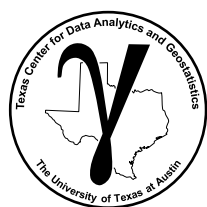


PGE 383

Tuning Hyperparameters

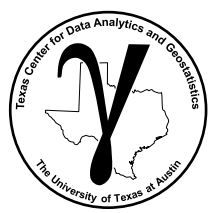
- **Training and Testing**
- **Model Goodness Metrics**
- **Cross Validation Workflows**

Michael Pyrcz, The University of Texas at Austin



Motivation

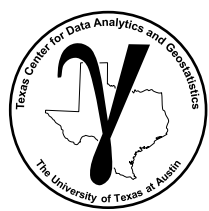
Let's formalize concepts and define terms for hyperparameter tuning.



PGE 383

Ridge Regression

- **Training and Testing**



Model Parameters Definition

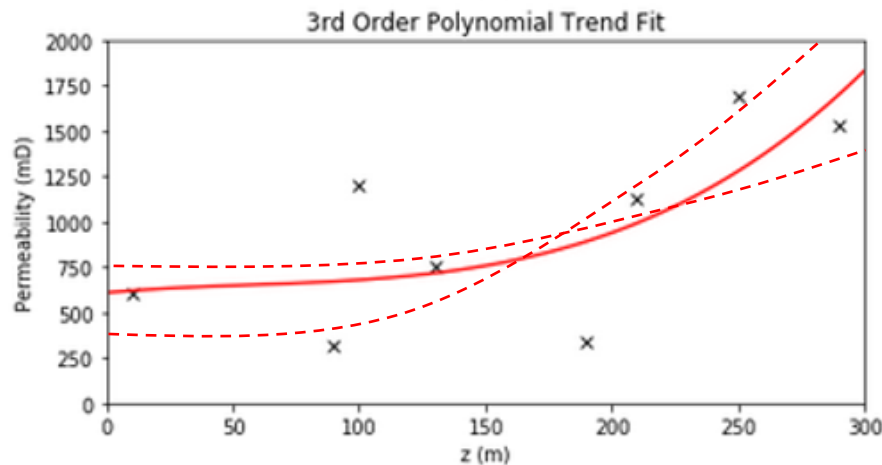
Model Parameters

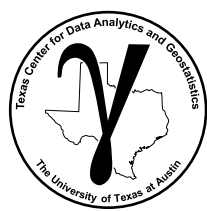
Derived during training phase to fit the model to the training data

