

Regression and generalization

Machine Learning

Hamid R Rabiee – Zahra Dehghanian

Spring 2025



Sharif University
of Technology

Topics

- Beyond linear regression models
- Evaluation & model selection
- Regularization



Recall: Linear regression (squared loss)

- ▶ Linear regression functions

$$f : \mathbb{R} \rightarrow \mathbb{R} \quad f(x; \mathbf{w}) = w_0 + w_1 x$$

$$f : \mathbb{R}^d \rightarrow \mathbb{R} \quad f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$ are the parameters we need to set.

- ▶ Minimizing the squared loss for linear regression

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

- ▶ We obtain $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Beyond linear regression

- How to extend the linear regression to non-linear functions?
 - Transform the data using basis functions
 - Learn a linear regression on the new feature vectors (obtained by basis functions)



Beyond linear regression

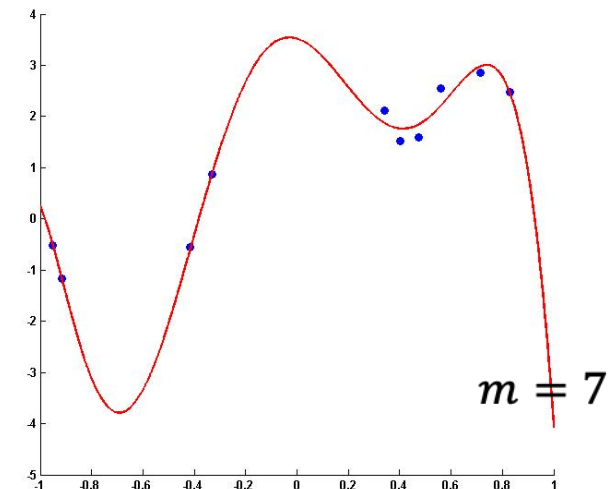
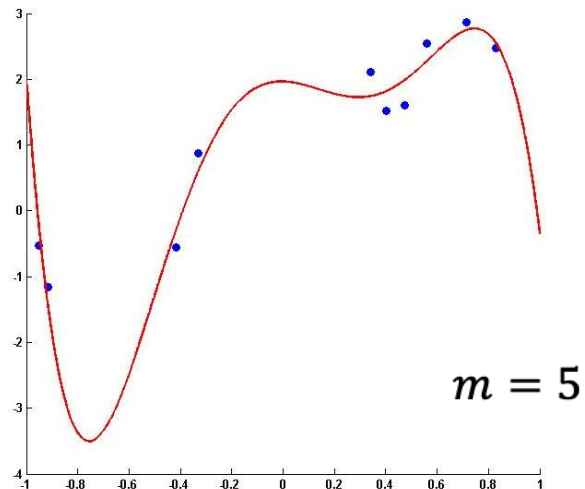
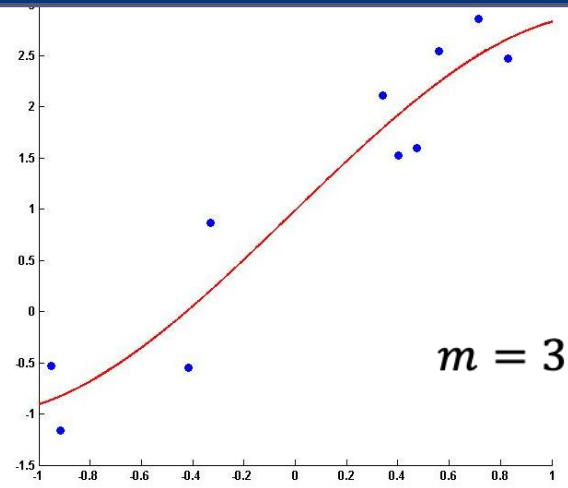
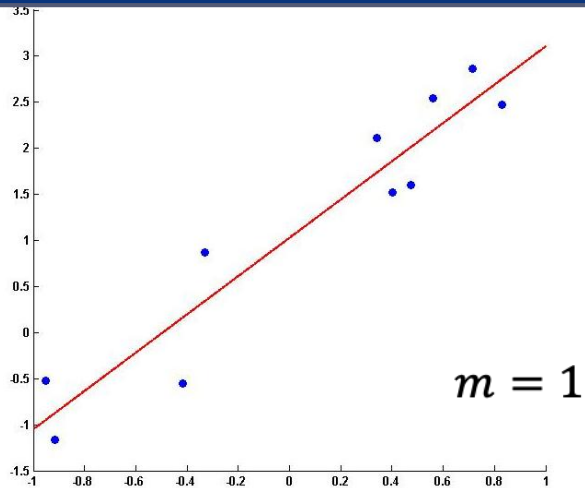
- m^{th} order polynomial regression (univariate $f : \mathbb{R} \rightarrow \mathbb{R}$)

$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{m-1}x^{m-1} + w_mx^m$$

- Solution: $\hat{\mathbf{w}} = (\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T \mathbf{y}$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X}' = \begin{bmatrix} 1 & x^{(1)1} & x^{(1)2} & \dots & x^{(1)m} \\ 1 & x^{(2)1} & x^{(2)2} & \dots & x^{(2)m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x^{(n)1} & x^{(n)2} & \dots & x^{(n)m} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} \hat{w}_0 \\ \hat{w}_1 \\ \vdots \\ \hat{w}_m \end{bmatrix}$$

Polynomial regression: example



Generalized linear

- Linear combination of fixed non-linear function of the input vector

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_m\phi_m(\mathbf{x})$$

$\{\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})\}$: set of basis functions (or features)

$$\phi_i(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$$



Basis functions: examples

- Linear

If $m = d$, $\phi_i(\mathbf{x}) = x_i$, $i = 1, \dots, d$, then

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + \dots + w_dx_d$$

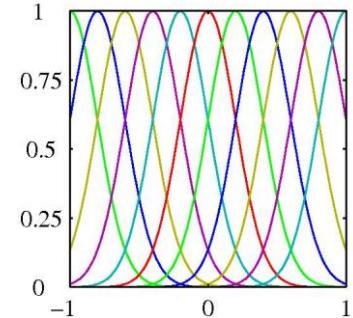
- Polynomial (univariate)

If $\phi_i(x) = x^i$, $i = 1, \dots, m$, then

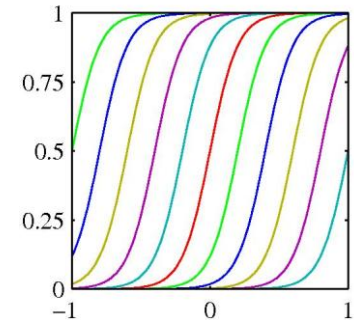
$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{m-1}x^{m-1} + w_mx^m$$

