

Linear Regression. Overfitting and regularization

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.

Linear regression cost function (residual sum of squares)

- Input: $\mathbf{x} \in \mathbb{R}$
- Output: $y \in \mathbb{R}$
- Hypotheses/Model: h_{θ} , with $h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 = \underline{\underline{\theta^T \mathbf{x}}}$

Residual sum of squares (RSS)

$$\underset{\theta}{\operatorname{argmin}} \ J(\underline{\underline{\theta}}) = \sum_n \underbrace{[y_n - \underline{\underline{h_{\theta}(\mathbf{x}_n)}}]}_{\text{residual}}^2 = \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2$$

Outline

- 1 Probabilistic interpretation
- 2 Multivariate solution
- 3 Nonlinear hypotheses
- 4 Basic ideas to overcome overfitting

Why is minimizing J sensible?

Probabilistic interpretation

- Noisy observation model

$$\begin{aligned} y_1 &= \theta_0 + \theta_1 x_1 + \eta_1 \rightarrow (x_1, y_1) \\ y_2 &= \theta_0 + \theta_1 x_2 + \eta_2 \rightarrow (x_2, y_2) \\ &\vdots \\ y_N &= \theta_0 + \theta_1 x_N + \eta_N \rightarrow (x_N, y_N) \end{aligned}$$

Known (assumption!)

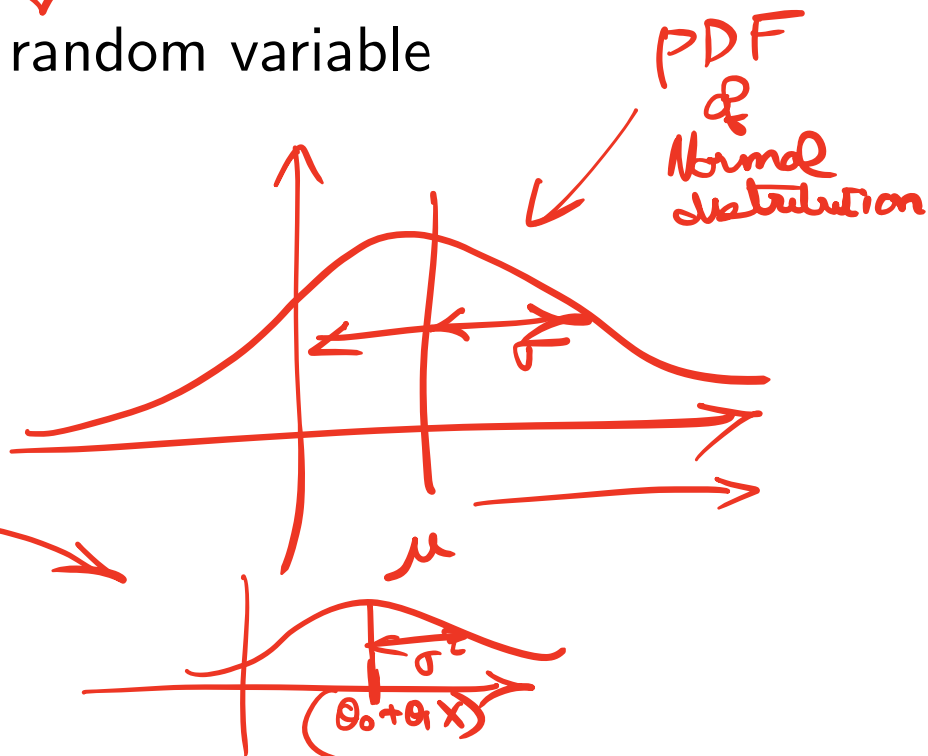
$$Y = \theta_0 + \theta_1 X + \eta \Rightarrow \mathcal{N}(\theta_0, \theta_1, \sigma^2) = P(\text{Data} | \theta_0, \theta_1, \sigma^2)$$

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

Mean \rightarrow Variance

$$y = \theta_0 + \theta_1 x + \eta$$

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



Why is minimizing J sensible?

Probabilistic interpretation

- Noisy observation model

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

$$x \sim \mathcal{N}(\mu, \sigma^2)$$
$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

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

- Likelihood of one training sample (x_n, y_n)

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

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

Probabilistic interpretation (cont'd)

Log-likelihood of the training data \mathcal{D} (assuming i.i.d)

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

Probabilistic interpretation (cont'd)

Log-likelihood of the training data \mathcal{D} (assuming i.i.d)

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

Probabilistic interpretation (cont'd)

Log-likelihood of the training data \mathcal{D} (assuming i.i.d)

$$\mathcal{LL}(\boldsymbol{\theta}) = \log P(\mathcal{D})$$

$$= \log \prod_{n=1}^N p(y_n | x_n) = \sum_n \log p(y_n | x_n)$$

$$= \sum_n \left\{ \frac{[y_n - (\theta_0 + \theta_1 x_n)]^2}{2\sigma^2} - \log \sqrt{2\pi} \sigma \right\}$$

$$= \left(-\frac{1}{2\sigma^2} \right) \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 - \frac{N}{2} \log \sigma^2 - N \log \sqrt{2\pi}$$

$$\begin{aligned} & \log \sqrt{2\pi\sigma^2} \\ &= \frac{1}{2} \log(2\pi\sigma^2) \\ &= \left(\frac{1}{2} \log 2\pi + \frac{1}{2} \log \sigma^2 \right) \end{aligned}$$

Probabilistic interpretation (cont'd)

Log-likelihood of the training data \mathcal{D} (assuming i.i.d)

$$\underset{\theta}{\text{argmax}} \mathcal{L}(\theta) = \log P(\mathcal{D})$$

$$\underset{\theta}{\text{argmin}} J(\theta) = \frac{\sum_n (y_n - (\theta_0 + \theta_1 x_n))^2}{2}$$

$$= \log \prod_{n=1}^N p(y_n | x_n) = \sum_n \log p(y_n | x_n)$$

$$= \sum_n \left\{ -\frac{[y_n - (\theta_0 + \theta_1 x_n)]^2}{2\sigma^2} - \log \sqrt{2\pi}\sigma \right\}$$

$$= -\frac{1}{2\sigma^2} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 - \frac{N}{2} \log \sigma^2 - \underline{\underline{N \log \sqrt{2\pi}}}$$

$$= -\frac{1}{2} \left\{ \frac{1}{\sigma^2} \left(\sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 \right) + N \log \sigma^2 \right\} + \underline{\underline{\text{const}}}$$

$$\underset{\theta}{\text{argmin}} \left\{ \frac{1}{\sigma^2} \sum_n (y_n - (\theta_0 + \theta_1 x_n))^2 + N \log \sigma^2 \right\}$$

