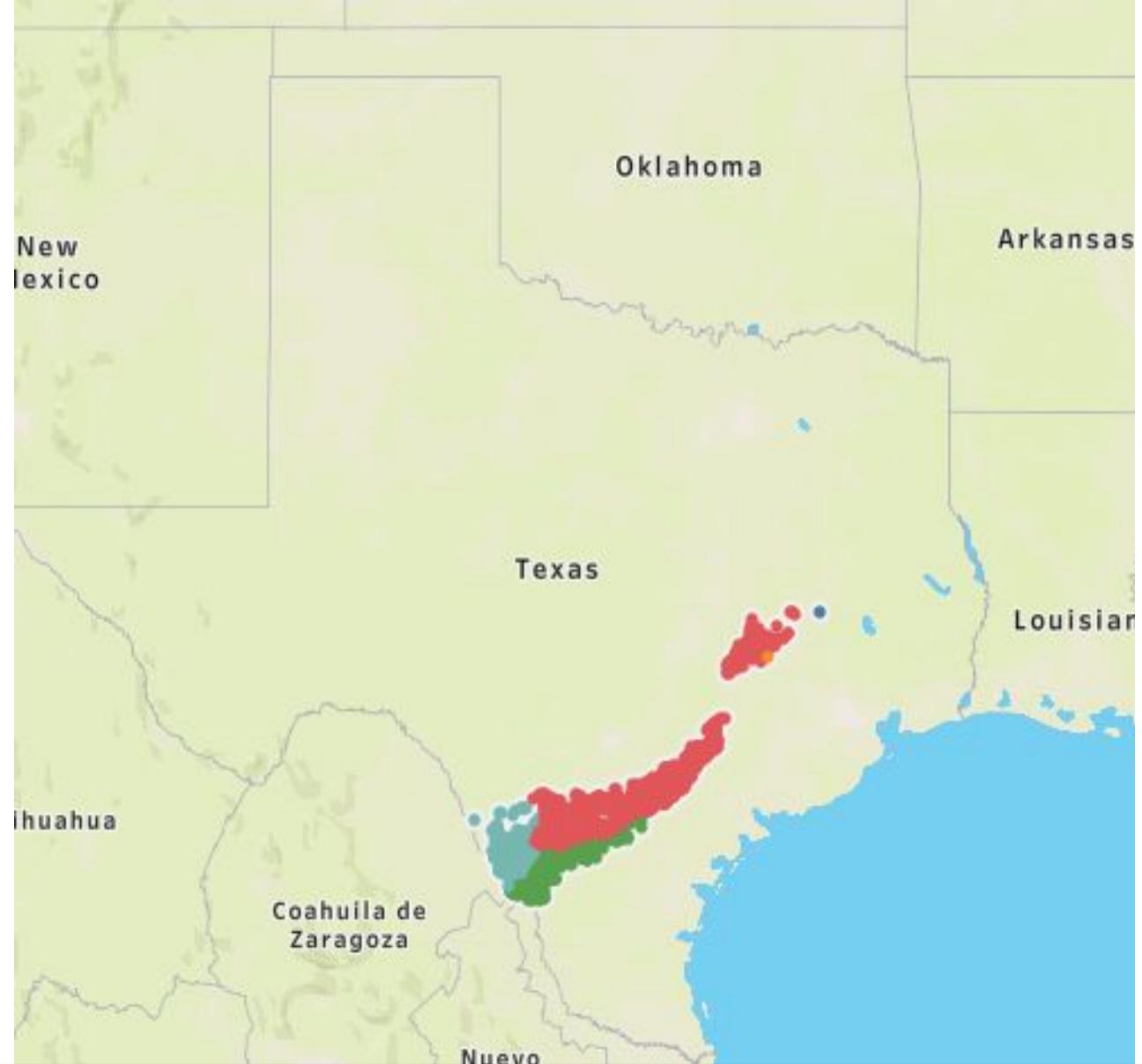


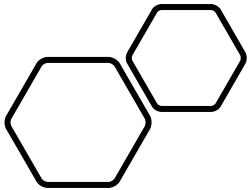
Note: Crude oil includes lease condensate.

Source: U.S. Energy Information Administration, *Petroleum Supply Monthly*, February 2020, preliminary data

The Eagle Ford Shale

- \$60 billion impact since initial development in 2006
- Production Forecasting Challenge
 - Limited production history
 - Unique areal variation





Case Study Overview

Data from 448 oil wells in EFS

- *Dynamic* production Data
 - Time-series oil, gas, water production
- *Static* well Information
 - Well design measurements, locational, PVT data.

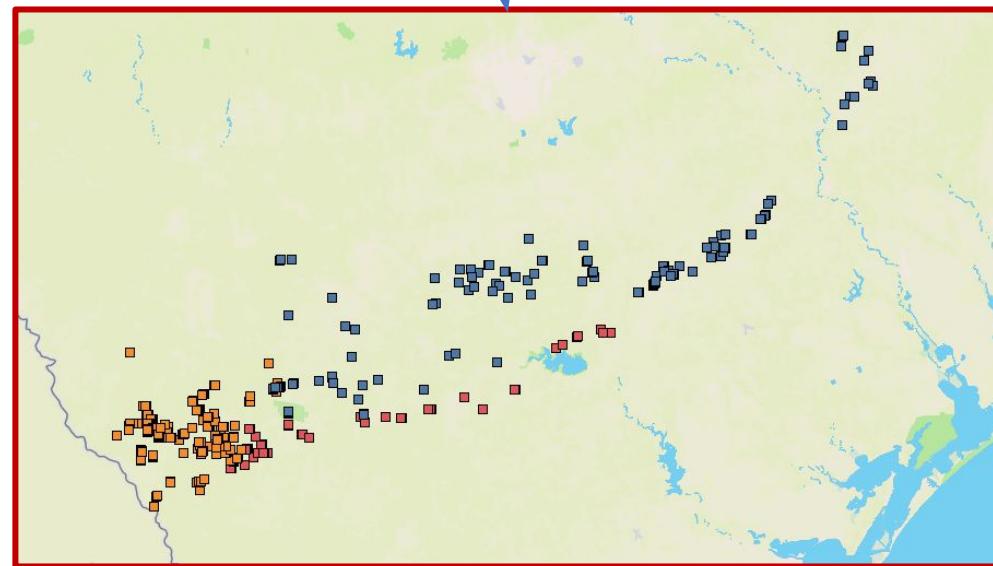


Scenarios for input data

- (1) 6-months
- (2) 12-months
- (3) 24-months

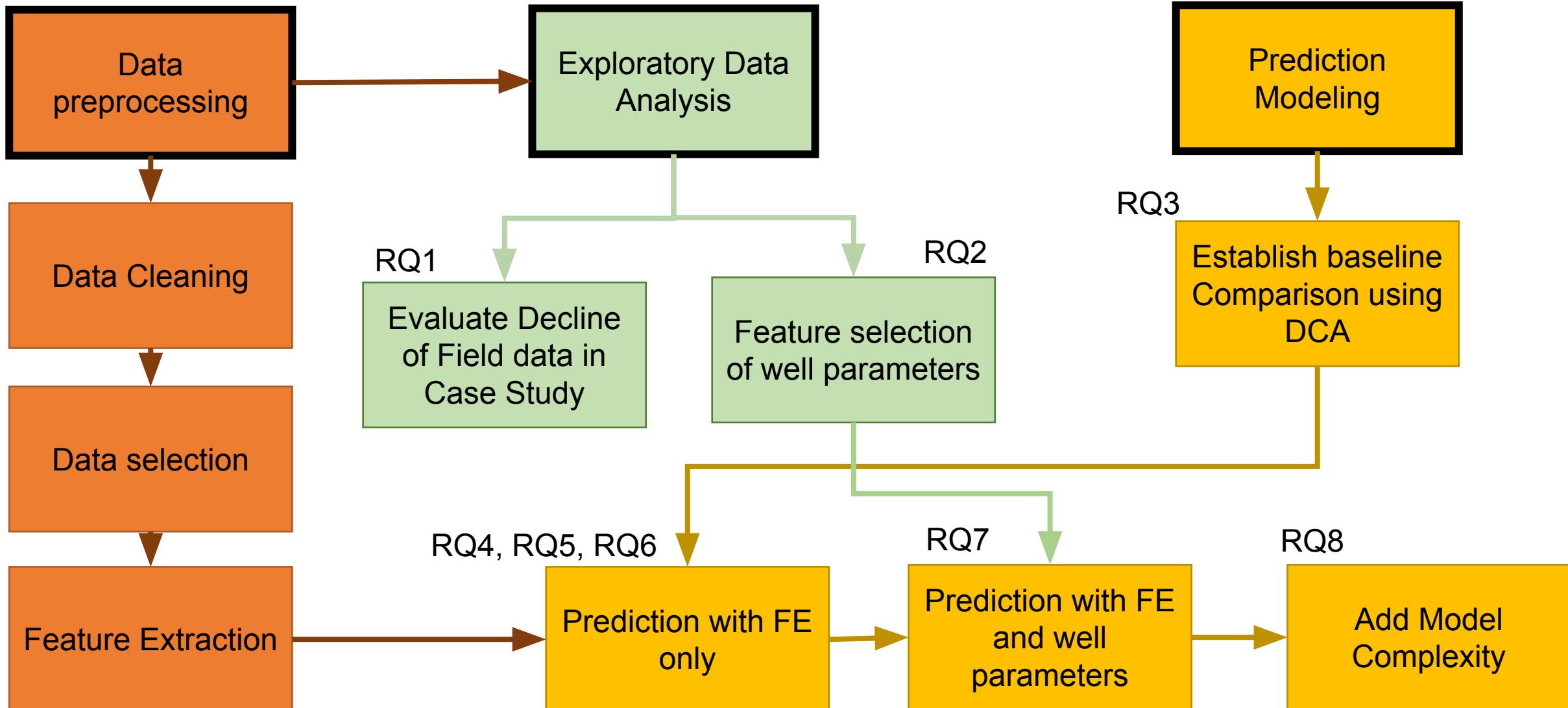
Target variable of interest

- three-year oil cumulative oil

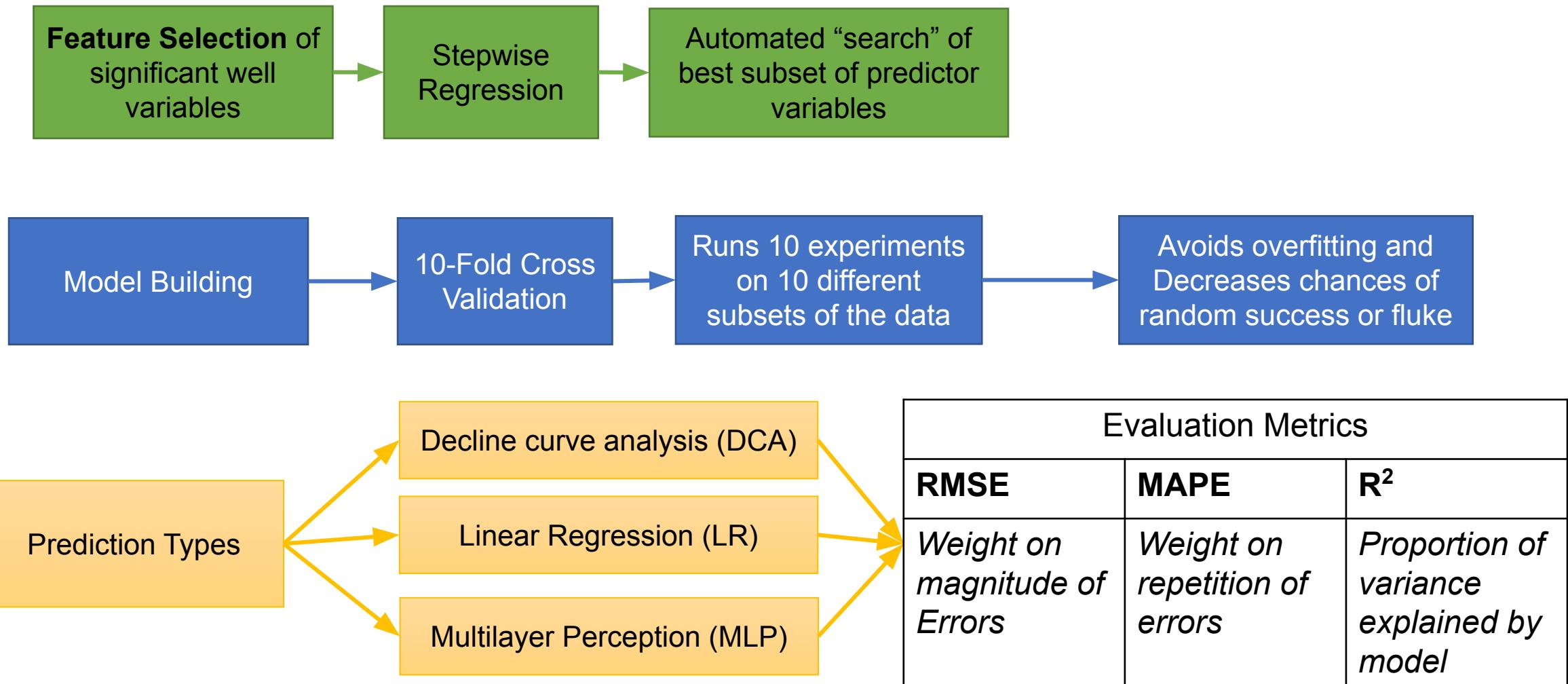


Research Framework

*FE = Feature Extraction
*DCA = decline curve analysis

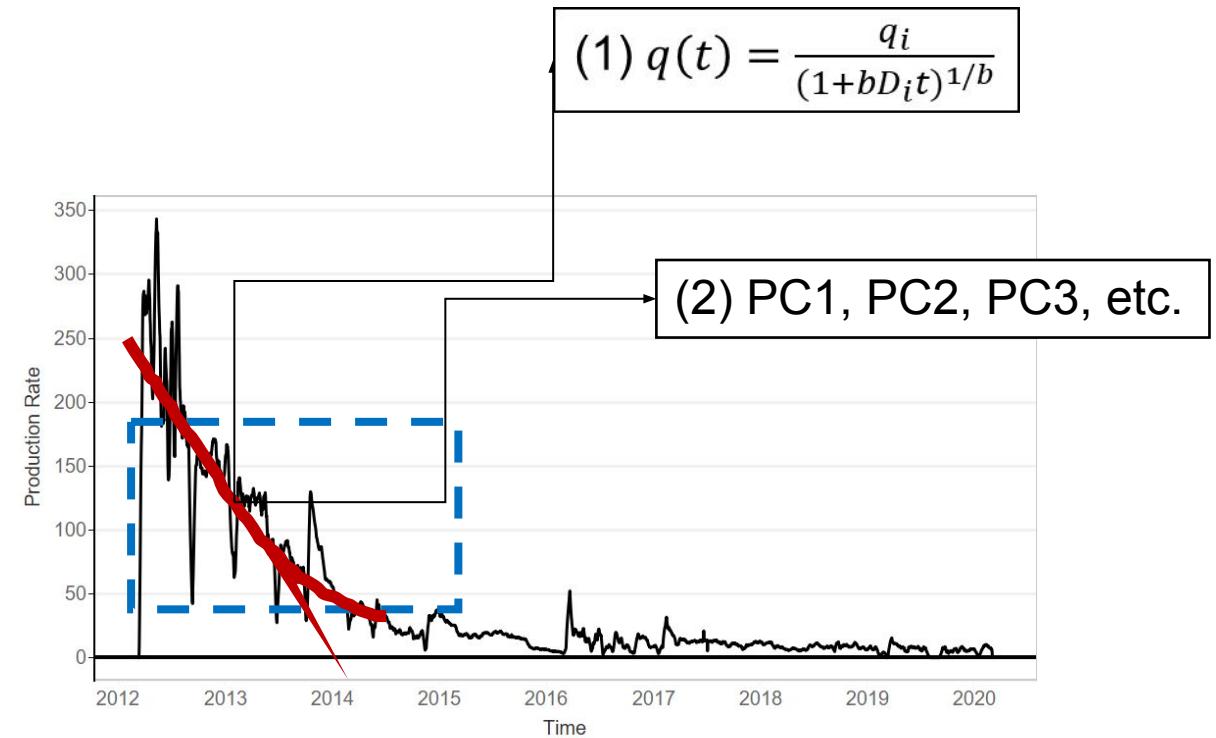
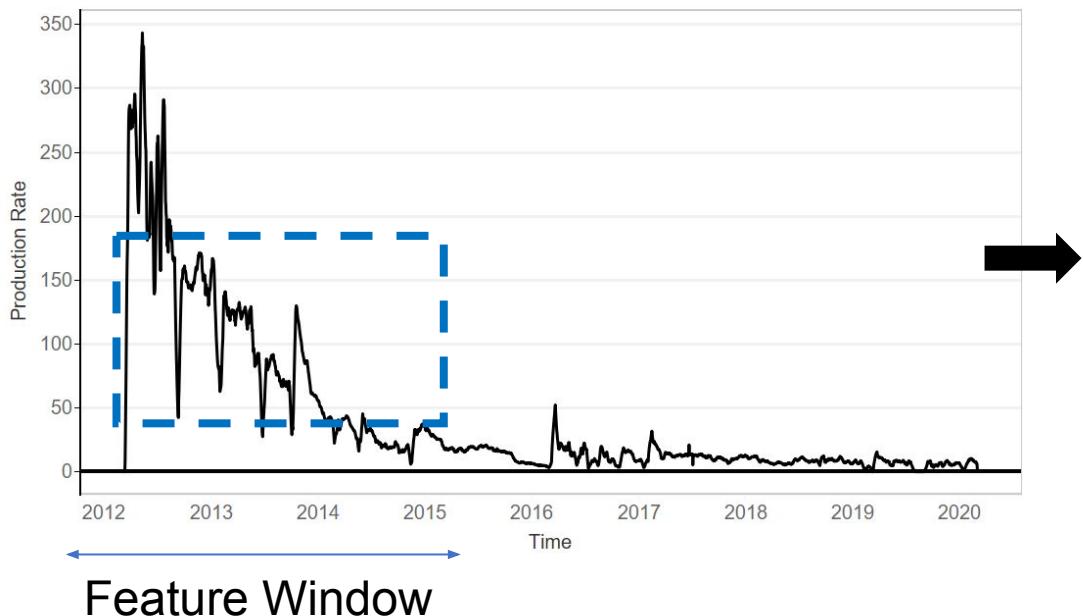