Basis functions: examples

- ▶ Gaussian: $\phi_j(\mathbf{x}) = \exp\left\{-\frac{(\mathbf{x}-\mathbf{c}_j)^2}{2\sigma_j^2}\right\}$



- ▶ Sigmoid: $\phi_j(\mathbf{x}) = \sigma\left(\frac{\|\mathbf{x}-\mathbf{c}_j\|}{\sigma_j}\right)$ $\sigma(a) = \frac{1}{1+\exp(-a)}$



Radial Basis Functions: prototypes

- Predictions based on similarity to “prototypes”:

$$\phi_j(\mathbf{x}) = \exp\left\{-\frac{1}{2\sigma_j^2}\|\mathbf{x} - \mathbf{c}_j\|^2\right\}$$

- Measuring the similarity to the prototypes $\mathbf{c}_1, \dots, \mathbf{c}_m$
 - σ^2 controls how quickly it vanishes as a function of the distance to the prototype.
 - Training examples themselves could serve as prototypes

Model complexity and overfitting

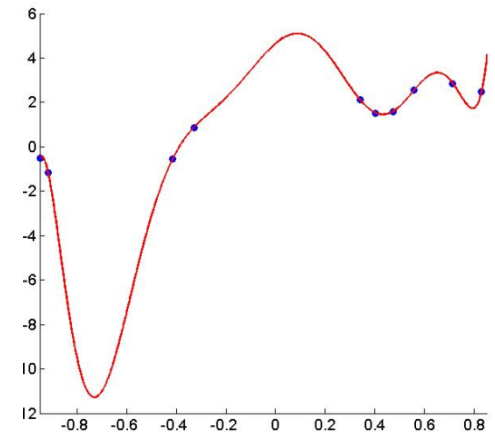
- With limited training data, models may achieve zero training error but a large test error.

Training
(empirical) loss

$$\frac{1}{n} \sum_{i=1}^n (y^{(i)} - f(x^{(i)}; \theta))^2 \approx 0$$

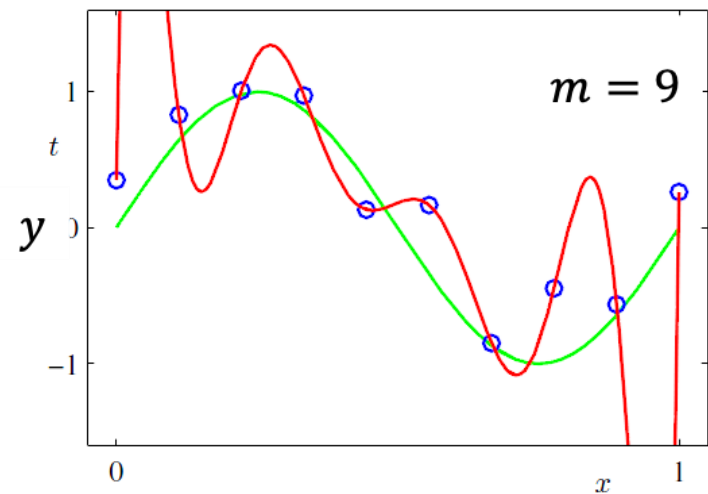
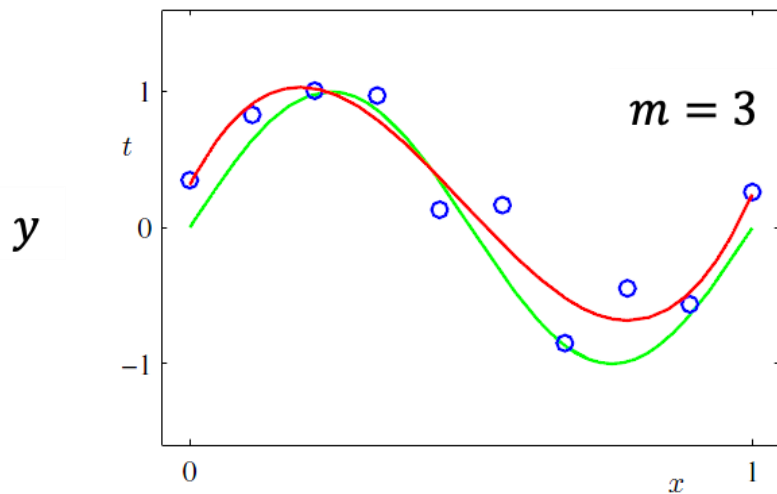
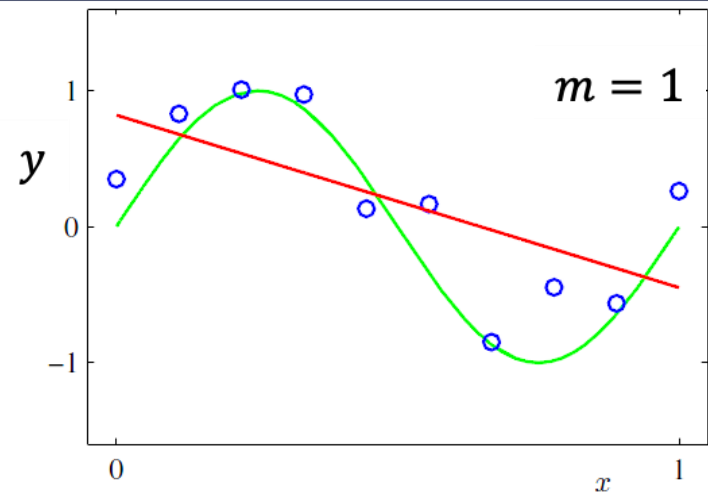
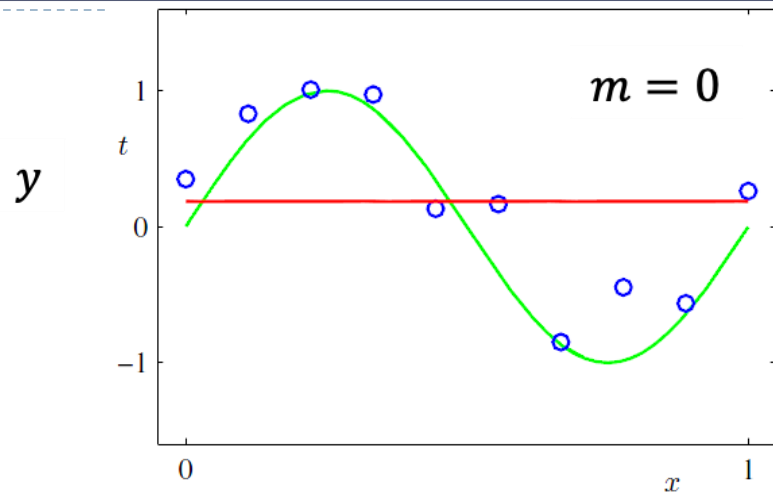
Expected
(true) loss

$$E_{\mathbf{x}, y} \left\{ (y - f(\mathbf{x}; \theta))^2 \right\} \gg 0$$



- Over-fitting: when the training loss no longer bears any relation to the test (generalization) loss.
 - Fails to generalize to unseen examples.

