

Overfitting, Bias / Variance Analysis

Professor Ameet Talwalkar

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Basic ideas to overcome overfitting
- 4 Bias/Variance Analysis

Announcements

- HW2 will be returned in section on Friday
- HW3 due in class next Monday
- Midterm is next Wednesday

Midterm

- Next Wednesday in class from 10am - 11:50am
- Completely closed-book (no notes allowed)
- Will include roughly 6 short answer questions and 3 long questions
 - ▶ Short questions should take 5 minutes on average
 - ▶ Long questions should take 15 minutes each
- Covers all material through (and including) today's lecture
 - ▶ Goal is to test conceptual understanding of the course material
 - ▶ Suggestion: carefully review lecture notes and problem sets
- Office hours / Section (see timing on course website)

Outline

- 1 Administration
- 2 Review of last lecture
 - Linear Regression
 - Ridge Regression for Numerical Purposes
 - Non-linear Basis
- 3 Basic ideas to overcome overfitting
- 4 Bias/Variance Analysis

Linear regression

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$
 - ▶ $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$: *weights, parameters, or parameter vector*
 - ▶ w_0 is called *bias*
 - ▶ We also sometimes call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ parameters too
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

Linear regression

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$
 - ▶ $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$: *weights, parameters, or parameter vector*
 - ▶ w_0 is called *bias*
 - ▶ We also sometimes call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ parameters too
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

Least Mean Squares (LMS) Objective: Minimize squared difference on training data (or residual sum of squares)

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2$$

Linear regression

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$
 - ▶ $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$: *weights, parameters, or parameter vector*
 - ▶ w_0 is called *bias*
 - ▶ We also sometimes call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ parameters too
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

Least Mean Squares (LMS) Objective: Minimize squared difference on training data (or residual sum of squares)

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2$$

1D Solution: Identify stationary points by taking derivative with respect to parameters and setting to zero, yielding ‘normal equations’

LMS when x is D-dimensional

$RSS(\tilde{w})$ in matrix form

$$RSS(\tilde{w}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2$$

LMS when \mathbf{x} is D-dimensional

$RSS(\tilde{\mathbf{w}})$ in matrix form

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

LMS when x is D-dimensional

$RSS(\tilde{w})$ in matrix form

$$RSS(\tilde{w}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{w}^T \tilde{x}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{x} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{w} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$RSS(\tilde{w}) = \sum_n (y_n - \tilde{w}^T \tilde{x}_n)(y_n - \tilde{x}_n^T \tilde{w})$$

LMS when \mathbf{x} is D-dimensional

$RSS(\tilde{\mathbf{w}})$ in matrix form

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$\begin{aligned} RSS(\tilde{\mathbf{w}}) &= \sum_n (y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)(y_n - \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}) \\ &= \sum_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}} - 2y_n \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}} + \text{const.} \end{aligned}$$

LMS when \mathbf{x} is D-dimensional

$RSS(\tilde{\mathbf{w}})$ in matrix form

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$\begin{aligned} RSS(\tilde{\mathbf{w}}) &= \sum_n (y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)(y_n - \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}) \\ &= \sum_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}} - 2y_n \tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}} + \text{const.} \\ &= \left\{ \tilde{\mathbf{w}}^T \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} \right\} + \text{const.} \end{aligned}$$

$RSS(\tilde{\mathbf{w}})$ in new notations

From previous slide

$$RSS(\tilde{\mathbf{w}}) = \left\{ \tilde{\mathbf{w}}^T \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} \right\} + \text{const.}$$

$RSS(\tilde{\mathbf{w}})$ in new notations

From previous slide

$$RSS(\tilde{\mathbf{w}}) = \left\{ \tilde{\mathbf{w}}^T \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} \right\} + \text{const.}$$

Design matrix and target vector

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

$RSS(\tilde{\mathbf{w}})$ in new notations

From previous slide

$$RSS(\tilde{\mathbf{w}}) = \left\{ \tilde{\mathbf{w}}^T \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^T \right) \tilde{\mathbf{w}} \right\} + \text{const.}$$