Parameters

$$k = b_3 z^3 + b_2 z^2 + b_1 z + c$$

b_3, b_2, b_1 and c





Model Hyperparameter Definition

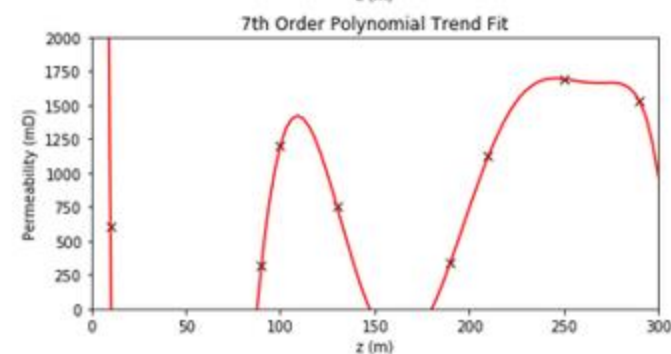
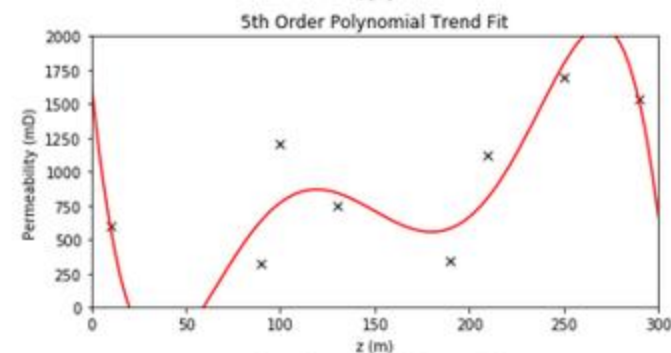
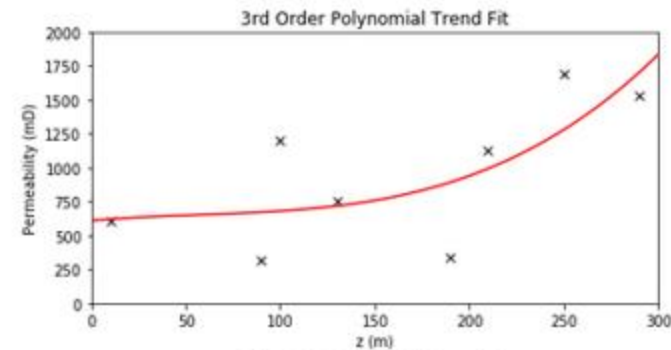
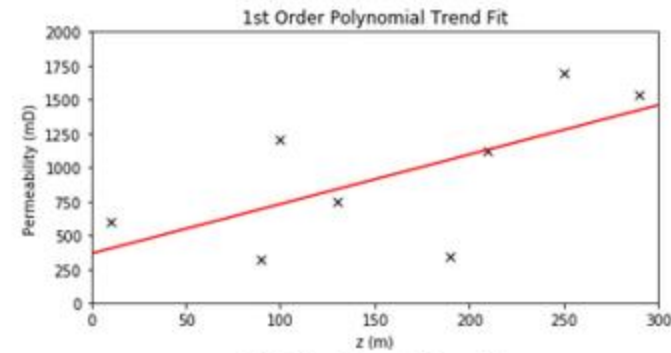
Model Hyperparameters

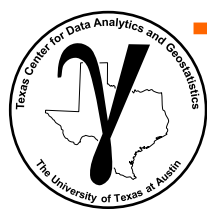
Set prior to learning from the data. Impact the form of the model and often the complexity.