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

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

Maximum likelihood estimation

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

- Maximize over θ_0 and θ_1

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

- Maximize over $s = \sigma^2$ (we could estimate σ directly)

$$\frac{\partial \log P(\mathcal{D})}{\partial s} = -\frac{1}{2} \left\{ -\frac{1}{s^2} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 + N \frac{1}{s} \right\} = 0$$

$$s^2 \times \frac{1}{s^2} \sum_n (y_n - (\theta_0 + \theta_1 x_n))^2 = \frac{N}{s} \times s^2$$

Maximum likelihood estimation

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

- Maximize over θ_0 and θ_1

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

- Maximize over $s = \sigma^2$ (we could estimate σ directly)

$$\frac{\partial \log P(\mathcal{D})}{\partial s} = -\frac{1}{2} \left\{ -\frac{1}{s^2} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2 + N \frac{1}{s} \right\} = 0$$

$$\rightarrow \sigma^{*2} = s^* = \frac{1}{N} \sum_n [y_n - (\theta_0 + \theta_1 x_n)]^2$$

$$(\theta_0 + \theta_1 x_n) + \eta$$

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.

$$P(x; \theta)$$

$$\sum_n (y_n - (\theta_0 + \theta_1 x_n))^2$$

$$y_n = \theta_0 + \theta_1 x_n + \eta$$

Independence

$$\eta \sim \mathcal{N}(0, \sigma^2)$$

y_{n+1}

Why does the probabilistic interpretation help us?

Assumptions underlying linear regression

y_n are

- Independent
- Normally distributed
- Mean is a linear function of $x_n \rightarrow \theta_0 + \theta_1 x_n$
- Constant variance σ^2

Outline

- 1 Probabilistic interpretation
- 2 Multivariate solution
 - D-dimensional inputs
 - Computational and numerical optimization
- 3 Nonlinear hypotheses
- 4 Basic ideas to overcome overfitting

Linear regression when x is D-dimensional

$J(\theta)$ in matrix form

$$\underset{\theta}{\operatorname{argmin}} J(\theta) = \sum_n [y_n - (\theta_0 + \sum_d \theta_d x_{nd})]^2 = \sum_n [y_n - \theta^T x_n]^2$$

where we have redefined some variables (by augmenting)

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

$J(\theta)$ in new notations

$$x_1 = \begin{pmatrix} x_{1,0} \\ x_{1,1} \\ \vdots \\ x_{1,D} \end{pmatrix}$$

Design matrix and target vector

$$\underline{\underline{X}} = \begin{pmatrix} \rightarrow x_1^T \\ \rightarrow x_2^T \\ \vdots \\ \rightarrow x_N^T \end{pmatrix} \in \mathbb{R}^{N \times (D+1)}, \quad \underline{\underline{y}} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

$(N \times 1)$

Vector of predictions for a given θ

$$X\theta = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} \theta = \begin{pmatrix} x_1^T \theta \\ x_2^T \theta \\ \vdots \\ x_N^T \theta \end{pmatrix}$$

$$= \begin{pmatrix} x_1^T \theta \\ x_2^T \theta \\ \vdots \\ x_N^T \theta \end{pmatrix} = \underline{\underline{X\theta}}$$

$J(\theta)$ in new notations

Vector of errors for a given θ

$$\begin{aligned} \underline{\underline{y}} - \underline{\underline{X\theta}} &= \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} - \begin{pmatrix} x_1^T \theta \\ x_2^T \theta \\ \vdots \\ x_N^T \theta \end{pmatrix} = \begin{pmatrix} y_1 - x_1^T \theta \\ y_2 - x_2^T \theta \\ \vdots \\ y_N - x_N^T \theta \end{pmatrix} \end{aligned}$$

(N x 1)

error on datapoint 1
error on datapoint 2
...

Expression for $J(\theta)$

$$\begin{aligned} \underline{\underline{\|y - X\theta\|_2^2}} &= (\underline{\underline{y - X\theta}})^T (\underline{\underline{y - X\theta}}) \\ &= (y_1 - x_1^T \theta \quad \cdots \quad y_N - x_N^T \theta) \begin{pmatrix} y_1 - x_1^T \theta \\ \vdots \\ y_N - x_N^T \theta \end{pmatrix} \\ &= \sum_n [y_n - x_n^T \theta]^2 \leftarrow \text{That is } \underline{\underline{J(\theta)}} \end{aligned}$$

$\|x\|_2^2 = (x^T x) = \sum_{d=1}^D x_d^2$

$\nabla J(\theta) = 0$

$J(\theta)$ in new notations

$$f(x) = \underline{b^T x} = \underline{b_1 x_1 + b_2 x_2 + \dots + b_n x_n}$$

$x \in \mathbb{R}^D$

$$\nabla f(x) = \underline{b}$$

Compact expression

$$\nabla_{\theta} J(\theta)$$

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

$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_D \end{pmatrix} = b$

$y^T X\theta = (X\theta)^T y = -2y^T X\theta$

$\nabla(-2(X^T y)^T \theta) = -2X^T y$

Solution in matrix form

Compact expression

$$\nabla (\theta^T \underline{x^T x} \theta) = 2(x^T x) \theta$$

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

Gradients of Linear and Quadratic Functions

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

$$f(x) = x^T \underline{Ax}$$

$x \in \mathbb{R}^D$

$(D \times D)$

$$\nabla f(x) = 2Ax$$

$$\frac{d}{dx} (Ax^2) = 2Ax$$

Solution in matrix form

Compact expression

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

Gradients of Linear and Quadratic Functions

- $\nabla b^T x = b$
- $\nabla x^T A x = 2Ax$ (symmetric A)

Normal equations

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

Handwritten notes below the equation:

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

$\mathbf{A} \boldsymbol{\theta} = \mathbf{b}$

Arrows indicate the mapping from the matrix equation to the standard linear system form.

Solution in matrix form

$$(X^T X)^{-1} \stackrel{?}{=} X^{-1} X^{-T}$$

Compact expression

$$\begin{matrix} X \\ N \times (D+1) \end{matrix} \quad \begin{matrix} (X^T X) \\ (D+1) \times (D+1) \end{matrix}$$

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

Gradients of Linear and Quadratic Functions

- $\nabla b^T x = b$
- $\nabla x^T A x = 2Ax$ (symmetric A)

Normal equations

$$\nabla J(\theta) = 2X^T X \theta - 2X^T y = 0$$

This leads to the linear regression solution

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

$$\begin{aligned} (X^T X)^{-1} &\downarrow \\ (X^T X) \theta &= X^T y \\ \theta &= (X^T X)^{-1} X^T y \end{aligned}$$

Computational complexity

Bottleneck of computing the solution?

