

# Generalization

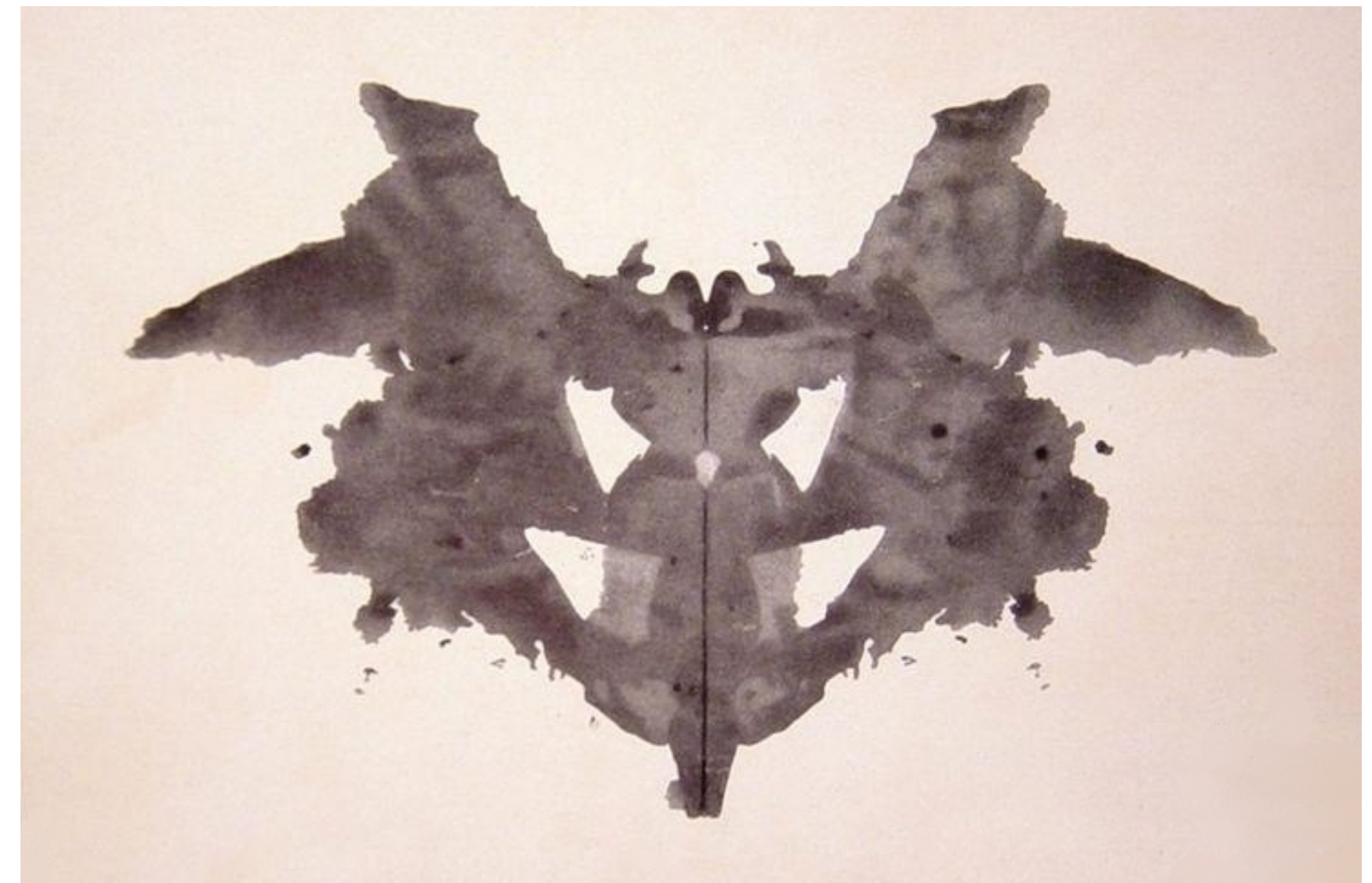
Bert Huang  
Data Intensive Studies Center  
Tufts University

