

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Jan 27, 2021

1 / 28

Logistics

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Linear regression with nonlinear basis
- 4 Overfitting and preventing overfitting

3 / 28

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Linear regression with nonlinear basis
- 4 Overfitting and preventing overfitting

2 / 28

Logistics

Logistics

- The schedule of lectures is available at <https://sirisharambhatla.com/CSCI567/index.html>

4 / 28

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Linear regression with nonlinear basis
- 4 Overfitting and preventing overfitting

5 / 28

Regression

Predicting a continuous outcome variable using past observations

- temperature, amount of rainfall, house price, etc.

Key difference from classification

- continuous vs discrete
- measure *prediction errors* differently.
- lead to quite different learning algorithms.

Linear Regression: regression with linear models: $f(x) = w^T x$

6 / 28

Least square solution

$$\begin{aligned}
 w^* &= \underset{w}{\operatorname{argmin}} \operatorname{RSS}(w) \\
 &= \underset{w}{\operatorname{argmin}} \|Xw - y\|_2^2 \\
 &= (X^T X)^{-1} X^T y
 \end{aligned}$$

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Two approaches to find the minimum:

- find **stationary points** by setting gradient = 0
- “**complete the square**”

7 / 28

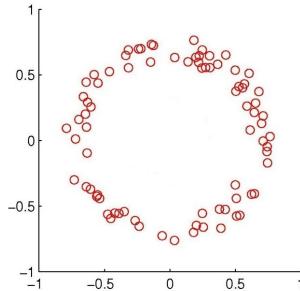
Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Linear regression with nonlinear basis
- 4 Overfitting and preventing overfitting

8 / 28

What if linear model is not a good fit?

Example: a straight line is a bad fit for the following data



9 / 28

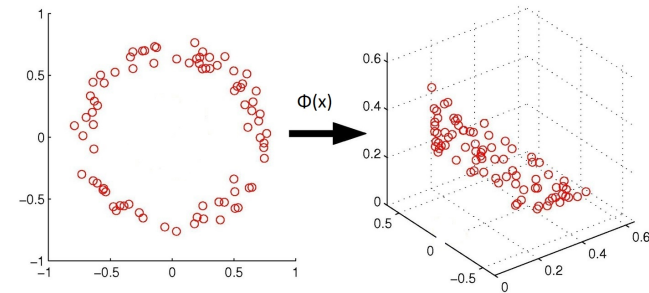
Solution: nonlinearly transformed features

1. Use a nonlinear mapping

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

to transform the data to a more complicated feature space

2. Then apply linear regression (hope: linear model is a better fit for the new feature space).



10 / 28

Regression with nonlinear basis

Model: $f(x) = w^T \phi(x)$ where $w \in \mathbb{R}^M$

Objective:

$$\text{RSS}(w) = \sum_n (w^T \phi(x_n) - y_n)^2$$

Similar least square solution:

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y \quad \text{where} \quad \Phi = \begin{pmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

11 / 28

Example

Polynomial basis functions for $D = 1$

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

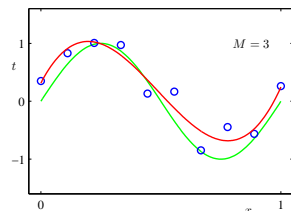
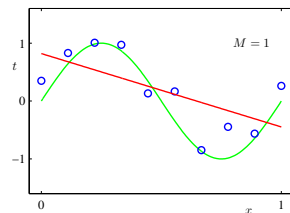
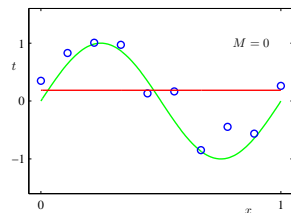
Learning a linear model in the new space

= learning an *M-degree polynomial model* in the original space

12 / 28

Example

Fitting a noisy sine function with a polynomial ($M = 0, 1, \text{ or } 3$):



13 / 28

Why nonlinear?