3rd Order: $k = b_3 z^3 + b_2 z^2 + b_1 z + c$

2nd Order: $k = b_2 z^2 + b_1 z + c$

1st Order: $k = b_1 z + c$

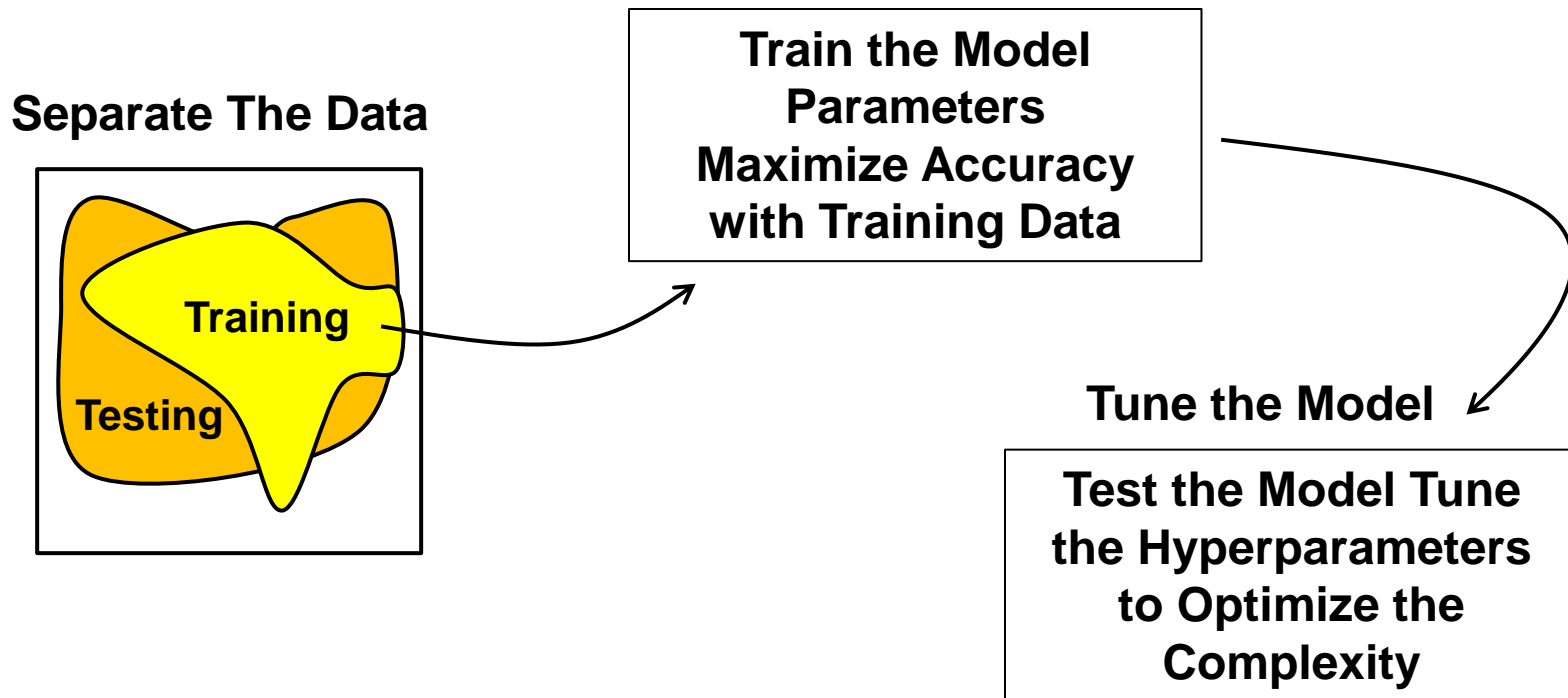




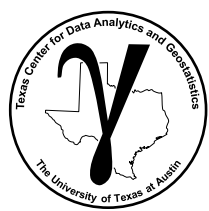
Training and Testing

The Training and Testing Workflow

- establish a subset of the data for fair testing of the model



We avoid the overfit problem.