Design of Experiments Overview

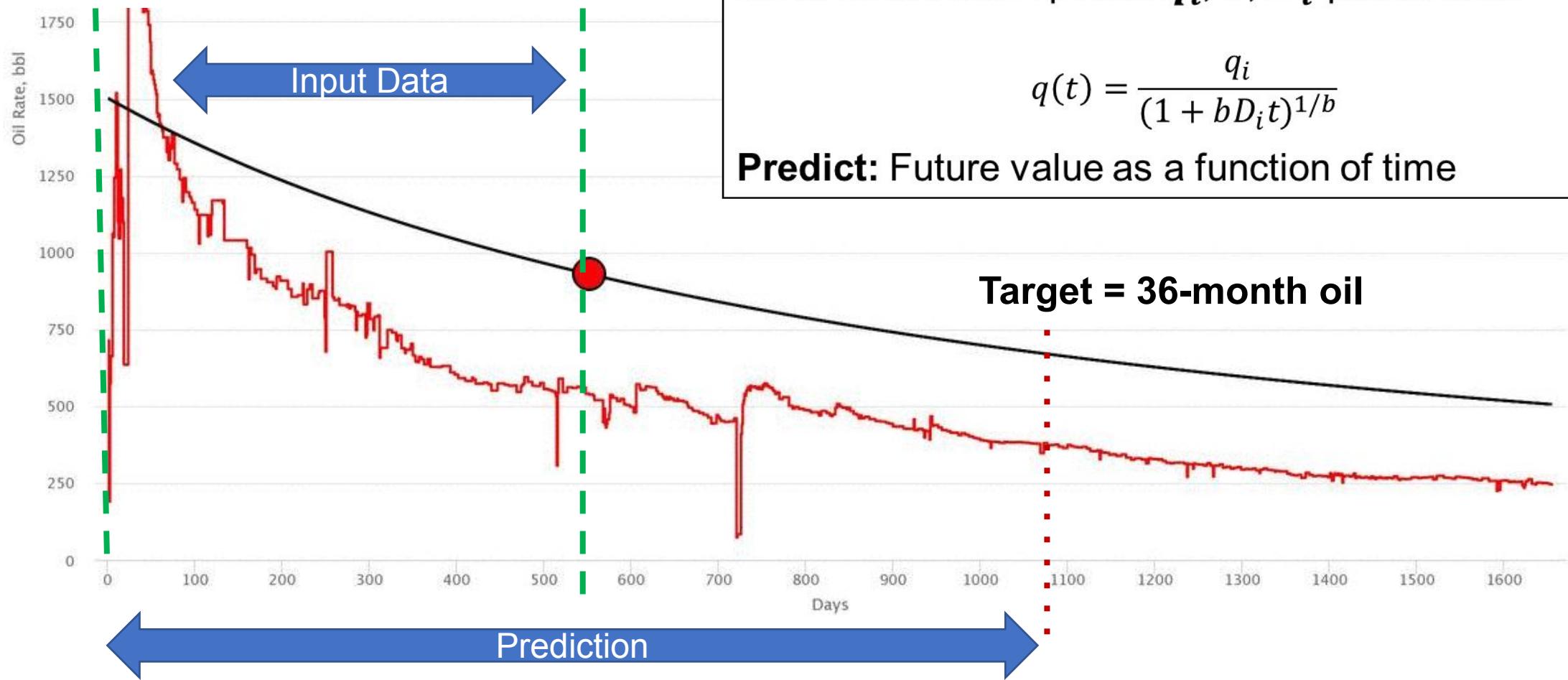


Feature Extraction Methods: A quick review

- (1) Arps : optimal q_i, b, D_i parameters
- (2) PCA: principal components of water, gas, oil data
- (3) Hybrid: combine Arps + PCA



Prediction Type 1: Decline Curve Analysis (DCA)



Prediction Type 2: Linear Regression

Linear Regression using extracted features as predictor variables

RQ4: parametric $\mathbf{Y} = \beta_0 + \beta_1 \mathbf{q_i} + \beta_2 \mathbf{b} + \beta_n \mathbf{d_i}$

RQ5: nonparametric $\mathbf{Y} = \beta_0 + \beta_1 \mathbf{PC_1} + \beta_2 \mathbf{PC_2} + \dots + \beta_n \mathbf{PC_R}$

RQ6: hybrid $\mathbf{Y} = \beta_0 + \beta_1 \mathbf{q_i} + \beta_2 \mathbf{b} + \beta_3 \mathbf{d_i} + \beta_4 \mathbf{PC1} + \beta_5 \mathbf{PC2} + \dots + \beta_n \mathbf{PC_R}$

*Compare these results with DCA from RQ3

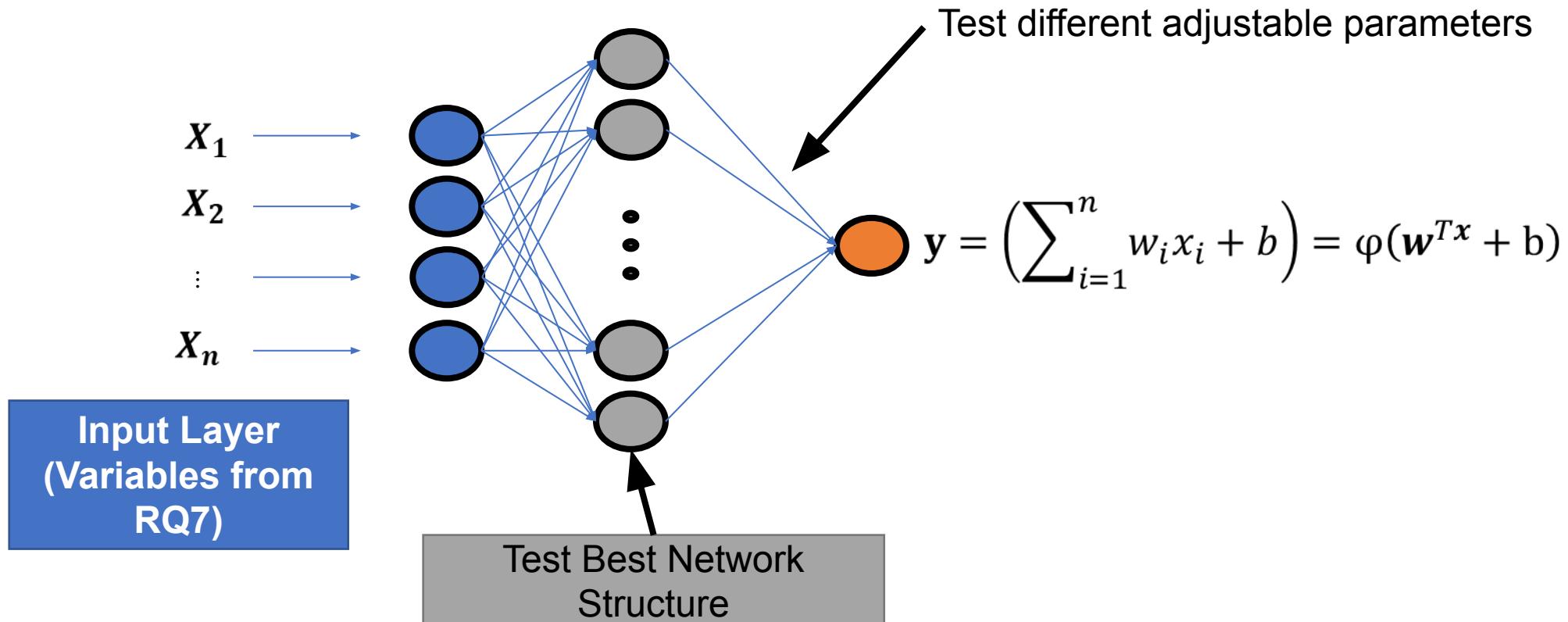
RQ7: Linear Regression using extracted features (FE) + static well features

RQ7: $\mathbf{Y} = \beta_0 + \beta_1 \mathbf{FE_1} + \beta_2 \mathbf{FE_2} + \dots + \beta_n \mathbf{X_1} + \beta_{n+1} \mathbf{X_2} + \dots + \beta_k \mathbf{X_k}$

Prediction Type 3: MLP

RQ8: Adding model complexity through MLP using existing input variables as RQ7

*Compare results with RQ7

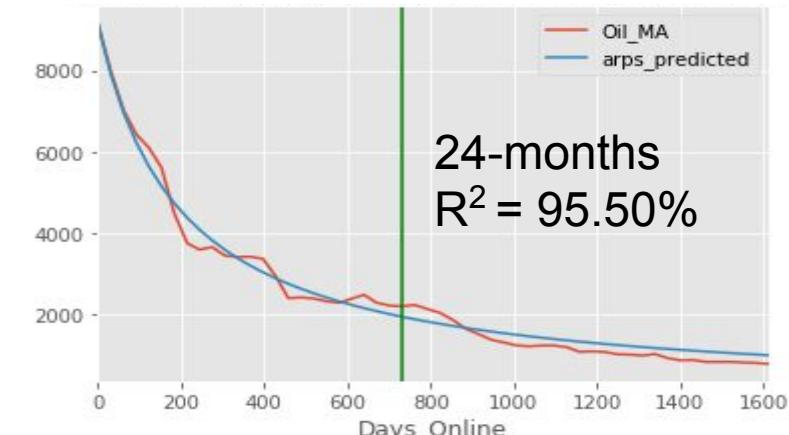
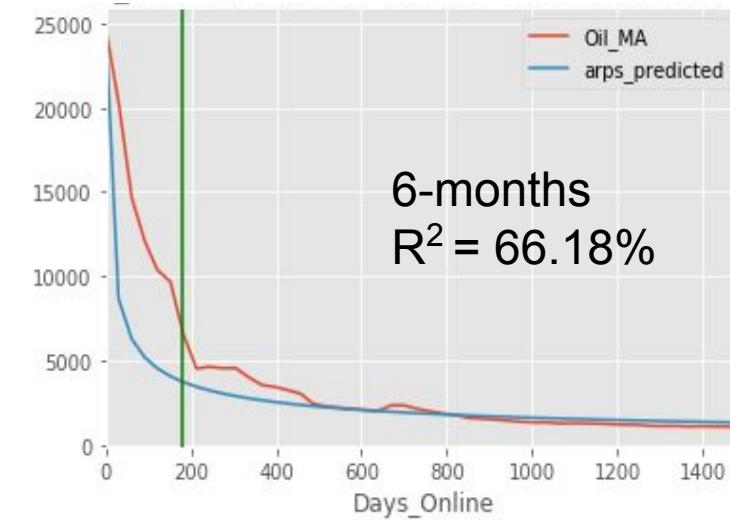


Results: Evaluating Arps Model

Answering research question RQ1

Goodness-of-fit metrics using Arps Model

Scenario	RMSE	MAPE	R ²
6 Months	14253	16.26%	66.18%
12 Months	4989	8.53%	91.63%
24 Months	3659	6.31%	95.50%



Results: Feature Selection of Static Well Variables

Answering research question RQ2

Parameter Type	Parameter	Rank
PVT	Oil Gravity	1
Completion	True Vertical Depth	2
Location	Surface Latitude	3
Location	BH Longitude	4
Completion	Proppant to Fluid Ratio	5
Completion	CLL	6
Completion	Total Fluid	7

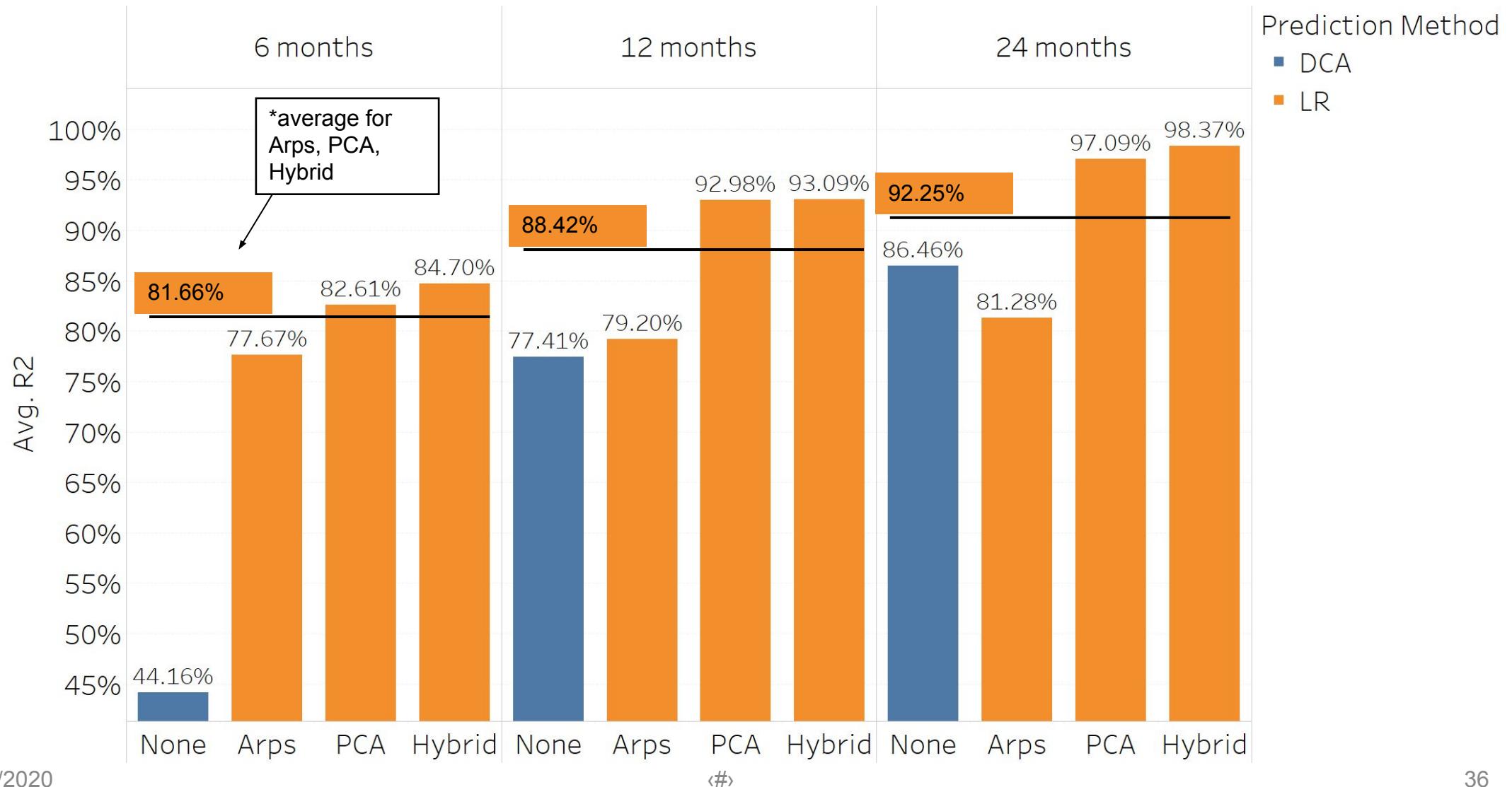
Summary

- Three types of well variables were found significant
- 51.11% proportion of variance explained by regression model

RMSE	MAPE	R ²
80942	21.18%	<u>52.11%</u>

Results: Can regression prediction outperform DCA?