Can I use a fancy **linear feature map**?

$$\phi(x) = \begin{bmatrix} x_1 - x_2 \\ 3x_4 - x_3 \\ 2x_1 + x_4 + x_5 \\ \vdots \end{bmatrix} = \mathbf{A}x \quad \text{for some } \mathbf{A} \in \mathbb{R}^{M \times D}$$

No, it basically *does nothing* since

$$\min_{w \in \mathbb{R}^M} \sum_n (w^T \mathbf{A}x_n - y_n)^2 = \min_{w' \in \text{Im}(\mathbf{A}^T) \subset \mathbb{R}^D} \sum_n (w'^T x_n - y_n)^2$$

We will see more nonlinear mappings soon.

14 / 28

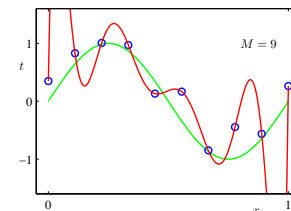
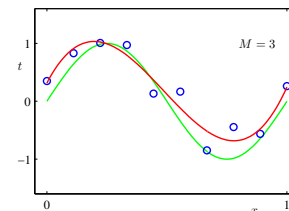
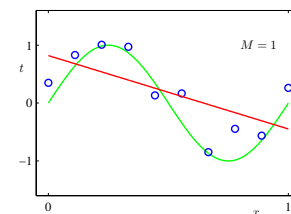
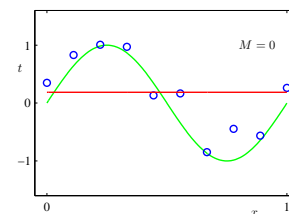
Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Linear regression with nonlinear basis
- 4 Overfitting and preventing overfitting

15 / 28

Should we use a very complicated mapping?

Ex: fitting a noisy sine function with a polynomial:



16 / 28

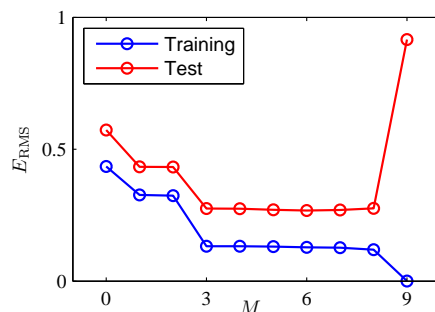
Underfitting and Overfitting

$M \leq 2$ is **underfitting** the data

- large training error
- large test error

$M \geq 9$ is **overfitting** the data

- small training error
- **large test error**



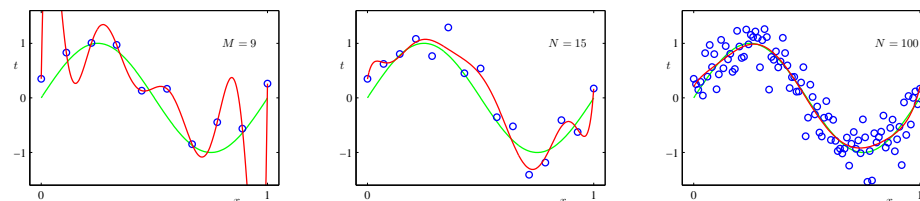
More complicated models \Rightarrow larger gap between training and test error

How to prevent overfitting?

17 / 28

Method 1: use more training data

The more, the merrier



More data \Rightarrow smaller gap between training and test error

18 / 28

Method 2: control the model complexity

For polynomial basis, the **degree M** clearly controls the complexity

- use cross-validation to pick hyper-parameter M

When M or in general Φ is fixed, are there still other ways to control complexity?

19 / 28

Magnitude of weights

Least square solution for the polynomial example:

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
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

Intuitively, **large weights \Rightarrow more complex model**

20 / 28

How to make w small?