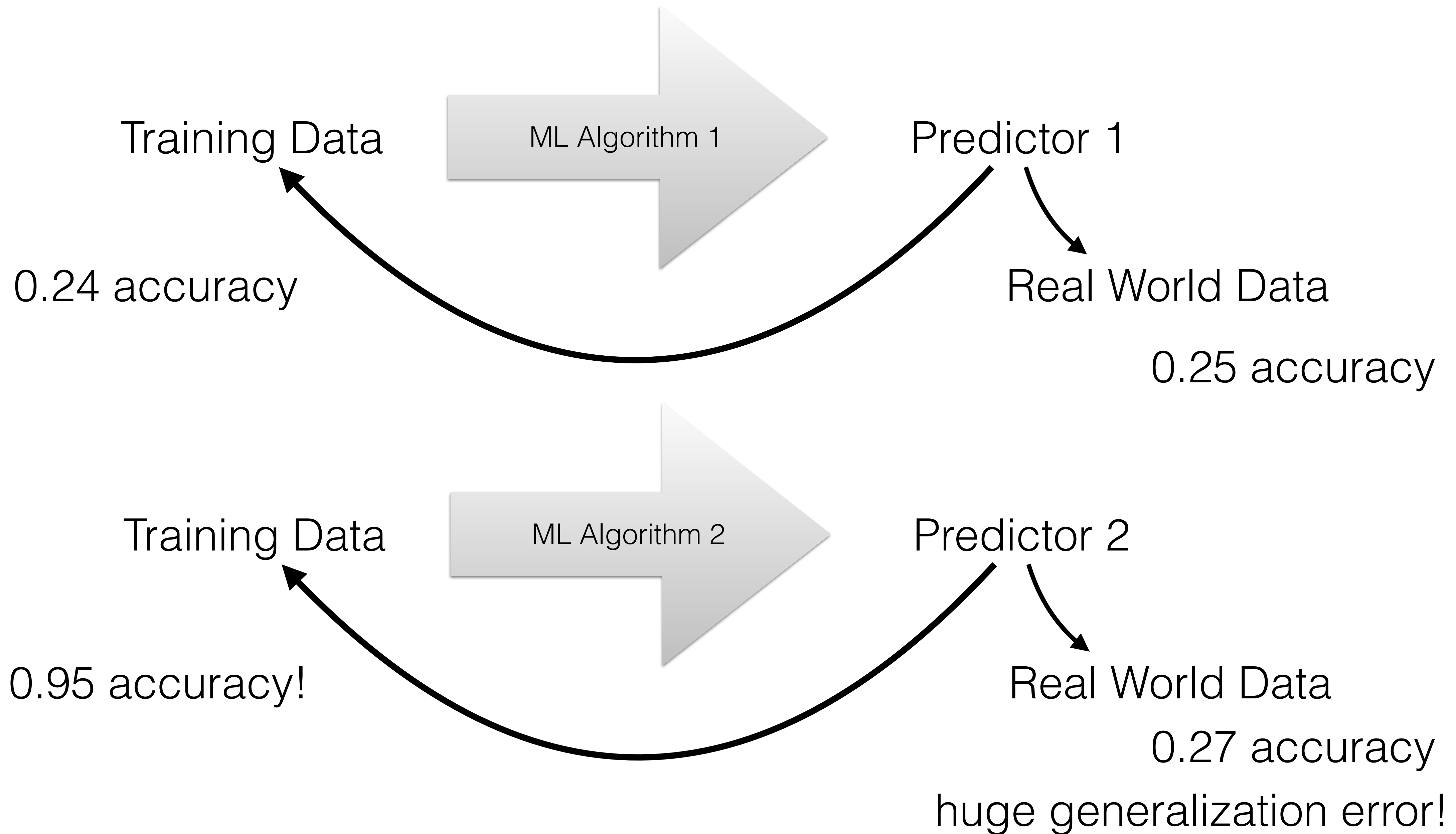
# Outline

- Introduction
- Bias and Variance (interactive website)
- Cross-Validation
- Relationship to my research

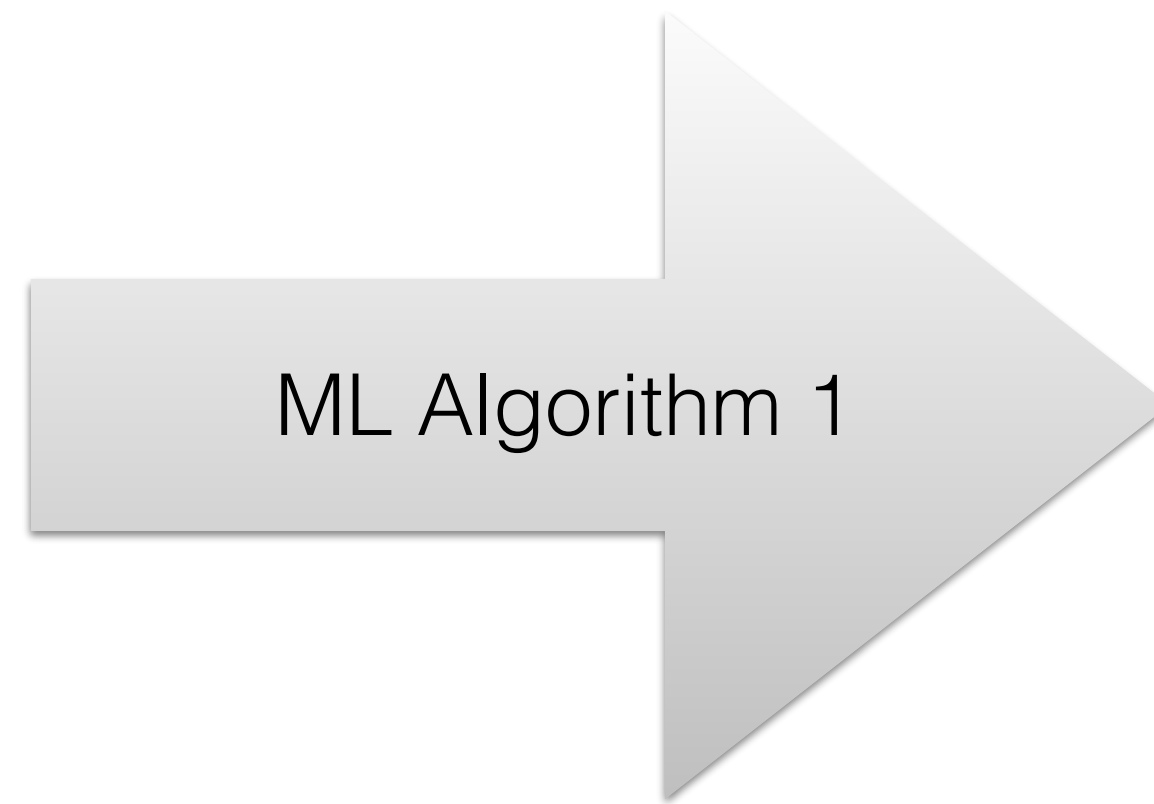
# Complexity in Machine Learning

- Relates to issues in philosophy of science, induction
  - Occam's Razor, etc.
- Baseball (sports) statistics