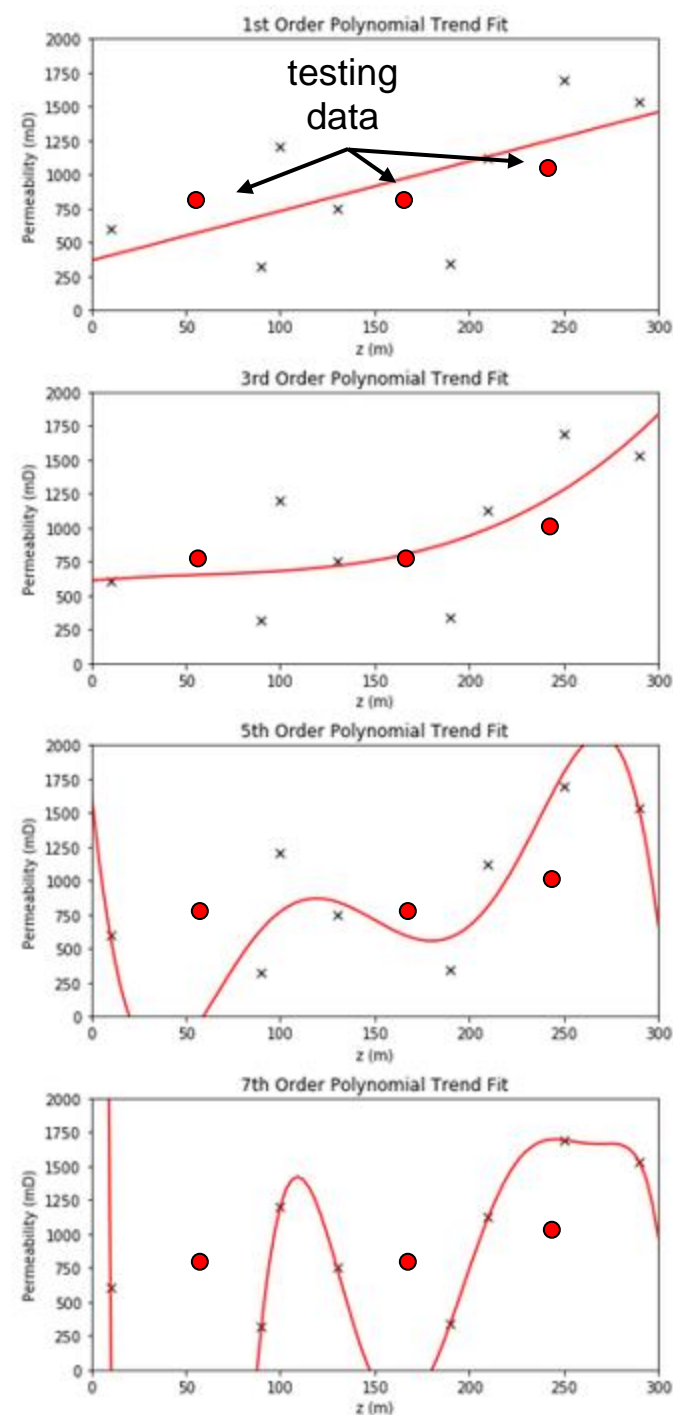


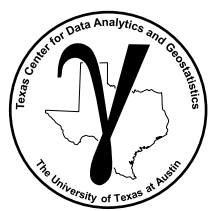
Hyperparameter Tuning

What do we have?

- A suite of models of variable level of complexity, and other decisions informed by a range of hyperparameters.
- We have a set of testing data withheld from the training of the model parameters.

Suite of Models with Increasing Complexity

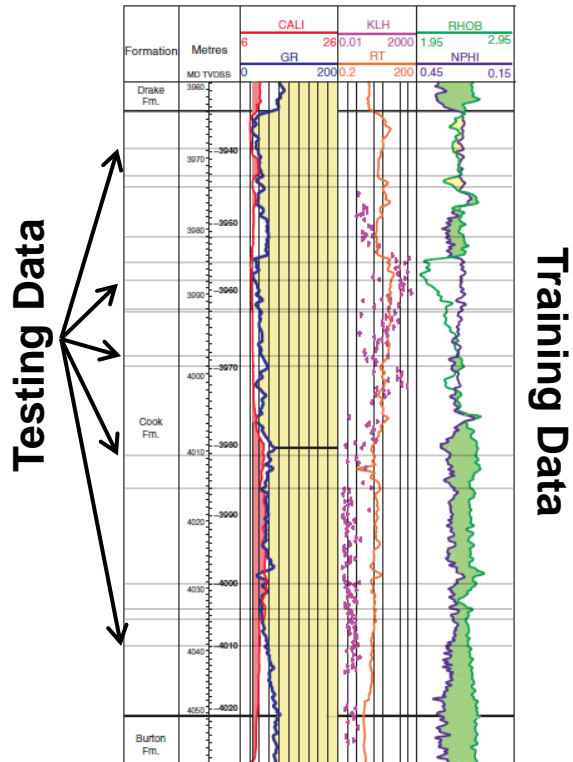




Training and Testing Split

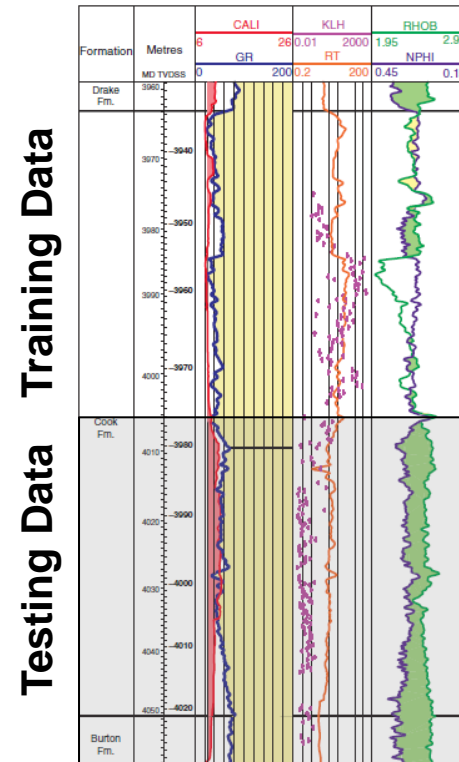
Fair Testing in Spatial / Temporal Settings

Too Easy



Predictions only at ½ ft offsets

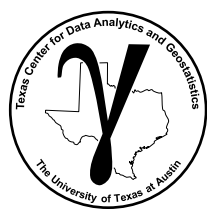
Too Hard



Predictions in a different rock.

The Train and Test Split Should Be Fair

The prediction difficulty (interpolation, extrapolation) should be similar to the planned real-world use of the prediction model.

