

innovate

achieve

lead



BITS Pilani
Pilani Campus

Machine Learning

DSECL ZG565

Dr. Monali Mavani

Contents



- Learning Model Parameters linear regression– Closed Form Solution (using vectorization)
- Linear basis function models (3.1 Bishop)
- Over fitting, Under fitting, Regularization (1.1 Bishop, Andrew Ng notes)
- Bias-Variance Decomposition (Bishop)

Learning Model Parameters – Closed Form Solution (using vectorization)

Vectorization

- Benefits of vectorization
 - More compact equations
 - Faster code (using optimized matrix libraries)

- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

- Can write the model in vectorized form as $h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$

Vectorization

- Consider our model for n instances:

$$h\left(\mathbf{x}^{(i)}\right)=\sum_{j=0}^d \theta_j x_j^{(i)}$$

- Let

$$\boldsymbol{\theta}=\left[\begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{array}\right] \quad \mathbf{X}=\left[\begin{array}{cccc} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{array}\right]$$

$\mathbb{R}^{(d+1) \times 1}$
 $\mathbb{R}^{n \times (d+1)}$

- Can write the model in vectorized form as $h_{\boldsymbol{\theta}}(\mathbf{x})=\mathbf{X} \boldsymbol{\theta}$

Vectorization

- For the linear regression cost function:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)^2$$

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

$$= \frac{1}{2n} \underbrace{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T}_{\mathbb{R}^{1 \times n}} \underbrace{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})}_{\mathbb{R}^{n \times 1}}$$

Annotations for dimensions:

- $\mathbb{R}^{n \times (d+1)}$ points to \mathbf{X}
- $\mathbb{R}^{(d+1) \times 1}$ points to $\boldsymbol{\theta}$
- $\mathbb{R}^{1 \times n}$ points to $(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T$
- $\mathbb{R}^{n \times 1}$ points to $(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$

Let:

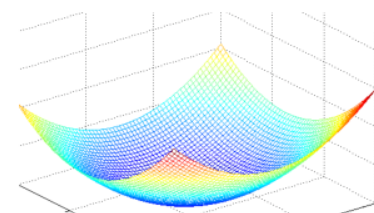
$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Closed Form Solution



- Instead of using GD, solve for optimal θ analytically

– Notice that the solution is when $\frac{\partial}{\partial \theta} J(\theta) = 0$



- Derivation:

$$\begin{aligned} \mathcal{J}(\theta) &= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - \mathbf{y}^\top \mathbf{X} \theta - \theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \end{aligned}$$

A blue arrow points from the 1×1 label to the $\mathbf{y}^\top \mathbf{X} \theta$ term in the second line of the derivation.

Take derivative and set equal to 0, then solve for θ :

$$\frac{\partial}{\partial \theta} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \cancel{\mathbf{y}^\top \mathbf{y}}) = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y} = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta = \mathbf{X}^\top \mathbf{y}$$

Closed Form Solution:

$$\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

46

Closed Form Solution

- Can obtain θ by simply plugging X and y into

$$\theta = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- If $X^T X$ is not invertible (i.e., singular), may need to:
 - Use pseudo-inverse instead of the inverse
 - In python, `numpy.linalg.pinv(a)`
 - Remove redundant (not linearly independent) features
 - Remove extra features to ensure that $d \leq n$

Gradient Descent vs Closed Form



Gradient Descent

- Requires multiple iterations
- Need to choose α
- Works well when n is large
- Can support incremental learning

Closed Form Solution

- Non-iterative
- No need for α
- Slow if n is large
 - Computing $(X^T X)^{-1}$ is roughly $O(n^3)$

Linear Basis Function Models

Extending Linear Regression to More Complex Models



- The inputs \mathbf{X} for linear regression can be:
 - Original quantitative inputs
 - Transformation of quantitative inputs
 - e.g. log, exp, square root, square, etc.
 - Polynomial transformation
 - example: $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot x^3$
 - Basis expansions
 - Dummy coding of categorical inputs
 - Interactions between variables
 - example: $x_3 = x_1 \cdot x_2$

This allows use of linear regression techniques to fit non-linear datasets.

Linear Basis Function Models

- Basic Linear Model:
$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$
- Generalized Linear Model:
$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j \phi_j(\mathbf{x})$$
- Once we have replaced the data by the outputs of the basis functions, fitting the generalized model is exactly the same problem as fitting the basic model
 - Unless we use the kernel trick – more on that when we cover support vector machines
 - Therefore, there is no point in cluttering the math with basis functions

Linear Basis Function Models

- Generally,

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j \underbrace{\phi_j(\mathbf{x})}_{\text{basis function}}$$

- Typically, $\phi_0(\mathbf{x}) = 1$ so that θ_0 acts as a bias
- In the simplest case, we use linear basis functions :

$$\phi_j(\mathbf{x}) = x_j$$

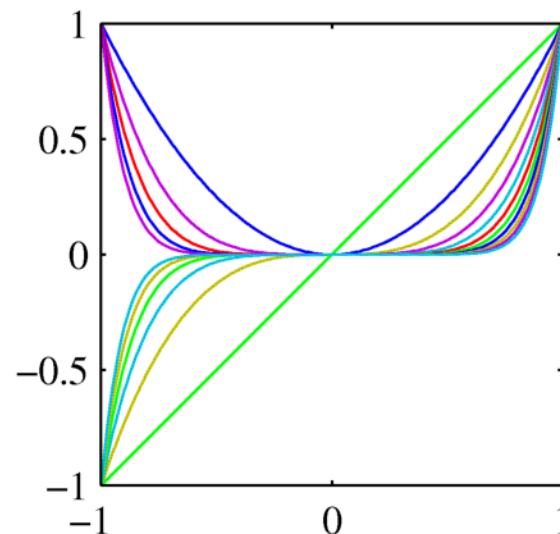
Linear Basis Function Models



- Polynomial basis functions:

$$\phi_j(x) = x^j$$

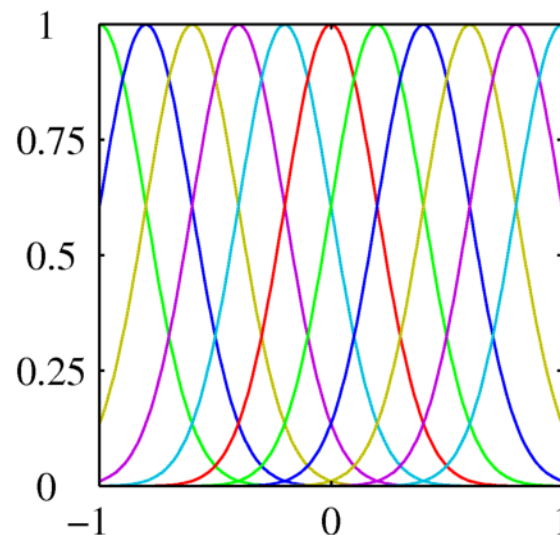
- These are global; a small change in x affects all basis functions



- Gaussian basis functions:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

- These are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (width).