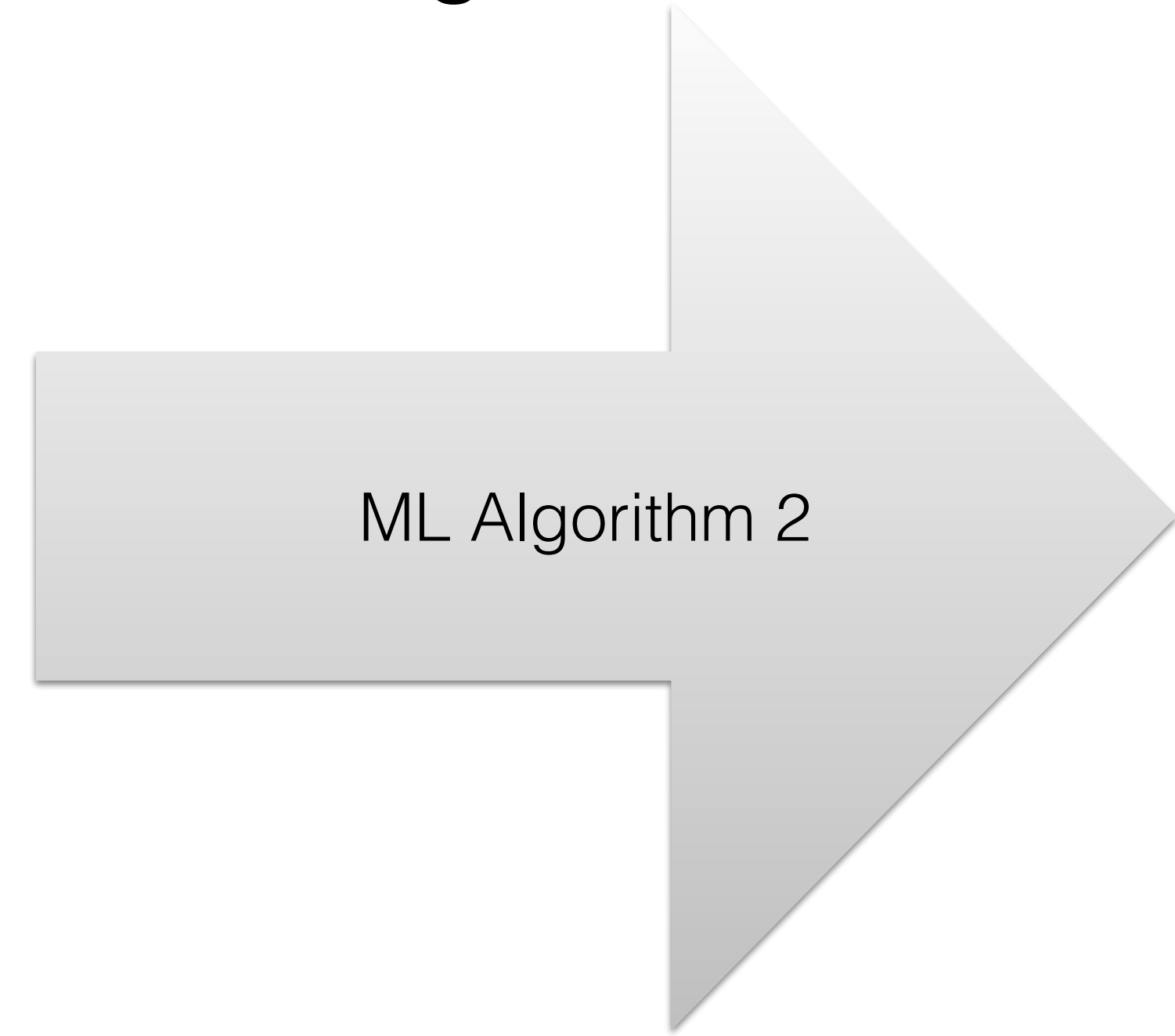


## Underfitting

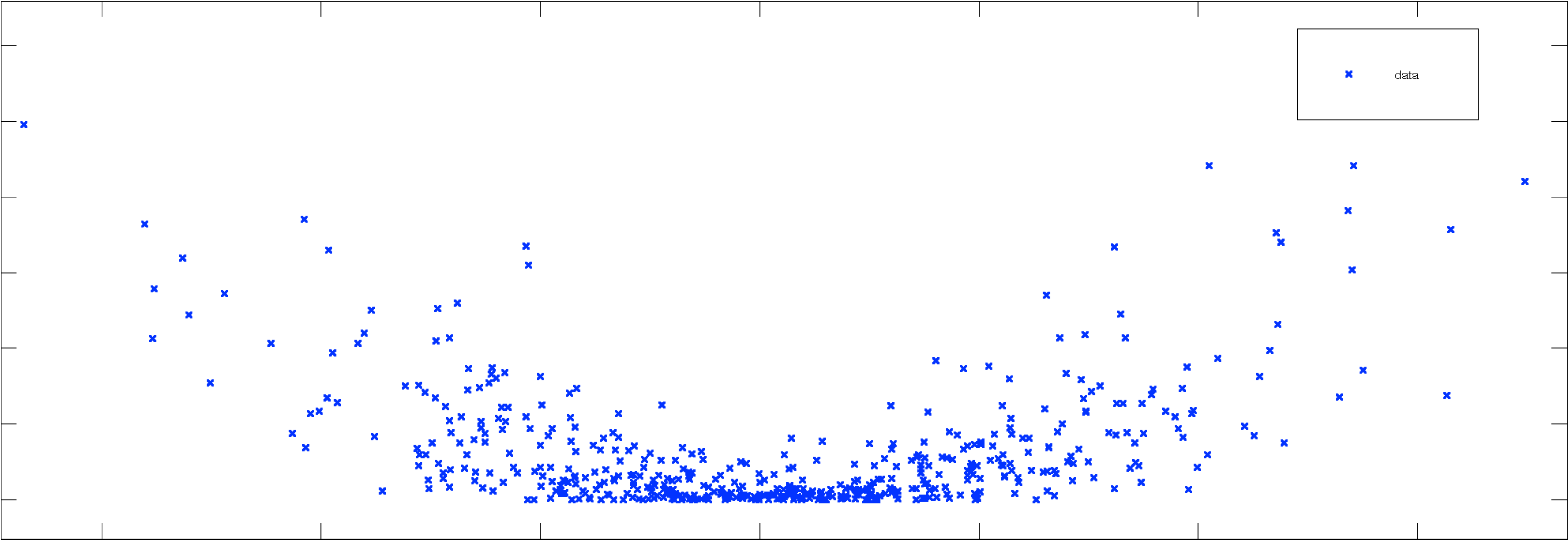


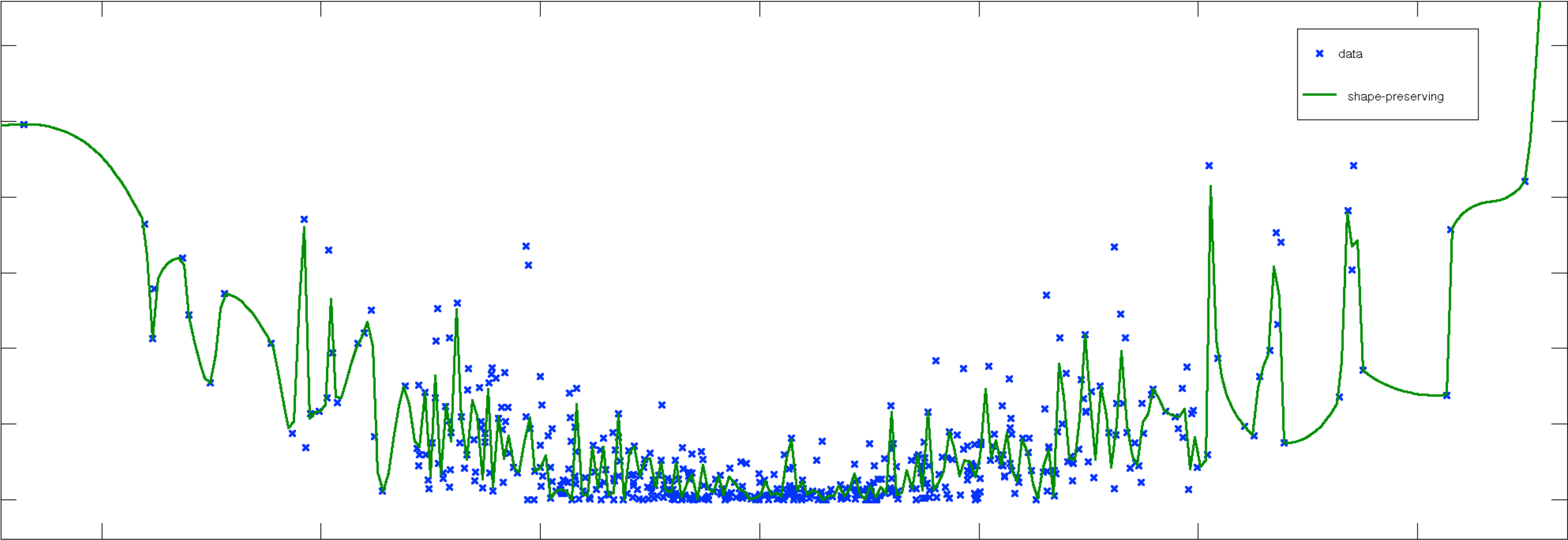
- Low dimensional
- Heavily regularized
- Bad modeling assumptions

