

Overfitting, Bias / Variance Analysis

Professor Ameet Talwalkar

Slide Credit: Professor Fei Sha

Outline

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

Announcements

- HW2 (and remaining HW1) will be returned today after class (see Nikos and Amogh)
- HW3 and HW4 due on Thursday
- No new HW assigned on Thursday (to give you time to work on project proposal)
- Project proposal guideline posted on course website last week

CS260: Project Proposal

Professor Ameet Talwalkar

Due: 11/5/15

Project proposals should be 1-2 pages long and contain the following information:

- **Motivation:** What high-level problem are you studying? Why is it important?
- **Background:** What previous work exists on your topic? How does your proposed project differ from these existing works? Note: If you have prior experience working on this problem, please clearly describe how your proposed work for this class project differs from your prior work.
- **Proposed Work:** What do you plan to do? What methods will you use and why will you use them?
- **Timeline:** What is your timeline for performing this work?
- **Deliverables / Evaluation:** What is the expected outcome of your work? How will you evaluate the quality of your work?
- **Data:** If your project involves data, what data are you using? How did you obtain this data? Why is this data interesting / relevant to your problem? What computing resources will you be using to analyze this data?
- **Software Tools / Libraries:** What (if any) software will be used for this project? Please describe any third-party software libraries you plan to use.
- **Team:** Who is working on this project (at most two students can work on a project)?
- **Prior Discussion:** You are required to discuss this project with Professor Talwalkar or one of the TAs prior to submitting your proposal. Please list who you spoke with and when. Note: You are also encouraged to speak with us as you continue working on the project if you have questions.

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 1D

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$

We also sometimes call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ parameters too!

Linear regression 1D

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Training data: $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$

We also sometimes call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ parameters too!

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 \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

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

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

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

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

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

Stochastic gradient descent

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

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); 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
$$\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_t$$

Stochastic gradient descent

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

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); 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 $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); 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 $O(\text{D})$ for SGD

What if $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ 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.

Ridge regression

For $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ that is not invertible

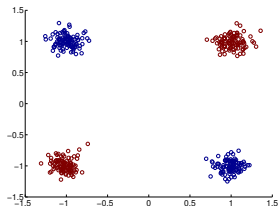
$$\tilde{\mathbf{w}} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

This is equivalent to adding an extra term to $RSS(\tilde{\mathbf{w}})$

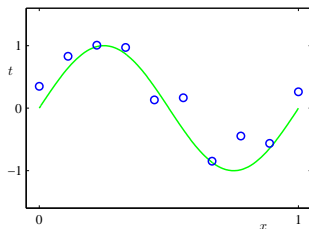
$$\overbrace{\frac{1}{2} \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\}}^{RSS(\tilde{\mathbf{w}})} + \underbrace{\frac{1}{2} \lambda \|\tilde{\mathbf{w}}\|_2^2}_{\text{regularization}}$$

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

Example of nonlinear classification



Example of nonlinear regression



General nonlinear basis functions

We can use a nonlinear mapping

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

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 $\mathbf{w}^T \phi(\mathbf{x})$
- 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

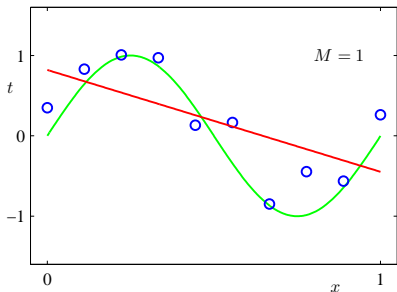
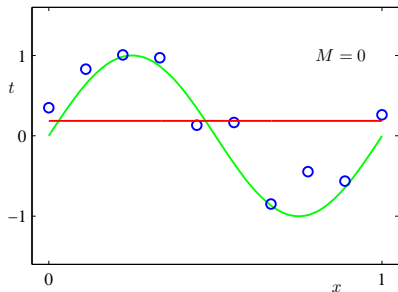
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Example with regression

Polynomial basis functions

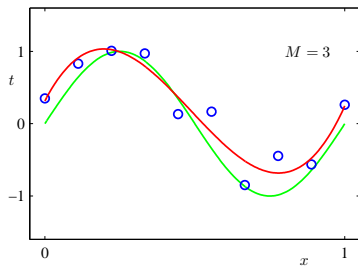
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

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



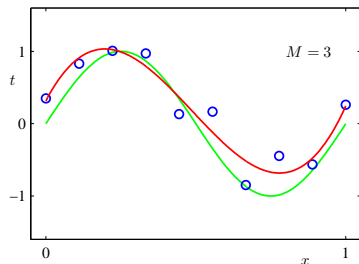
Adding high-order terms

$M=3$

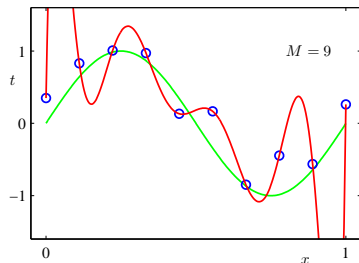


Adding high-order terms

$M=3$



$M=9$: *overfitting*



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

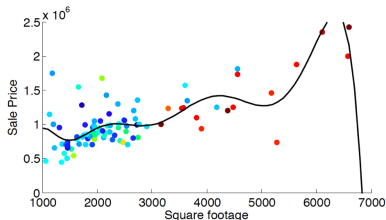
Overfitting

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 $M = 3$



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

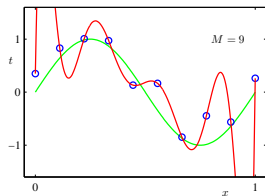
How might we prevent overfitting?