Linear Basis Function Models



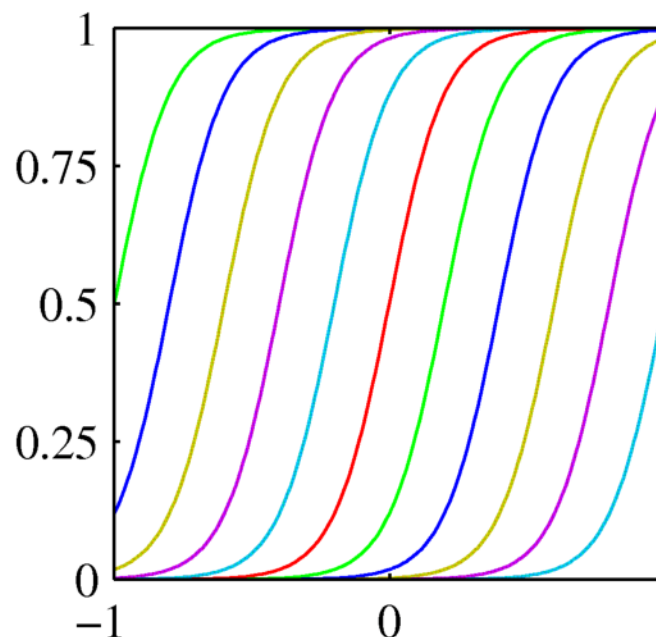
- Sigmoidal basis functions:

$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

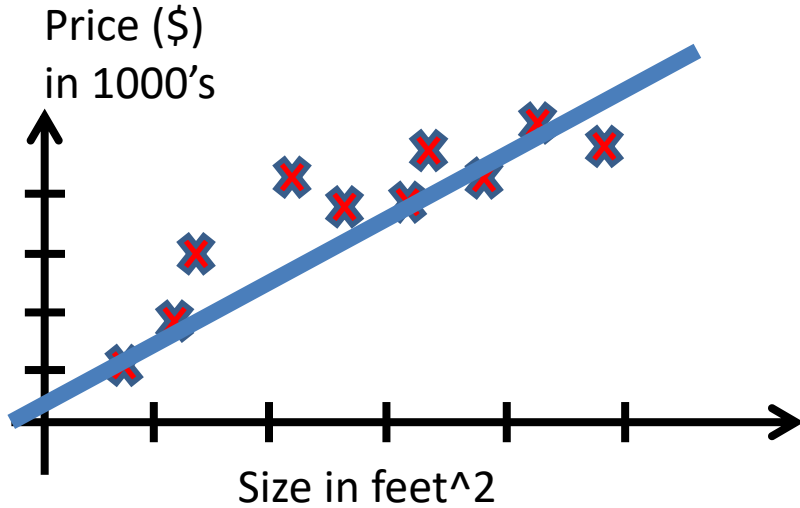
- These are also local; a small change in x only affects nearby basis functions. μ_j and s control location and scale (slope).



By using nonlinear basis functions, we allow the function $y(\mathbf{x}, \mathbf{w})$ to be a nonlinear function of the input vector \mathbf{x} . They are called linear models because this function is linear in \mathbf{w} . It is this

