

Linear Regression (continued)

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Announcements

- Mid-term on Feb 10 (Monday).
 - ▶ In-class
 - ▶ Closed notes, closed book.
 - ▶ Will cover material from lectures and problem sets.
 - ▶ All material till and including lecture on Wednesday.
 - ▶ Please come to office hours with any questions.
 - ▶ Encouraged to go to discussion section on Friday. TAs will discuss some practice problems.

Outline

- 1 Linear regression
 - Univariate solution (review)
 - Probabilistic interpretation
- 2 Multivariate solution
- 3 Nonlinear hypotheses

Linear regression (ordinary least squares)

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Hypotheses/Model:

$$h_{\mathbf{w},b}, \text{ with } h_{\mathbf{w},b}(\mathbf{x}) = b + \sum_d w_d x_d = b + \mathbf{w}^T \mathbf{x} = \boldsymbol{\theta}^T \mathbf{x}$$

► $\mathbf{w} = [w_1 \ w_2 \ \cdots \ w_D]^T$: **weights**

b is called the **bias or offset or intercept term**.

$$\boldsymbol{\theta} = [b \ w_1 \ w_2 \ \cdots \ w_D]^T$$

- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

How do we learn parameters?

Minimize prediction error on training data

Minimize the sum of squared errors (also called residual sum of squares *RSS*):

cost function for linear regression.

$$J(\boldsymbol{\theta}) = \sum_n [y_n - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]^2$$

How do we minimize the cost function (residual sum of squares)?

Numerical optimization

- Gradient descent

Analytical solution

- Can compute minimum in closed form for linear regression!

A simple case: x is just one-dimensional ($D=1$)

Residual sum of squares (RSS)

$$J(\boldsymbol{\theta}) = \sum_n [y_n - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]^2 = \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2$$

Identify stationary points by taking derivative with respect to parameters and setting to zero

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_0} = 0 \Rightarrow -2 \sum_n [y_n - (\theta_0 + \theta_1 x_n)] = 0$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_1} = 0 \Rightarrow -2 \sum_n [y_n - (\theta_0 + \theta_1 x_n)] x_n = 0$$

Simplify these expressions to get “Normal Equations”

$$\sum y_n = N\theta_0 + \theta_1 \sum x_n$$

$$\sum x_n y_n = \theta_0 \sum x_n + \theta_1 \sum x_n^2$$

We have two equations and two unknowns! Do some algebra to get:

$$\theta_1 = \frac{\sum (x_n - \bar{x})(y_n - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{and} \quad \theta_0 = \bar{y} - \theta_1 \bar{x}$$

where $\bar{x} = \frac{1}{n} \sum_n x_n$ and $\bar{y} = \frac{1}{n} \sum_n y_n$.

Why is minimizing J sensible?

Probabilistic interpretation

- Noisy observation model

$$Y = \theta_0 + \theta_1 X + \eta$$

where $\eta \sim \mathcal{N}(0, \sigma^2)$ is a Gaussian or normal random variable

Use properties of Gaussian random variables

$$Z \sim \mathcal{N}(0, 1), W = a + bZ \quad \text{then} \quad W \sim \mathcal{N}(a, b^2)$$

$$Y \sim \mathcal{N}(\theta_0 + \theta_1 X, \sigma^2)$$

- Distribution of Y determined by parameters $(\theta_0, \theta_1, \sigma^2)$.
- Likelihood of one training sample (x_n, y_n)

$$p(y_n | x_n; \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\theta_0 + \theta_1 x_n, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{[y_n - (\theta_0 + \theta_1 x_n)]^2}{2\sigma^2}}$$

Probabilistic interpretation

Log-likelihood of the training data \mathcal{D} (assuming independence)

$$\begin{aligned}\mathcal{LL}(\theta, \sigma^2) &= \log P(\mathcal{D}) \\&= \log \prod_{n=1}^N p(y_n | x_n) = \sum_n \log p(y_n | x_n) \\&= \sum_n \left\{ -\frac{[y_n - (\theta_0 + \theta_1 x_n)]^2}{2\sigma^2} - \log \sqrt{2\pi\sigma^2} \right\} \\&= -\frac{1}{2\sigma^2} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 - \frac{N}{2} \log 2\pi\sigma^2 \\&= -\frac{1}{2} \left\{ \frac{1}{\sigma^2} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 + N \log \sigma^2 \right\} + \text{const}\end{aligned}$$

What is the relationship between minimizing J and maximizing the log-likelihood?