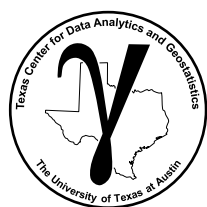


Training and Testing Split

How Much Data Should be Used in Testing?

The proportion in testing is recommended by various sources from 30% - 15% of the total dataset.

- Data withheld for testing reduces the data available for training; therefore, reduces the accuracy of the model.
- Data withheld for testing improves the accuracy of the assessment of the model performance.
- Various authors have experimented on a variety of training and testing ratios and have recommended splits for their applications:
 - The optimum ratio of training and testing split depends on problem setting
 - Could consider the difficulty in training (e.g. the number of parameters) and the difficulty in testing (e.g. number of hyperparameters).



Alternative Testing Workflows

Training, Validation and Testing

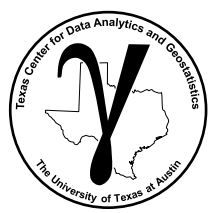
There is a more complete workflow commonly applied.

Note: to avoid confusion in our class we will use the train and test approach only.

- **Train with training data.** Models sees and learns from this data to train the model parameters.
- **Validate with validation data.** Unbiased evaluation of model fit to tune the model hyperparameters.
- **Test with testing data.** Data withheld until the model is complete to provide a final evaluation. Commonly applied to compare multiple competing models. This data had no role in building the model.



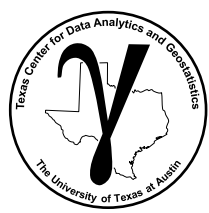
A visualisation of the splits



PGE 383

Ridge Regression

- **Model Goodness Metrics**

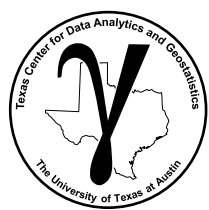


Testing Metrics

Testing Metrics

To evaluate model performance we need to assess the goodness of the suite of models with respect to the testing data.

- There are a variety of metrics that are applied for this assessment.
- They depend primarily on classification vs. regression.
- We will cover regression metrics first.



Regression Testing Metrics

Regression Testing Metrics

Mean Square Error (MSE)

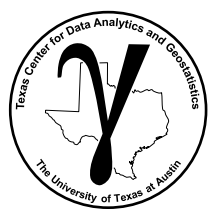
- L^2 norm – sensitive to large errors

$$\text{Test MSE} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (y_i - \hat{y}_i)^2 = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (\Delta y_i)^2$$

Mean Absolute Error (MAE)

- L^1 norm – less sensitive to large errors

$$\text{Test MAE} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} |y_i - \hat{y}_i| = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} |\Delta y_i|$$



Regression Testing Metrics

Regression Testing Metrics

Variance Explained – note, **we only use it for linear models**

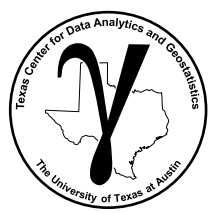
- Proportion of variance of the response feature captured by the model
- Takes advantage of the additivity of variance
 - Total Variance = Variance Explained + Variance Not Explained

$$\sigma_{explained}^2 = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} (\hat{y}_i - \bar{y})^2 \quad \sigma_{not\ explained}^2 = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} (y_i - \hat{y}_i)^2$$

$$r^2 = \frac{\sigma_{explained}^2}{\sigma_{explained}^2 + \sigma_{not\ explained}^2} = \frac{\sigma_{explained}^2}{\sigma_{total}^2}$$

Issues

- recall $r^2 = (\rho)^2$, suffers the same issues of correlation coefficients, linearity, do not account for redundancy, sensitive to outliers.
- e.g. adding outliers and 2 populations make your model look better!



Classification Testing Metrics

Classification Testing Metrics

Confusion Matrix

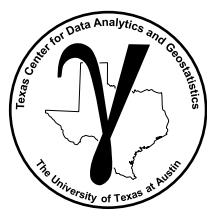
- Matrix with the categorical truth vs. predicted
- Visualize and diagnose all the combinations of correct and misclassification with the classification model.

Truth	$C_{k=1}$	$C_{k=2}$	$C_{k=3}$
	15	15	2
	5	22	9
	7	15	4
	$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$
	Predicted		

Model says category 3,
Data value is category 1.