Over fitting, Under fitting, Regularization

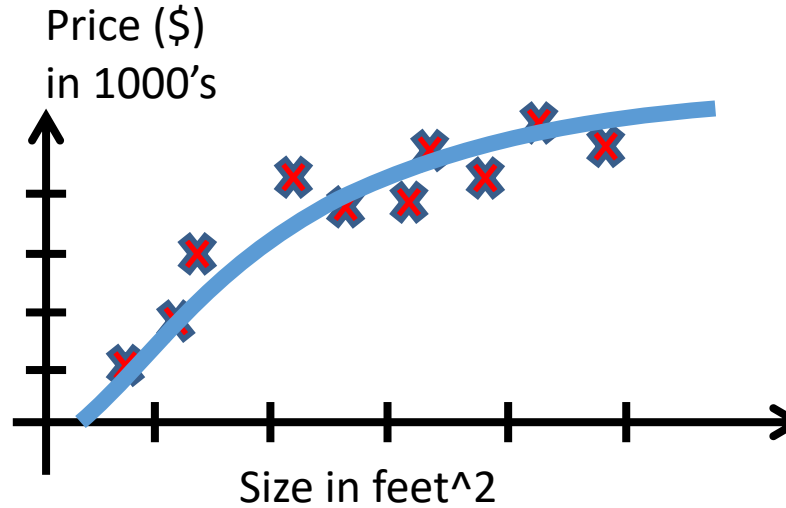
Linear regression



$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

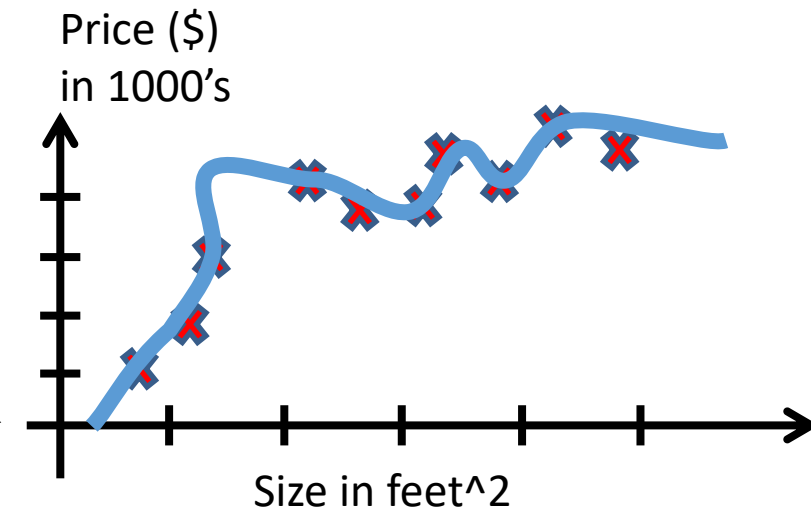
Underfitting

High bias



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

Just right

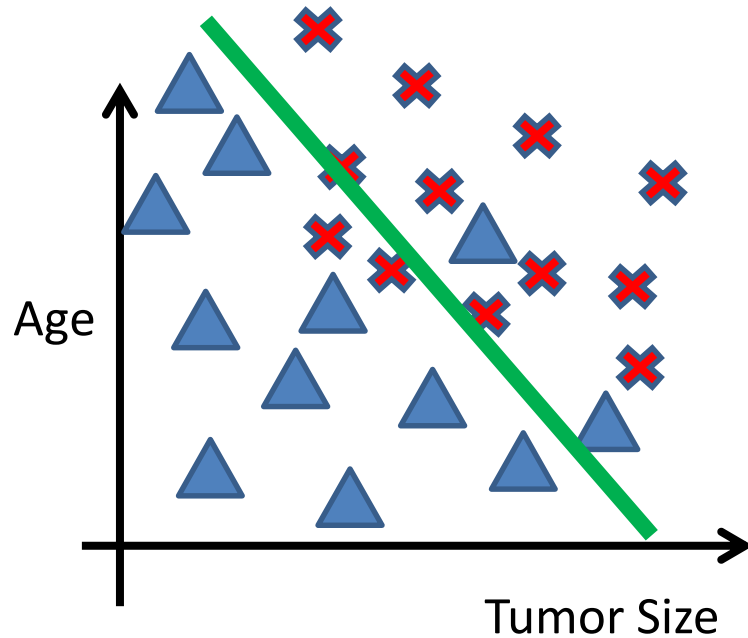


$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \dots$$

Overfitting

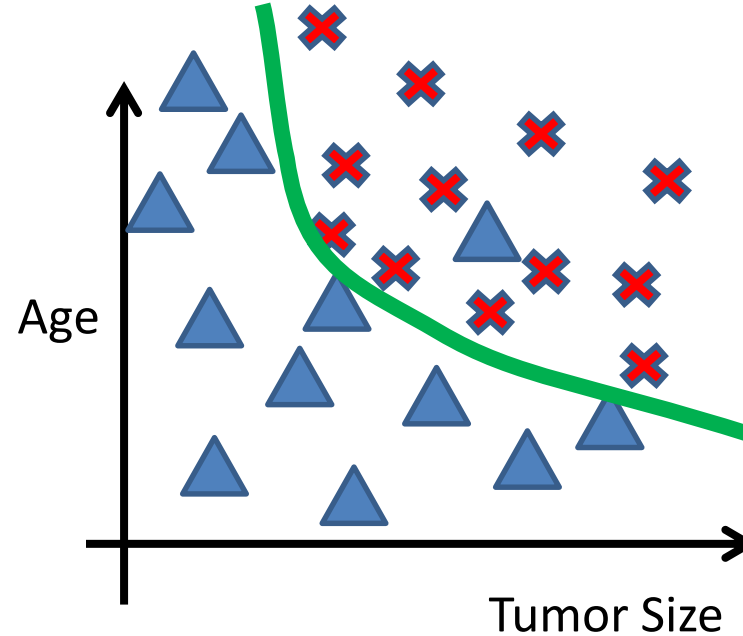
High variance

Logistic regression

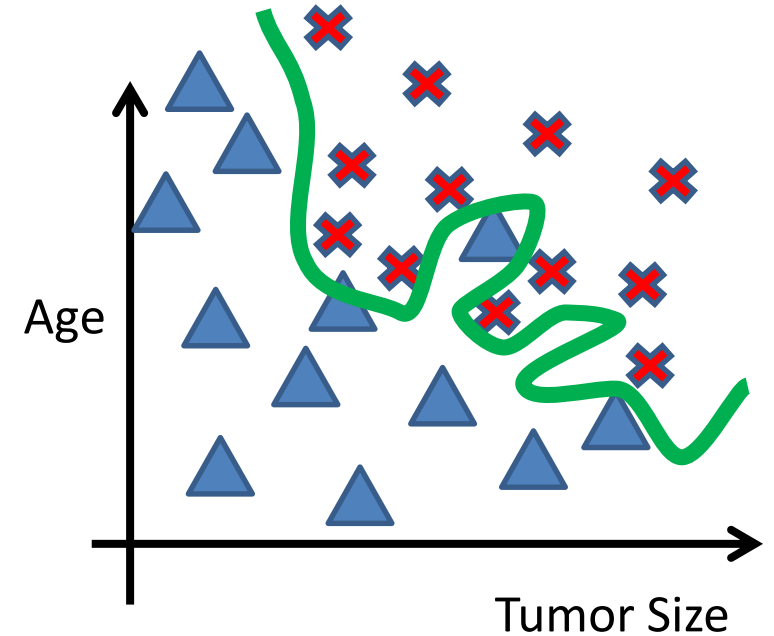


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2)$$

Underfitting



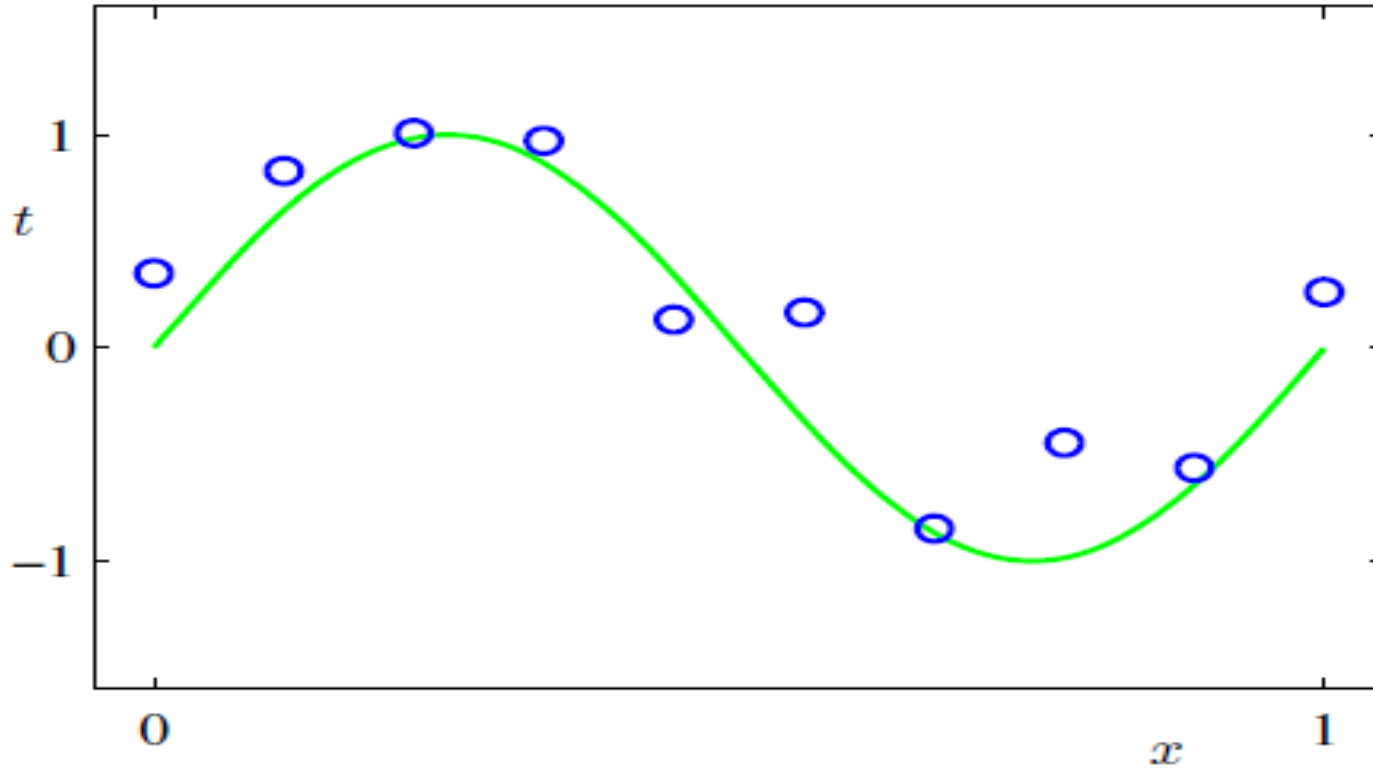
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \theta_6 x_1^3 x_2 + \theta_7 x_1 x_2^3 + \dots)$$