~~$N \times D$~~

y
 $N \times 1$

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

$O(D^2 N) \Rightarrow X^T X$

$(X^T X)$
 $(D \times D)$

$(X^T X)^{-1}$
 $O(D^3)$

$X^T y$
 $(D \times 1)$
 $O(ND)$

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?

Computational complexity

Bottleneck of computing the solution?

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

Matrix multiply of $X^T X \in \mathbb{R}^{(D+1) \times (D+1)}$
Inverting the matrix $\underline{\underline{X^T X}}$

Intercept $\rightarrow \theta_0 + \underline{\theta_1} x_{n1} + \theta_2 x_{n2} + \dots + \underline{\theta_D} x_{nD}$

$x_n = (1, x_{n1}, x_{n2}, \dots, x_{nD})$
 $\theta = (\theta_0, \theta_1, \dots, \theta_D)$

$\theta^T x_n$

How many operations do we need?

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

$\begin{matrix} \uparrow & \uparrow \\ A & B = I \\ D \times D & D \times D \end{matrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1, \dots \end{pmatrix}$

$(D \times D) \times D \rightarrow \begin{matrix} A b_1 = I_1 \\ A b_2 = I_2 \dots A b_D = I_D \end{matrix}$

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: $\nabla J(\theta^{(t)}) = \mathbf{X}^T \mathbf{X} \theta^{(t)} - \mathbf{X}^T \mathbf{y} = \sum_n (\underbrace{\mathbf{x}_n^T \theta^{(t)}}_{O(D)} - y_n) \underbrace{\mathbf{x}_n}_{O(D)} \Rightarrow O(ND)$

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

Hessian matrix of $J(\theta)$ is positive-semi-definite

Stochastic gradient descent

Update parameters using one example at a time

$O(ND)$

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: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \underline{\eta g_t}$
 - 7: $t \leftarrow t + 1$
 - 8: **until** convergence
 - 9: Return final value of θ
-

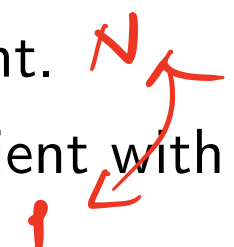
$O(D)$

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 $X^T X$ is not invertible

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

Can you think of any reasons why that could happen?

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

$$\begin{array}{c} 1, 1 \\ \theta \rightarrow 100, -100 \\ \quad \rightarrow 10, -10 \end{array}$$

What if $X^T X$ 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 $X^T X$ 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.

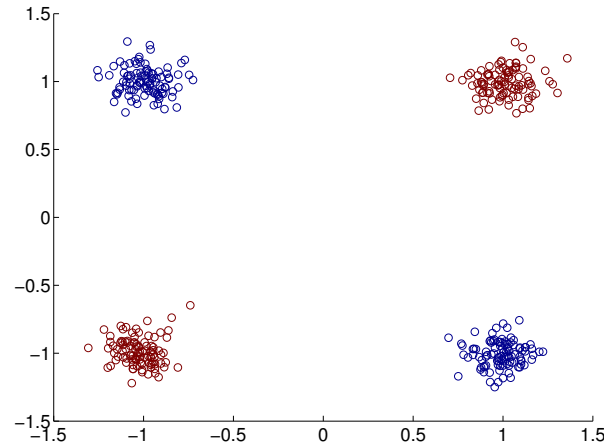
Answer 2: X columns are not linearly independent. Intuitively, there are two features that are perfectly correlated. In this case, solution is not unique.

Outline

- 1 Probabilistic interpretation
- 2 Multivariate solution
- 3 Nonlinear hypotheses**
- 4 Basic ideas to overcome overfitting

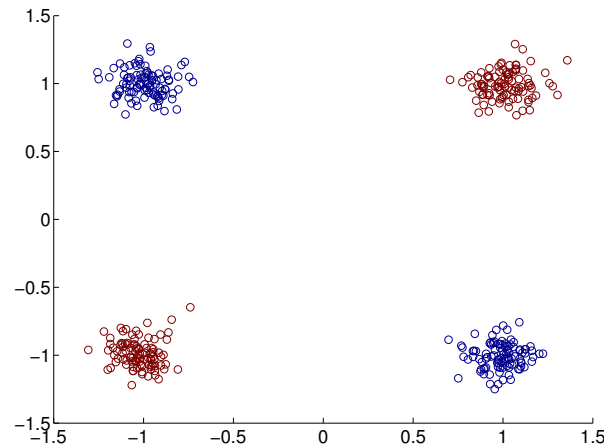
What if data is not linearly separable or fits to a line