- Perfect accuracy is number of each class in the training data on the diagonal.

Truth	n_1	0	0
	0	n_2	0
	0	0	n_3
	Predicted		



Classification Testing Metrics

Precision – for group k , the ratio of true positive over all positives.

$$precision_k = \frac{n_k \text{ true positive}}{n_k \text{ true positive} + n_k \text{ false positive}}$$

$precision_k = \text{Probability}(k \text{ is happening} \mid \text{the model says } k \text{ is happening})$

$k = 1$

	true positive	false positive	
Truth			
$C_{k=1}$	15	15	2
$C_{k=2}$	5	22	9
$C_{k=3}$	7	15	4
	$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

$$\frac{15}{15 + (5 + 7)} = \frac{15}{27} = 0.56$$

$k = 2$

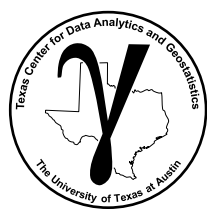
	true positive	false positive	
Truth			
$C_{k=1}$	15	15	2
$C_{k=2}$	5	22	9
$C_{k=3}$	7	15	4
	$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

$$\frac{22}{22 + (15 + 15)} = \frac{22}{52} = 0.42$$

$k = 3$

	true positive	false positive	
Truth			
$C_{k=1}$	15	15	2
$C_{k=2}$	5	22	9
$C_{k=3}$	7	15	4
	$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

$$\frac{4}{4 + (2 + 9)} = \frac{4}{15} = 0.27$$



Classification Testing Metrics

Recall – for group k , the ratio of true positive over all cases of k .

$$Recall_k = \frac{n_k \text{ true positive}}{n_k}$$

How many of group k did we catch? Does not account for false positives.

$k = 1$

true positive

Truth	$C_{k=1}$	15	15	2
	$C_{k=2}$	5	22	9
	$C_{k=3}$	7	15	4
		$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

$$\frac{15}{15 + (15 + 2)} = \frac{15}{32} = 0.47$$

$k = 2$

true positive

Truth	$C_{k=1}$	15	15	2
	$C_{k=2}$	5	22	9
	$C_{k=3}$	7	15	4
		$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

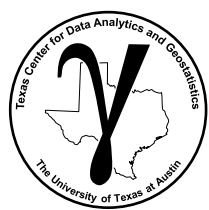
$$\frac{22}{22 + (5 + 9)} = \frac{22}{36} = 0.61$$

$k = 3$

true positive

Truth	$C_{k=1}$	15	15	2
	$C_{k=2}$	5	22	9
	$C_{k=3}$	7	15	4
		$\hat{C}_{k=1}$	$\hat{C}_{k=2}$	$\hat{C}_{k=3}$

$$\frac{4}{4 + (7 + 15)} = \frac{4}{26} = 0.15$$



Classification Testing Metrics

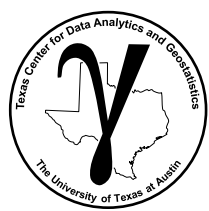
	Precision	Recall	f1-score
$k = 1$	0.56	0.47	0.51
$k = 2$	0.42	0.61	0.49
$k = 3$	0.27	0.15	0.19

Precision and Recall measure 2 components of categorical accuracy. Let's combine them into one measure.

$$f1 - score_k = \frac{2}{\frac{1}{Precision_k} + \frac{1}{Recall_k}}$$

$f1 - score$ is the Harmonic mean of precision and recall for k .

- Sensitive the to lowest score, good performance in one score cannot make up for bad performance in the other!