Overfitting

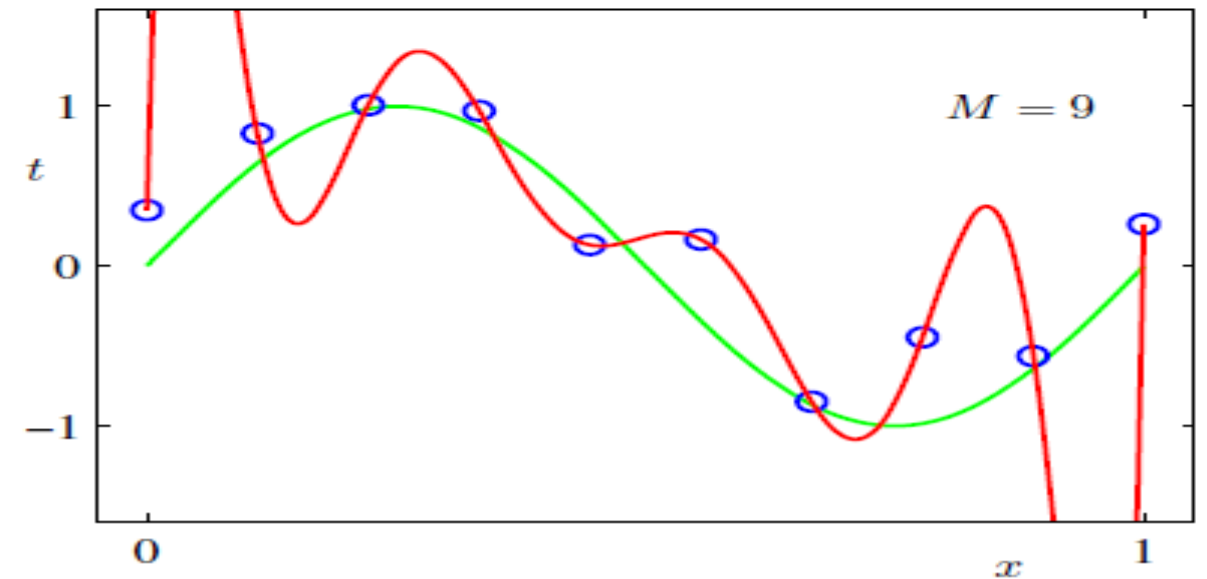
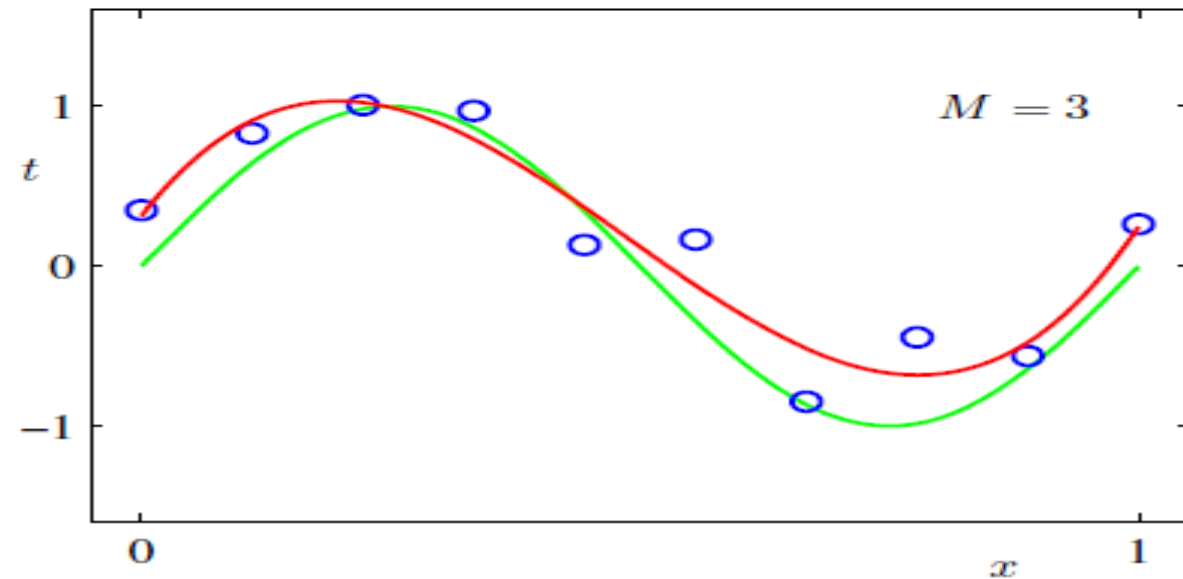
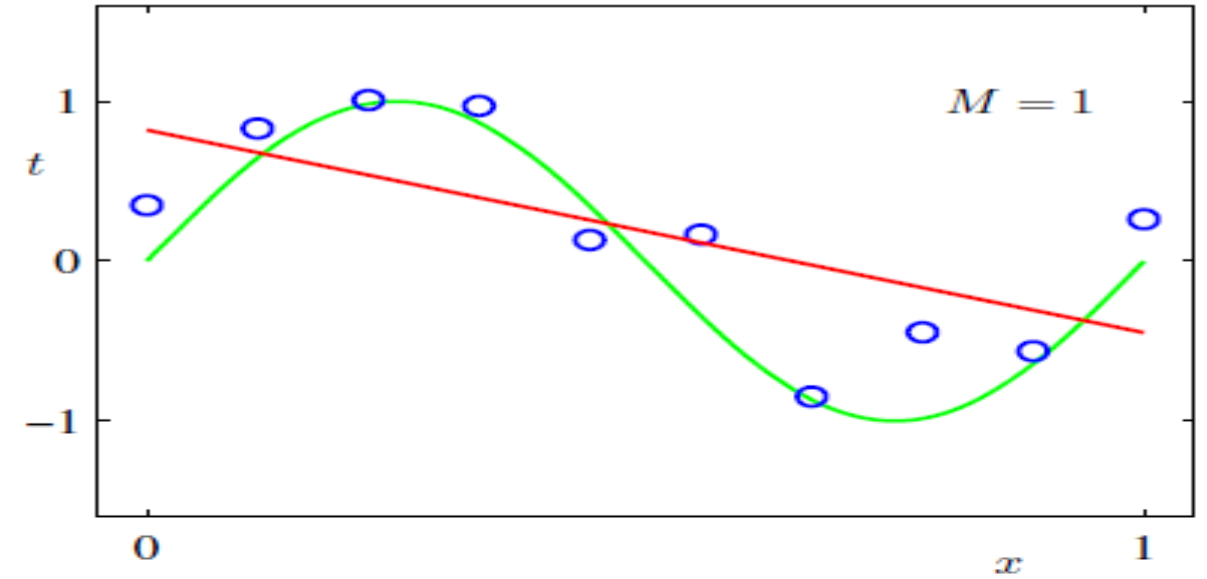
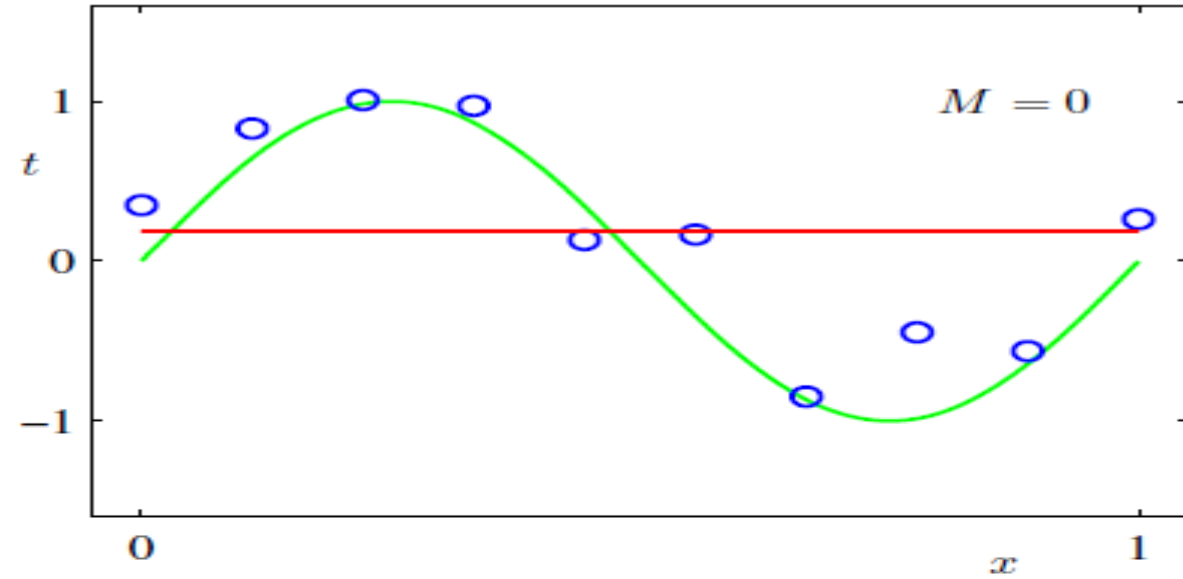
Example: Curve fitting problem