Example of nonlinear classification

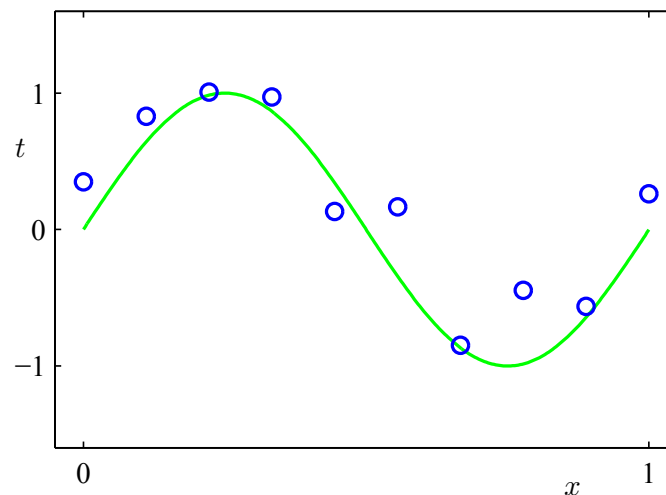


What if data is not linearly separable or fits to a line

Example of nonlinear classification



Example of nonlinear regression



Nonlinear basis for classification

Transform the input/feature

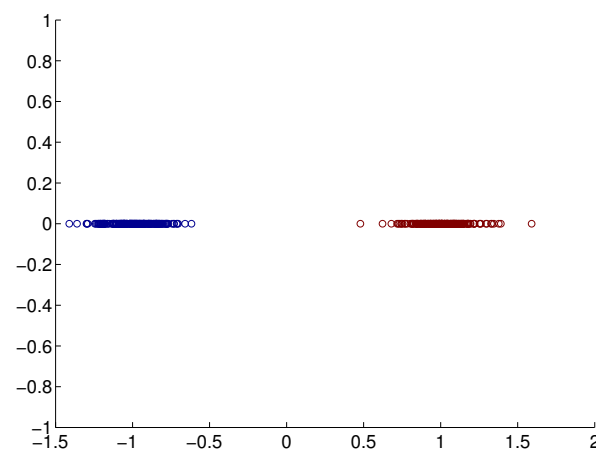
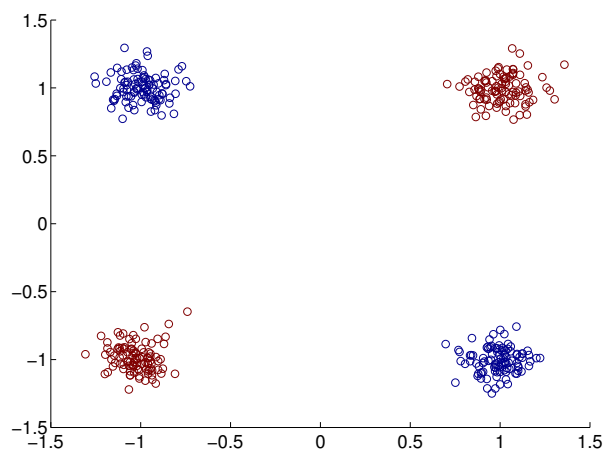
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Nonlinear basis for classification

Transform the input/feature

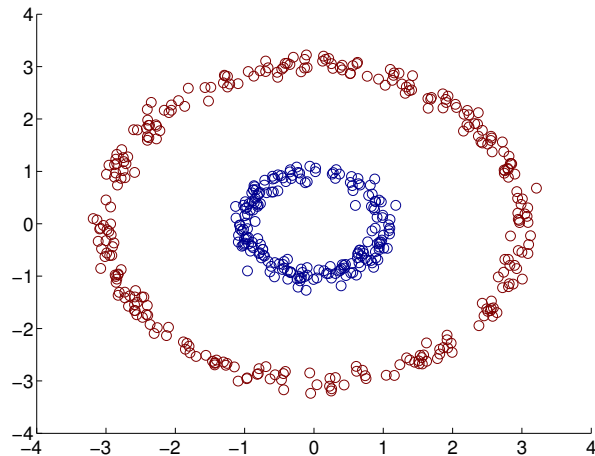
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Transformed training data: linearly separable!



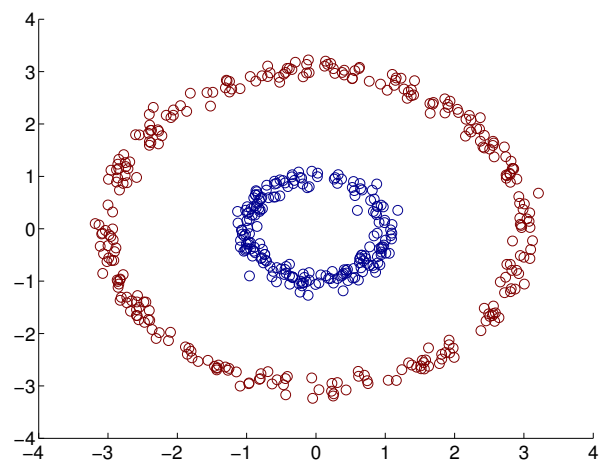
Another example

How to transform the input/feature?



Another example

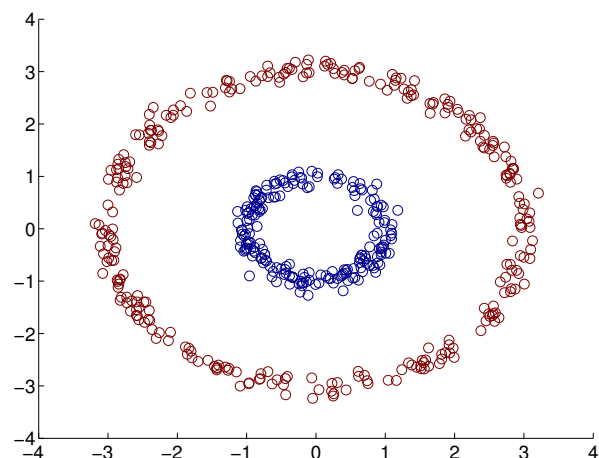
How to transform the input/feature?



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

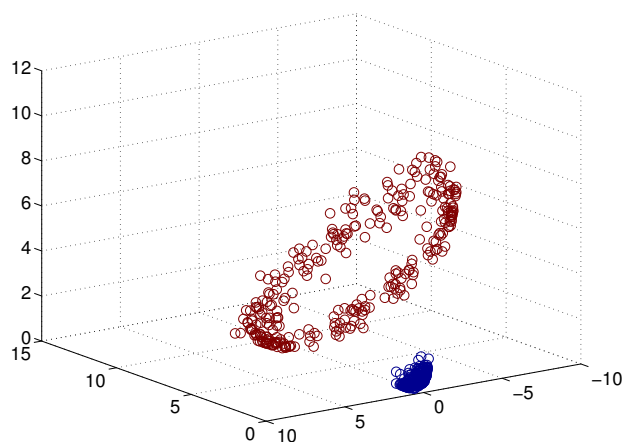
Another example

How to transform the input/feature?