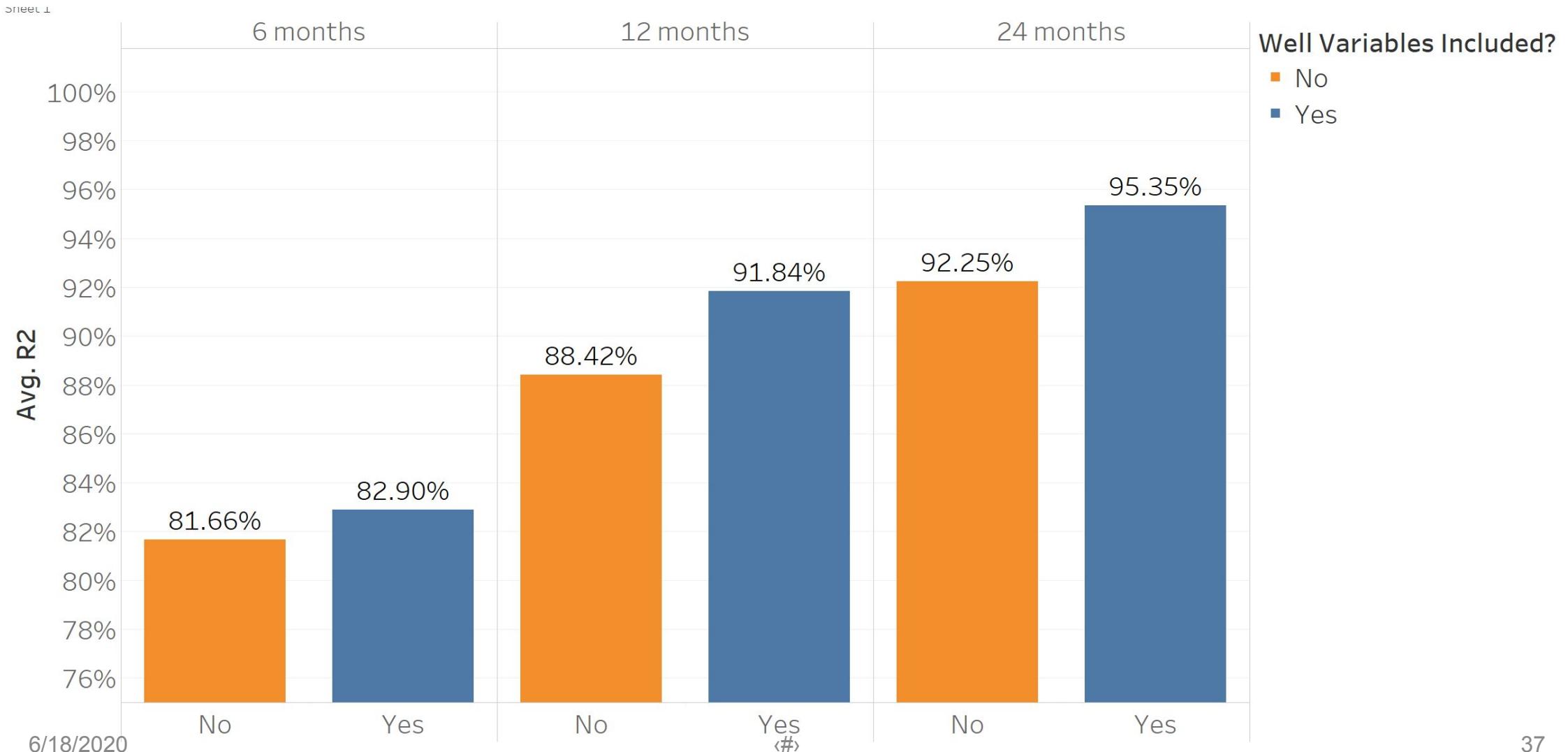
Answering Research Questions RQ3 – RQ6



Results: Does Including Well Variables Aid the prediction?

Answers Research Question RQ7

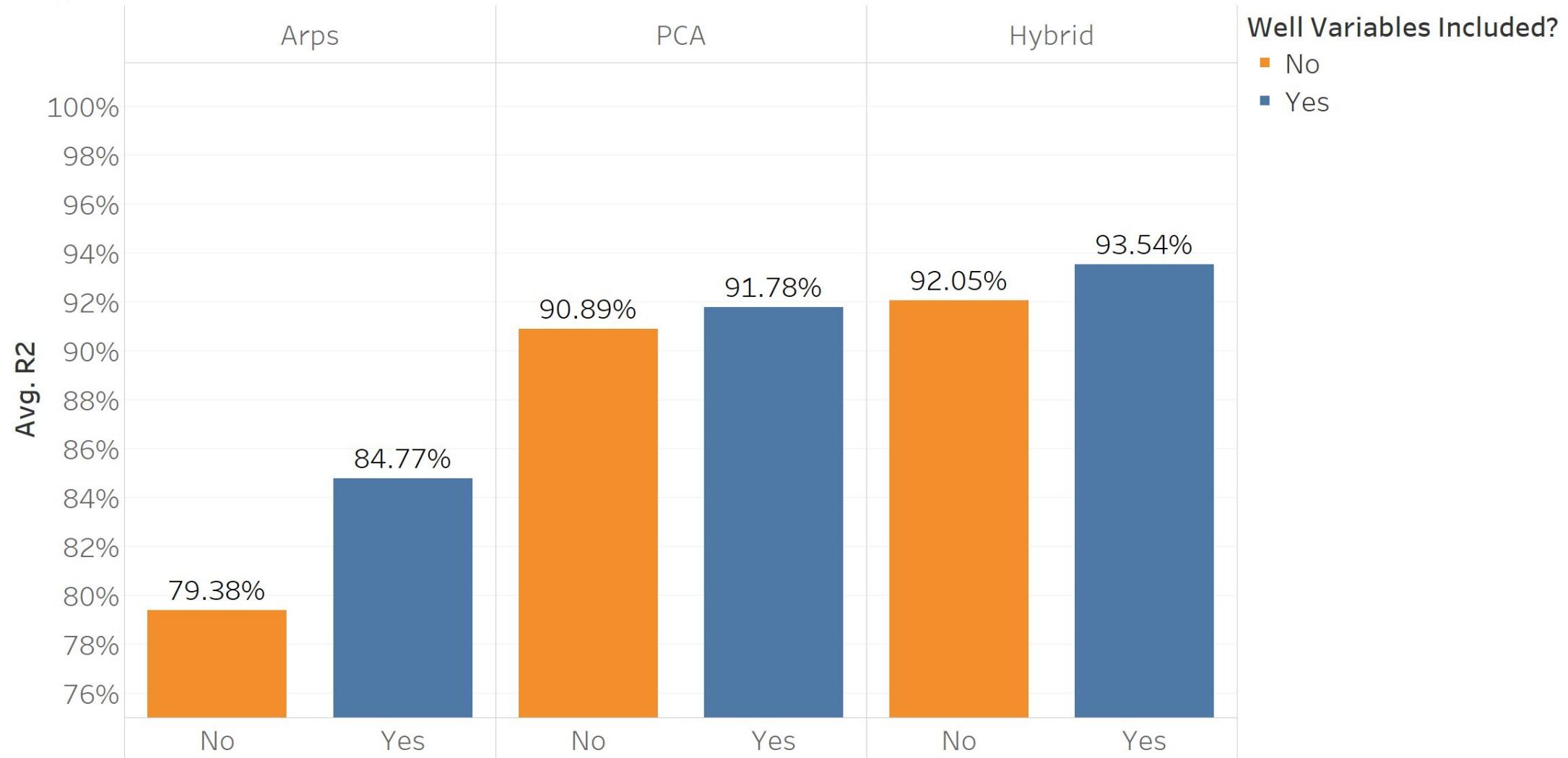
Average R² by prediction scenario when well variables included in the model



Results: Does Including Well Variables Aid the prediction?

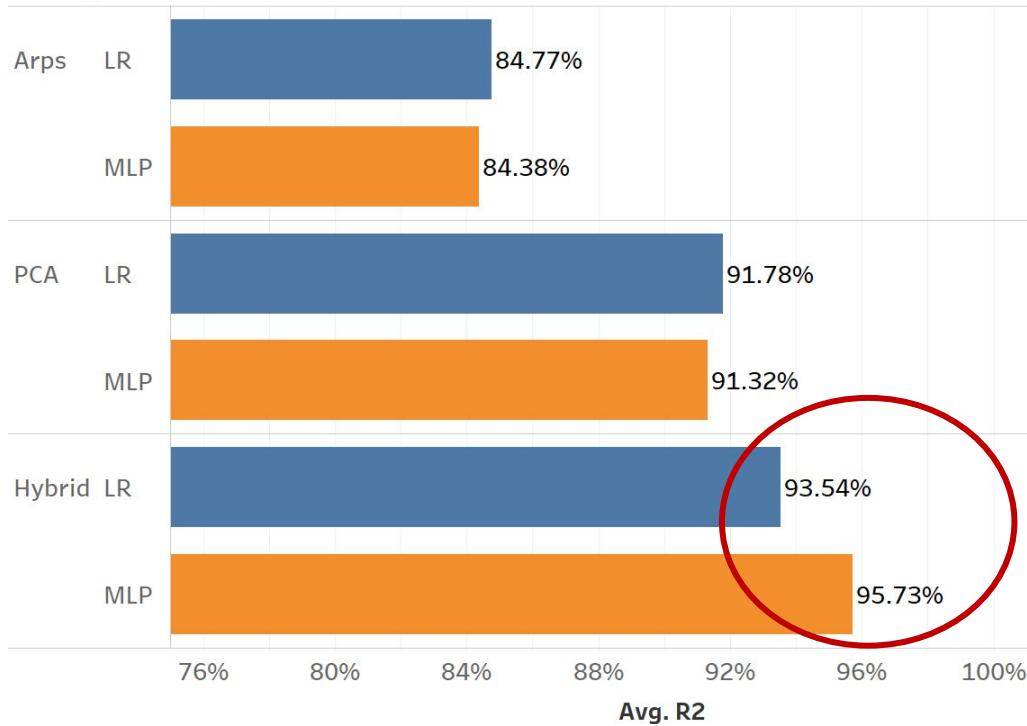
Answers Research Question RQ7

Average R^2 by feature extraction method when well variables included in the model



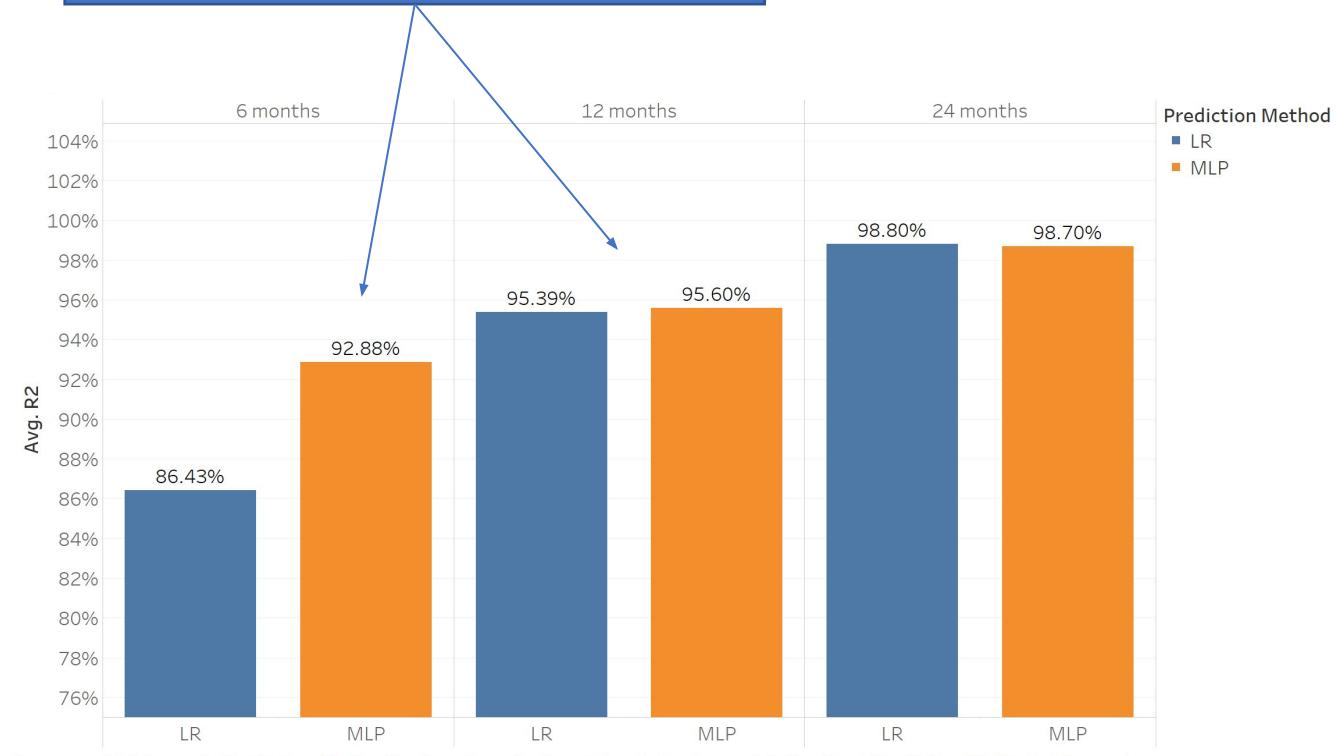
Results: Does adding model complexity aid the prediction results?