Maximum likelihood estimation

Estimating σ , θ_0 and θ_1 can be done in two steps

- Maximize over θ_0 and θ_1 for a fixed σ

$$\max \log P(\mathcal{D}) \Leftrightarrow \min \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 \leftarrow \text{That is } J(\boldsymbol{\theta})!$$

- Maximize over σ

Why does the probabilistic interpretation help us?

Gives us a template for modeling

- Probabilistic model $P(x; \theta)$
 - ▶ “Story” for generating the data
- Given data x and probabilistic model $P(x; \theta)$, we can find a “good” value for θ .
 - ▶ Maximize the likelihood
- We have seen this for logistic regression, linear regression (minimizing the RSS) but this principle is very general.
- Makes clear what the assumptions are.

Why does the probabilistic interpretation help us?

Assumptions underlying linear regression

y_n are

- Independent
- Normally distributed
- Mean is a linear function of x_n
- Constant variance σ^2

Outline

1 Linear regression

2 Multivariate solution

- Input is D-dimensional
- Computational and numerical optimization

3 Nonlinear hypotheses

Linear regression when \mathbf{x} is D-dimensional

$J(\boldsymbol{\theta})$ in matrix form

$$J(\boldsymbol{\theta}) = \sum_n [y_n - (\theta_0 + \sum_d \theta_d x_{nd})]^2 = \sum_n [y_n - \boldsymbol{\theta}^T \mathbf{x}_n]^2$$

where we have redefined some variables (by augmenting)

$$\mathbf{x} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \boldsymbol{\theta} \leftarrow [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_D]^T$$

$J(\boldsymbol{\theta})$ in new notations

Design matrix and target vector

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times (D+1)}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

$J(\boldsymbol{\theta})$ in new notations

Design matrix and target vector

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times (D+1)}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Vector of predictions for a given $\boldsymbol{\theta}$

$$\mathbf{X}\boldsymbol{\theta} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \boldsymbol{\theta} = \begin{pmatrix} \mathbf{x}_1^T \boldsymbol{\theta} \\ \mathbf{x}_2^T \boldsymbol{\theta} \\ \vdots \\ \mathbf{x}_N^T \boldsymbol{\theta} \end{pmatrix}$$

$J(\boldsymbol{\theta})$ in new notations

Vector of errors for a given $\boldsymbol{\theta}$

$$\mathbf{y} - \mathbf{X}\boldsymbol{\theta} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \boldsymbol{\theta} \\ \mathbf{x}_2^T \boldsymbol{\theta} \\ \vdots \\ \mathbf{x}_N^T \boldsymbol{\theta} \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \boldsymbol{\theta} \\ y_2 - \mathbf{x}_2^T \boldsymbol{\theta} \\ \vdots \\ y_N - \mathbf{x}_N^T \boldsymbol{\theta} \end{pmatrix}$$

$J(\theta)$ in new notations

Vector of errors for a given θ

$$\mathbf{y} - \mathbf{X}\theta = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \theta \\ \mathbf{x}_2^T \theta \\ \vdots \\ \mathbf{x}_N^T \theta \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \theta \\ y_2 - \mathbf{x}_2^T \theta \\ \vdots \\ y_N - \mathbf{x}_N^T \theta \end{pmatrix}$$

Expression for $J(\theta)$

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 &= (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \\ &= (y_1 - \mathbf{x}_1^T \theta \quad \cdots \quad y_N - \mathbf{x}_N^T \theta) \begin{pmatrix} y_1 - \mathbf{x}_1^T \theta \\ \vdots \\ y_N - \mathbf{x}_N^T \theta \end{pmatrix} \\ &= \sum_n [y_n - \mathbf{x}_n^T \theta]^2 \leftarrow \text{That is } J(\theta) \end{aligned}$$

$J(\boldsymbol{\theta})$ in new notations

Compact expression

$$\begin{aligned} J(\boldsymbol{\theta}) &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= (\mathbf{y}^\top - \{\mathbf{X}\boldsymbol{\theta}\}^\top) (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= (\mathbf{y}^\top - \boldsymbol{\theta}^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= (\mathbf{y}^\top) (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} - \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \\ &= \text{constant} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \end{aligned}$$

Solution in matrix form

Compact expression

$$J(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \left\{ \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2 (\mathbf{X}^T \mathbf{y})^T \boldsymbol{\theta} \right\} + \text{constant}$$

Gradients of Linear and Quadratic Functions

- $\nabla \mathbf{b}^T \mathbf{x} = \mathbf{b}$
- $\nabla \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$ (symmetric \mathbf{A})