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

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

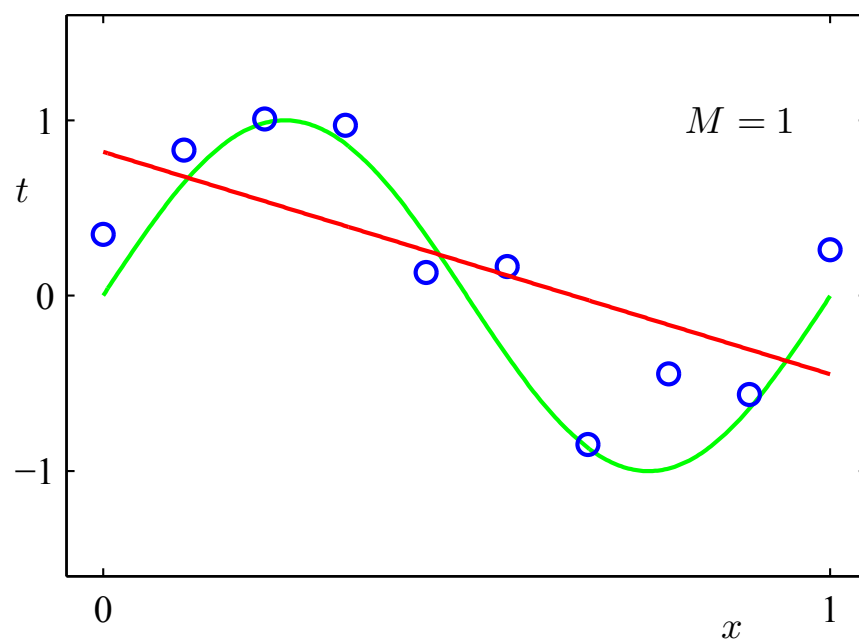
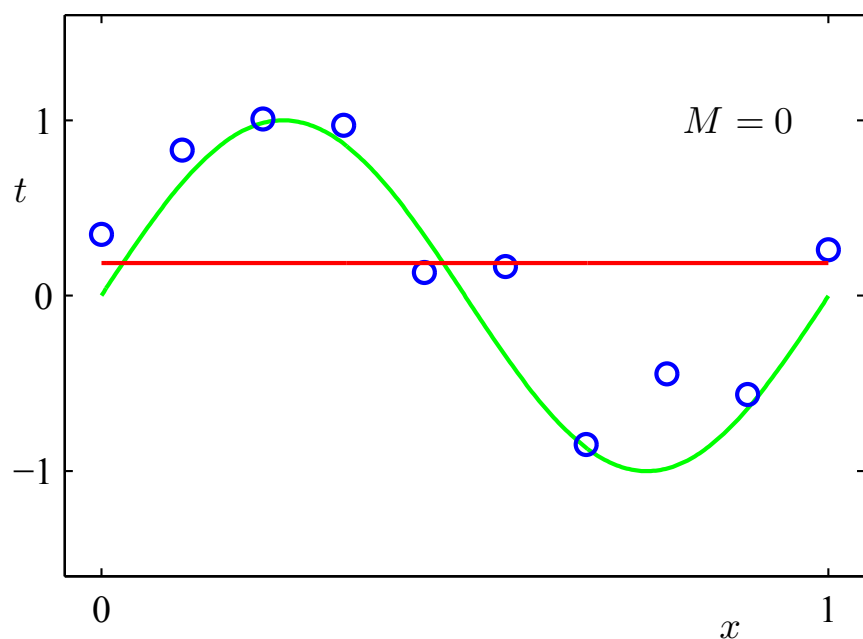
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow h(x) = \theta_0 + \sum_{m=1}^M \theta_m x^m$$

Example with regression

Polynomial basis functions

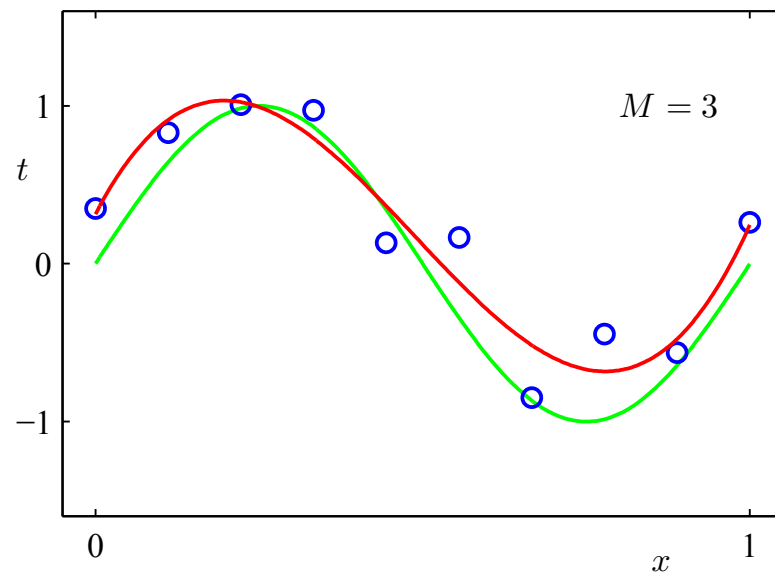
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow h(x) = \theta_0 + \sum_{m=1}^M \theta_m x^m$$

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



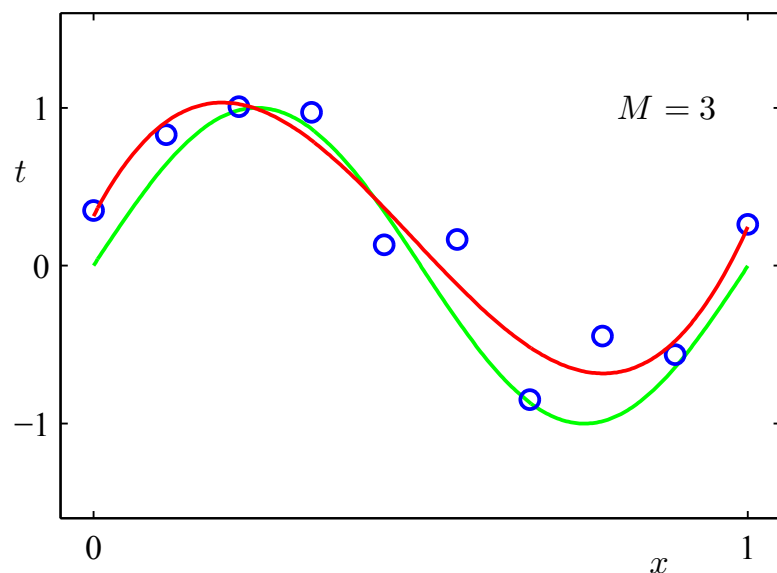
Adding high-order terms

M=3

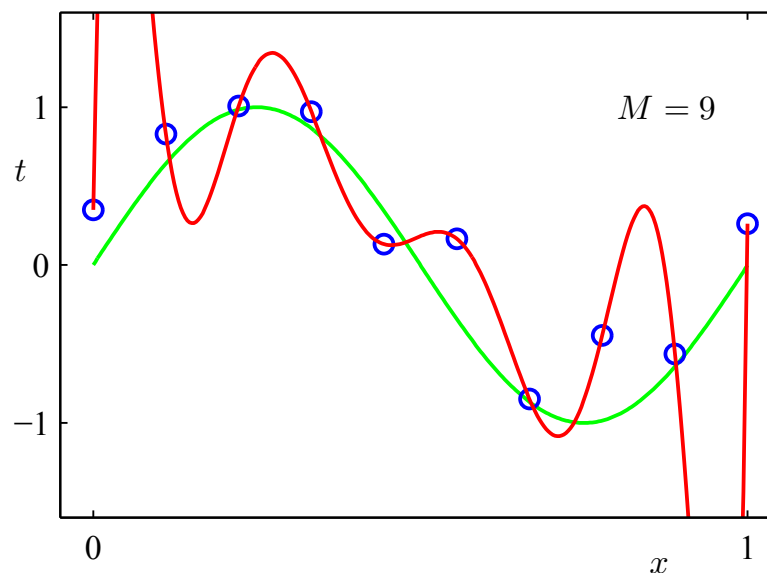


Adding high-order terms

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!