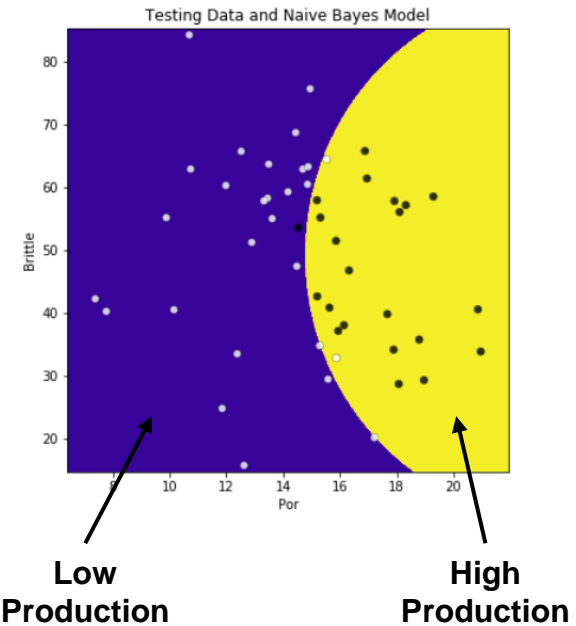


Classification Testing Metrics

Classification Testing Metrics

Another example from the naïve Bayes workflow.

Truth Low High
 High [1 21]
 Low High
 Predicted

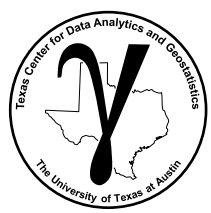


$$\bullet \text{ Precision}_k = \frac{n_k \text{ true positive}}{n_k \text{ true positive} + n_k \text{ false positive}} = \frac{26}{26+1}$$

$$\bullet \text{ Recall}_k = \frac{n_k \text{ true positive}}{n_k} = \frac{26}{28}$$

	precision	recall	f1-score	support
Low	0.96	0.93	0.95	28
High	0.91	0.95	0.93	22
avg / total	0.94	0.94	0.94	50

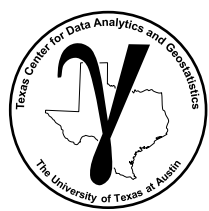
$$\bullet \text{ f1-score}_k = \frac{2}{\frac{1}{\text{Precision}_k} + \frac{1}{\text{Recall}_k}}$$



PGE 383

Tuning Hyperparameters

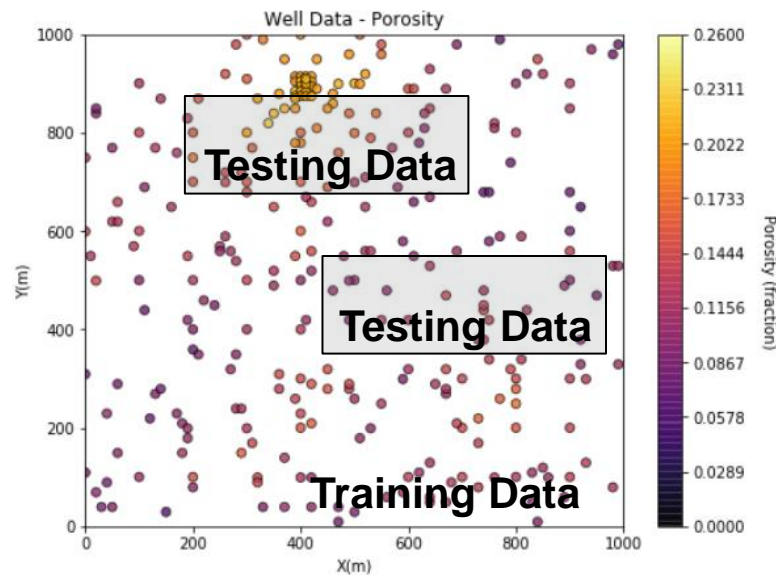
- **Cross Validation Workflows**



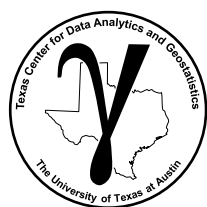
Model Cross Validation Workflows

Cross Validation

- Withhold subset of the data during model training
- Then testing the trained model with withheld subset dataset
- Must make sure cross validation is fair
- Training data set (used for training), Testing data set (withheld for testing)



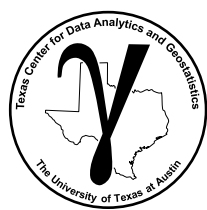
Predictions in a distinctly different range of reservoir values.