Solution in matrix form

Compact expression

$$J(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \left\{ \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2 (\mathbf{X}^T \mathbf{y})^T \boldsymbol{\theta} \right\} + \text{constant}$$

Gradients of Linear and Quadratic Functions

- $\nabla \mathbf{b}^T \mathbf{x} = \mathbf{b}$
- $\nabla \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$ (symmetric \mathbf{A})

Normal equations

$$\nabla J(\boldsymbol{\theta}) = 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y} = 0$$

Solution in matrix form

Compact expression

$$J(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \left\{ \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2 (\mathbf{X}^T \mathbf{y})^T \boldsymbol{\theta} \right\} + \text{constant}$$

Gradients of Linear and Quadratic Functions

- $\nabla \mathbf{b}^T \mathbf{x} = \mathbf{b}$
- $\nabla \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$ (symmetric \mathbf{A})

Normal equations

$$\nabla J(\boldsymbol{\theta}) = 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y} = 0$$

This leads to the linear regression solution

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

Mini-Summary

- Linear regression is the linear combination of features
 $h : \mathbf{x} \rightarrow y$, with $h(\mathbf{x}) = \theta_0 + \sum_d \theta_d x_d = \boldsymbol{\theta}^T \mathbf{x}$
- If we minimize residual sum of squares as our learning objective, we get a closed-form solution of parameters
- Probabilistic interpretation: maximum likelihood if assuming the output is Gaussian distributed

Computational complexity

Bottleneck of computing the solution?

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

Computational complexity

Bottleneck of computing the solution?

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

Matrix multiply of $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$

Inverting the matrix $\mathbf{X}^T \mathbf{X}$

How many operations do we need?

Computational complexity

Bottleneck of computing the solution?

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

Matrix multiply of $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$

Inverting the matrix $\mathbf{X}^T \mathbf{X}$

How many operations do we need?

- $O(ND^2)$ for matrix multiplication
- $O(D^3)$ for matrix inversion
- Impractical for very large D or N

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

Algorithm 1 Gradient Descent (J)

- 1: $t \leftarrow 0$
 - 2: Initialize $\theta^{(0)}$
 - 3: **repeat**
 - 4: $\nabla J(\theta^{(t)}) = \mathbf{X}^T \mathbf{X} \theta^{(t)} - \mathbf{X}^T \mathbf{y} = \sum_n (\mathbf{x}_n^T \theta^{(t)} - y_n) \mathbf{x}_n$
 - 5: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla J(\theta^{(t)})$
 - 6: $t \leftarrow t + 1$
 - 7: **until** convergence
 - 8: Return final value of θ
-

What is the complexity of each iteration?

Why would this work?

If gradient descent converges, it will converge to the same solution as using matrix inversion.

This is because $J(\theta)$ is a convex function in its parameters θ

Stochastic gradient descent

Update parameters using one example at a time

Algorithm 2 Stochastic Gradient Descent (J)

- 1: $t \leftarrow 0$
 - 2: Initialize $\theta^{(0)}$
 - 3: **repeat**
 - 4: Randomly choose a training a sample x_t
 - 5: Compute its contribution to the gradient $g_t = (x_t^T \theta^{(t)} - y_t)x_t$
 - 6: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta g_t$
 - 7: $t \leftarrow t + 1$
 - 8: **until** convergence
 - 9: Return final value of θ
-

How does the complexity per iteration of stochastic gradient descent (SGD) compare with gradient descent (GD)?

How does the complexity per iteration of stochastic gradient descent (SGD) compare with gradient descent (GD)?

- $O(ND)$ for GD versus $O(D)$ for SGD

Mini-summary

- Batch gradient descent computes the exact gradient.
- Stochastic gradient descent approximates the gradient with a single data point;
- Mini-batch variant: trade-off between accuracy of estimating gradient and computational cost
- Similar ideas extend to other ML optimization problems.
 - ▶ For large-scale problems, stochastic gradient descent often works well.

What if $\mathbf{X}^T \mathbf{X}$ is not invertible

Can you think of any reasons why that could happen?

What if $\mathbf{X}^T \mathbf{X}$ is not invertible

Can you think of any reasons why that could happen?

Answer 1: $N < D$. Intuitively, not enough data to estimate all the parameters.

What if $\mathbf{X}^T \mathbf{X}$ is not invertible

Can you think of any reasons why that could happen?

Answer 1: $N < D$. Intuitively, not enough data to estimate all the parameters.