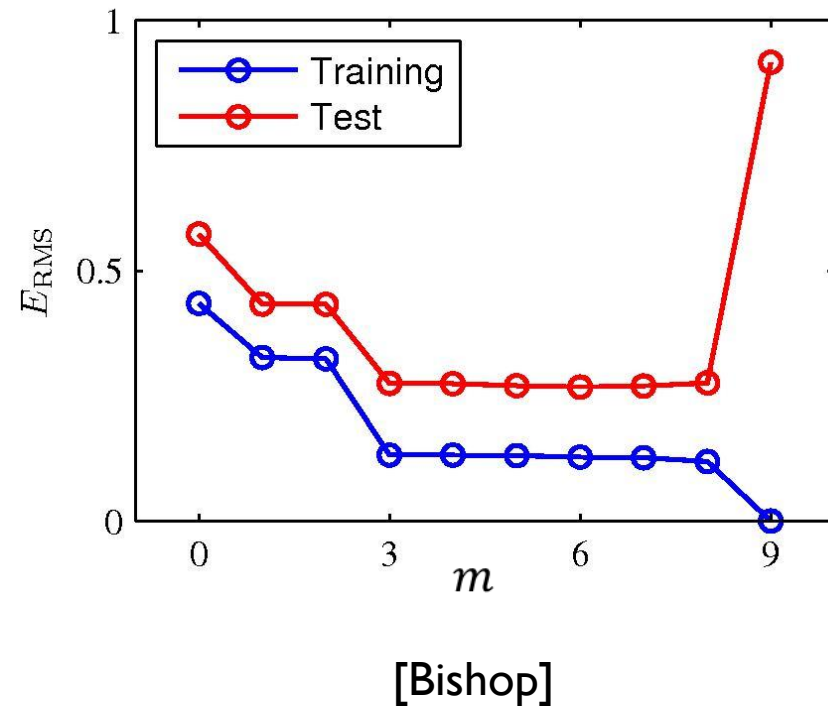
Polynomial regression



[Bishop]

Polynomial regression: training and test error

$$RMSE = \sqrt{\frac{\sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \right)^2}{n}}$$



Over-fitting causes

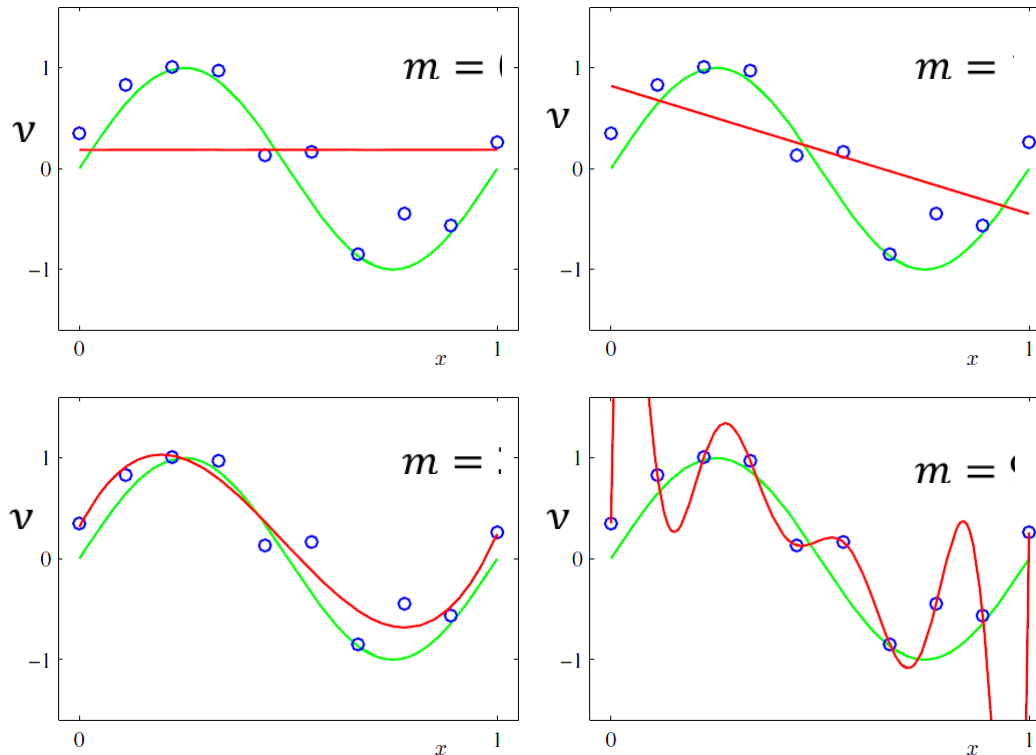
- Model complexity
 - E.g., Model with a large number of parameters (degrees of freedom)
- Low number of training data
 - Small data size compared to the complexity of the model



Model complexity

- **Example:**

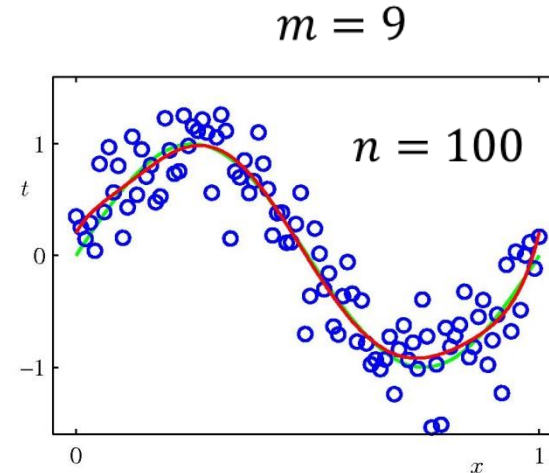
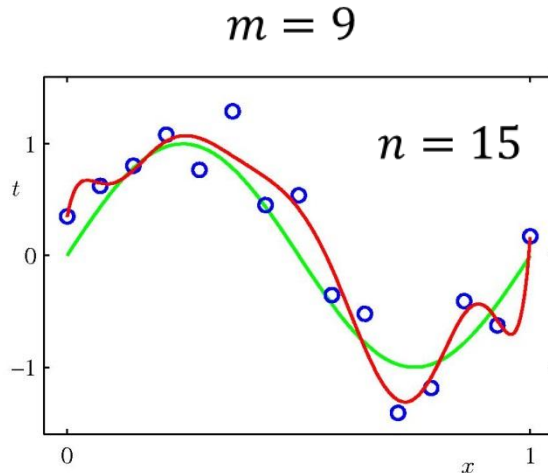
- Polynomials with larger m are becoming increasingly tuned to the random noise on the target values.