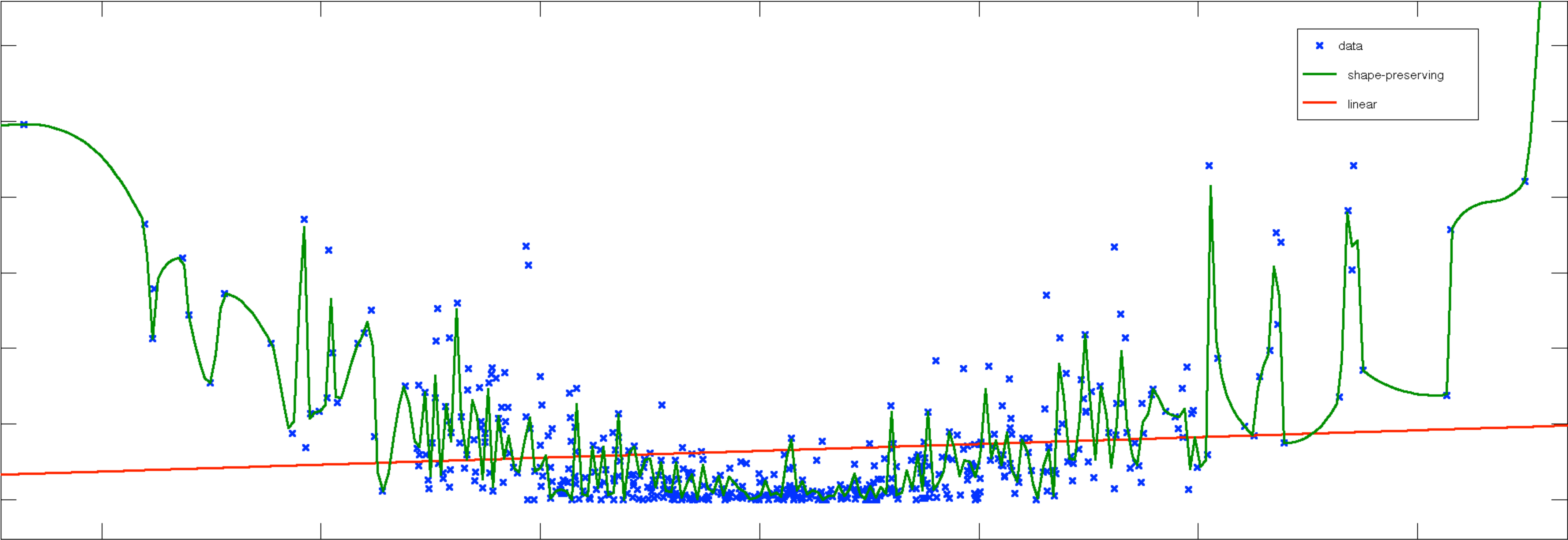
## Overfitting



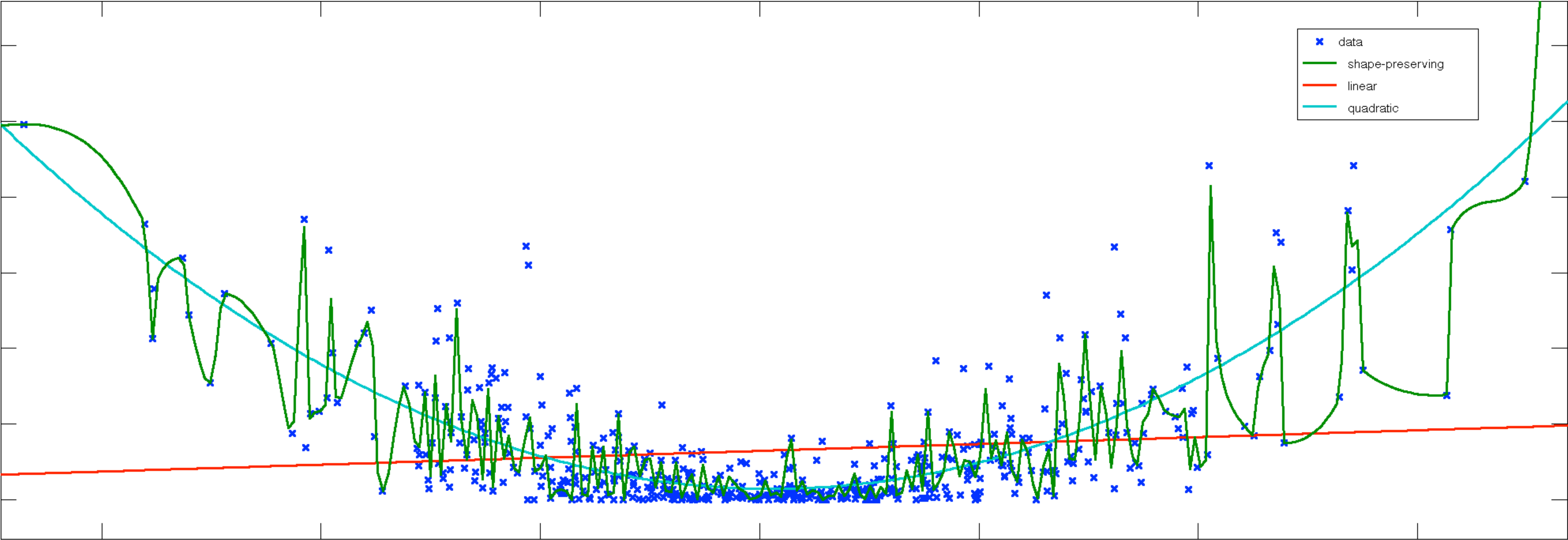
- High dimensional or non-parametric
- Weakly regularized
- Not enough modeling assumptions
- Not enough data











# Overfitting and Underfitting

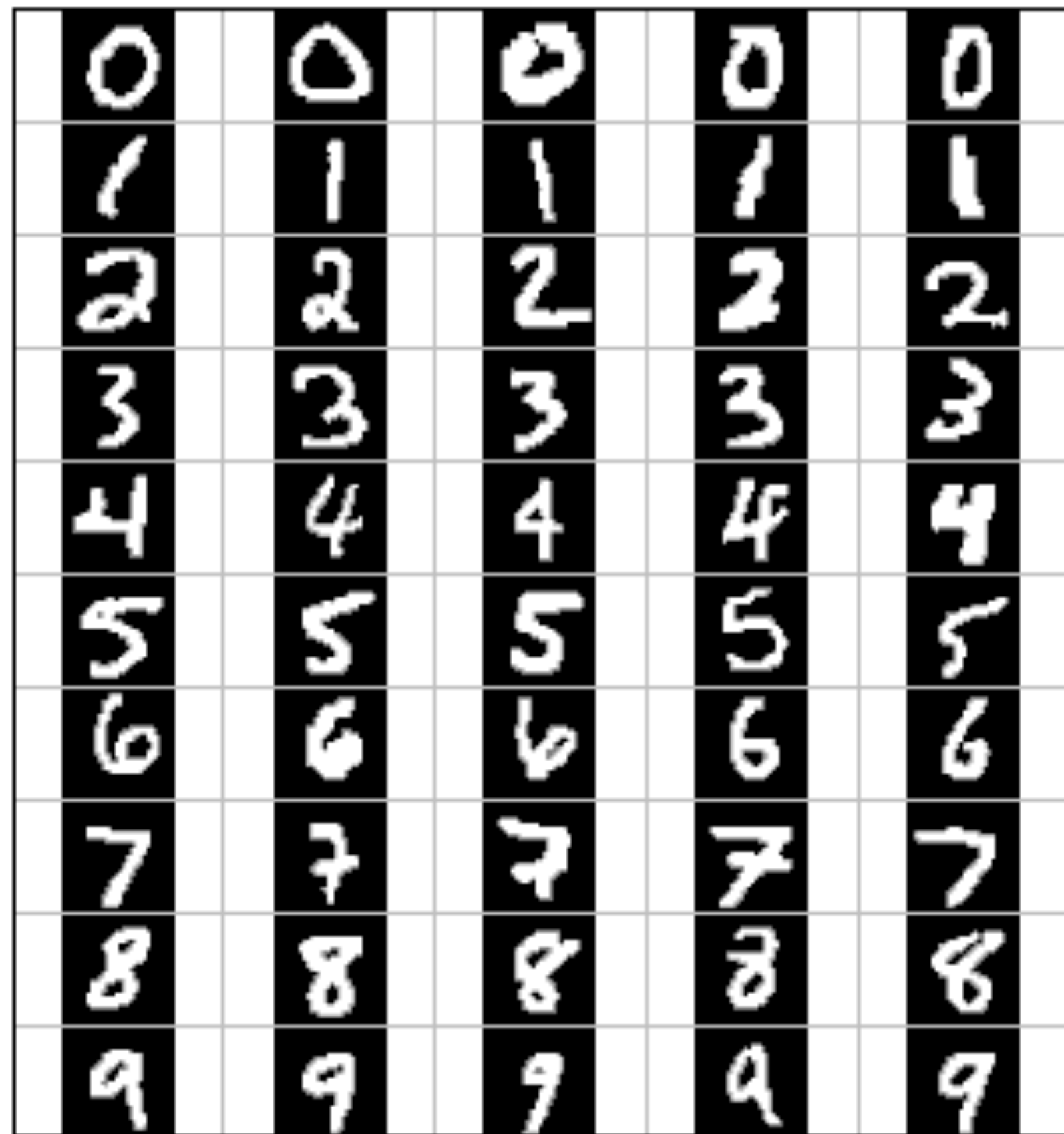
- Training models too complex can cause overfitting
- Training models too simple (or wrong) can cause underfitting

# Outline

- Overfitting and underfitting
- Bias and variance
- Validation for model selection

<http://www.r2d3.us/visual-intro-to-machine-learning-part-2/>

# Nearest-Neighbor Classifiers



classifier = {

 : 0,

 : 0,

 : 0,

 : 0,

 : 0,

 : 1,

 : 1,

...