Design matrix and target vector

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

Compact expression

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

Solution in matrix form

Compact expression

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

Solution in matrix form

Compact expression

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

Gradients of Linear and Quadratic Functions

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

Solution in matrix form

Compact expression

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

Gradients of Linear and Quadratic Functions

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

Normal equation

$$\nabla_{\tilde{\mathbf{w}}} RSS(\tilde{\mathbf{w}}) \propto \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \tilde{\mathbf{X}}^T \mathbf{y} = 0$$

Solution in matrix form

Compact expression

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

Gradients of Linear and Quadratic Functions

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

Normal equation

$$\nabla_{\tilde{\mathbf{w}}} RSS(\tilde{\mathbf{w}}) \propto \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \tilde{\mathbf{X}}^T \mathbf{y} = 0$$

This leads to the least-mean-square (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

Practical concerns

Bottleneck of computing the LMS solution

$$\mathbf{w} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}} \mathbf{y}$$

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

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

Scalable methods

- Batch gradient descent
- Stochastic gradient descent

(Batch) Gradient Descent

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (e.g., randomly); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 Compute the gradient
$$\nabla RSS(\tilde{\mathbf{w}}) = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{X}}^T \mathbf{y}$$
 - 2 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \nabla RSS(\tilde{\mathbf{w}})$$
 - 3 $t \leftarrow t + 1$

What is the complexity of each iteration?

(Batch) Gradient Descent

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (e.g., randomly); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 Compute the gradient
$$\nabla RSS(\tilde{\mathbf{w}}) = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{X}}^T \mathbf{y}$$
 - 2 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \nabla RSS(\tilde{\mathbf{w}})$$
 - 3 $t \leftarrow t + 1$

What is the complexity of each iteration? $O(\text{ND})$

Why does this work?

(Batch) Gradient Descent

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (e.g., randomly); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 Compute the gradient
$$\nabla RSS(\tilde{\mathbf{w}}) = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{X}}^T \mathbf{y}$$
 - 2 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \nabla RSS(\tilde{\mathbf{w}})$$
 - 3 $t \leftarrow t + 1$

What is the complexity of each iteration? $O(\text{ND})$

Why does this work? $RSS(\tilde{\mathbf{w}})$ is convex (Hessian is PSD)

Stochastic gradient descent

Widrow-Hoff rule: update parameters using one example at a time

Stochastic gradient descent

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to some $\tilde{\mathbf{w}}^{(0)}$; set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_t
 - 2 Compute its contribution to the gradient

$$\mathbf{g}_t = (\tilde{\mathbf{x}}_t^T \tilde{\mathbf{w}}^{(t)} - y_t) \tilde{\mathbf{x}}_t$$

- 3 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_t$$
- 4 $t \leftarrow t + 1$

Stochastic gradient descent

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to some $\tilde{\mathbf{w}}^{(0)}$; set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_t
 - 2 Compute its contribution to the gradient

$$\mathbf{g}_t = (\tilde{\mathbf{x}}_t^T \tilde{\mathbf{w}}^{(t)} - y_t) \tilde{\mathbf{x}}_t$$

- 3 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_t$$
- 4 $t \leftarrow t + 1$

How does the complexity per iteration compare with gradient descent?

Stochastic gradient descent

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to some $\tilde{\mathbf{w}}^{(0)}$; set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_t
 - 2 Compute its contribution to the gradient

$$\mathbf{g}_t = (\tilde{\mathbf{x}}_t^T \tilde{\mathbf{w}}^{(t)} - y_t) \tilde{\mathbf{x}}_t$$

- 3 Update the parameters
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_t$$
- 4 $t \leftarrow t + 1$

How does the complexity per iteration compare with gradient descent?

- $O(\text{ND})$ for gradient descent versus $O(D)$ for SGD

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

Why might this happen?

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

Why might this happen?

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

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

Why might this happen?

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

Answer 2: Columns of \mathbf{X} are not linearly independent, e.g., some features are perfectly correlated. In this case, solution is not unique.