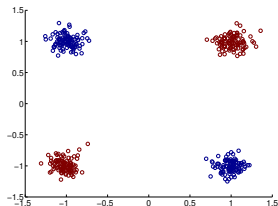
Answer 2: \mathbf{X} columns are not linearly independent. Intuitively, there are two features that are perfectly correlated. In this case, solution is not unique.

Outline

- 1 Linear regression
- 2 Multivariate solution
- 3 Nonlinear hypotheses**

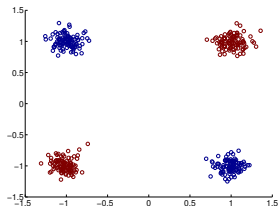
What if data is not linearly separable or fits to a line

Example of nonlinear classification

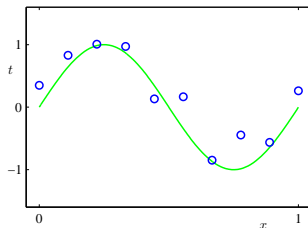


What if data is not linearly separable or fits to a line

Example of nonlinear classification



Example of nonlinear regression



Nonlinear basis for classification

Transform the input/feature

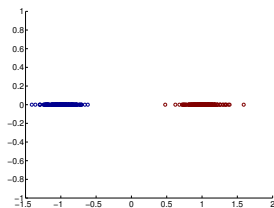
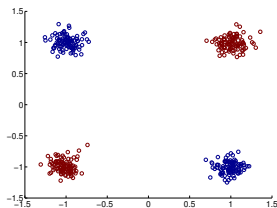
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Nonlinear basis for classification

Transform the input/feature

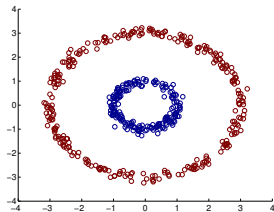
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Transformed training data: linearly separable!



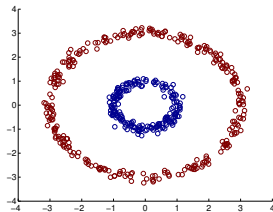
Another example

How to transform the input/feature?



Another example

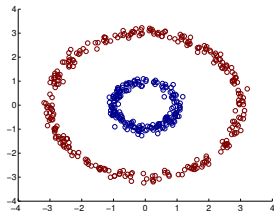
How to transform the input/feature?



$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow \mathbf{z} = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

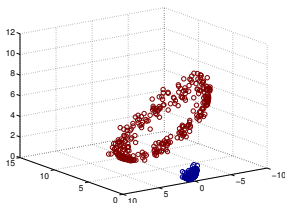
Another example

How to transform the input/feature?



$$\phi(x) : x \in \mathbb{R}^2 \rightarrow z = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

Transformed training data: linearly separable



General nonlinear basis functions

We can use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

where M is the dimensionality of the new feature/input \mathbf{z} (or $\phi(\mathbf{x})$). Note that M could be either greater than D or less than or the same.

General nonlinear basis functions

We can use a nonlinear mapping

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$$

where M is the dimensionality of the new feature/input \mathbf{z} (or $\phi(\mathbf{x})$). Note that M could be either greater than D or less than or the same.

With the new features, we can apply our learning techniques to minimize our errors on the transformed training data

- linear methods: prediction is based on $\theta^T \phi(\mathbf{x})$
- other methods: nearest neighbors, decision trees, etc

Regression with nonlinear basis

Residual sum squares

$$\sum_n [\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n]^2$$

where $\boldsymbol{\theta} \in \mathbb{R}^M$, the same dimensionality as the transformed features $\boldsymbol{\phi}(\mathbf{x})$.

Regression with nonlinear basis

Residual sum squares

$$\sum_n [\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}_n) - y_n]^2$$

where $\boldsymbol{\theta} \in \mathbb{R}^M$, the same dimensionality as the transformed features $\boldsymbol{\phi}(\mathbf{x})$.