[Bishop]

Number of training data & overfitting

- Over-fitting problem becomes less severe as the size of training data increases.



[Bishop]

How to evaluate the learner's performance?

- Generalization error: true (or expected) error that we would like to optimize
- Two ways to assess the generalization error are:
 - Practical: Use a separate data set to test the model
 - Theoretical: Law of Large numbers
 - Bias-variance decomposition of out-of-sample error
 - statistical bounds on the difference between training and expected errors



Avoiding over-fitting

- Determine a suitable value for model complexity (Model Selection)
 - **Simple hold-out method**
 - **Cross-validation**
- Regularization (Occam's Razor)
 - Explicit preference towards simple models
 - Penalize for the model complexity in the objective function
- Bayesian approach



Avoiding over-fitting

- Determine a suitable value for model complexity (Model Selection)
 - **Simple hold-out method**
 - **Cross-validation**
- Regularization (Occam's Razor)
 - Explicit preference towards simple models
 - Penalize for the model complexity in the objective function
- Bayesian approach



Evaluation and model selection

- **Evaluation:**
 - We need to measure how well the learned function can predict the target for unseen examples
- **Model selection:**
 - Most of the time we need to select among a set of models
 - Example: polynomials with different degree m
 - and thus we need to evaluate these models first

Model Selection

- **Learning algorithm** defines the data-driven search over the hypothesis space
 - search for good parameters
- **Hyper-parameters** are the tunable aspects of the model, that the learning algorithm does *not* select

This slide has been adopted from CMU ML course:
<http://www.cs.cmu.edu/~mgormley/courses/10601-s18/>



Model Selection

- **Model selection** is the process by which we choose the “best” model among a set of candidates
 - assume access to a function capable of measuring the quality of a model
 - typically done “outside” the main training algorithm
- Model selection / hyper-parameter optimization is just another form of learning



Simple hold-out: model selection

- **Steps:**
 - Divide training data into training and validation set v_set
 - Use only the training set to train a set of models
 - Evaluate each learned model on the validation set
 - $J_v(\mathbf{w}) = \frac{1}{|v_set|} \sum_{i \in v_set} \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$
 - Choose the best model based on the validation set error

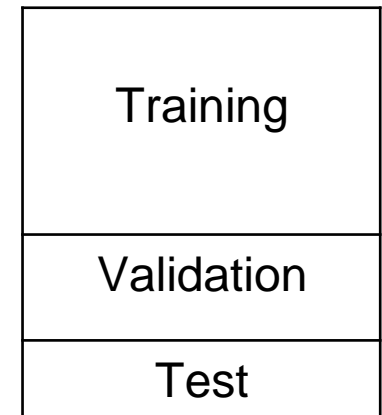


Simple hold-out: model selection

- **Steps:**
 - Divide training data into training and validation set v_set
 - Use only the training set to train a set of models
 - Evaluate each learned model on the validation set
 - $J_v(\mathbf{w}) = \frac{1}{|v_set|} \sum_{i \in v_set} \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$
 - Choose the best model based on the validation set error
- **Usually, too wasteful of valuable training data**
 - Training data may be limited.
 - On the other hand, small validation set obtains a relatively noisy estimate of performance.

Simple hold out: training, validation, and test sets

- Simple hold-out chooses the model that minimizes error on validation set.
- $J_v(\hat{\mathbf{w}})$ is likely to be an optimistic estimate of generalization error.
 - extra parameter (e.g., degree of polynomial) is fit to this set.
- Estimate generalization error for the test set
 - performance of the selected model is finally evaluated on the test set



Avoiding over-fitting

- Determine a suitable value for model complexity (Model Selection)
 - **Simple hold-out method**
 - **Cross-validation**
- **Regularization (Occam's Razor)**
 - Explicit preference towards simple models
 - Penalize for the model complexity in the objective function



Regularization

- Adding a penalty term in the cost function to discourage the coefficients from reaching large values.



Regularization

- ▶ Adding a penalty term in the cost function to discourage the coefficients from reaching large values.
- ▶ Ridge regression (weight decay):

$$J(\mathbf{w}) = \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w}$$

Regularization

- Adding a penalty term in the cost function to discourage the coefficients from reaching large values.
- Ridge regression (weight decay):

$$J(\mathbf{w}) = \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w}$$

$$\hat{\mathbf{w}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \boldsymbol{\Phi} = \begin{bmatrix} 1 & \phi_1(\mathbf{x}^{(1)}) & \cdots & \phi_m(\mathbf{x}^{(1)}) \\ 1 & \phi_1(\mathbf{x}^{(2)}) & \cdots & \phi_m(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}^{(n)}) & \cdots & \phi_m(\mathbf{x}^{(n)}) \end{bmatrix} \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

Regularization

- Adding a penalty term in the cost function to discourage the coefficients from reaching large values.
- Ridge regression (weight decay):

$$J(\mathbf{w}) = \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w}$$

$$\hat{\mathbf{w}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \boldsymbol{\Phi} = \begin{bmatrix} 1 & \phi_1(\mathbf{x}^{(1)}) & \cdots & \phi_m(\mathbf{x}^{(1)}) \\ 1 & \phi_1(\mathbf{x}^{(2)}) & \cdots & \phi_m(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}^{(n)}) & \cdots & \phi_m(\mathbf{x}^{(n)}) \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

Polynomial order

- Polynomials with larger m are becoming increasingly tuned to the random noise on the target values.
 - magnitude of the coefficients typically gets larger by increasing m .