Ridge regression

What can we do when $\tilde{X}^T \tilde{X}$ is not invertible?

Ridge regression

What can we do when $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is not invertible?

- Add regularizer so that all singular values are at least $\lambda > 0$!
- This is equivalent to adding an extra term to $RSS(\tilde{\mathbf{w}})$

$$\overbrace{\frac{1}{2} \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\}}^{RSS(\tilde{\mathbf{w}})} + \underbrace{\frac{1}{2} \lambda \|\tilde{\mathbf{w}}\|_2^2}_{\text{regularization}}$$

Ridge regression

What can we do when $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is not invertible?

- Add regularizer so that all singular values are at least $\lambda > 0$!
- This is equivalent to adding an extra term to $RSS(\tilde{\mathbf{w}})$

$$\overbrace{\frac{1}{2} \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\}}^{RSS(\tilde{\mathbf{w}})} + \underbrace{\frac{1}{2} \lambda \|\tilde{\mathbf{w}}\|_2^2}_{\text{regularization}}$$

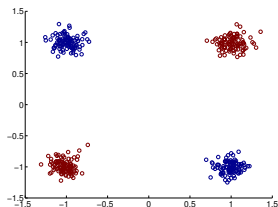
Solution

- Can derive normal equations as before
- Solution is of the form:

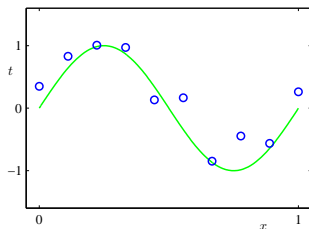
$$\tilde{\mathbf{w}} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

Is a linear modeling assumption always a good idea?

Example of nonlinear classification



Example of nonlinear regression



General nonlinear basis functions

We can use a nonlinear mapping

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

- M is dimensionality of new features \mathbf{z} (or $\phi(\mathbf{x})$)
- M could be greater than, less than, or equal to D

General nonlinear basis functions

We can use a nonlinear mapping

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

- M is dimensionality of new features \mathbf{z} (or $\phi(\mathbf{x})$)
- M could be greater than, less than, or equal to D

We can apply existing learning methods on the transformed data

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

Regression with nonlinear basis

Residual sum squares

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

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

The LMS solution can be formulated with the new design matrix

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

Example with regression

Polynomial basis functions

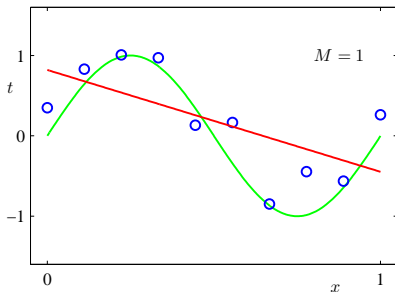
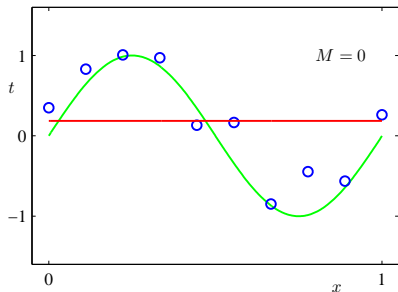
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Example with regression

Polynomial basis functions

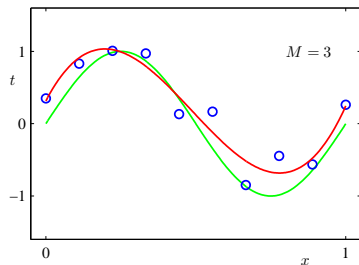
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