Regularized linear regression: new objective

$$\mathcal{E}(w) = \text{RSS}(w) + \lambda R(w)$$

Goal: find $w^* = \text{argmin}_w \mathcal{E}(w)$

- $R: \mathbb{R}^D \rightarrow \mathbb{R}^+$ is the **regularizer**
 - measure how complex the model w is
 - common choices: $\|w\|_2^2$, $\|w\|_1$, etc.
- $\lambda > 0$ is the **regularization coefficient**
 - $\lambda = 0$, no regularization
 - $\lambda \rightarrow +\infty$, $w \rightarrow \text{argmin}_w R(w)$
 - i.e. control **trade-off** between training error and complexity

21 / 28

The effect of λ

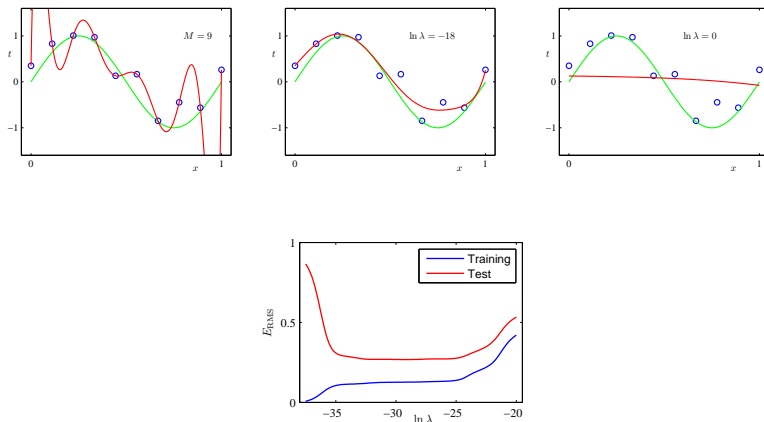
when we increase regularization coefficient λ

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
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

22 / 28

The trade-off

when we increase regularization coefficient λ



23 / 28

How to solve the new objective?

Simple for $R(w) = \|w\|_2^2$:

$$\mathcal{E}(w) = \text{RSS}(w) + \lambda \|w\|_2^2 = \|\Phi w - y\|_2^2 + \lambda \|w\|_2^2$$

$$\nabla \mathcal{E}(w) = 2(\Phi^T \Phi w - \Phi^T y) + 2\lambda w = 0$$

$$\Rightarrow (\Phi^T \Phi + \lambda I) w = \Phi^T y$$

$$\Rightarrow w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

Note the same form as in the fix when $X^T X$ is not invertible!

For other regularizers, as long as it's **convex**, standard optimization algorithms can be applied.

24 / 28

Equivalent form

Regularization is also sometimes formulated as

$$\underset{w}{\operatorname{argmin}} \operatorname{RSS}(w) \quad \text{subject to } R(w) \leq \beta$$

where β is some hyper-parameter.

Finding the solution becomes a *constrained optimization problem*.

Choosing either λ or β can be done by cross-validation.

25 / 28

Summary

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

Important to understand the derivation than remembering the formula

Overfitting: small training error but large test error

Preventing Overfitting: more data + regularization

26 / 28

Recall the question

Typical steps of developing a machine learning system:

- Collect data, split into training, development, and test sets.
- *Train a model with a machine learning algorithm.* Most often we apply cross-validation to tune hyper-parameters.
- Evaluate using the test data and report performance.
- Use the model to predict future/make decisions.

How to do the *red part* exactly?

27 / 28

General idea to derive ML algorithms

1. Pick a set of **models** \mathcal{F}
 - e.g. $\mathcal{F} = \{f(x) = w^T x \mid w \in \mathbb{R}^D\}$
 - e.g. $\mathcal{F} = \{f(x) = w^T \Phi(x) \mid w \in \mathbb{R}^M\}$
2. Define **error/loss** $L(y', y)$
3. Find **empirical risk minimizer (ERM)**:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{n=1}^N L(f(x_n), y_n)$$

or **regularized empirical risk minimizer**:

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{n=1}^N L(f(x_n), y_n) + \lambda R(f)$$

ML becomes optimization

28 / 28