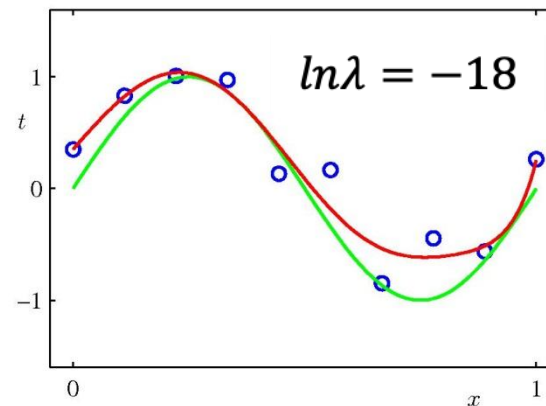
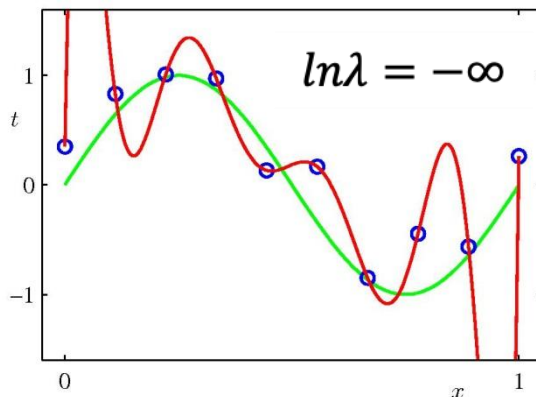
	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

[Bishop]

Regularization parameter

	$m = 9$		
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
\hat{w}_0	0.35	0.35	0.13
\hat{w}_1	232.37	4.74	-0.05
\hat{w}_2	-5321.83	-0.77	-0.06
\hat{w}_3	48568.31	-31.97	-0.05
\hat{w}_4	-231639.30	-3.89	-0.03
\hat{w}_5	640042.26	55.28	-0.02
\hat{w}_6	-1061800.52	41.32	-0.01
\hat{w}_7	1042400.18	-45.95	-0.00
\hat{w}_8	-557682.99	-91.53	0.00
\hat{w}_9	125201.43	72.68	0.01

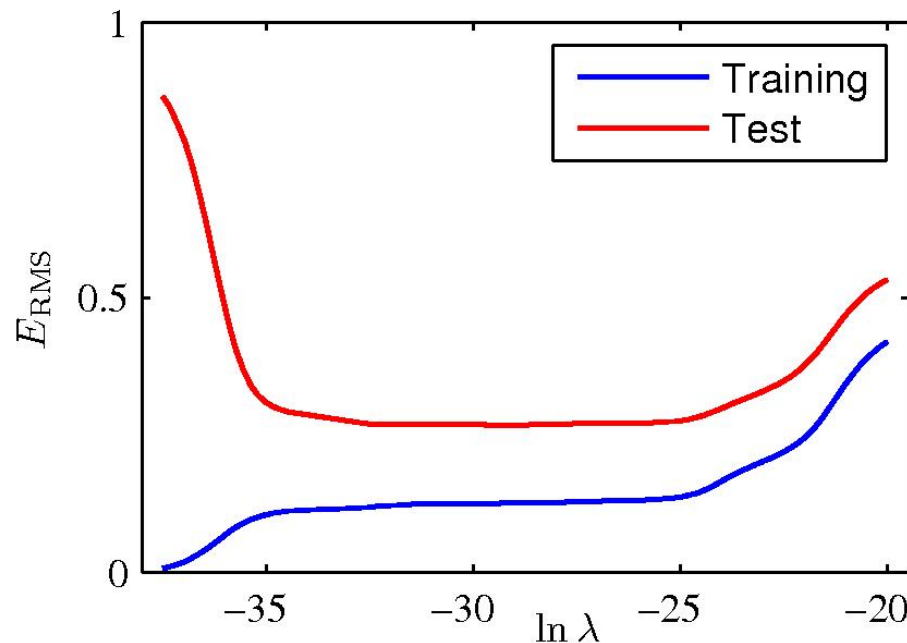
[Bishop]



Regularization parameter

Generalization

- ▶ λ now controls the effective complexity of the model and hence determines the degree of over-fitting



[Bishop]

Choosing the regularization parameter

- A set of models with different values of λ .
- Find $\hat{\mathbf{w}}$ for each model based on training data
- Find $J_v(\hat{\mathbf{w}})$ (or $J_{cv}(\hat{\mathbf{w}})$) for each model
 - $J_v(\mathbf{w}) = \frac{1}{n_v} \sum_{i \in v_set} \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$
- Select the model with the best $J_v(\hat{\mathbf{w}})$ (or $J_{cv}(\hat{\mathbf{w}})$)



The approximation-generalization trade-off

- - ▶ Small true error shows good approximation of f out of sample
 - ▶ More complex $\mathcal{H} \Rightarrow$ better chance of approximating f
 - ▶ Less complex $\mathcal{H} \Rightarrow$ better chance of generalization out of f