The green curve shows the function $\sin(2\pi x)$ used to generate the data.

$$F(x) = \sin(2\pi x) + \text{noise (Gaussian)}$$

Curve fitting problem: Model Selection



Overfitting



- If we have too many features (i.e. complex model), the learned hypothesis may fit the training set very well

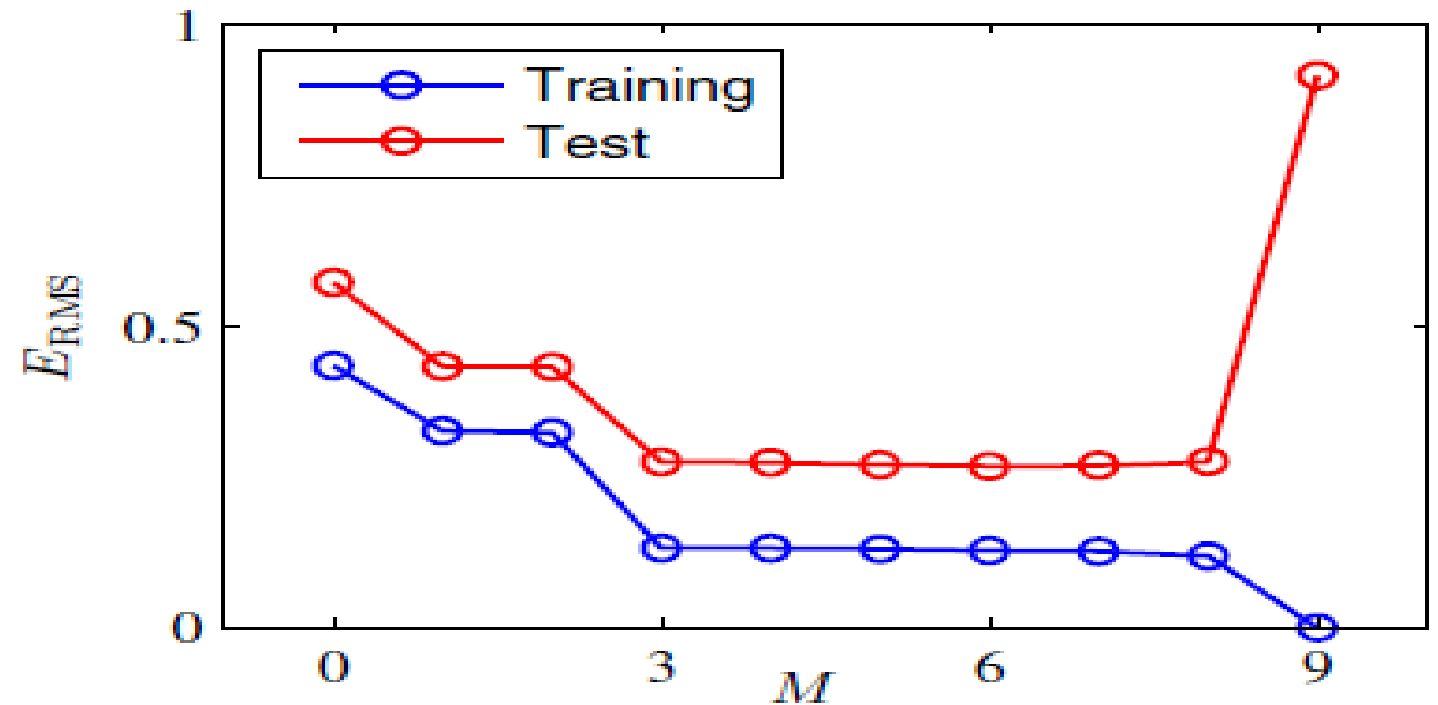
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$$

but fail to generalize to new examples

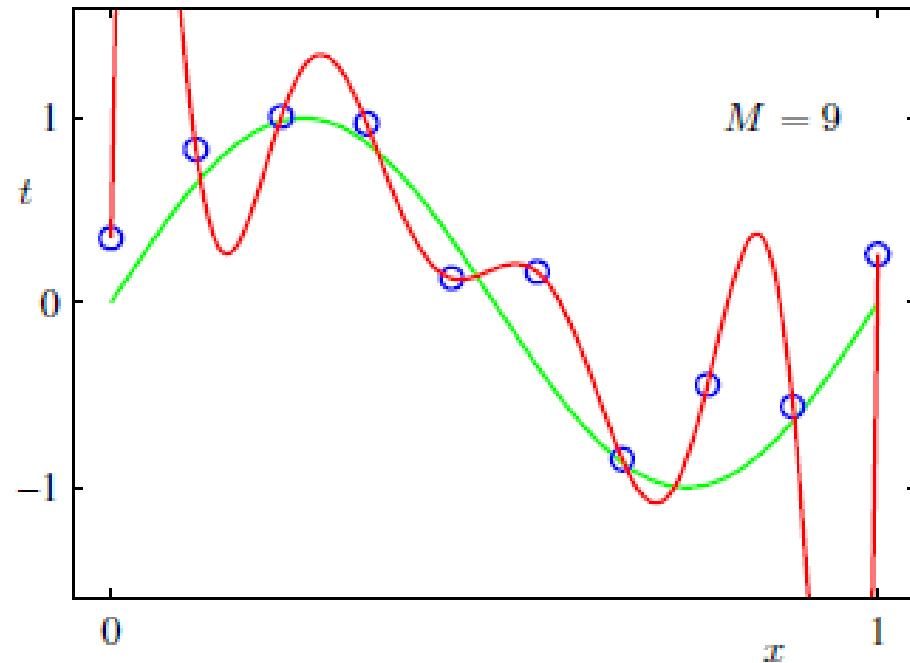
Model Selection



- Goal is to achieve good generalization.
- Generalization performance can be evaluated using **test data set**.
- Values of M in the range 3 to 8 give small values for the test set error, and these also give reasonable representations of the generating function $\sin(2\pi x)$
- For $M = 9$, polynomial contains 10 degrees of freedom corresponding to the 10 coefficients w_0, \dots, w_9 , can be tuned exactly to the 10 data points