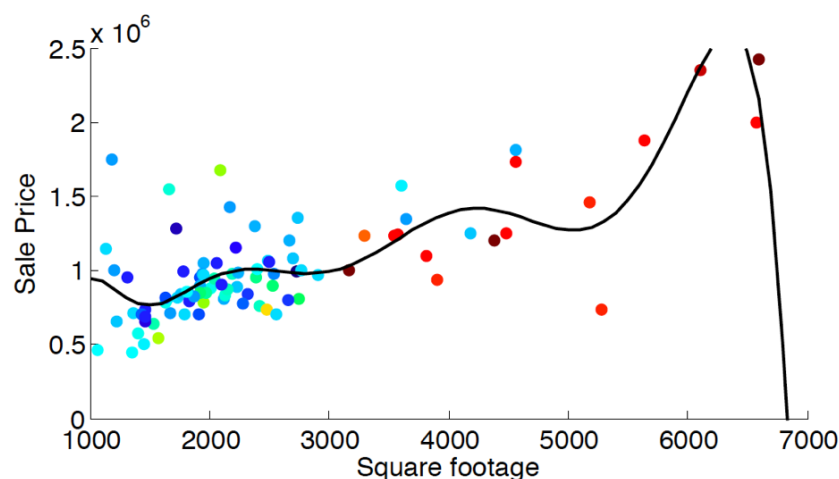
Overfitting

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 $M = 3$



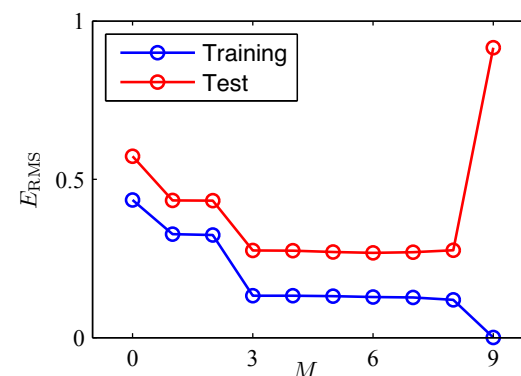
Note that the price would go to zero (or negative) if you buy bigger ones!
This is called poor generalization/overfitting.

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.

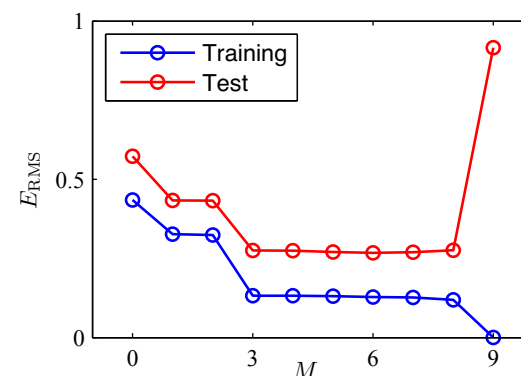
- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.



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.



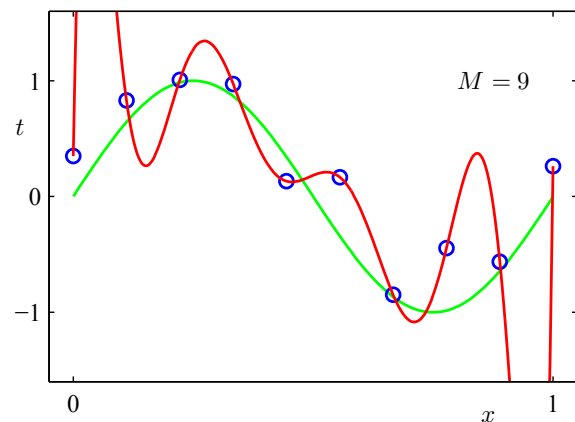
- Horizontal axis: *measure of model complexity*; in this example complexity defined by order of the polynomial basis functions.
- Vertical axis:
 - 1 For regression, residual sum of squares or residual mean squared (squared root of RSS)
 - 2 For classification, classification error rate.

Outline

- 1 Probabilistic interpretation
- 2 Multivariate solution
- 3 Nonlinear hypotheses
- 4 Basic ideas to overcome overfitting
 - Use more training data
 - Regularization methods
 - Regularized classification

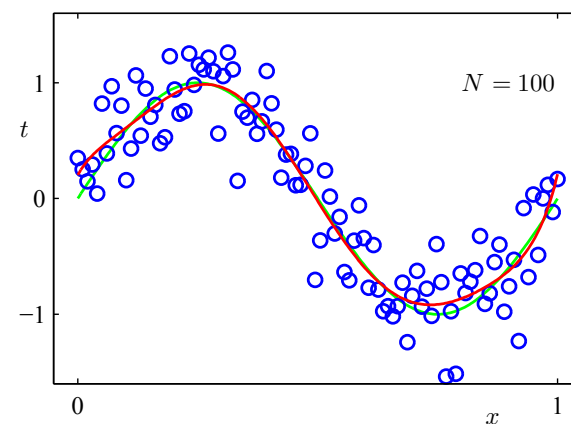
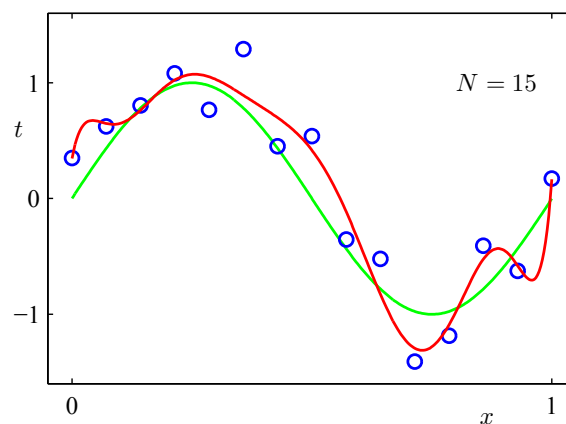
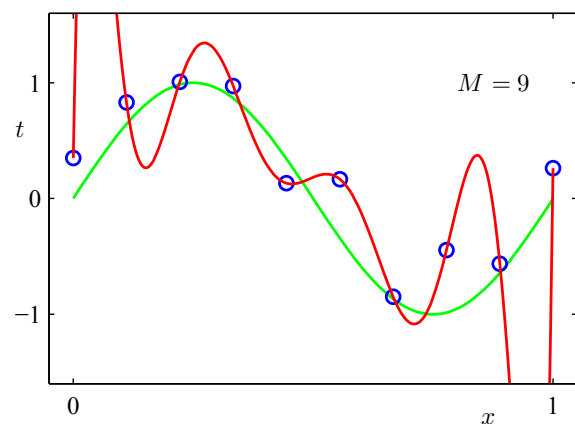
Use more training data to prevent over fitting