Complexity of Hypothesis Space: Example

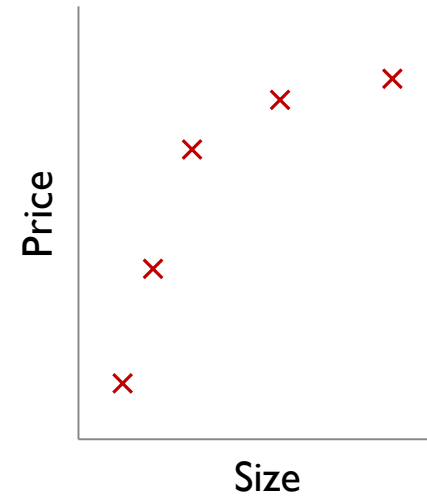


$$w_0 + w_1x$$

Less complex \mathcal{H}



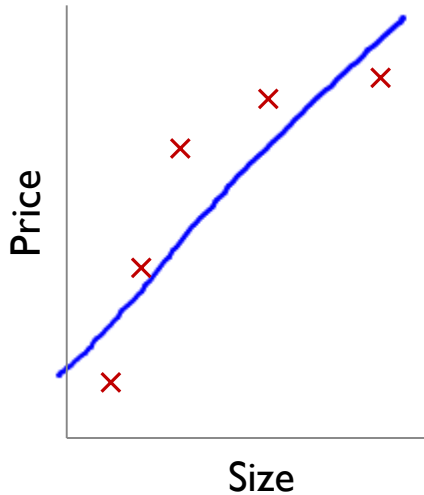
$$w_0 + w_1x + w_2x^2$$



$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

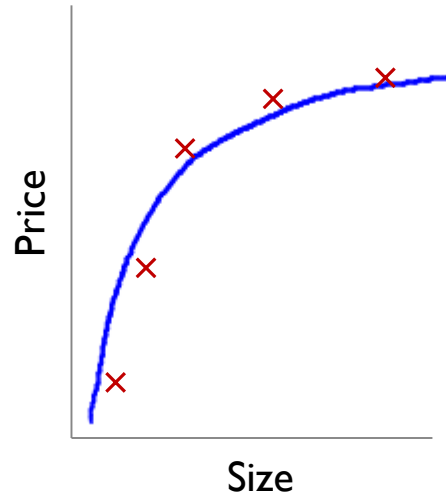
More complex \mathcal{H}

Complexity of Hypothesis Space: Example

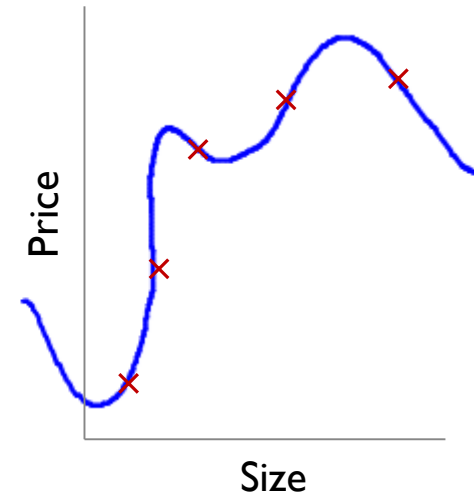


$$w_0 + w_1x$$

Underfitting



$$w_0 + w_1x + w_2x^2$$



$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

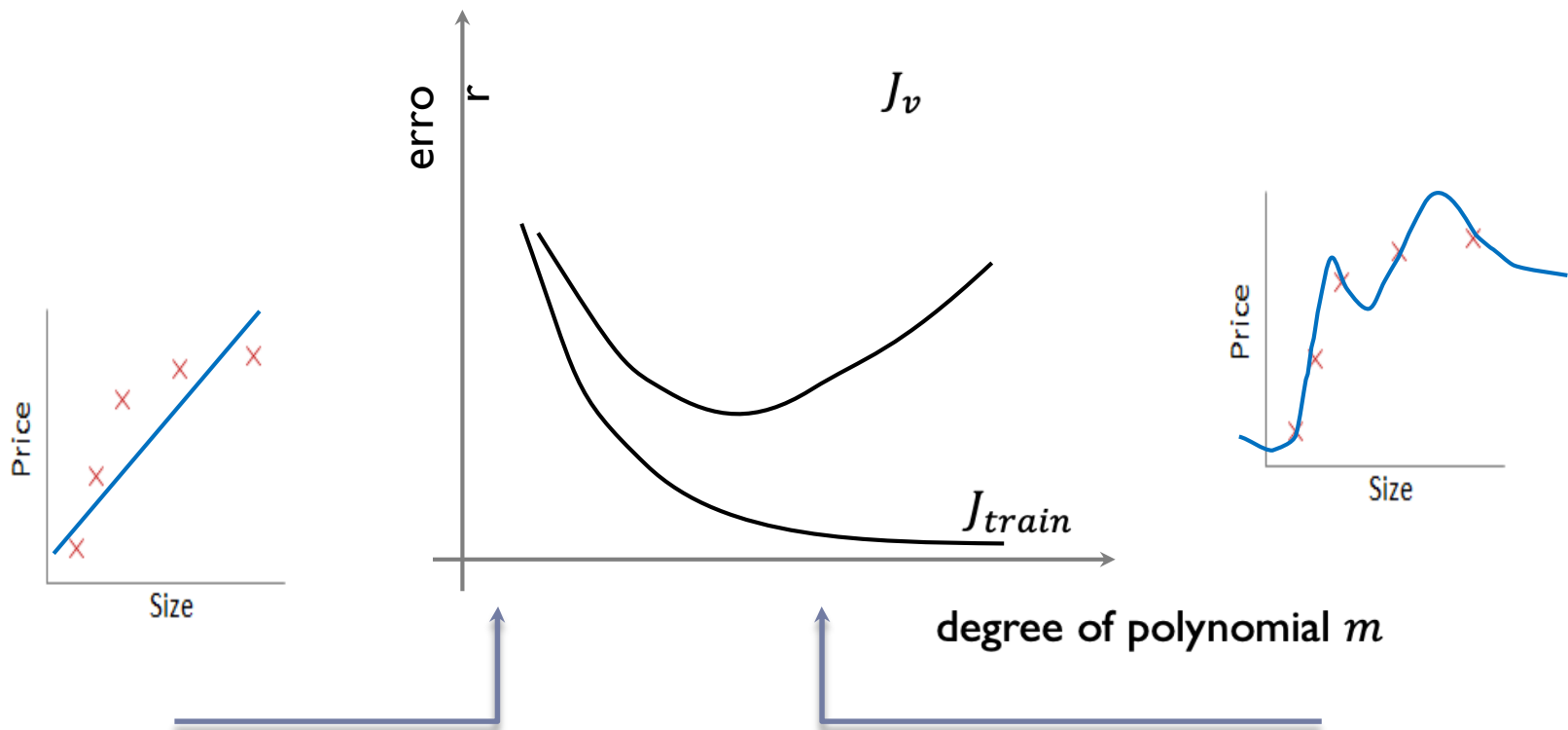
Overfitting

This example has been adapted from: Prof. Andrew Ng's slides

Complexity of Hypothesis Space: Example

$$J_v(\mathbf{w}) = \frac{1}{n_{val}} \sum_{i \in \text{eval_set}} (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2$$

$$J_{train}(\mathbf{w}) = \frac{1}{n_{train}} \sum_{i \in \text{train_set}} (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2$$



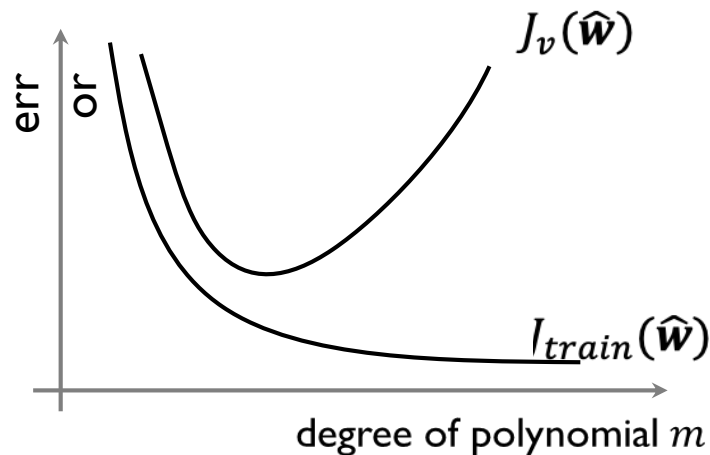
Complexity of Hypothesis Space

▶ Less complex \mathcal{H} :