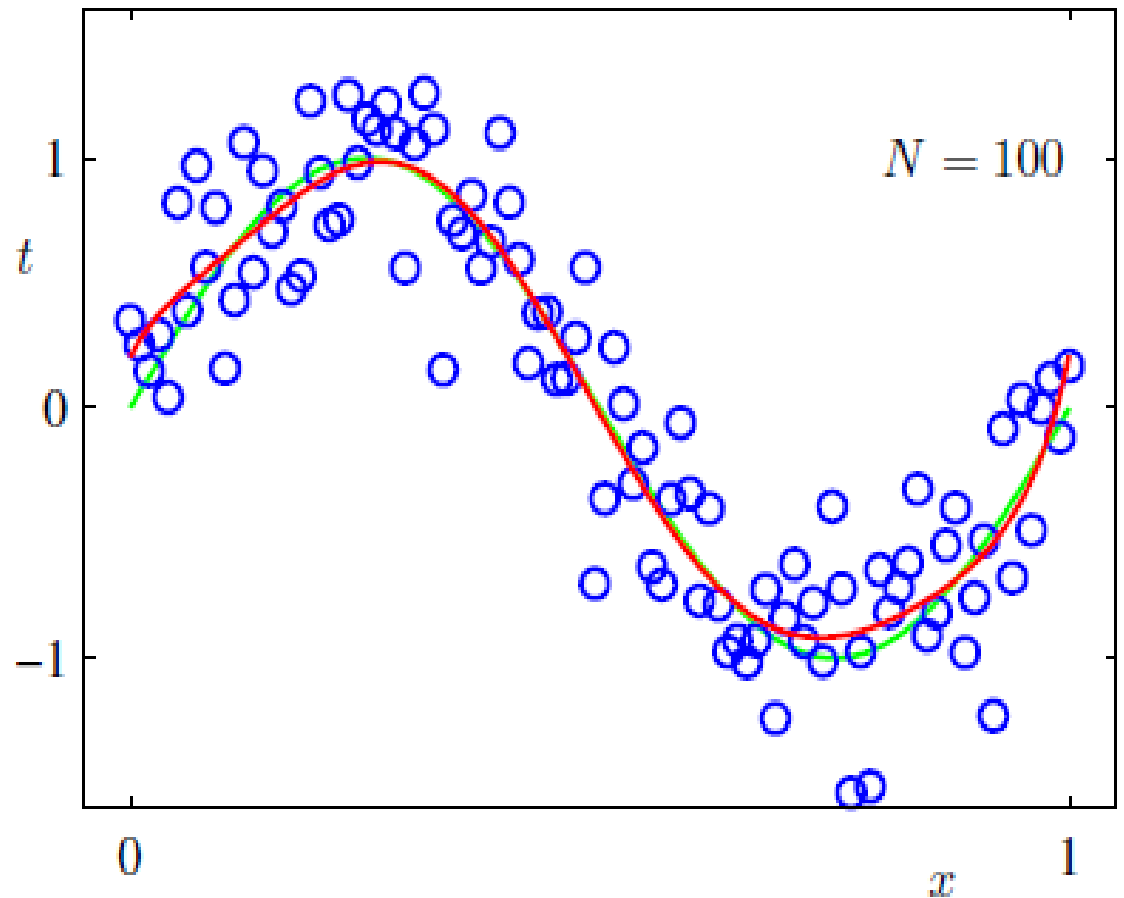
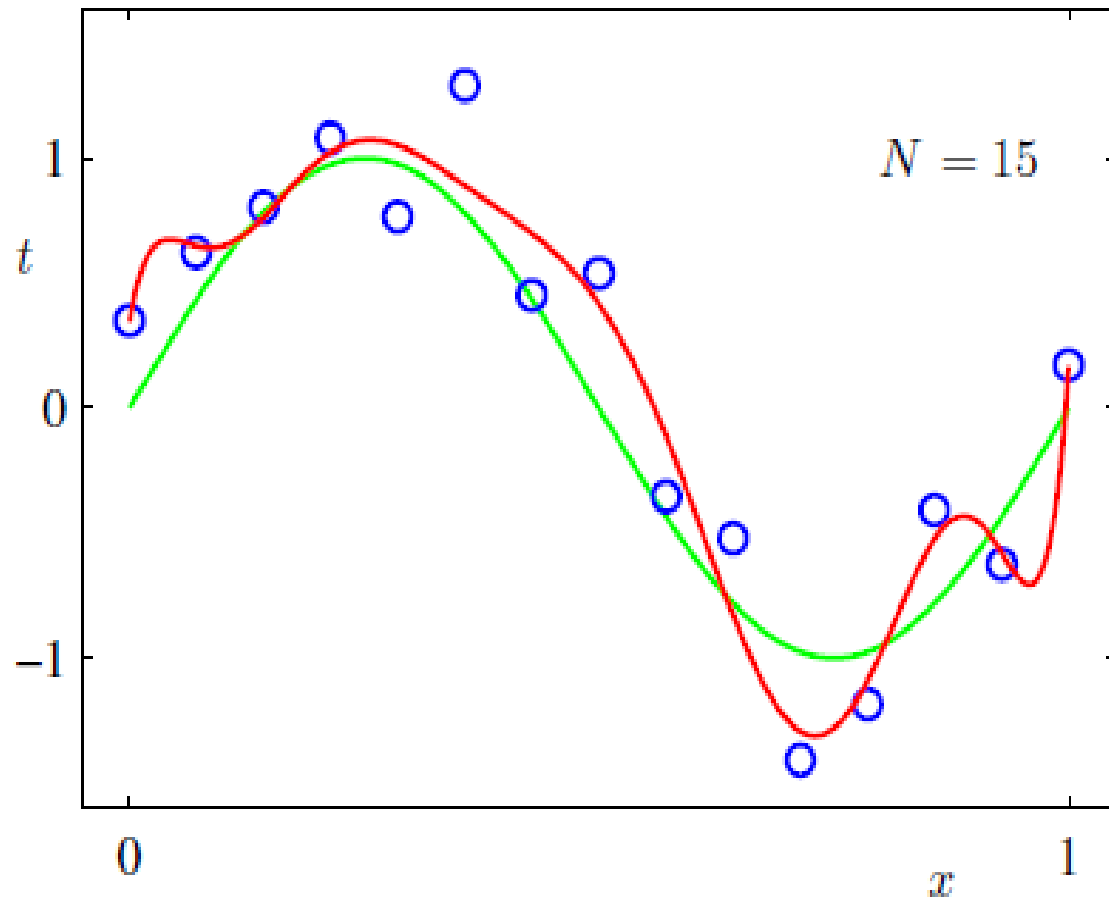
Coefficients for polynomials of various order



	$M = 0$	$M = 1$	$M = 6$	$M = 9$
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

Coefficients have become finely tuned to the data by developing large positive and negative values, matches each of the data points exactly

Size of the data set reduces the over-fitting problem



Addressing overfitting



- **Reduce number of features**
 - Manually select which features to keep.
 - One rough heuristic says that data points should be multiple (5 to 10) parameters.
 - It seems more reasonable to take number of parameters according to the problem complexity rather than given data points.
- **Any Solution?**
 - **Regularization**
 - Adding a penalty term to the error function in order to discourage the coefficients from reaching large values. E.g a sum of squares of all of the coefficients,
 - Keep all the features, but reduce magnitude/values of parameters θ_j .
 - Works well when we have a lot of features, each of which contributes a bit to predicting y .
 - ***Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.***

Regularization

- Linear regression objective function

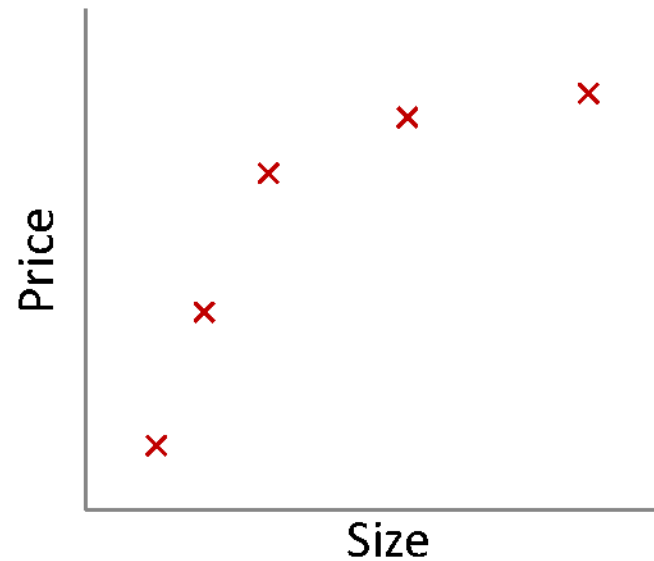
$$J(\boldsymbol{\theta}) = \underbrace{\frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)^2}_{\text{model fit to data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2}_{\text{regularization}}$$

- λ is the regularization parameter ($\lambda \geq 0$)
- No regularization on θ_0 !

Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- What happens if we set λ to be huge (e.g., 10^{10})?

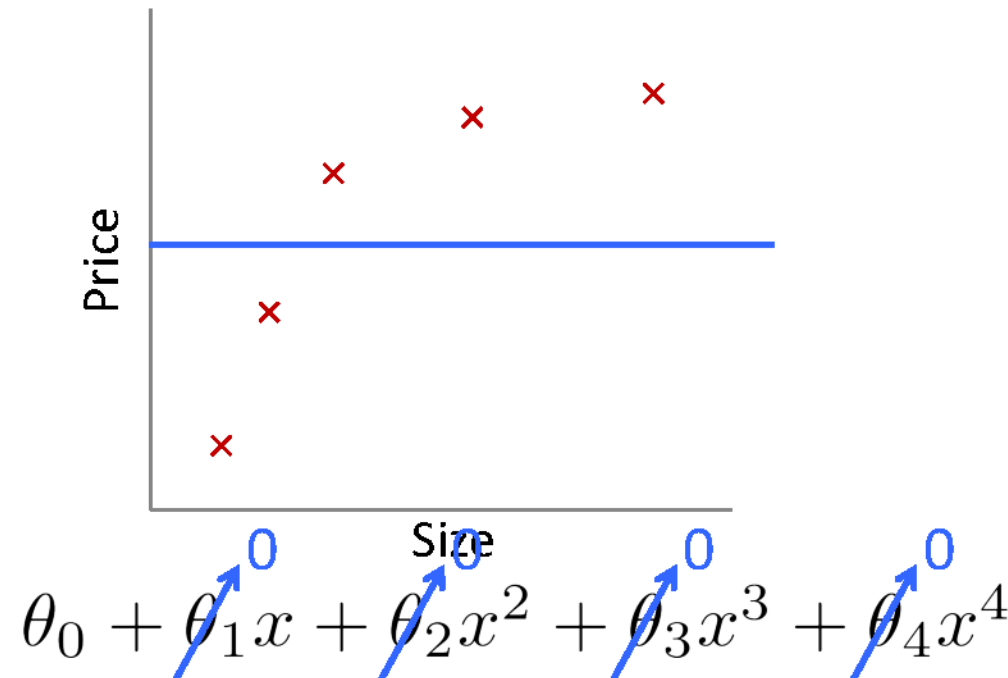


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- What happens if we set λ to be huge (e.g., 10^{10})?



Regularized linear regression

- Cost Function

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- Fit by solving , Gradient update $\min_{\theta} J(\theta)$

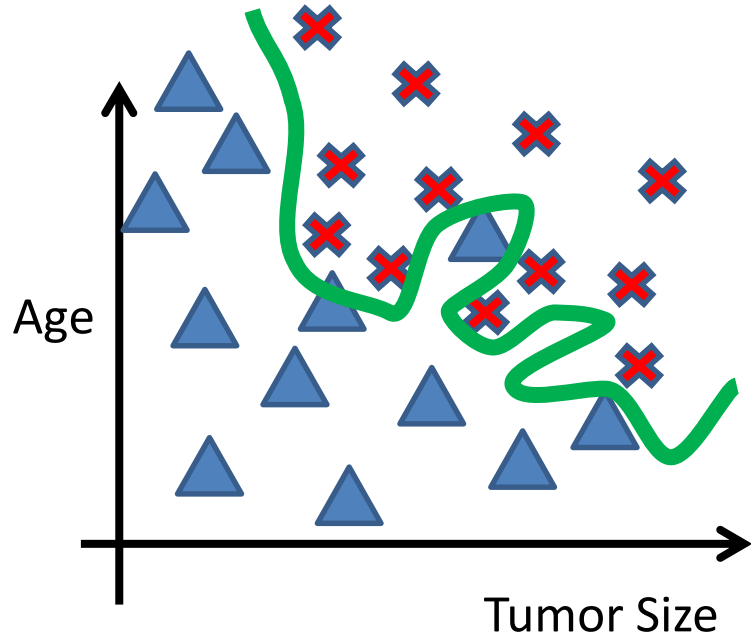
$$\begin{aligned} \frac{\partial}{\partial \theta_0} J(\theta) \quad \theta_0 &\leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) \\ \frac{\partial}{\partial \theta_j} J(\theta) \quad \theta_j &\leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \lambda \theta_j \end{aligned}$$

regularization

- Can be rearranged

$$\theta_j \leftarrow \theta_j (1 - \alpha \lambda) - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

Regularized logistic regression



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \theta_6 x_1^3 x_2 + \theta_7 x_1 x_2^3 + \dots)$$

- Cost function: $J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$

Regularized logistic regression

Regularized cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient descent with regularized cost function

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[\left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\}$$

}

Just like with linear regression, we will want to **separately** update θ_0 and the rest of the parameters because we do not want to regularize θ_0 .

$|\theta|_1$: Lasso regularization



$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j| \right]$$

LASSO: Least Absolute Shrinkage and Selection Operator

Regularization function

Name

$$\|\theta\|_2^2 = \sum_{j=1}^n \theta_j^2$$

Tikhonov regularization
Ridge regression

$$\|\theta\|_1 = \sum_{j=1}^n |\theta_j|$$

LASSO regression

$$\alpha \|\theta\|_1 + (1 - \alpha) \|\theta\|_2^2$$

Elastic net regularization

α is the mixing parameter between ridge ($\alpha = 0$) and lasso ($\alpha = 1$).

Bias-Variance Decomposition

Expectation of a random variable

- Let X be the random variable
- $E[X]$ = **mean** of a large number of observations of the random variable
- The expectation of any constant c : $E[c]=c$
- $E[E[X]] = E[X]$ because the expectation of a random variable is a constant.
- $E[X+Y] = E[X] + E[Y]$
- $E[cX] = E[c] E[X] = cE[X]$