Answers research question RQ8

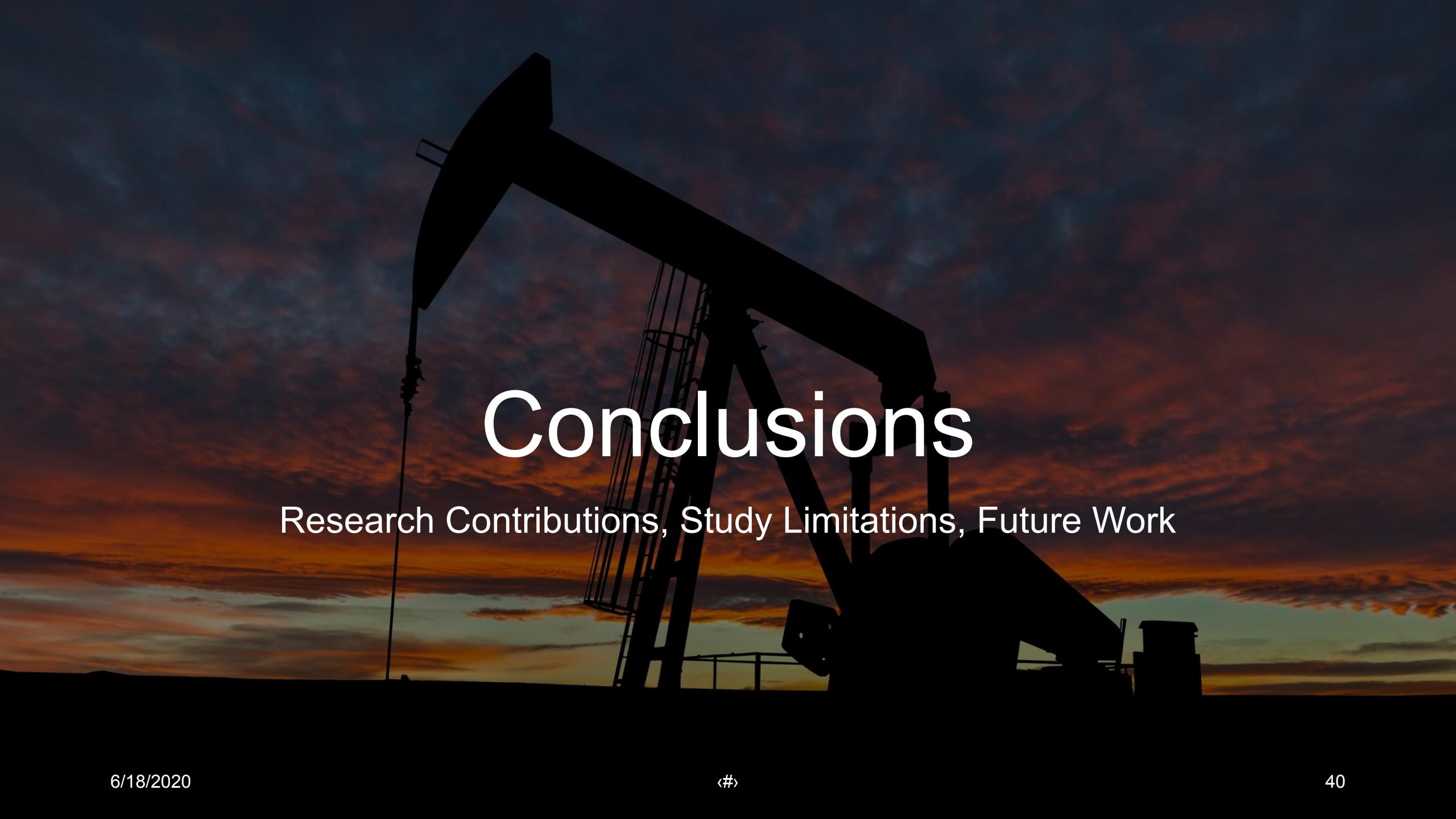


Average across all data scenarios reveals
Hybrid as best feature extraction

Recall: 24-month DCA method
yielded 86.46% R²



Comparison of hybrid methods only

A silhouette of an oil pumpjack is positioned in the center-left of the slide, facing right. The background is a dramatic sunset or sunrise with orange, red, and blue hues in the clouds.

Conclusions

Research Contributions, Study Limitations, Future Work

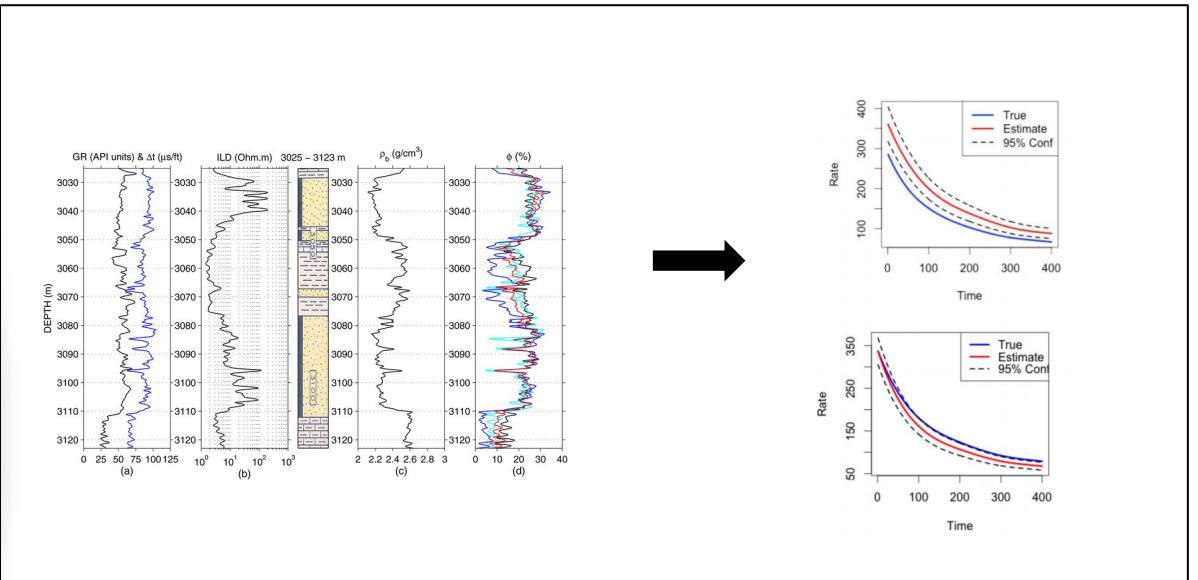
Conclusions and Contributions

- **Regression using feature extraction** is more effective than DCA
 - PCA of production data was effective in explaining statistical variance
 - Combination of DCA and PCA most effective in prediction
- **Inclusion of static well variables** can make valuable predictions about future cumulative oil production
 - Improved regression prediction results for all scenarios
- **Model complexity** not needed when richer production history available
 - Valuable in case of 6-months of data; predicted better with less data than DCA

Future Work

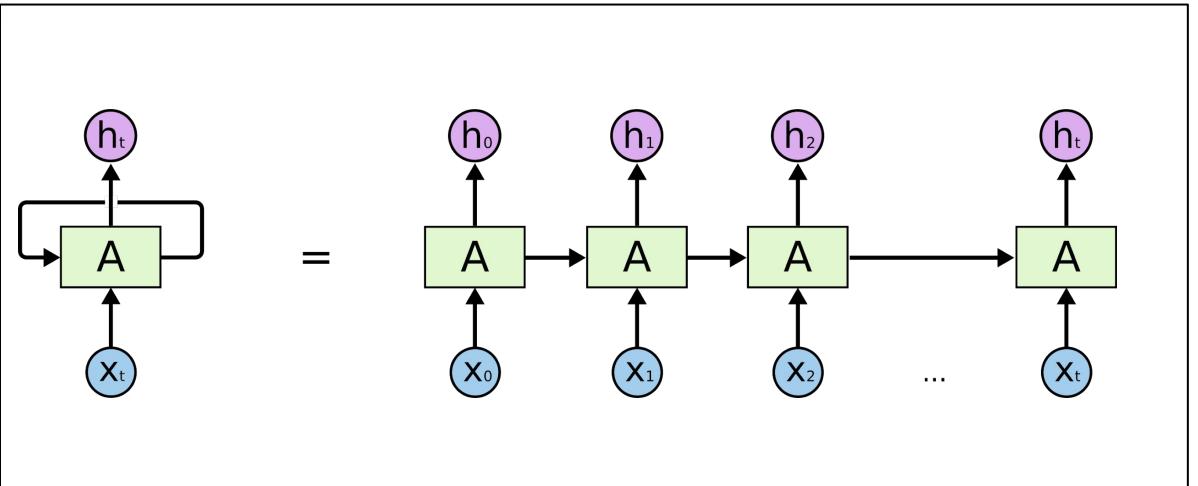
Include richer data

- Geologic, well logs, seismic, and reservoir pressures
- Multivariate prediction



New Prediction Scenarios

- Time-series neural network for month-to-month rate prediction
- Estimated Ultimate Recovery (EUR)





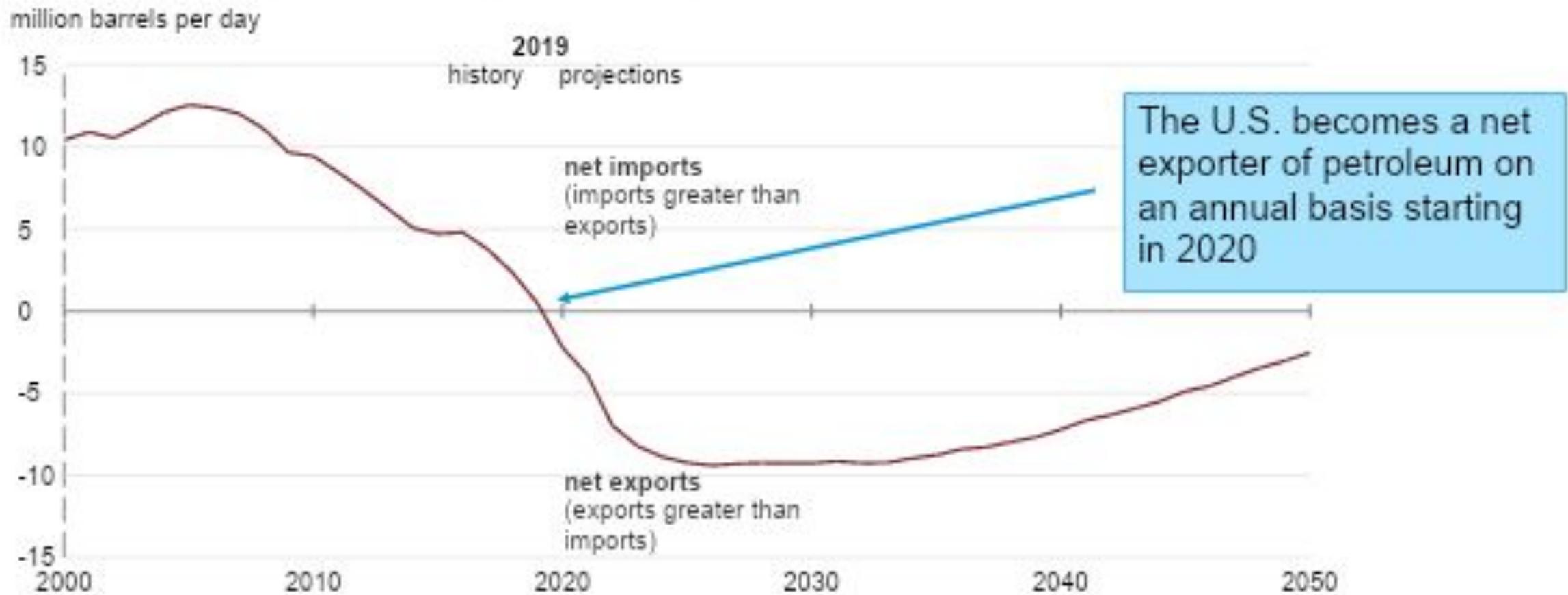
Questions

Appendix

Appendix: Oil Industry Background

The United States exports more petroleum than it imports from 2020 to 2050

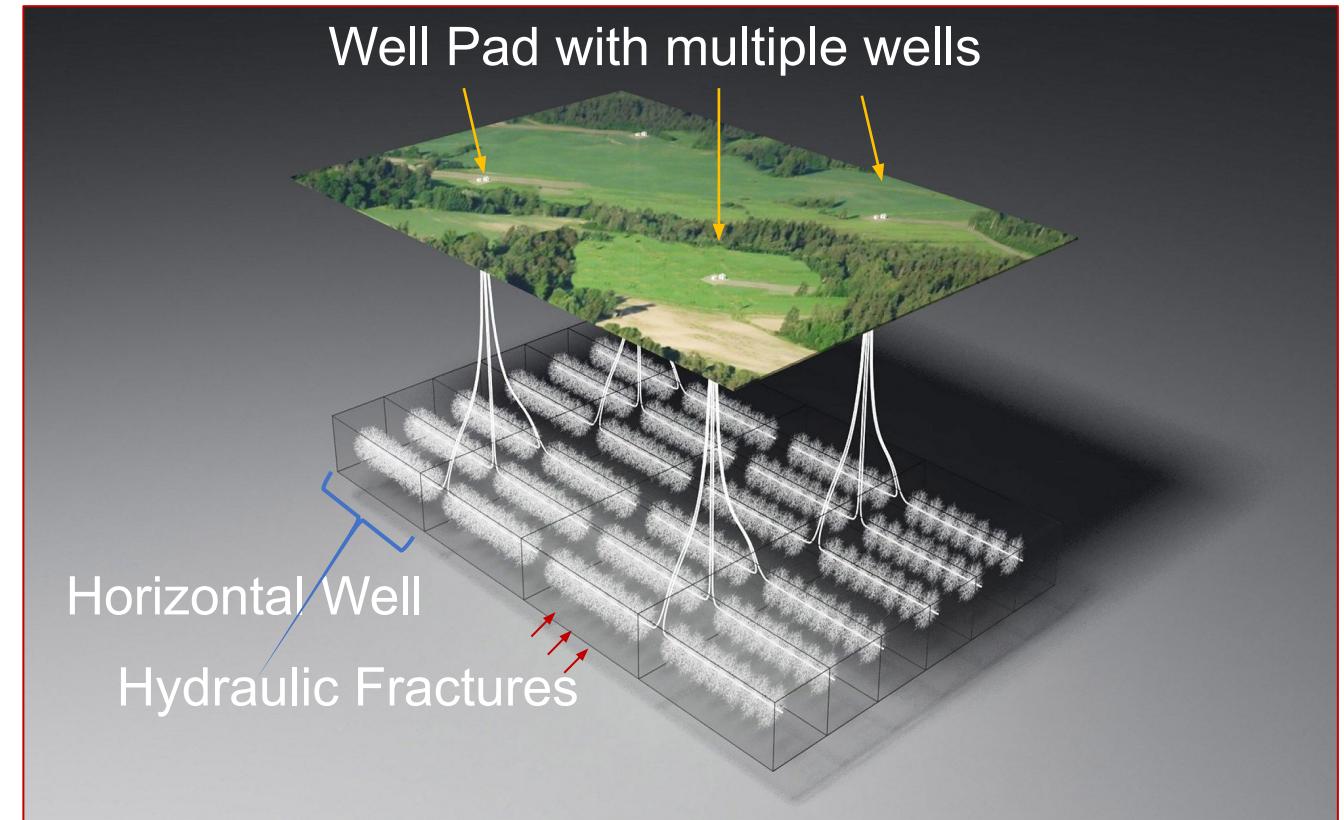
Crude Oil and petroleum liquids net exports



The U.S. becomes a net exporter of petroleum on an annual basis starting in 2020

Unconventional Reservoirs

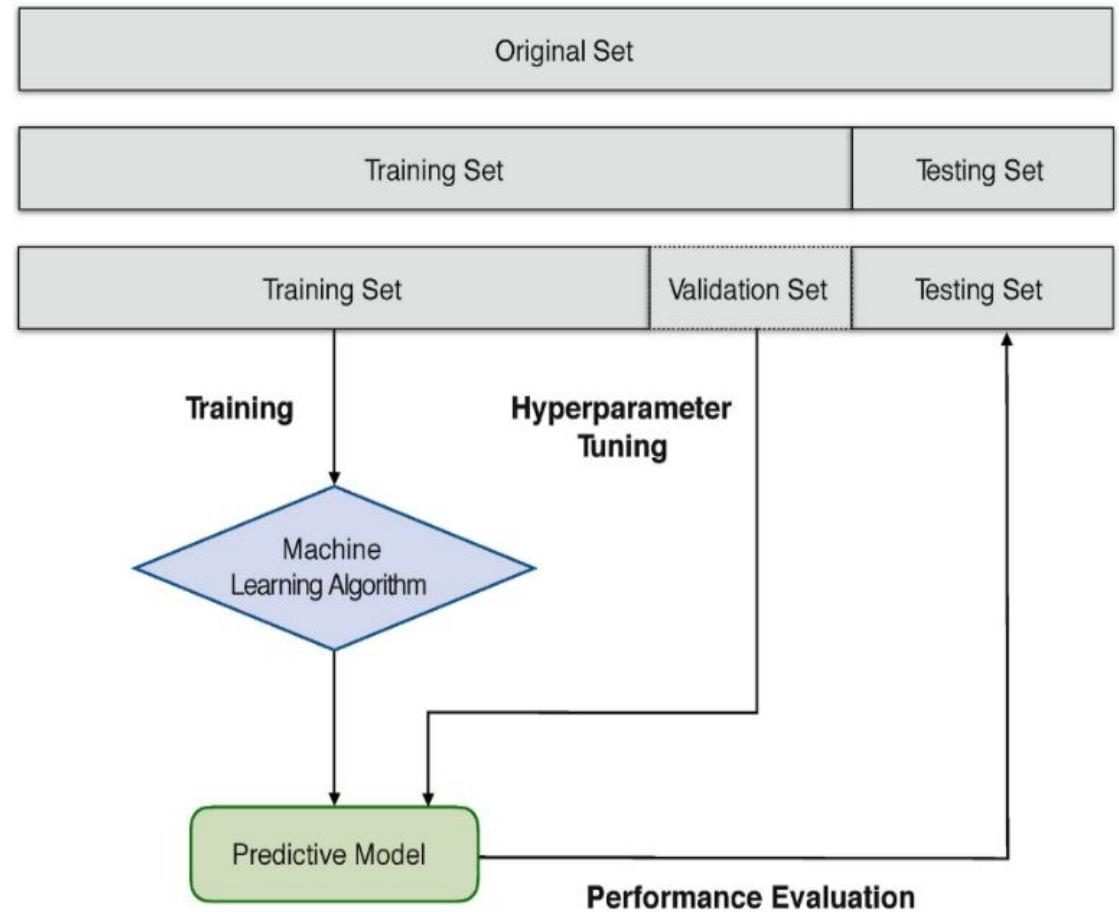
- Very low permeability
- *Horizontal* wells
- *Multistage* fracturing
- Well pad with multiple wells