Outline

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

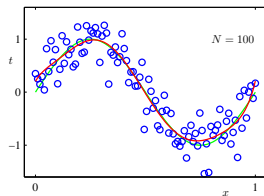
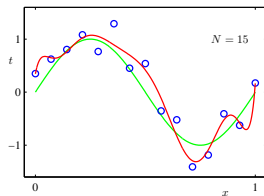
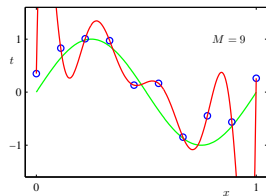
Use more training data to prevent over fitting

The more, the merrier



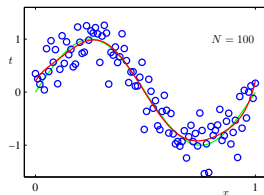
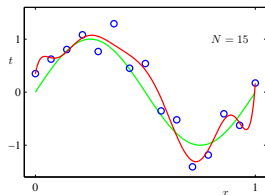
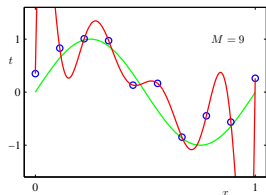
Use more training data to prevent over fitting

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: For a linear model for regression

$$\mathbf{w}^T \mathbf{x}$$

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

$$\mathbf{w}^T \mathbf{x}$$

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

Assumption We can place a *prior* on our weights, assuming that w_d is centered around zero.

With this reasoning, we will interpret w as a random variable and we will use the observed data \mathcal{D} to update our prior belief on w

Review: Probabilistic interpretation for LMS

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)
- We define $p(y|\mathbf{x}, \mathbf{w}, \sigma_0^2)$ as the sampling distribution given fixed values for the parameters \mathbf{w}, σ_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)
- We define $p(y|\mathbf{x}, \mathbf{w}, \sigma_0^2)$ as the sampling distribution given fixed values for the parameters \mathbf{w}, σ_0^2
- The likelihood function maps parameters to probabilities

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

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)
- We define $p(y|\mathbf{x}, \mathbf{w}, \sigma_0^2)$ as the sampling distribution given fixed values for the parameters \mathbf{w}, σ_0^2
- The likelihood function maps parameters to probabilities

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

- 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

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 with a prior distribution (*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 with a prior distribution (*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 with a prior distribution (*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:

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

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 with a prior distribution (*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:

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

- 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 with a prior distribution (*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:

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

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

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

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

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

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

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

If we have to use numerical procedure, the gradients and the Hessian matrix would 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

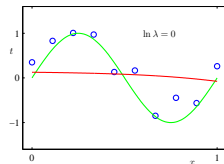
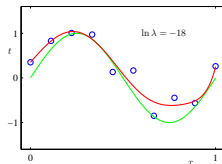
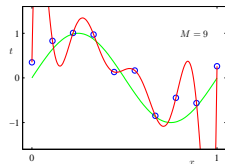
$$y = \sum_{m=1}^M w_m x^m$$

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

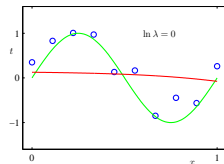
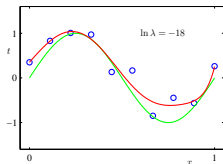
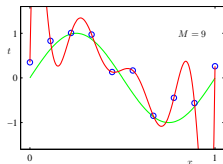
Overfitting in terms of λ

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

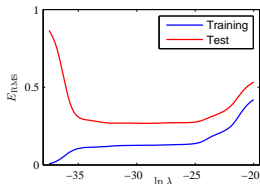


Overfitting in terms of λ

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



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 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/holdout dataset independent of training and testing dataset.
- We can do leave-one-out (LOO)

The procedure is similar to choose 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/holdout dataset independent of training and testing dataset.
- We can do leave-one-out (LOO)

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

For different λ , we get w^{MAP} and evaluate the model on the development/holdout dataset (or, the samples being left in LOO).

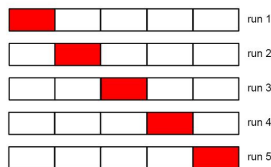
We then plot the curve λ versus prediction error (accuracy, classification error) and find the place that the performance on the holdout/LOO is the best.

Use cross-validation to choose λ

Procedure

- Randomly partition training data into K *disjoint* parts
Normally, K is chosen to be 10, 5, etc.
- For each possible value of λ
 - 1 Use one part as holdout; use other $(K - 1)$ parts as training
 - 2 Evaluate the model on the holdout
 - 3 Do this K times, and average the performance on the holdouts
- Choose the λ with the best performance

When $K = N$ (the number of training examples), this becomes LOO.



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

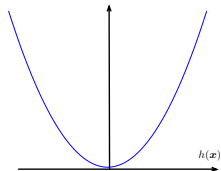
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

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

Example: quadratic loss function for regression when y is continuous

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

Ex: when $y = 0$

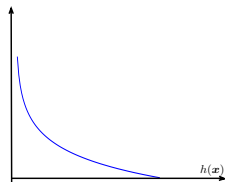


Other types of loss functions

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: assume we know 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: assume we know 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: assume we know 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}

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

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

- 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.
- Let's try to understand 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...