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



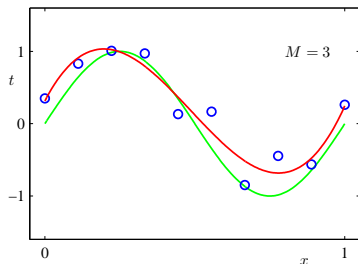
Adding high-order terms

$M=3$

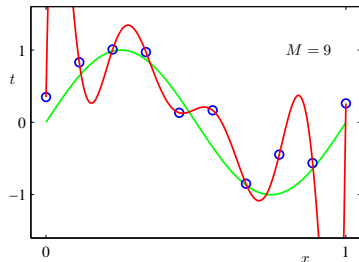


Adding high-order terms

$M=3$



$M=9$: *overfitting*



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

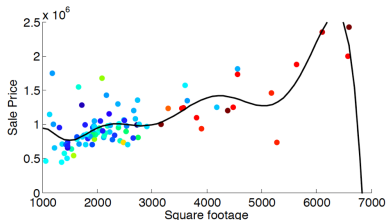
Overfitting

Parameters for higher-order polynomials are very large

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
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 can be quite disastrous

Fitting the housing price data with large M



Predicted price goes to zero (and is ultimately negative) if you buy a big enough house!

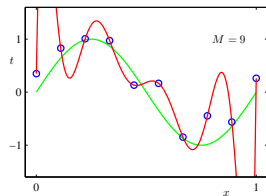
How might we prevent overfitting?

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Basic ideas to overcome overfitting**
 - Use more training data
 - Regularization methods
- 4 Bias/Variance Analysis

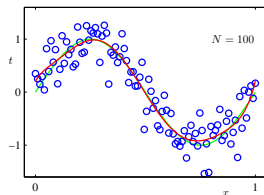
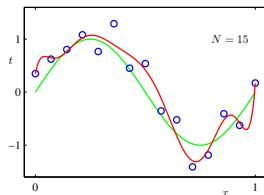
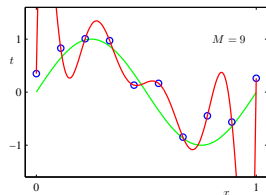
Use more training data to prevent over fitting

The more, the merrier



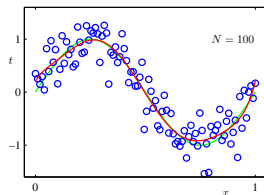
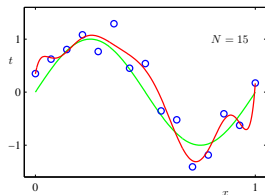
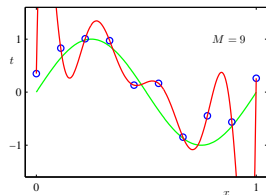
Use more training data to prevent over fitting

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: Give preference to 'simpler' models

- How do we define a simple linear regression model — $w^T x$?

Regularization methods

Intuition: Give preference to 'simpler' models

- How do we define a simple linear regression model — $w^T x$?

Our Strategy: Place a *prior* on our weights

- Interpret w as a random variable
- Assume that each w_d is centered around zero
- Use observed data \mathcal{D} to update our prior belief on w

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $\eta \sim N(0, \sigma_0^2)$ is a Gaussian random variable
 - ▶ Thus, $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $\eta \sim N(0, \sigma_0^2)$ is a Gaussian random variable
 - ▶ Thus, $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$
- We assume that \mathbf{w} is fixed (*Frequentist* interpretation)

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $\eta \sim N(0, \sigma_0^2)$ is a Gaussian random variable
 - ▶ Thus, $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$
- We assume that \mathbf{w} is fixed (*Frequentist* interpretation)
- The likelihood function maps parameters to probabilities

$$L : \mathbf{w}, \sigma_0^2 \mapsto p(\mathcal{D} | \mathbf{w}, \sigma_0^2)$$

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $\eta \sim N(0, \sigma_0^2)$ is a Gaussian random variable
 - ▶ Thus, $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$
- We assume that \mathbf{w} is fixed (*Frequentist* interpretation)
- The likelihood function maps parameters to probabilities

$$L : \mathbf{w}, \sigma_0^2 \mapsto p(\mathcal{D} | \mathbf{w}, \sigma_0^2) = p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_0^2) = \prod_n p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma_0^2)$$