Appendix: Machine Learning Introduction

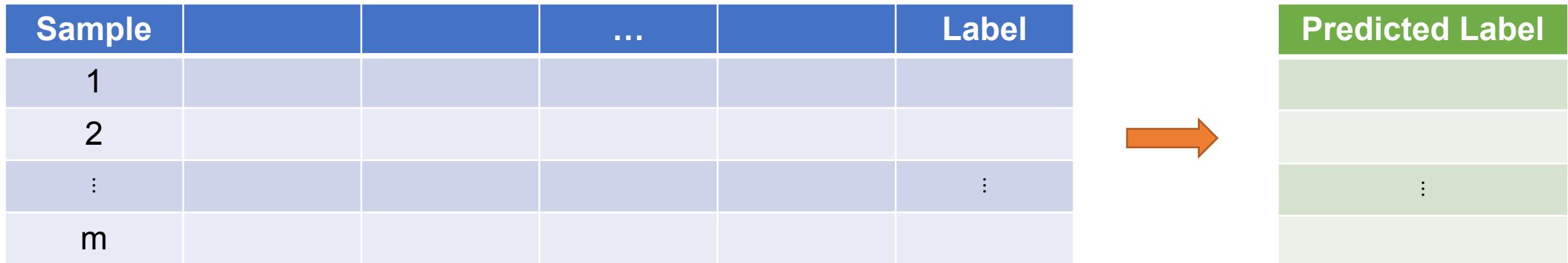
Supervised Learning

- Train set
 - ML algorithm “learns” data
- Validation set
 - ML algorithm “sees” data
- Testing set
 - ML algorithm “predicts”



Supervised Learning

- **Purpose:** Optimize a performance measure
 - Inclusion of a class label or “target” during calculations
 - Examples: Regression, Classification

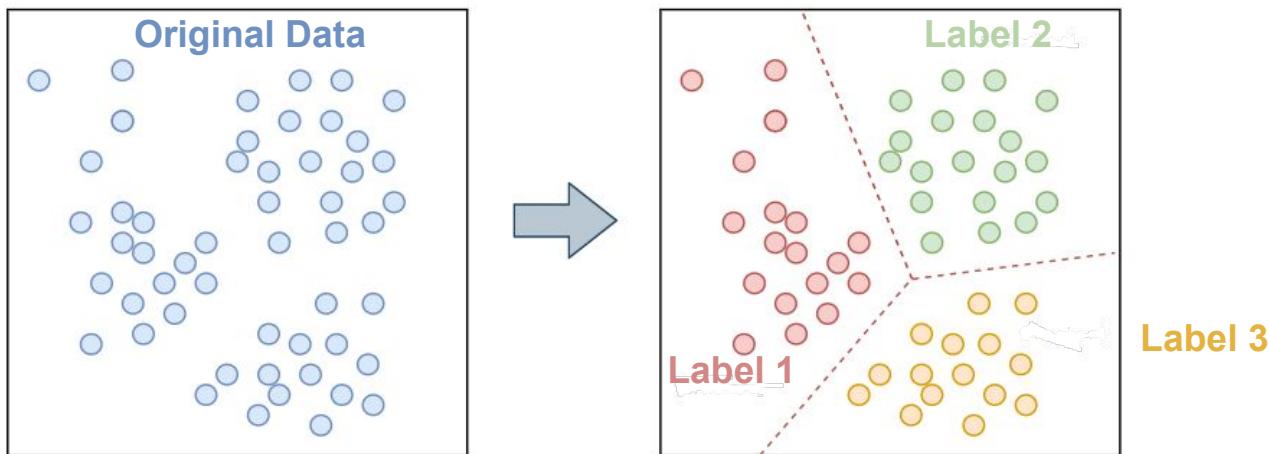


Historical/observed oil production represents target variable in this study

Unsupervised Learning

Purpose: discover inherent patterns in the data

- Typically excludes class label or “target” during calculations but not always
- Examples: Clustering and Dimension Reduction

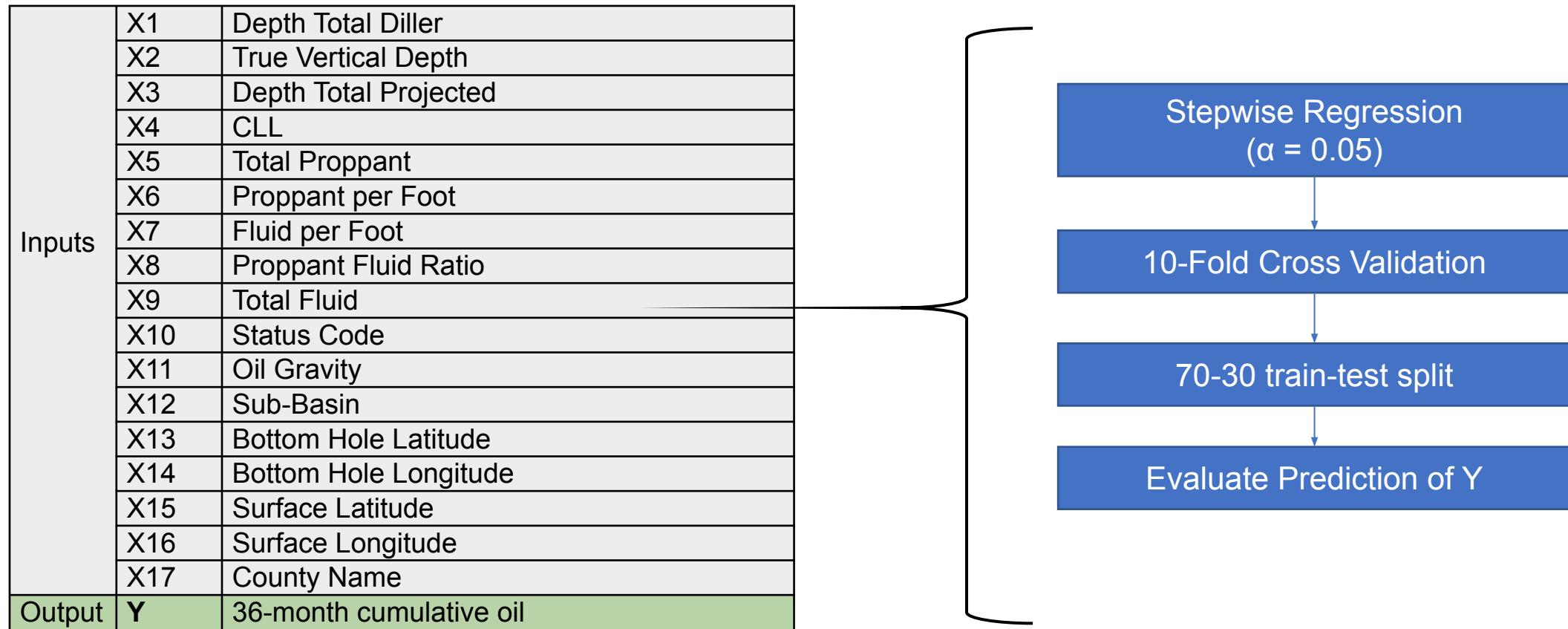


Appendix: Design of Experiments

Research Questions

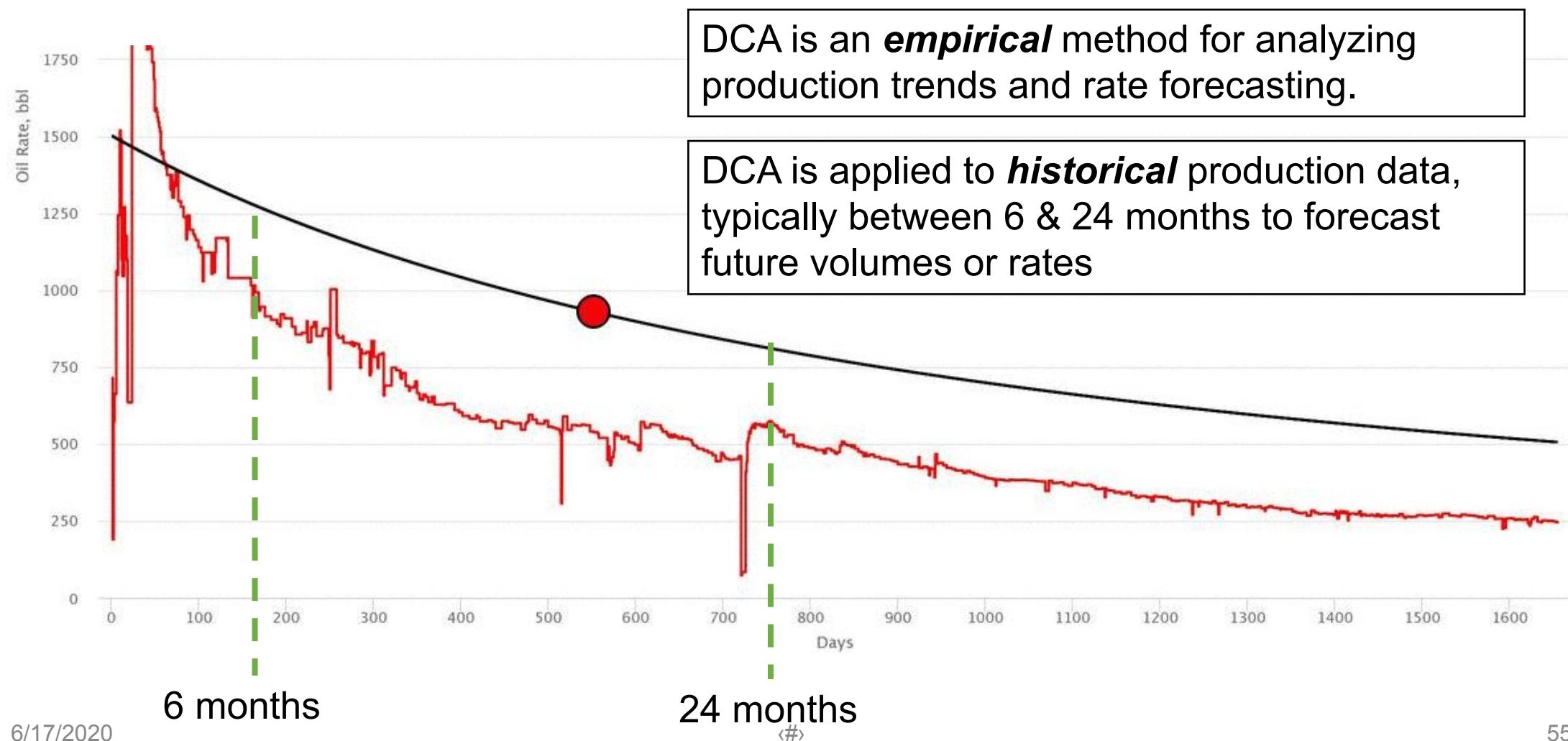
- **RQ1:** How effective is Arps Model in explaining the variance of the field data used in this study?
- **RQ2:** Is there a statistical relationship between static well variables and cumulative production?
- **RQ3:** How effective is Arps Model in **predicting** the target variable
- **RQ4:** Can Arps Variables act as an effective feature extraction method to predict cumulative oil production?
- **RQ5:** Can PCA as an effective feature extraction method to predict cumulative oil production?
- **RQ6:** Can the combination of Arps variables and PCA be used as an effective feature extraction method to predict cumulative oil production?
- **RQ7:** Does including the well variables aid the prediction performance of the proposed algorithms?
- **RQ8:** Using the same prediction framework as RQ7, does using MLP instead of LR help the prediction accuracy?

RQ2: Feature Selection of Static Well Variables



$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_{17} X_{17}$$

Decline Curve Analysis (DCA)

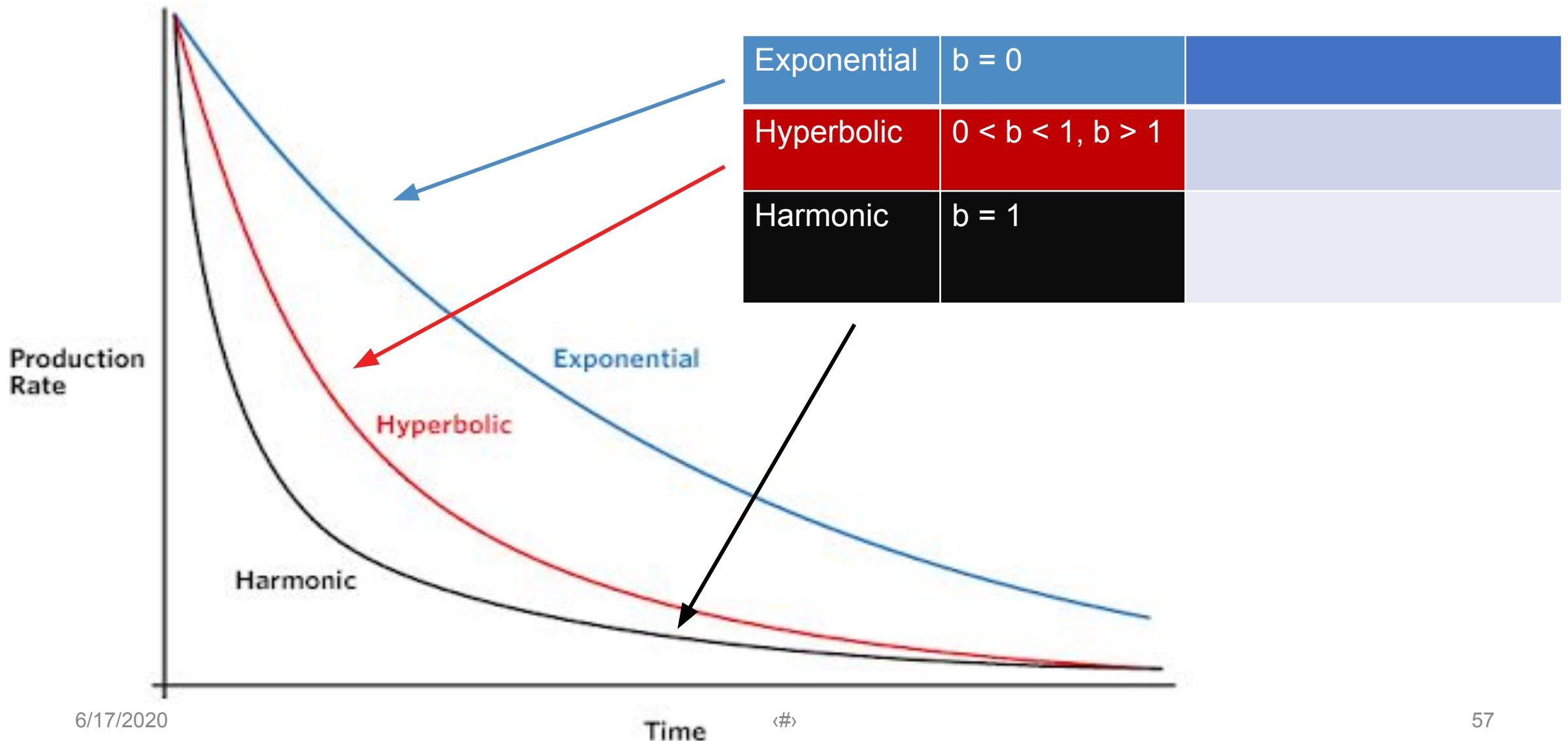


Arps DCA Model

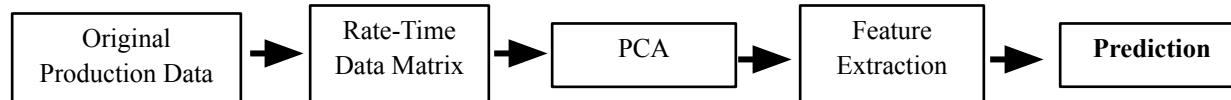
- $$q(t) = \frac{q_i}{(1 + bD_i t)^{1/b}}$$
- $q(t)$ = production rate at time t
- q_i = initial production rate
- D_i = decline constant or “nominal decline rate”
- b = decline exponent or “degree of curvature”