The linear regression solution can be formulated with the new design matrix

$$\boldsymbol{\Phi} = \begin{pmatrix} \boldsymbol{\phi}(\mathbf{x}_1)^T \\ \boldsymbol{\phi}(\mathbf{x}_2)^T \\ \vdots \\ \boldsymbol{\phi}(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}, \quad \hat{\boldsymbol{\theta}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

Example with regression

Polynomial basis functions

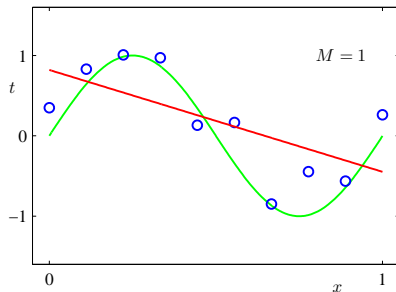
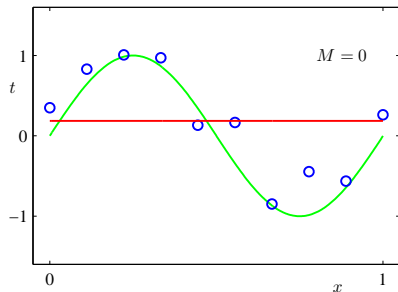
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow h(x) = \theta_0 + \sum_{m=1}^M \theta_m x^m$$

Example with regression

Polynomial basis functions

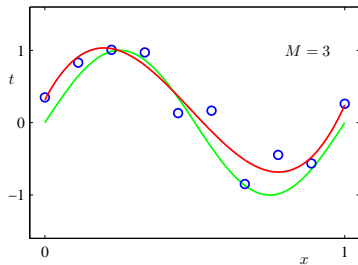
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow h(x) = \theta_0 + \sum_{m=1}^M \theta_m x^m$$

Fitting samples from a sine function: *underrfitting* as $h(x)$ is too simple



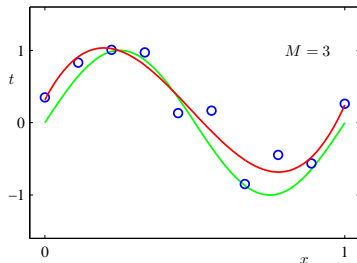
Adding high-order terms

$M=3$

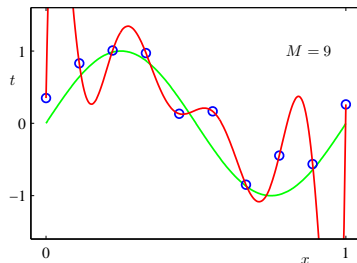


Adding high-order terms

M=3



M=9: *overfitting*



More complex features lead to better results on the training data, but potentially worse results on new data, e.g., test data!

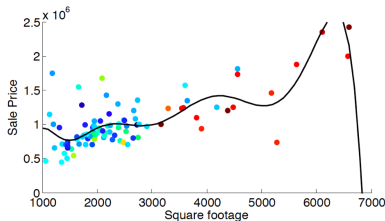
Overfitting

Parameters for higher-order polynomials are very large

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
θ_0	0.19	0.82	0.31	0.35
θ_1		-1.27	7.99	232.37
θ_2			-25.43	-5321.83
θ_3			17.37	48568.31
θ_4				-231639.30
θ_5				640042.26
θ_6				-1061800.52
θ_7				1042400.18
θ_8				-557682.99
θ_9				125201.43

Overfitting can be quite disastrous

Fitting the housing price data with $M = 3$



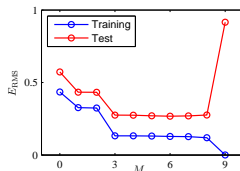
Note that the price would go to zero (or negative) if you buy bigger ones!
This is called poor generalization/overfitting.

Detecting overfitting

Plot model complexity versus objective function

As a model increases in complexity, performance on training data keeps improving while performance on test data may first improve but eventually deteriorate.

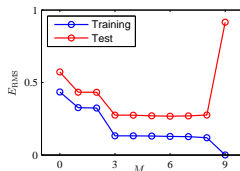
- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.



Detecting overfitting

Plot model complexity versus objective function

As a model increases in complexity, performance on training data keeps improving while performance on test data may first improve but eventually deteriorate.



- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.
- Vertical axis:
 - 1 For regression, residual sum of squares or residual mean squared (squared root of RSS)
 - 2 For classification, classification error rate.

Summary

- Hypotheses/models: linear functions of features.
- Objective: choose a function (*i.e.*, parameters for the function) that minimize a cost function.
 - ▶ Cost function measures loss or error of the predictions made by a model/hypothesis on training set.
 - ▶ Cost function depends on the learning problem.
- Probabilistic interpretation
 - ▶ Minimizing cost function equivalent to maximizing likelihood.
 - ▶ Allows us to understand assumptions and to generalize our models.
- How do we minimize the cost function ?
 - ▶ Numerical methods: gradient and stochastic gradient descent.
 - ▶ Sometimes analytical solutions are available.
 - ▶ Computational considerations influence the choice.
- Can allow non-linear models/hypotheses by transforming features using non-linear functions.