▶ $J_{train}(\hat{\mathbf{w}}) \approx J_v(\hat{\mathbf{w}})$ and $J_{train}(\hat{\mathbf{w}})$ is very high

▶ More complex \mathcal{H} :

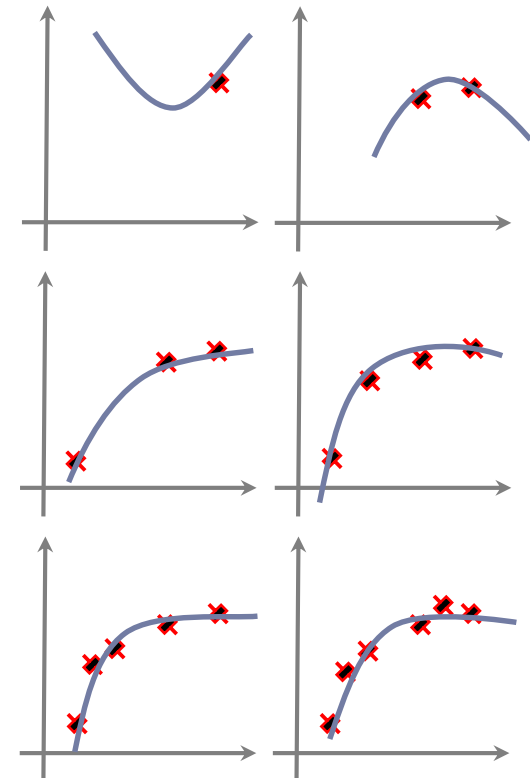
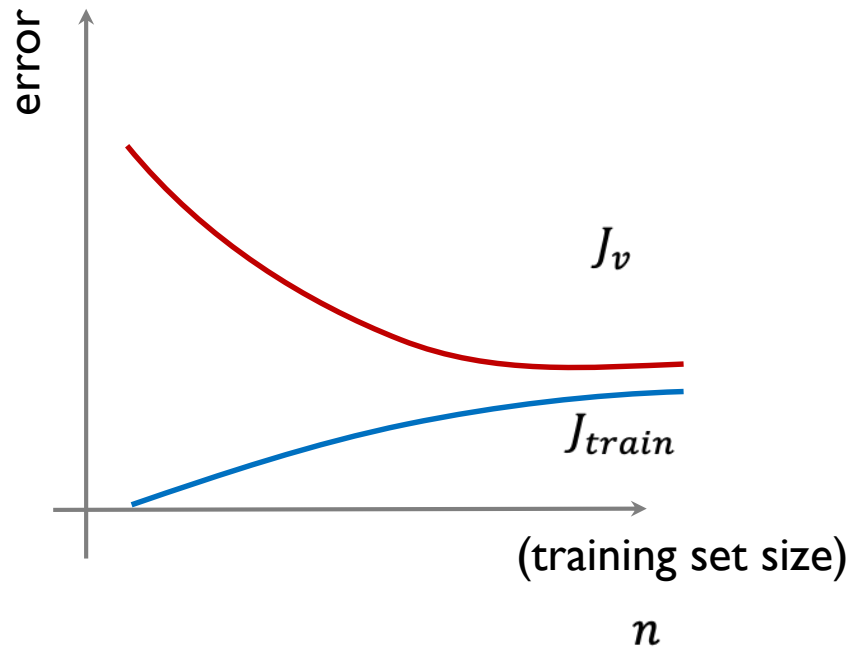
▶ $J_{train}(\hat{\mathbf{w}}) \ll J_v(\hat{\mathbf{w}})$ and $J_{train}(\hat{\mathbf{w}})$ is low



Size of training set

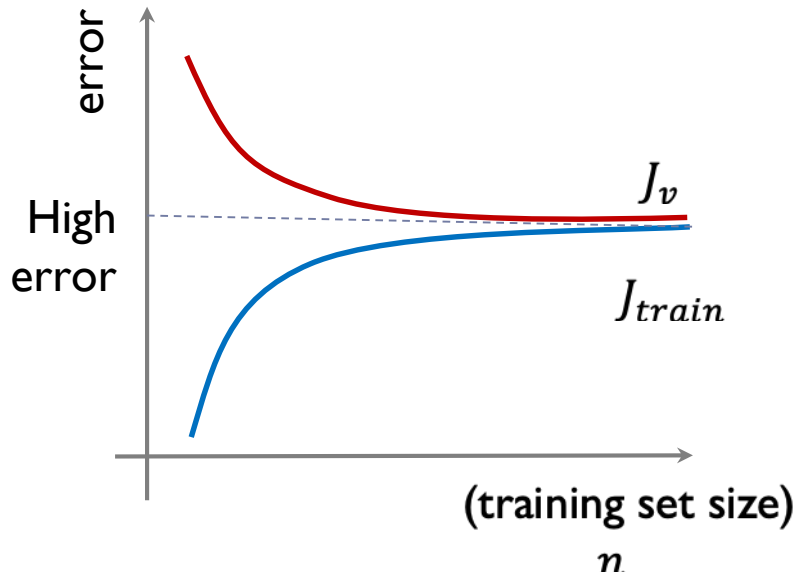
$$J_v(\mathbf{w}) = \frac{1}{n_{\text{val}}} \sum_{i \in \text{eval_set}} \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$
$$J_{\text{train}}(\mathbf{w}) = \frac{1}{n_{\text{train}}} \sum_{i \in \text{train_set}} \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$