“This natural declining trend is attributed to number of factors including fluid mechanics in porous media, changes in reservoir pressure, changing relative volumes of fluids through, and efficiency of vertical movement by drilling equipment” (Okoro, 2019)

Mathematics of DCA



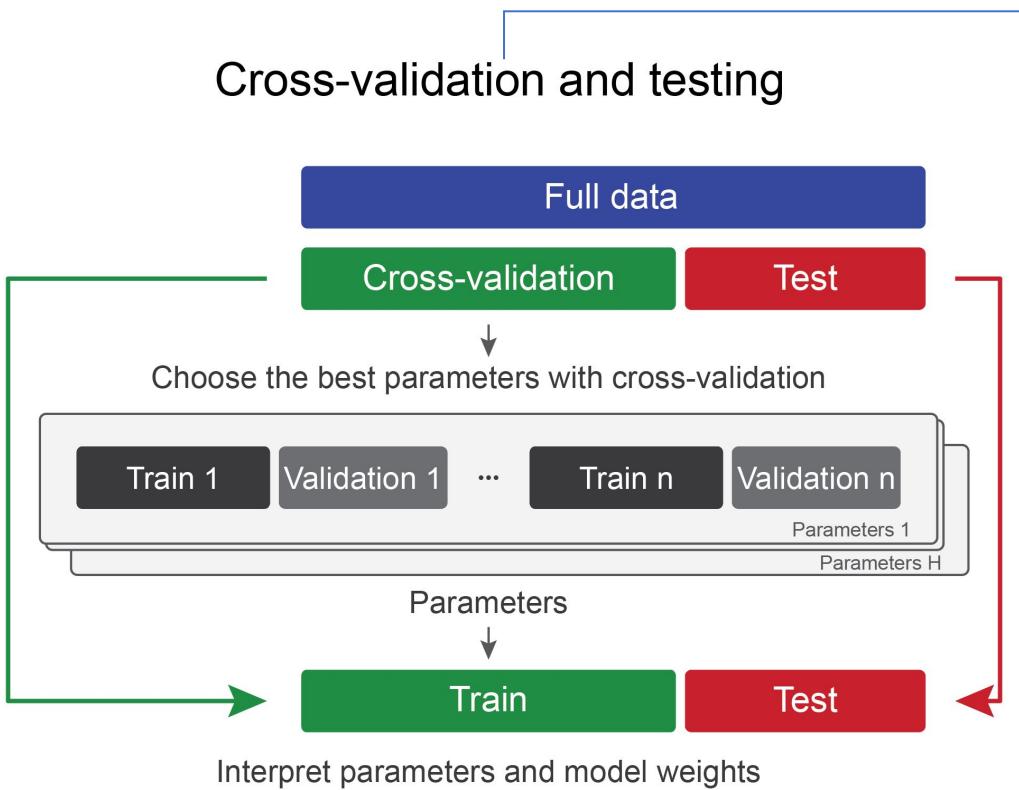
Dataset Transformation for PCA



Well ID	Production Month 1	Production Month 2	...	Production Month n
1				
2				
⋮				
448				

Well ID	PC1	PC2	...	PC n
1				
2				
⋮				
448				

Cross-validation: creating *reproducibility* in experiments



K-Fold Cross Validation Data Partition

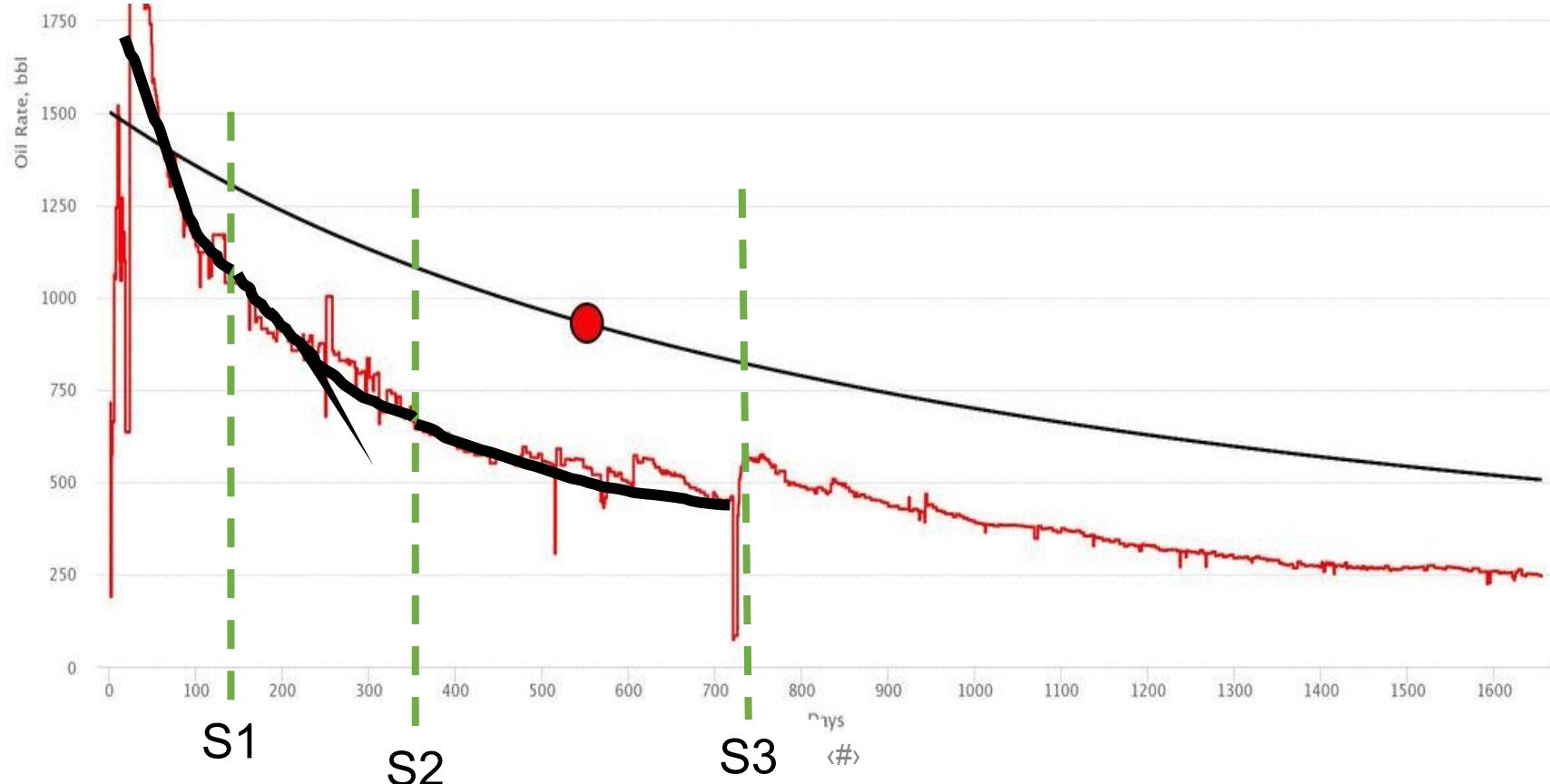
$K = 1$	Holdout	Train	Train	Train	Train
$K = 2$	Train	Holdout	Train	Train	Train
$K = 3$	Train	Train	Holdout	Train	Train
$K = 4$	Train	Train	Train	Holdout	Train
$K = 5$	Train	Train	Train	Train	Holdout

Evaluating Arps Model on Field Data

Fit Arps Model to each scenario of data

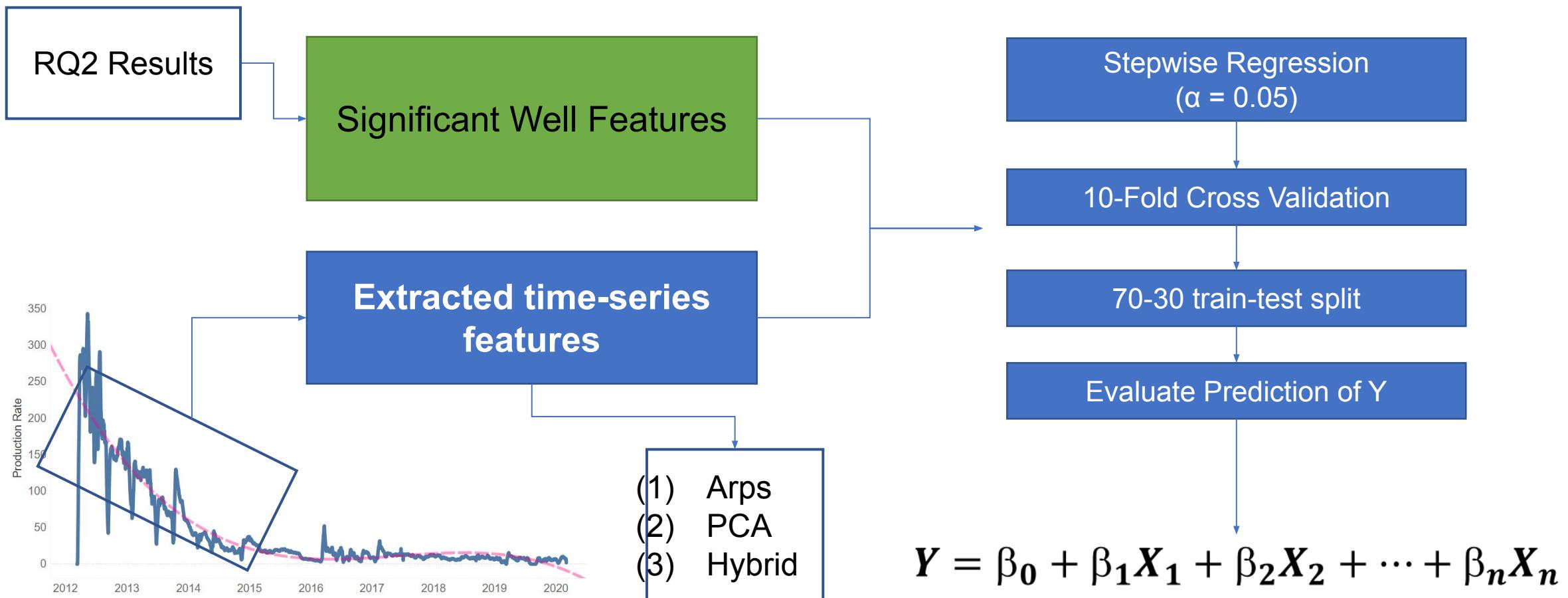
Evaluate Goodness-of-Fit

Assess as Feature Extraction



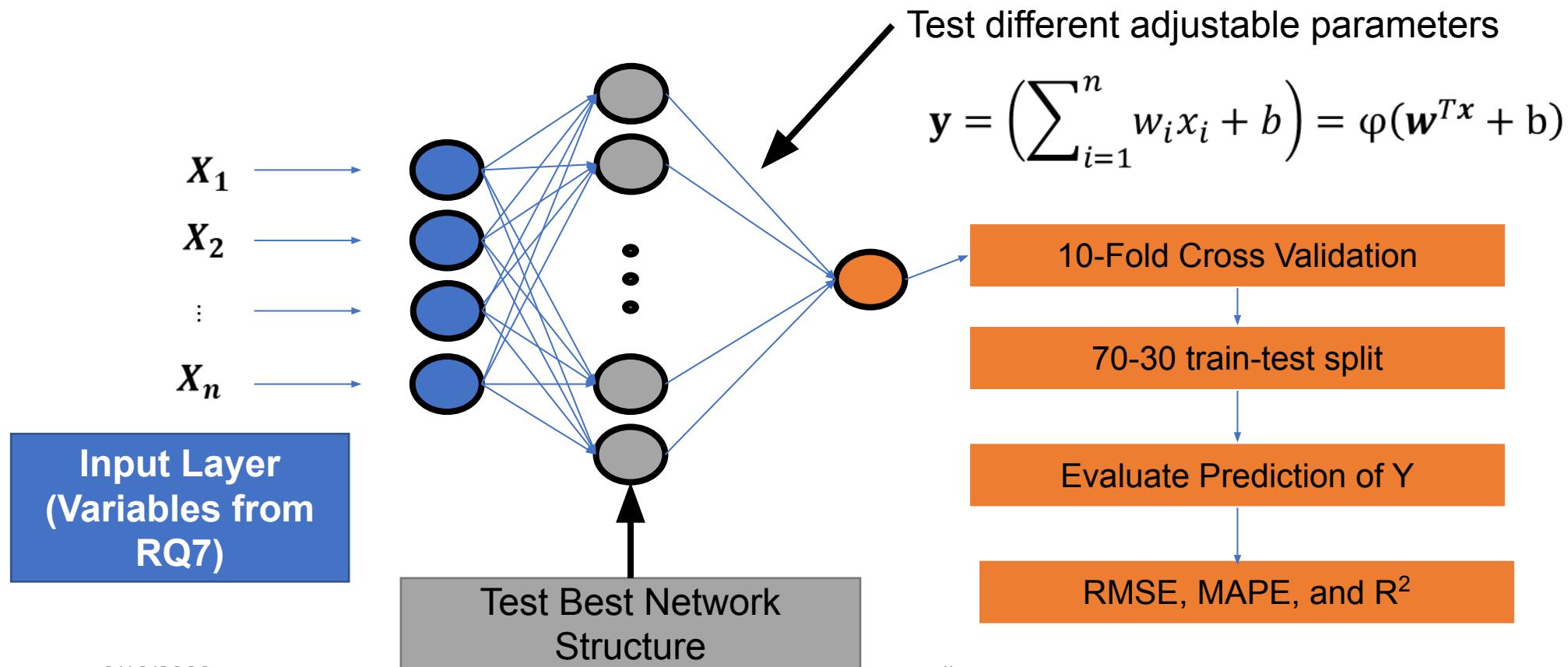
Design of Experiments: RQ7

- Including extracted features with well parameters



Design of Experiment: RQ8

Adding model complexity through MLP using existing input variables as RQ7:



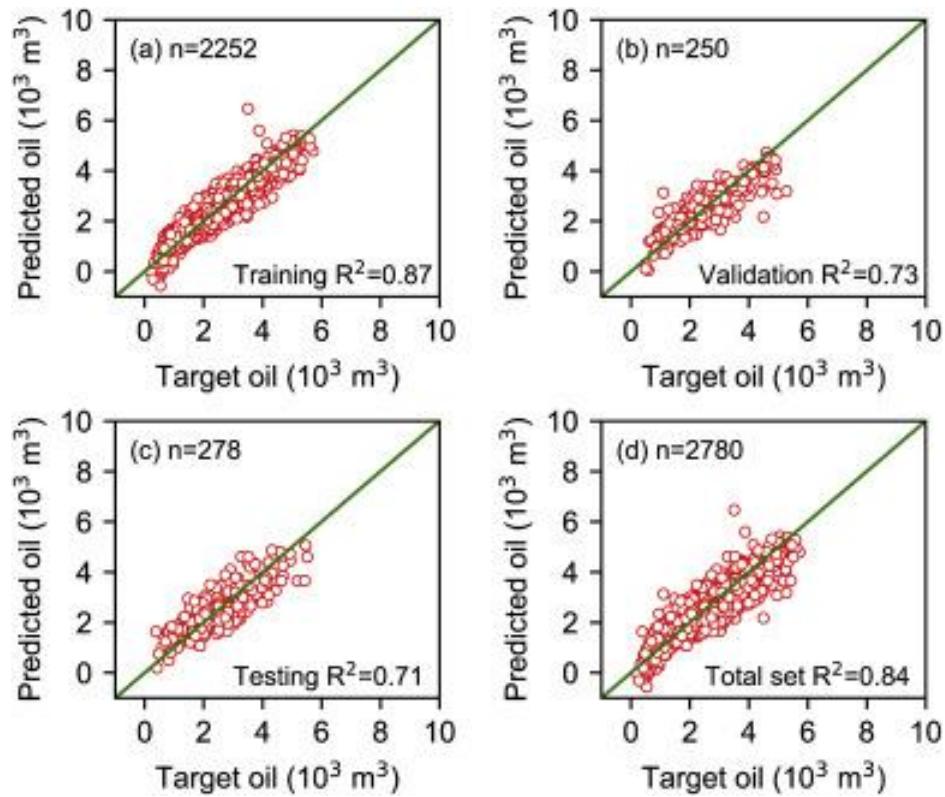
Appendix: Literature Review

Summary of Key Related Works

Source	Study Purpose	Methods	Evaluation Metrics	Outcomes	Relation to this study
[1], [2], [3]	Compare machine learning time-series predictions to empirical DCA methods	LSTM ANN	R ² , MAPE, RMSE, MSE	Successful time-series application in predicting oil and gas production in short-term intervals	Comparison to DCA
[4]	Evaluate static parameters to reconstruct dynamic profiles	PCA, MLP	R ² , MAPE, RMSE	Applied PCA to static well parameters as inputs to MLP to predict DCA model coefficients	Incorporate DCA parameters for model building
[5]	Prediction of initial flow rates and evaluate variable importance	MLP	MSE, R ²	Used pressure and time well log data to predict initial oil flow rates to 93% R ²	Study of initial flow rates (Qi in Arps Model)
[6], [7]. [8]	Prediction of cumulative oil production and evaluate variable importance of static well parameters	LR, MLP	R ² , MAPE, RMSE, MSE	Predicted cumulative oil production with resulting test set R ² 71% (6-month) and 72% (18-month)	Incorporates static well parameters including completions, geology, location, PVT data
[9]	Prediction of cumulative oil production and evaluate variable importance of static well parameters specific to EFS	LR and MLP	R ² , MAPE, RMSE	Predicted 6-month cumulative oil production with resulting test set R ² 39% (LR) and 62% (MLP)	EFS
[10]	Study efficacy of FPCA on oil production time-series	LR and PCA	R ² , MAPE, RMSE	Reconstruct production profiles to less than 5% MAPE using first 4 PCs	Feature Extraction using PCA
[11], [12]	Apply FPCA to match historical production and predict future values	PCA, LR	R ² , MAPE, RMSE	MAPE: 21.87% to match field data [11] MAPE: 10.3% to match simulated forecasts [12]	Feature Extraction using PCA for prediction

Appendix: Literature Review

- Shuhua Wang, Zan Chen, Shengnan Chen,
- Applicability of deep neural networks on production forecasting in Bakken shale reservoirs,



Variables	Parameters	Range	Mean	Standard deviation
Independent	Easting (10^5 m)	5.71–6.77	6.42	0.16
	Northing (10^6 m)	5.45–5.55	5.5	0.01
	True vertical depth (m)	728.90–2338.50	1545.3	93.34
	Well lateral length (m)	186.00–3040.60	1254.7	282.95
	Upper Bakken thickness (m)	1.13–6.82	2.84	0.84
	Middle Bakken thickness (m)	1.13–15.11	10.05	2.21
	Lower Bakken thickness (m)	2.56–11.19	5.93	1.29
	Number of fracture stages	1–50	15	8
	Fracture spacing (m)	17.44–716.56	117.88	68.68
	Average fluid pumped per stage (m ³)	0.20–221.00	33.6	21.25
	Load fluid recovery factor (%)	0.00–90.00	18.41	19.26
	Average proppant per stage (tons)	2.00–99.00	6.69	4.27
Dependent	6 months oil production (10^3 m^3)	0.01–8.22	2.51	1.45
	18 months oil production (10^3 m^3)	0.05–17.96	5.35	3.19

(Grujic, Da Silva, & Caers, 2015)

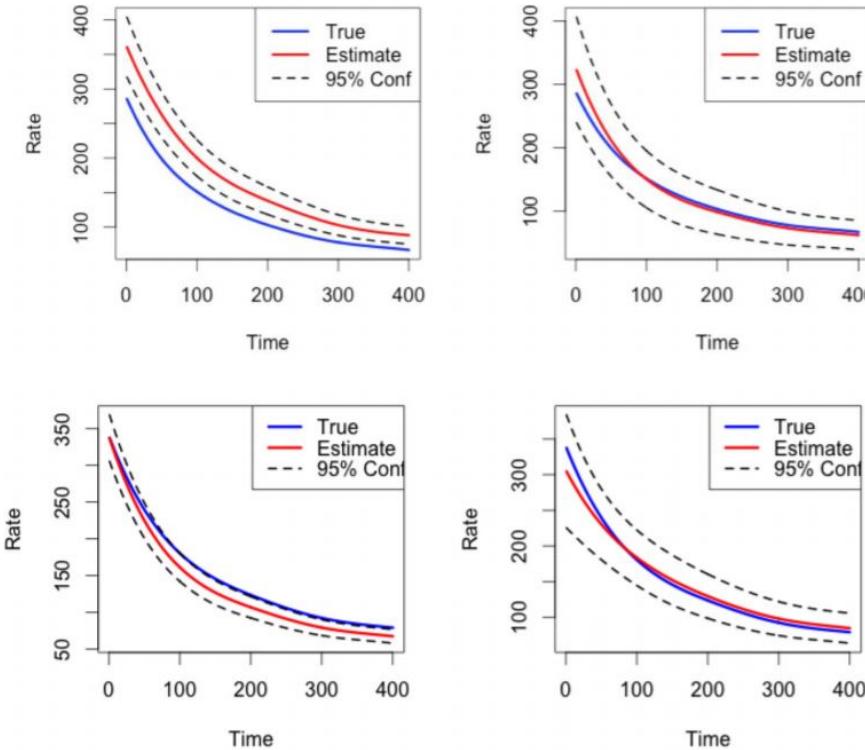


Table 1—Summary of all available data

Type	Parameter	Unit	Source	Kind	Abbreviation
OUTPUT	Oil Rates	stb/day	raw	function	-
	Gas Rates	stb/day	raw	function	-
	Water Rates	stb/day	raw	function	-
Completions	Number of Completion Stages	#	raw	scalar	CMP STAGES STIMULATED
	Total Amount of Injected Fluid	gal	raw	scalar	CMP TOTAL FLUID PUMPED GAL
	Total Amount of Injected Proppant	lbs	raw	scalar	CMP TOTAL PROPPANT USED
	Stimulated Lateral Length	ft	raw	scalar	CMP STIMULATED LATERAL LENGTH
	Total Amount of Slick Water	bbl	raw	scalar	CMP AMT SLICKWATER BBL
	Total Amount of Injected x-link fluid	bbl	raw	scalar	CMP AMT CROSSLINK BBL
	Completion Stage Interval	ft	raw	scalar	CompStageInterval
INPUT	Total amount of linear fluid	bbl	raw	scalar	CMP AMT LINEAR BBL
	X location	ft	raw	scalar	GeoIX Rel
	Y location	ft	raw	scalar	GeoY Rel
	Z location (depth)	ft	raw	scalar	GeoZ
Geophysical	Oil API Gravity	api units	raw	scalar	GeoAPIGrav
	Total organic content (TOC)	%	well log - interpretation	scalar	PetroTOC
	Clay Content	%	well log - interpretation	scalar	PetroVClay
	Water Saturation	%	well log - interpretation	scalar	PetroSwt
	Porosity	%	well log - interpretation	scalar	PetroPor
	Total Amount of Quartz	%	well log - interpretation	scalar	PetroVQtz
Petrophysical	Amount of Pyrite	%	well log - interpretation	scalar	PetroPyr

Method	Average MAPE	% Data within 95% CB
OLS Regression	21.87%	47%

Appendix: Results Tables

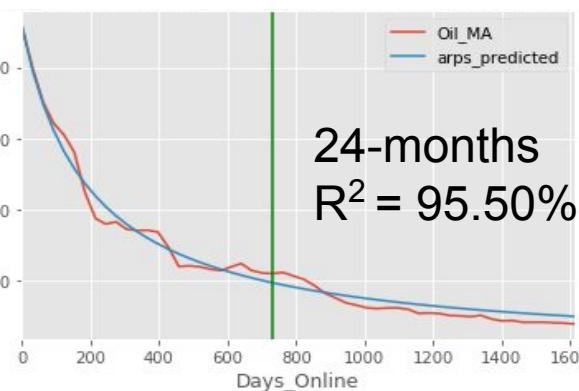
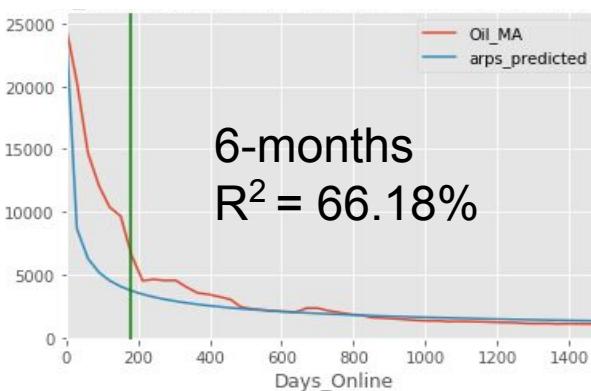
RQ1 Results: Evaluating Arps Model vs PCA on Field Data

Goodness-of-fit metrics using Arps Model

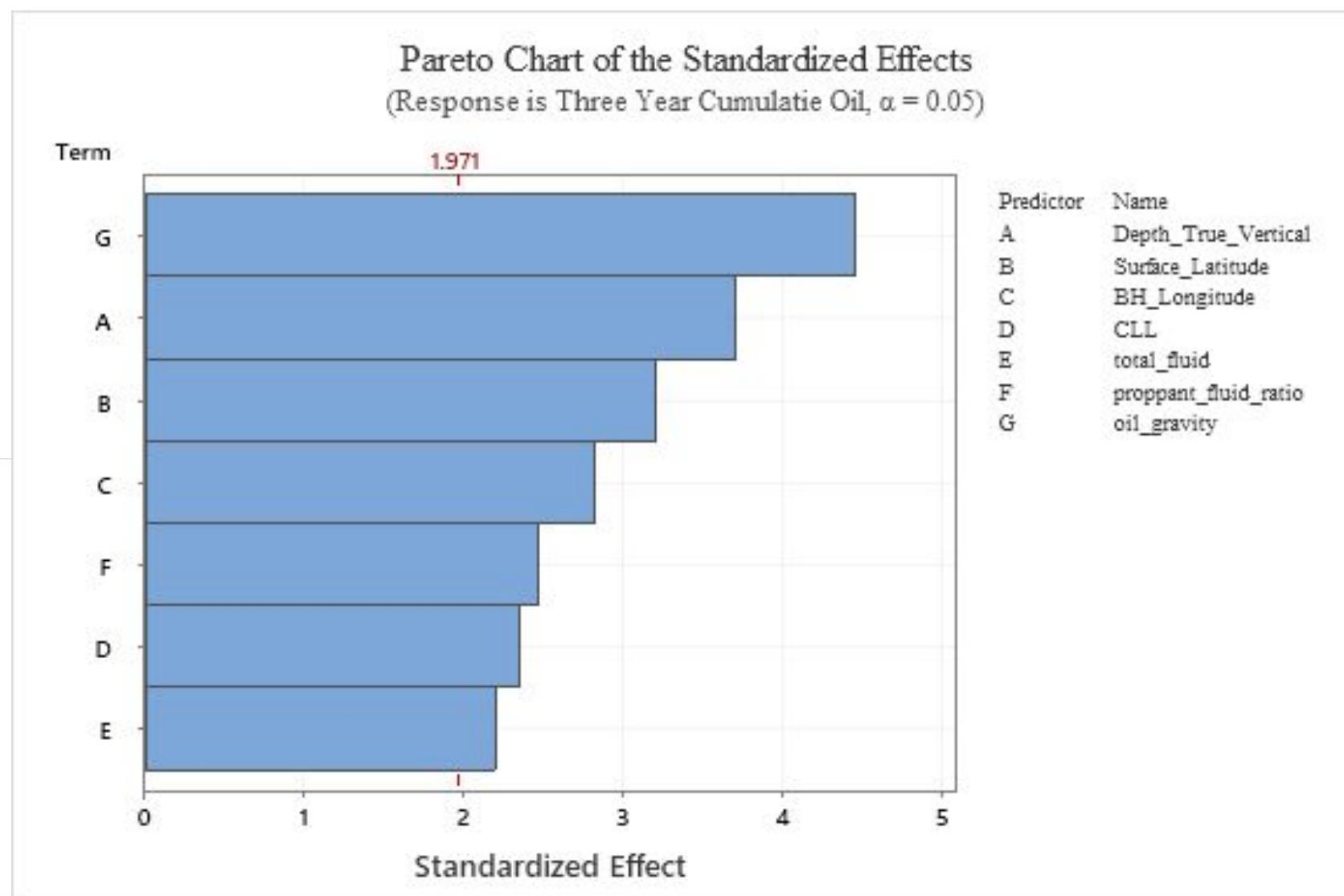
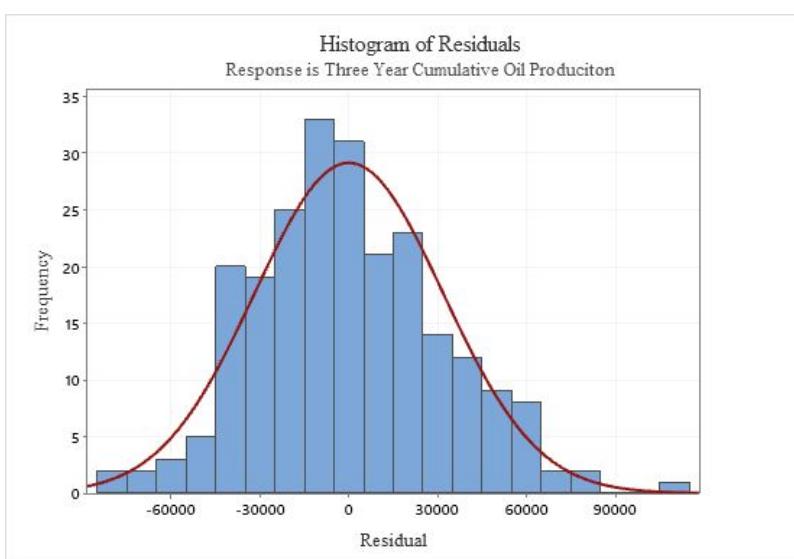
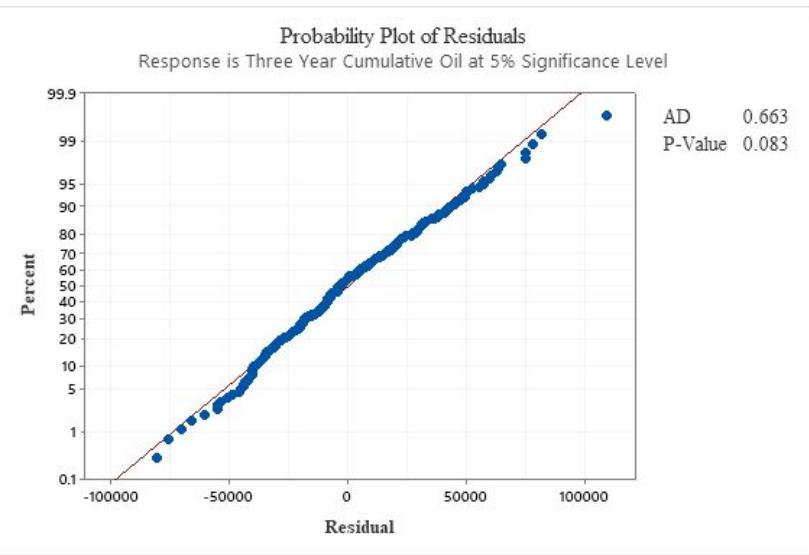
Scenario	RMSE	MAPE	R ²
6 Months	14253	16.26%	66.18%
12 Months	4989	8.53%	91.63%
24 Months	3659	6.31%	95.50%