Review: Probabilistic interpretation for LMS

- LMS model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $\eta \sim N(0, \sigma_0^2)$ is a Gaussian random variable
 - ▶ Thus, $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$
- We assume that \mathbf{w} is fixed (*Frequentist* interpretation)
- The likelihood function maps parameters to probabilities

$$L : \mathbf{w}, \sigma_0^2 \mapsto p(\mathcal{D} | \mathbf{w}, \sigma_0^2) = p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma_0^2) = \prod_n p(y_n | \mathbf{x}_n, \mathbf{w}, \sigma_0^2)$$

- Maximizing likelihood with respect to \mathbf{w} minimizes RSS and yields the LMS solution:

$$\mathbf{w}^{\text{LMS}} = \mathbf{w}^{\text{ML}} = \arg \max_{\mathbf{w}} L(\mathbf{w}, \sigma_0^2)$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) =$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

- What's the relationship between MAP and MLE?

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

- What's the relationship between MAP and MLE?
 - ▶ MAP reduces to MLE if we assume uniform prior for $p(\mathbf{w})$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$$

- What's the relationship between MAP and MLE?
 - ▶ MAP reduces to MLE if we assume uniform prior for $p(\mathbf{w})$
- Fully Bayesian treatment considers entire posterior, not just the mode

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Joint log likelihood

Plugging in Gaussian PDF, we get:

$$\begin{aligned} \log p(\mathcal{D}, \mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Joint log likelihood

Plugging in Gaussian PDF, we get:

$$\begin{aligned} \log p(\mathcal{D}, \mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

MAP estimate: $\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} \log p(\mathcal{D}, \mathbf{w})$

- As with LMS, set gradient equal to zero and solve (for \mathbf{w})

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

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

where $\lambda > 0$ is used to denote σ_0^2/σ^2 . This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

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

where $\lambda > 0$ is used to denote σ_0^2/σ^2 . This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

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

where $\lambda > 0$ is used to denote σ_0^2/σ^2 . This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then $\sigma_0^2 \gg \sigma^2$. That is, the variance of noise is far greater than what our prior model can allow for \mathbf{w} . In this case, our prior model on \mathbf{w} would be more accurate than what data can tell us. Thus, we are getting a *simple* model. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{0}$$

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

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

where $\lambda > 0$ is used to denote σ_0^2/σ^2 . This extra term $\|\mathbf{w}\|_2^2$ is called regularization/regularizer and controls the model complexity.

Intuitions

- If $\lambda \rightarrow +\infty$, then $\sigma_0^2 \gg \sigma^2$. That is, the variance of noise is far greater than what our prior model can allow for \mathbf{w} . In this case, our prior model on \mathbf{w} would be more accurate than what data can tell us. Thus, we are getting a *simple* model. Numerically,

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{0}$$

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

$$\mathbf{w}^{\text{MAP}} \rightarrow \mathbf{w}^{\text{LMS}} = \arg \min \sum_n (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the LMS solution

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

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

Closed-form solution

For regularized linear regression: the solution changes very little (in form) from the LMS solution

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

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

Gradients and Hessian change nominally too

$$\nabla \mathcal{E}(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{w}), \quad \mathbf{H} = 2(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$$

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

Example: fitting data with polynomials

Our regression model

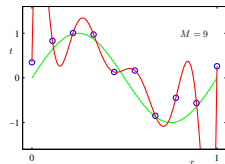
$$y = \sum_{m=1}^M w_m x^m$$

Regularization would discourage large parameter values as we saw with the LMS 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

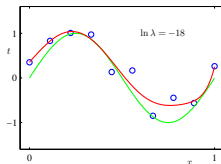
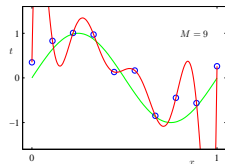
Overfitting in terms of λ