}

100% training accuracy!



53% testing accuracy...

# Held-out Validation

	0			0			0			0			0		
	1			1			1			1			1		
	2			2			2			2			2		
	3			3			3			3			3		
	4			4			4			4			4		
	5			5			5			5			5		
	6			6			6			6			6		
	7			7			7			7			7		
	8			8			8			8			8		
	9			9			9			9			9		

# Held-out Validation

0	0	0	0
1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4
5	5	5	5
6	6	6	6
7	7	7	7
8	8	8	8
9	9	9	9

training data

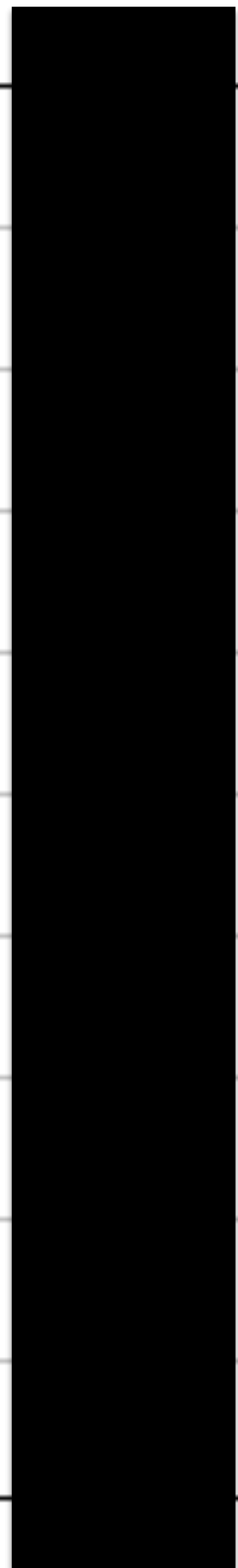
	Accuracy on training data	Accuracy on validation data
Simple	0.91	0.83
Medium	0.95	0.88
Complex	0.99	0.79
Super Complex	1.0	0.54

0
1
2
3
4
5
6
7
8
9

validation data

# Cross Validation

Fold 1



0	0	0	0
1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4
5	5	5	5
6	6	6	6
7	7	7	7
8	8	8	8
9	9	9	9

training data

0
1
2
3
4
5
6
7
8
9

validation data



# Cross Validation

Fold 2

0				0				0				0			
1				1				1				1			
2				2				2				2			
3				3				3				3			
4				4				4				4			
5				5				5				5			
6				6				6				6			
7				7				7				7			
8				8				8				8			
9				9				9				9			

training data

0
1
2
3
4
5
6
7
8
9

validation data

# Cross Validation

Fold 3

0	0				0	0
1	1				1	1
2	2				2	2
3	3				3	3
4	4				4	4
5	5				5	5
6	6				6	6
7	7				7	7
8	8				8	8
9	9				9	9

training data

0
1
2
3
4
5
6
7
8
9

validation data

# Cross Validation

Fold 4

0	0	0
1	1	1
2	2	2
3	3	3
4	4	4
5	5	5
6	6	6
7	7	7
8	8	8
9	9	9

training data

0
1
2
3
4
5
6
7
8
9

validation data

# Cross Validation

Fold 5

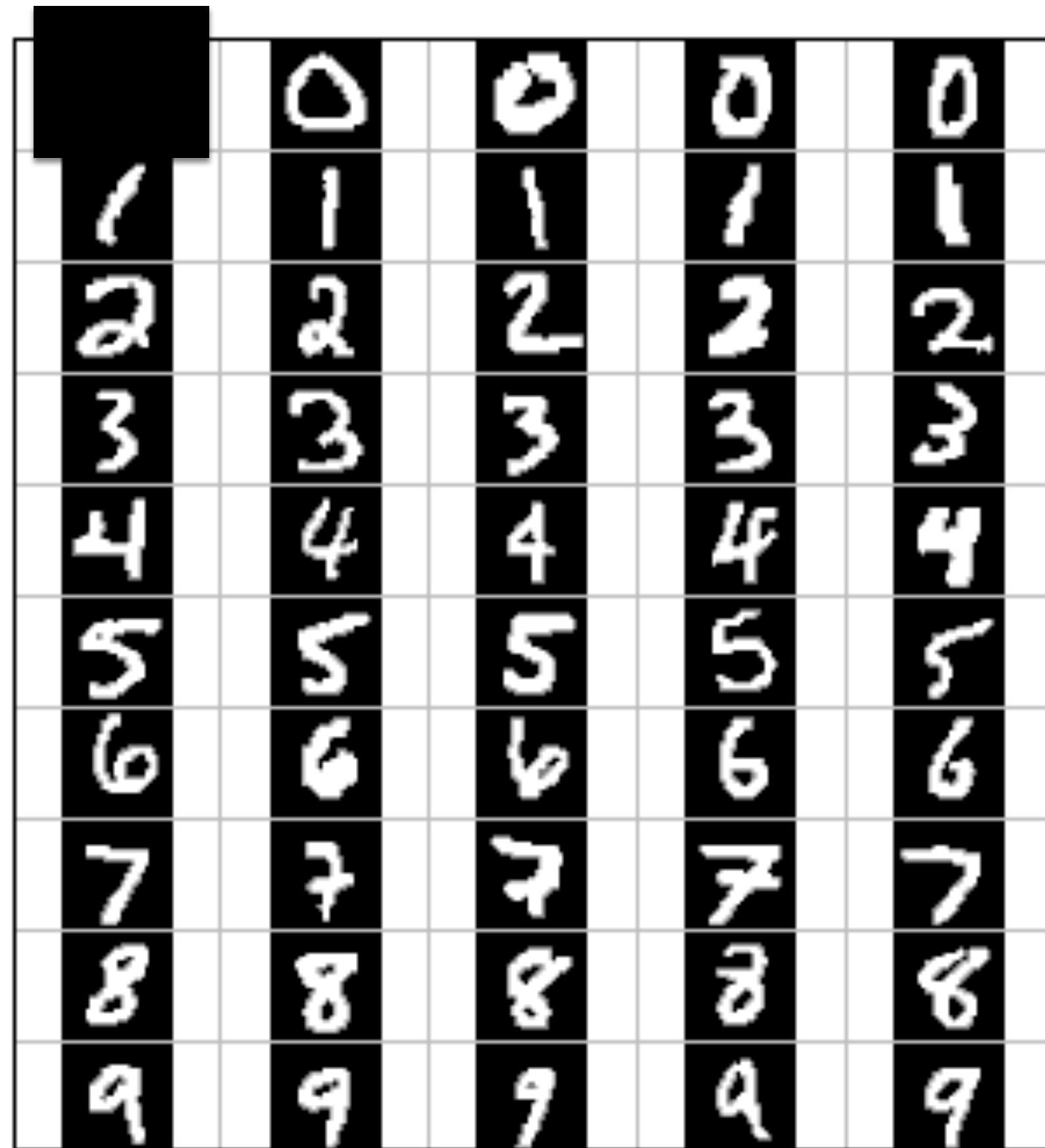
0	0	0	0
1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4
5	5	5	5
6	6	6	6
7	7	7	7
8	8	8	8
9	9	9	9