Cumulative Proportion of Variance captured by PCA

Scenario	PC1	PC2	PC3
6-months	89.6%	98.5%	99.7%
12-months	90.6%	97.7%	99.5%
24-months	89.1%	97.1%	99.4%



Results RQ2

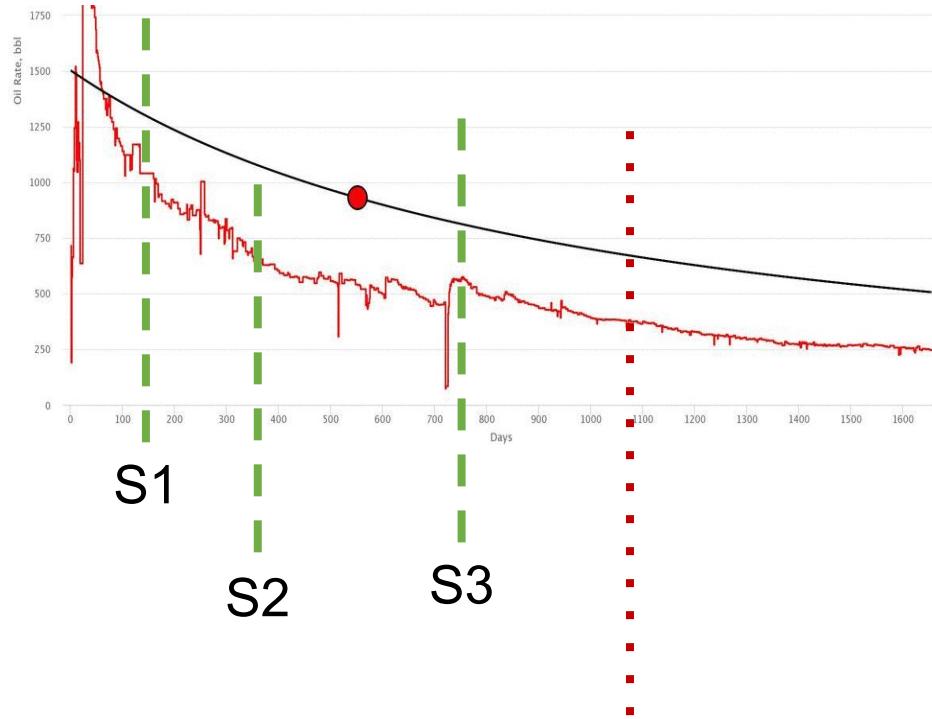


RQ3 Results

Fit Arps Model to each scenario of data

Predict target variable

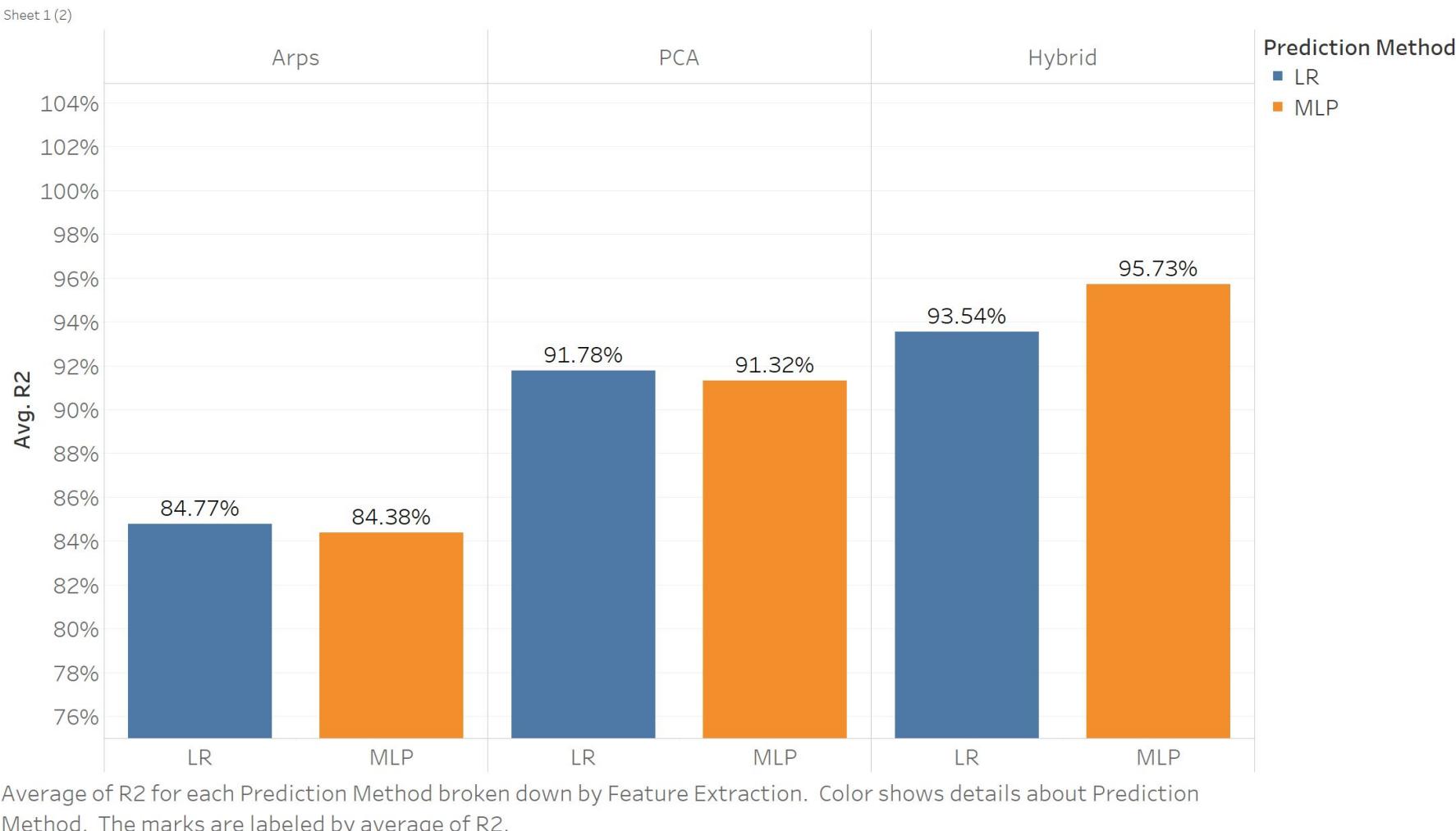
Assess Prediction quality



Scenario	RMSE	MAPE	R ²
S1: 6 Months	26084	28.3%	44.16%
S2: 12 Months	17130	16.23%	77.41%
S3: 24 Months	13263	8.81%	86.46%

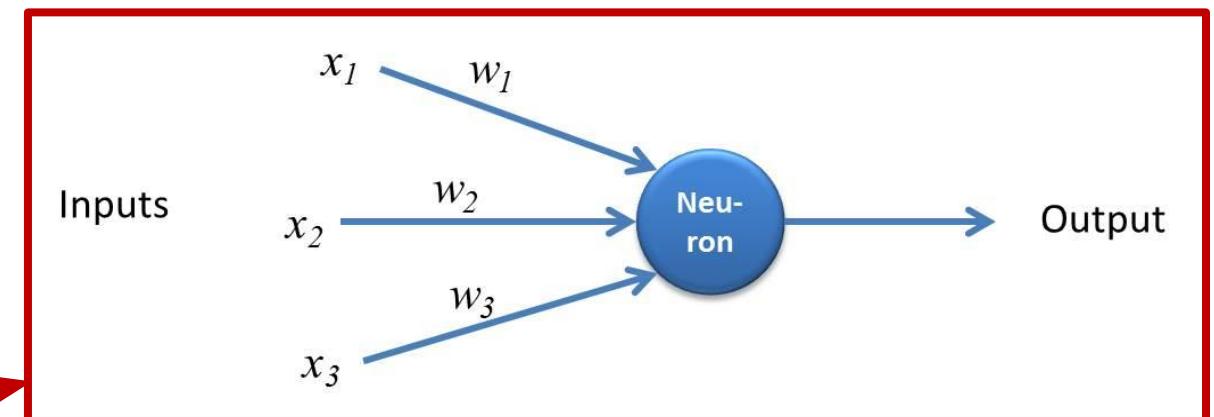
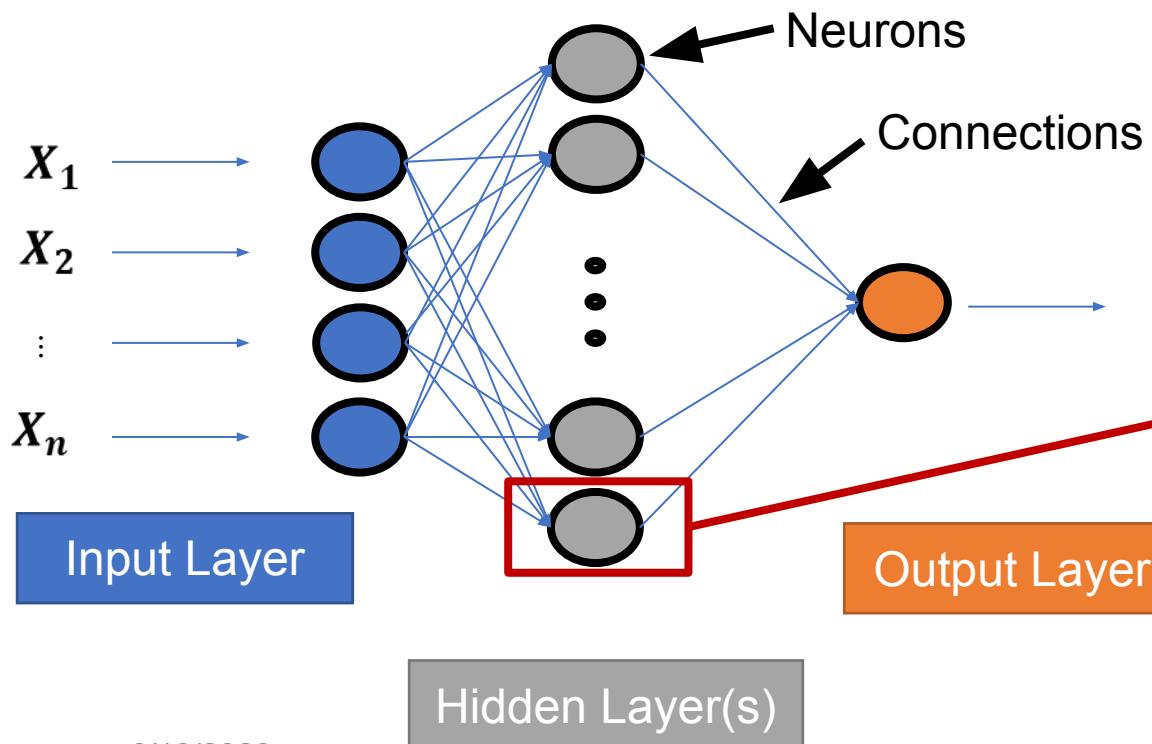
Target = 36-month oil

Appendix: RQ8 LR vs MLP by Feature Extraction



Artificial Neural Networks (ANN)

- Experimental branch of science to simulate human learning
 - Multilayer Perception (MLP) most famous and widely used ANN
 - Nonlinear approach to production forecasting

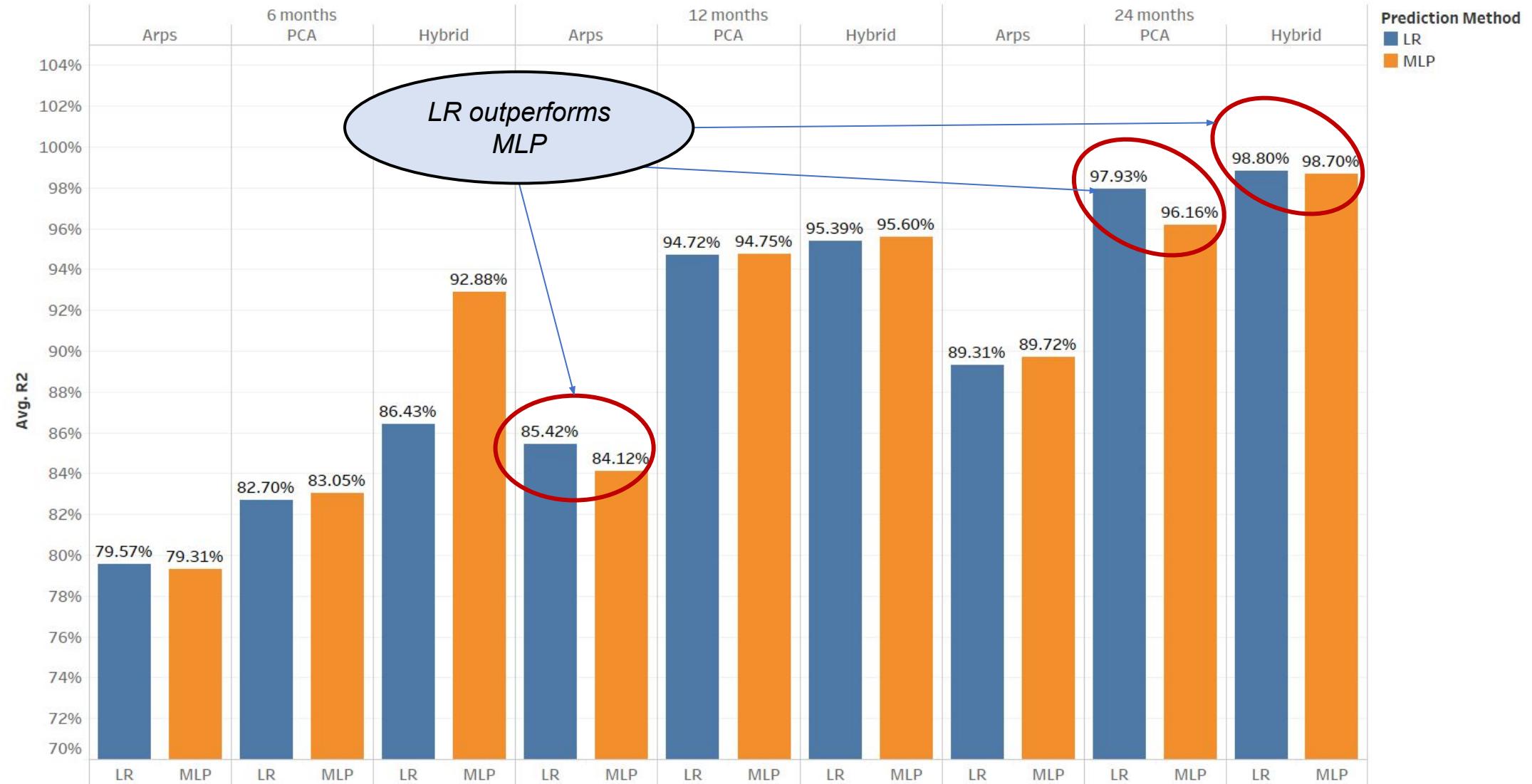


$$\mathbf{y} = \left(\sum_{i=1}^n w_i x_i + b \right) = \varphi(\mathbf{w}^T \mathbf{x} + b)$$

MLP Activation Function

RQ8 Results: Broken down by feature extraction and Scenario

A closer look at answering research question RQ8



LR Limitations

- Linear Regression limitations•
 - **Only Linear:**
 - Only capable of fitting the data to linear models.
 - – Non-linear modes might be able to explain and predict better.
 - **Sensitive to outliers:–**
 - The estimation of β_0 will be skewed significantly to account for a handful of outliers.
 - **Heavy design necessary:** The prediction performance heavily relies on the design of the regression model.
 - – Inclusion and exclusion of an attribute has significant impact on prediction performance.
 - **Limited prediction power:** The method is not designed for prediction.