Model Cross Validation Workflows

Jackknife

- Combines training and testing into one step
- Loop over all data, withhold that data
 - Train on $n - 1$ data and test on the withheld single data
 - Calculate model goodness metric
- Aggregate over all data, n
- Typically too easy of an estimation problem
- K-fold is a more general and robust approach



Model Cross Validation Workflows

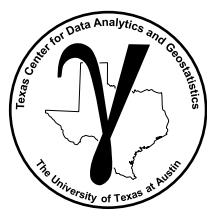
k-fold Approach

- Select k , integer number of folds
- Break data set into k subsets, equal size n/k
- Loop over k subsets:
 - use data outside the k subset to predict inside the k subset
 - calculate the model goodness metric
- Aggregate metric over all subsets, k

	1	Sample Data				n
k=1	Test					$Metric_1$
		Test				
			Test			
				Test		
k=5					Test	$Metric_5$

Cycle test over k folds, all other data are train.

Metric Aggregation



Model Cross Validation Workflows Limitations

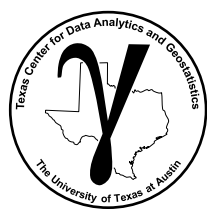
Methods in the Python scikit-learn package:

`sklearn.model_selection.train_test_split(X,y,train_size,random_state)`

- random selection for training and testing
- specify `train_size` or test size (other is defaulted as complement)
- use `random_state` for repeatability
- stratify will enforce whole sample statistics in the train and test sample subsets

`sklearn.cross_val_score(model_object,X,y,cv,scoring)`

- wrapper that performs k-fold cross validation with k random subsets
- k is the `cv` parameter
- scoring allows assignment of a custom model accuracy metric

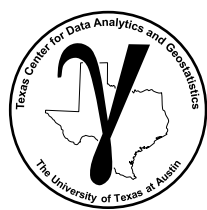


Model Cross Validation Workflows Limitations

Methods in the Python scikit-learn package:

`sklearn.cross_validate(model_object,X,y,cv,scoring)`

- wrapper that performs k-fold cross validation with k random subsets
- k is the cv parameter
- scoring allows assignment of a custom model in training accuracy metric
- scoring parameter may be a list of model accuracy metrics
- outputs include fit and scoring times, estimates and all scores in a dictionary object



Model Cross Validation Workflows Limitations

Issues with Cross Validation

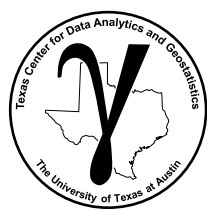
Peeking, information leakage – some information is transmitted from the withheld data into the model, some model decision(s) use all the data

Fair Train and Test Split – many practitioners use random selection for train and test split (we use it, it is built into scikit-learn) and this may be too easy of a prediction problem

Black Swans / Stationarity – the model cannot be tested for data events not available in the data

This is also known as the '**No Free Lunch Theorem**' in machine learning

'even after the observation of the frequent or constant conjunction of objects, we have no reason to draw any inference concerning any object beyond those of which we have had experience' - Hume (1739–1740)

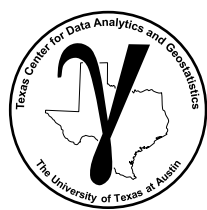


Model Cross Validation Workflows Limitations

Issues with Cross Validation

Subsurface Validation – it is not possible to validate open earth systems – Oreskes et al., 1994. Here's the abstract from their paper:

'Verification and validation of numerical models of natural systems is impossible. This is because natural systems are never closed and because model results are always nonunique. Models can be confirmed by the demonstration of agreement between observation and prediction, but confirmation is inherently partial. Complete confirmation is logically precluded by the fallacy of affirming the consequent and by incomplete access to natural phenomena. Models can only be evaluated in relative terms, and their predictive value is always open to question. The primary value of models is heuristic.'



Model Cross Validation Workflows Limitations

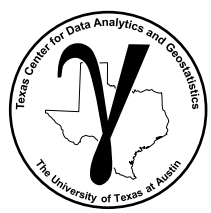
Issues with Cross Validation

‘All models are wrong, but some are useful’ – George Box

Parsimony – since all models are wrong, an economical description of the system. Occam’s Razor

Worrying Selectively – since all models are wrong, figure out what is most importantly wrong.

‘Be humble, the earth will surprise you!’ – me.



PGE 383

Tuning Hyperparameters

- **Training and Testing**
- **Model Goodness Metrics**
- **Cross Validation Workflows**

Michael Pyrcz, The University of Texas at Austin