The more, the merrier



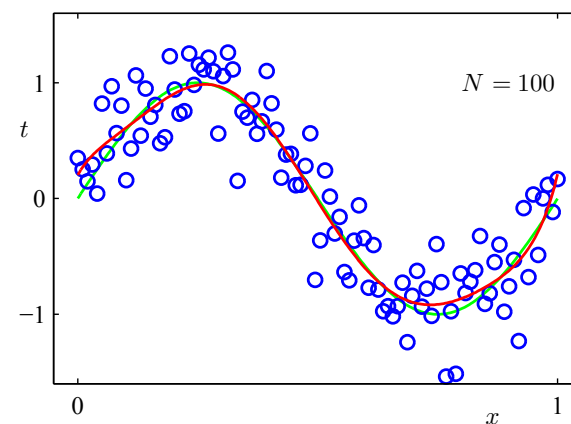
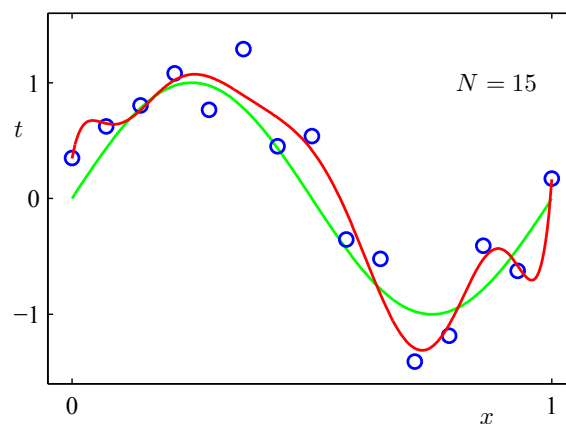
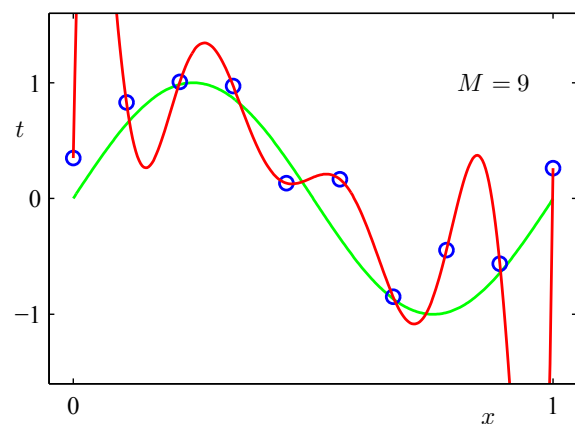
Use more training data to prevent over fitting

The more, the merrier



Use more training data to prevent over fitting

The more, the merrier



What if we do not have a lot of data?

Regularization methods

Intuition: For a linear model for regression

$$w^T x + b$$

we can try to identify ‘simpler’ models. But what does it mean for a model to *be simple*?

Regularization methods

Intuition: For a linear model for regression

$$w^T x + b$$

we can try to identify ‘simpler’ models. But what does it mean for a model to *be simple*?

Assumption (inductive bias)

A simpler model is one where most of the weights are zero.

A simpler model is one with smaller weights.

Why are smaller weights associated with simpler models?

Simpler functions are smoother, *i.e.*, nearby values of x have similar outputs \hat{y} .

Two values x and x' that differ in the first component by a small value ϵ .

Their predictions \hat{y} and \hat{y}' differ by ϵw_1 .

Smaller w_1 (closer to zero), more similar are the predictions.

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then

$$\hat{\mathbf{w}} \rightarrow \mathbf{0}$$

Regularized linear regression

A new cost function or error function to minimize

$$J(\mathbf{w}, b) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$. This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then

$$\hat{\mathbf{w}} \rightarrow \mathbf{0}$$

- If $\lambda \rightarrow 0$, then we trust our data more. Numerically,

$$\hat{\mathbf{w}} \rightarrow \arg \min \sum_n (\mathbf{w}^T \mathbf{x}_n + b - y_n)^2$$

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the OLS (Ordinary Least Squares) solution

$$\arg \min \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2 \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and reduces to the OLS solution when $\lambda = 0$, as expected.

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the OLS (Ordinary Least Squares) solution

$$\arg \min \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n - b)^2 + \lambda \|\mathbf{w}\|_2^2 \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and reduces to the OLS solution when $\lambda = 0$, as expected.

If we have to use numerical procedure, the gradient would change nominally too,

$$\nabla J(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{w})$$

As long as $\lambda \geq 0$, the optimization is convex.

Example: fitting data with polynomials

Our regression model

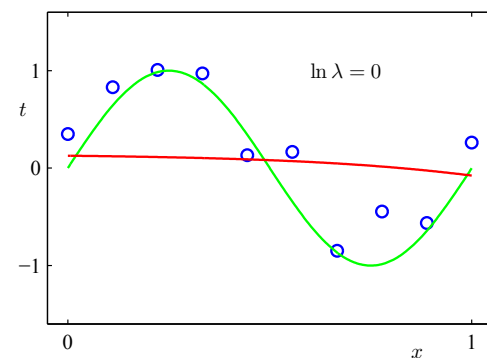
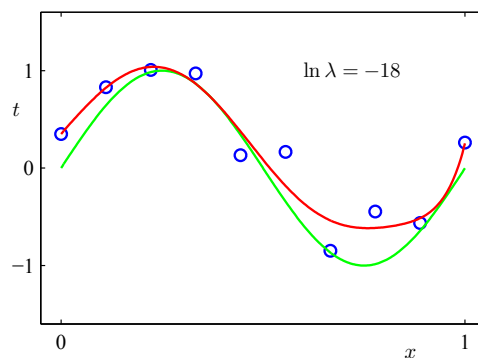
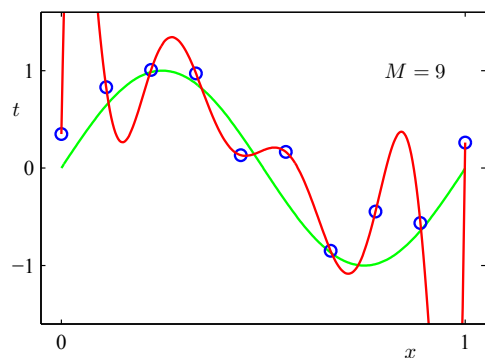
$$y = \sum_{m=1}^M w_m x^m$$