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



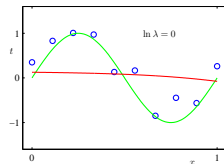
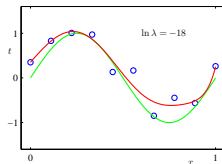
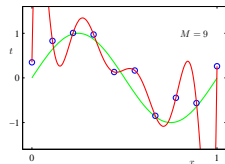
Overfitting in terms of λ

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



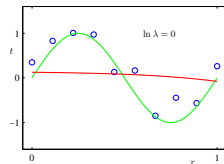
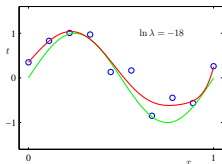
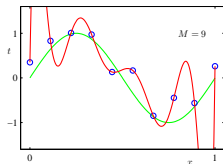
Overfitting in terms of λ

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

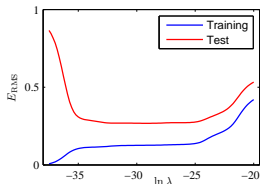


Overfitting in terms of λ

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



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

Regularized methods for classification

Adding regularizer to the cross-entropy functions used for binary and multinomial logistic regression

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\} + \lambda \|\mathbf{w}\|_2^2$$

$$\mathcal{E}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = - \sum_n \sum_k \log P(C_k | \mathbf{x}_n) + \lambda \sum_k \|\mathbf{w}_k\|_2^2$$

Regularized methods for classification

Adding regularizer to the cross-entropy functions used for binary and multinomial logistic regression

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\} + \lambda \|\mathbf{w}\|_2^2$$

$$\mathcal{E}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = - \sum_n \sum_k \log P(C_k | \mathbf{x}_n) + \lambda \sum_k \|\mathbf{w}_k\|_2^2$$

Numerical optimization

- Objective functions remain to be convex as long as $\lambda \geq 0$.
- Gradients and Hessians change 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 always set λ to zero, i.e., no regularization, defeating our intention of controlling model complexity

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

- We can use a validation set or do cross validation (as previously discussed)

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Basic ideas to overcome overfitting
- 4 Bias/Variance Analysis**

Basic and important machine learning concepts

Supervised learning

We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(h(\mathbf{x}), y)$$

Basic and important machine learning concepts

Supervised learning

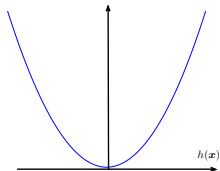
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(h(\mathbf{x}), y)$$

Example: quadratic loss function for regression when y is continuous

$$\ell(h(\mathbf{x}), y) = [h(\mathbf{x}) - y]^2$$

Ex: when $y = 0$

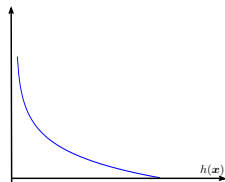


Other types of loss functions

For classification: cross-entropy loss (also called *logistic* loss)

$$\ell(h(\mathbf{x}), y) = -y \log h(\mathbf{x}) - (1-y) \log[1-h(\mathbf{x})]$$

Ex: when $y = 1$



Measure how good our predictor is

Risk: Given the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} d y$$

Measure how good our predictor is

Risk: Given the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

However, we cannot compute $R[h(\mathbf{x})]$, so we use *empirical risk*, given a training dataset \mathcal{D}

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Measure how good our predictor is

Risk: Given the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} d y$$

However, we cannot compute $R[h(\mathbf{x})]$, so we use *empirical risk*, given a training dataset \mathcal{D}

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as $N \rightarrow +\infty$,

$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow R[h(\mathbf{x})]$$

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

ERM might be problematic

How this relates to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

ERM might be problematic

- If $h(\mathbf{x})$ is complicated enough,

$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow 0$$

- But then $h(\mathbf{x})$ is unlikely to do well in predicting things out of the training dataset \mathcal{D}
- This is called *poor generalization* or *overfitting*. We have just discussed approaches to address this issue.
- We'll explore why regularization might work from the context of the bias-variance tradeoff, focusing on regression / squared loss

Bias/variance tradeoff (Looking ahead)

Error decomposes into 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

We will prove this result, and interpret what it means...