$$f(x; \mathbf{w}) = w_0 + w_1 x + w_2 x^2$$

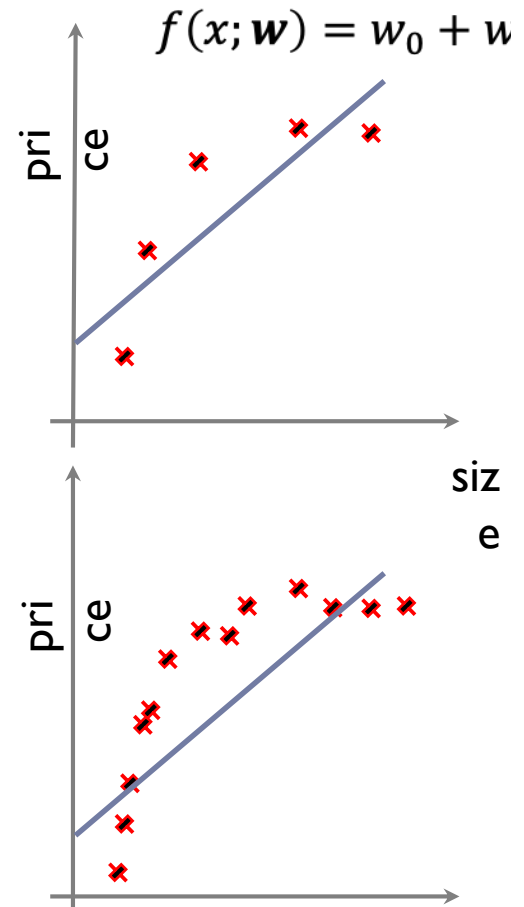


This slide has been adapted from: Prof. Andrew Ng's slides

Less complex \mathcal{H}

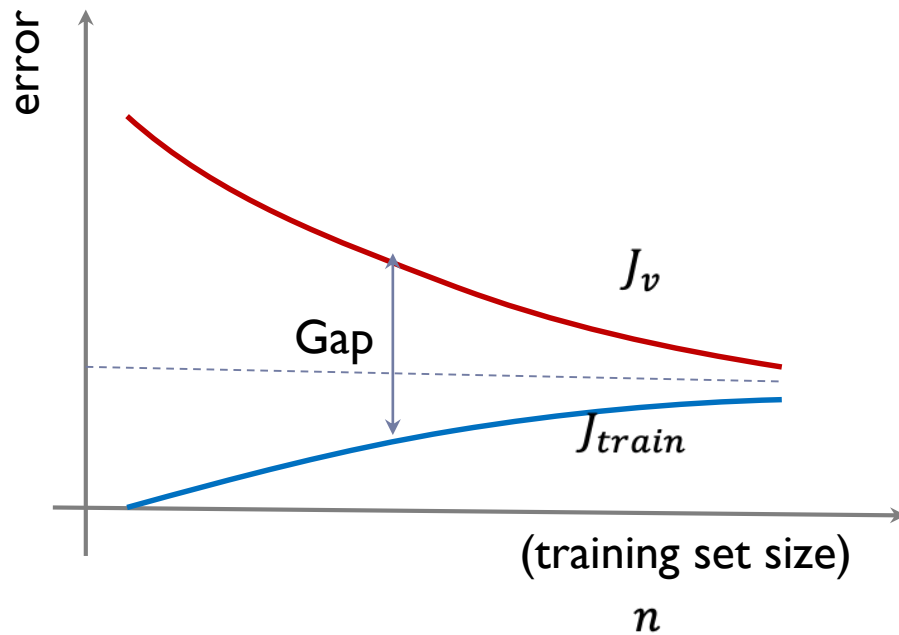


If model is very simple, getting more training data will not (by itself) help much.

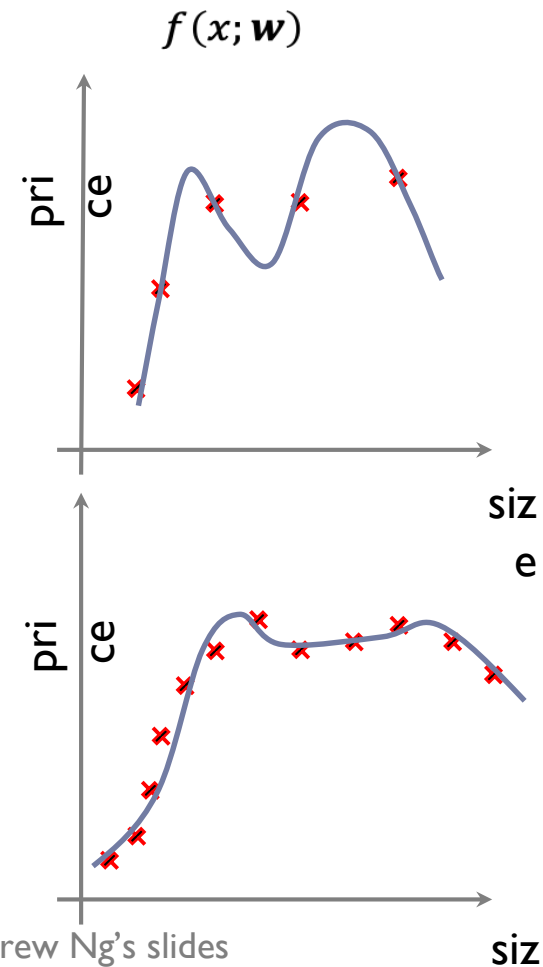


This slide has been adapted from: Prof. Andrew Ng's slides

More complex \mathcal{H}



For more complex models, getting more training data is usually helps.



This slide has been adapted from: Prof. Andrew Ng's slides

Regularization: Example

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

$$J(\mathbf{w}) = \frac{1}{n} \left(\sum_{i=1}^n (y^{(i)} - f(x^{(i)}; \mathbf{w}))^2 + \lambda \mathbf{w}^T \mathbf{w} \right)$$

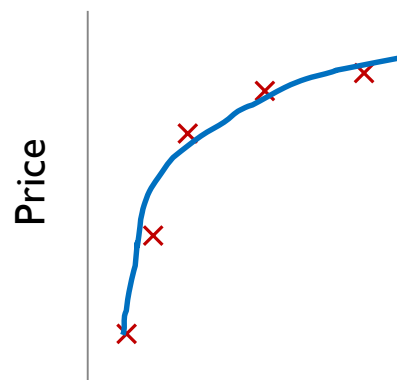


Size

Large λ

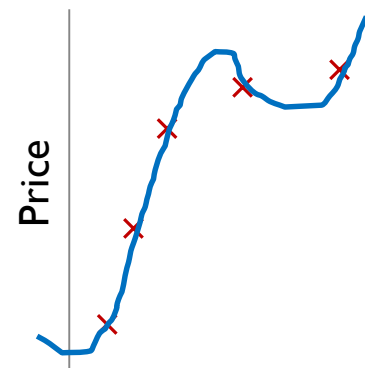
(Prefer to more simple models)

$$w_1 = w_2 \approx 0$$



Size

Intermediate λ



Size

Small λ

(Prefer to more complex models)

$$\lambda = 0$$

This example has been adapted from: Prof. Andrew Ng's slides

Resources

- C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 1.1,1.3, 3.1
- Course CE-717, Dr. M.Soleymani
- CMU ML course: [http:// www.cs.cmu.edu /~mgormley /courses /10601-s18/](http://www.cs.cmu.edu/~mgormley/courses/10601-s18/)