training data

0
1
2
3
4
5
6
7
8
9

validation data

# Leave-one-out Cross Validation



0	0	0	0	0
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

training data



validation data

# Leave-one-out Cross Validation

0				0			0			0
1				1			1			1
2				2			2			2
3				3			3			3
4				4			4			4
5				5			5			5
6				6			6			6
7				7			7			7
8				8			8			8
9				9			9			9

training data



validation data

# Leave-one-out Cross Validation

0	0				0	0	0
1	1				1	1	1
2	2				2	2	2
3	3				3	3	3
4	4				4	4	4
5	5				5	5	5
6	6				6	6	6
7	7				7	7	7
8	8				8	8	8
9	9				9	9	9

training data



validation data

# Leave-one-out Cross Validation

0	0	0	0	
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

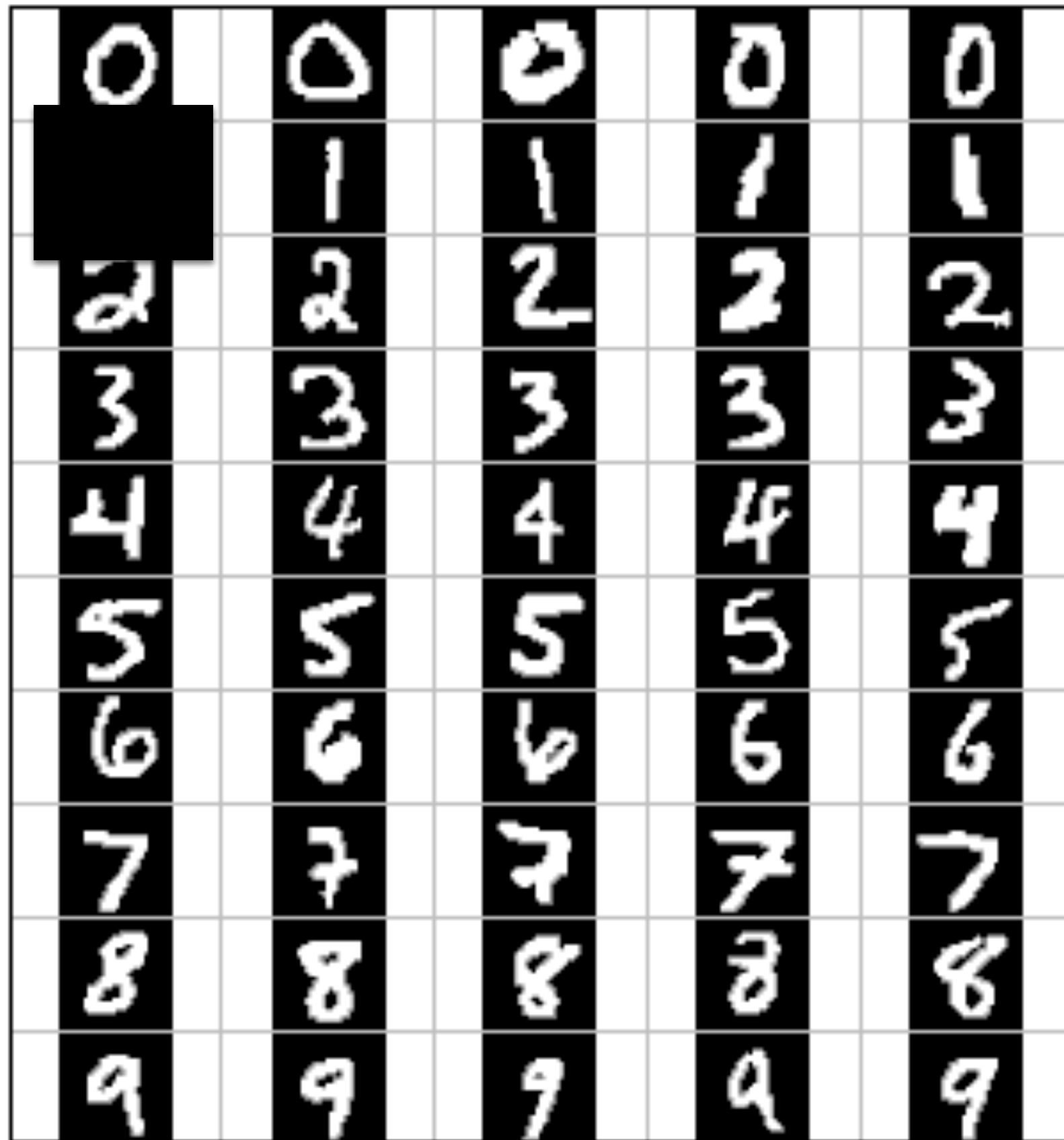
training data



validation data



# Leave-one-out Cross Validation



training data



validation data

# Leave-one-out Cross Validation

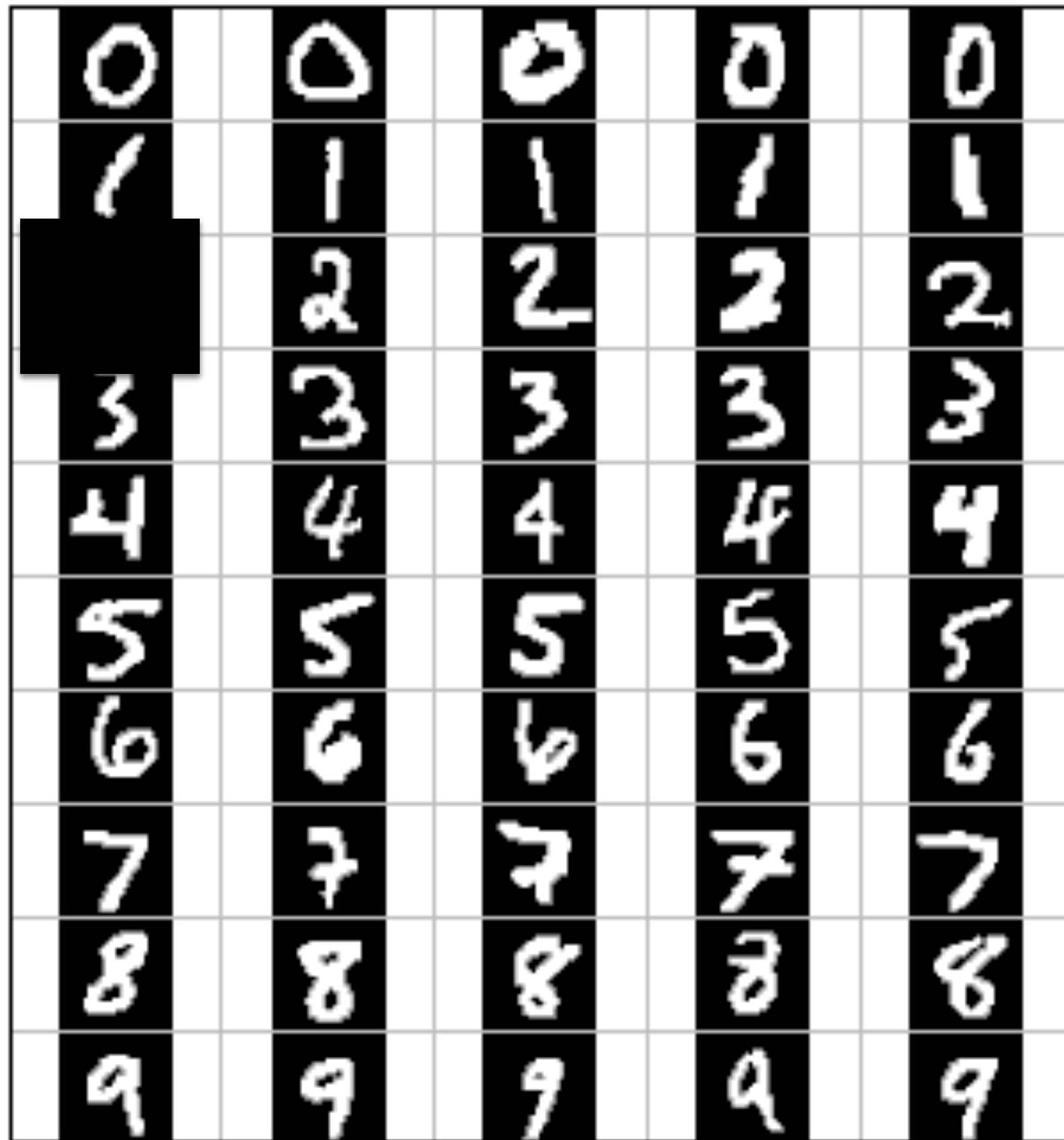
0	0	0	0	0
1	1	1	1	
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

training data



validation data

# Leave-one-out Cross Validation



training data



validation data

# How Many Folds?

- What are the pros and cons of leave-one-out cross-validation?
- We usually train on N-1 folds and test on 1 fold. What are pros and cons of doing the inverse: train on 1 fold and test on N-1 folds?