Regularization would discourage large parameter values as we saw with the OLS solution, thus potentially preventing overfitting.

| | $M = 0$ | $M = 1$ | $M = 3$ | $M = 9$ |
|-------|---------|---------|---------|-------------|
| 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 |

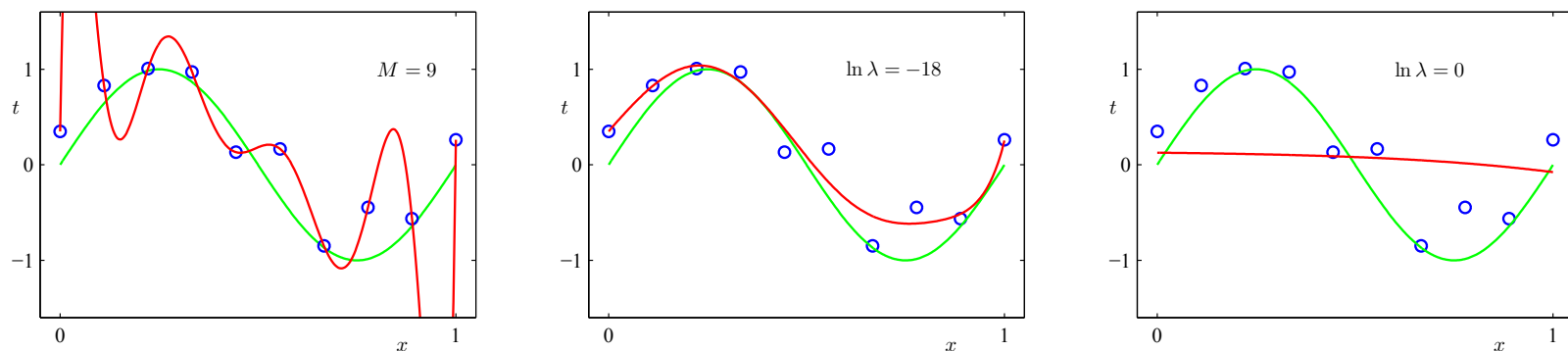
Overfitting in terms of λ

Overfitting is reduced from complex model to simpler one with the help of increasing regularizer

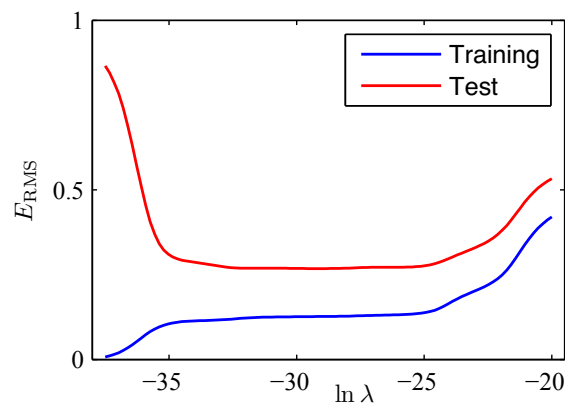


Overfitting in terms of λ

Overfitting is reduced from complex model to simpler one with the help of increasing regularizer



λ vs. residual error shows the difference of the model performance on training and testing dataset



The effect of λ

Large λ attenuates parameters towards 0

| | $\ln \lambda = -\infty$ | $\ln \lambda = -18$ | $\ln \lambda = 0$ |
|-------|-------------------------|---------------------|-------------------|
| w_0 | 0.35 | 0.35 | 0.13 |
| w_1 | 232.37 | 4.74 | -0.05 |
| w_2 | -5321.83 | -0.77 | -0.06 |
| w_3 | 48568.31 | -31.97 | -0.06 |
| w_4 | -231639.30 | -3.89 | -0.03 |
| w_5 | 640042.26 | 55.28 | -0.02 |
| w_6 | -1061800.52 | 41.32 | -0.01 |
| w_7 | 1042400.18 | -45.95 | -0.00 |
| w_8 | -557682.99 | -91.53 | 0.00 |
| w_9 | 125201.43 | 72.68 | 0.01 |

l_2 regularized logistic regression

Adding regularizer to the cost function for logistic regression

$$J(\mathbf{w}, b) = - \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

$$h_{\mathbf{w}, b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

l_2 regularized logistic regression

Adding regularizer to the cost function for logistic regression

$$J(\mathbf{w}, b) = - \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{regularization}}$$

$$h_{\mathbf{w}, b}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

Numerical optimization

- Objective functions remains convex as long as $\lambda \geq 0$.
- Gradients and Hessians are changed marginally and can be easily derived.

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

No: as this will set λ to zero, i.e., without regularization, defeating our intention to use it to control model complexity and to gain better generalization.

λ is thus a hyperparameter. To tune it,

- We can use a development/validation dataset independent of training and testing dataset.

The procedure is similar to choosing K in the nearest neighbor classifiers.

How to choose the right amount of regularization?

Can we tune λ on the training dataset?

No: as this will set λ to zero, i.e., without regularization, defeating our intention to use it to control model complexity and to gain better generalization.

λ is thus a **hyperparameter**. To tune it,

- We can use a development/validation dataset independent of training and testing dataset.

The procedure is similar to choosing K in the nearest neighbor classifiers.

For different λ , we get $\hat{\mathbf{w}}$ and evaluate the model on the development/validation dataset.

We then plot the curve λ versus prediction error (accuracy, classification error) and find the place that the performance on the validation is the best.

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.

Summary

- Regularization: preferring a simpler model by penalizing large weights or many non-zero weights.
- Helpful to reduce overfitting and to prevent numerical instability.
- Can be added to any linear model.
- Only a minor modification to the unregularized algorithms.
- New hyperparameter λ needs to be chosen using cross-validation.