Bias-Variance Tradeoff



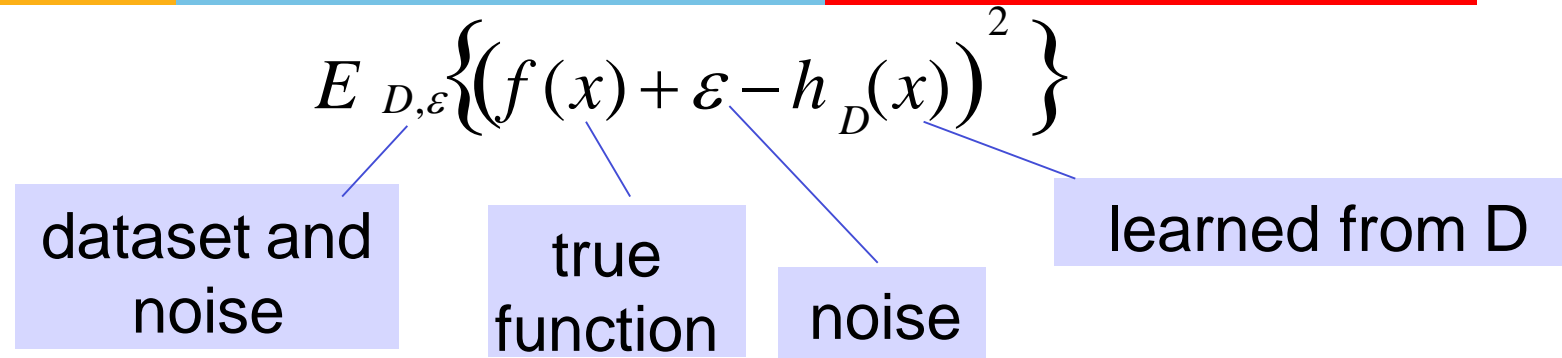
- **Bias:** difference between **what you expect to learn** and **truth**
 - Measures how well you expect to represent true solution
 - Decreases with more complex model
- **Variance:** difference between **what you expect to learn** and **what you learn from a particular dataset**
 - Measures how sensitive learner is to specific dataset
 - Increases with more complex model

generalisation error = Bias² + Variance + Irreducible Error

Bias – Variance decomposition of error

$$E_{D,\varepsilon} \left\{ \left(f(x) + \varepsilon - h_D(x) \right)^2 \right\}$$

dataset and noise true function noise learned from D



Fix test case x , then do this experiment:

1. Draw size n sample $D=(x_1, y_1), \dots, (x_n, y_n)$
2. Train linear regressor h_D using D
3. Draw one test example $(x, f(x)+\varepsilon)$
4. Measure squared error of h_D on that example x
5. What's the expected error?

Bias – Variance decomposition of error



Notation - to simplify this

$$f \equiv f(x) + \varepsilon \qquad \hat{y} = \hat{y}_D \equiv h_D(x)$$

$$E_{D,\varepsilon} \left\{ \left(\overbrace{f(x) + \varepsilon} - \overbrace{h_D(x)}^2 \right) \right\}$$

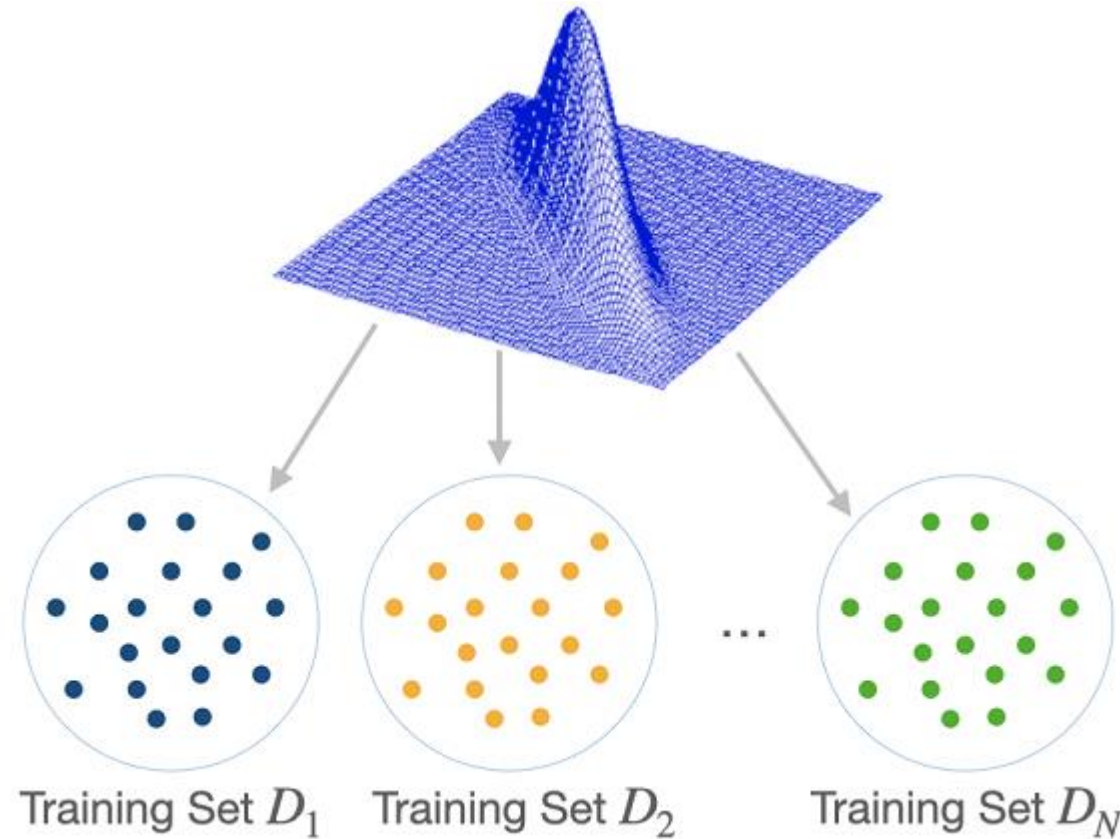
dataset and
noise

true
function

noise

learned from D

Bias – Variance decomposition of error



$$h \equiv E_D\{h_D(x)\}$$

long-term expectation of learner's prediction
on this x averaged over many data sets D

Bias – Variance decomposition of error

$$\begin{aligned}
 & E_{D,\varepsilon}\{(f-y)^2\} \\
 &= E\left\{([f-h]+[h-\hat{y}])^2\right\} \\
 &= E\left\{[f-h]^2+[h-\hat{y}]^2+2[f-h][h-\hat{y}]\right\} \\
 &= E\left\{[f-h]^2+[h-\hat{y}]^2+2[fh-f\hat{y}-h^2+h\hat{y}]\right\} \\
 &= E[(f-h)^2]+E[(h-\hat{y})^2]+2\left(\cancel{E[fh]}-\cancel{E[f\hat{y}]}-\cancel{E[h^2]}+\cancel{E[h\hat{y}]}\right)
 \end{aligned}$$

$h = E_D\{\hat{y}_D\}$
 $\hat{y} = \hat{y}_D \equiv h_D(x)$
 $f \equiv f(x) + \varepsilon$

Bias – Variance decomposition of error

$$E_{D,\varepsilon}\{(f - \hat{y})^2\} = E[(f - h)^2] + E[(h - \hat{y})^2]$$

$$h \equiv E_D\{h_D(x)\}$$

$$\hat{y} = \hat{y}_D \equiv h_D(x)$$

$$f \equiv f(x) + \varepsilon$$

BIAS²

Squared difference between best possible prediction for x , $f(x)$, and our “long-term” expectation for what the learner will do if we averaged over many datasets D , $E_D[h_D(x)]$

VARIANCE

Squared difference btwn our long-term expectation for the learners performance, $E_D[h_D(x)]$, and what we expect in a representative run on a dataset D (\hat{y})

Bayesian linear regression

- Bayesian analysis will show that
 - under certain assumptions any learning algorithm that minimizes the squared error between the output hypothesis predictions and the training data will output a maximum likelihood hypothesis

Training Linear Regression: MCLE, MLE



$$h_{ML} = \operatorname{argmax}_{h \in H} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

The first term in this expression is a constant independent of h , and can therefore be discarded, yielding

$$h_{ML} = \operatorname{argmax}_{h \in H} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

Maximizing this negative quantity is equivalent to minimizing the corresponding positive quantity.

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

Finally, we can again discard constants that are independent of h .

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2 \quad (6.6)$$

References



- <https://www.desmos.com/calculator/bgontvxotm>
- <https://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf>

Thanks