0	0	0	0	0
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

Training

0
1
2
3
4
5
6
7
8
9

Testing

# How Many Folds?

- What are the pros and cons of leave-one-out cross-validation?
- We usually train on  $N-1$  folds and test on 1 fold. What are pros and cons of doing the inverse: train on 1 fold and test on  $N-1$  folds?

0
1
2
3
4
5
6
7
8
9

Training

0	0	0	0	0
1	1	1	1	1
2	2	2	2	2
3	3	3	3	3
4	4	4	4	4
5	5	5	5	5
6	6	6	6	6
7	7	7	7	7
8	8	8	8	8
9	9	9	9	9

Testing

# Testing versus Validation

- Best practice for experiments:
  - Hold out test set completely hidden from training
  - Use validation on training data for model (or parameter) selection
  - Evaluate on held-out test data

# Model Selection via Validation

- Measure performance on **held-out** training data
  - Simulate testing environment
- Rotate **folds** of held-out subsets
- Can even hold out one at a time: **leave-one-out** validation
- Use (cross) validation performance to tune extra parameters

# Take-Away Points

- Overfitting and underfitting, bias and variance
  - bias -> modeling error, variance -> sampling error
  - Always have a mix of two
- Validation for model selection
- Reducing bias may need more complex models, but comes with challenges that researchers are working on solving