

CSCI-567: Machine Learning (Fall 2019)

Prof. Victor Adamchik

U of Southern California

Sep. 10, 2019

September 10, 2019 1 / 54

Outline

- 1 Linear regression
 - Classification and Regression
 - Motivation
 - Setup and Algorithm
 - Discussions
- 2 Linear regression with nonlinear basis
- 3 Overfitting and Preventing Overfitting

September 10, 2019 3 / 54

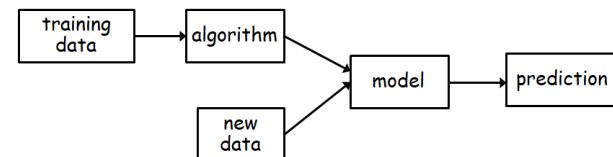
Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and Preventing Overfitting

September 10, 2019 2 / 54

Predictive modeling

Predictive modeling (i.e. supervised learning) is a process of creating a model using data to make a prediction on new data.



Predictive modeling is a problem of finding a mapping function f from training data ($x \in \mathbb{R}^D$) to output variables.

There are important differences between classification and regression problems.

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

September 10, 2019 4 / 54

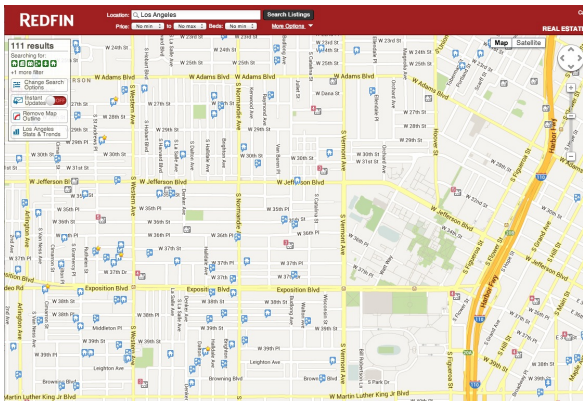
Classification

Classification is a problem of finding a mapping function f from training data $(x \in \mathbb{R}^D)$ to *discrete* output variables $(y \in C)$.

- The output variables are called labels or classes or categories.
- The mapping function predicts the class for a given observation.
- The classification *accuracy* is computed as the percentage of correctly classified examples out of all examples.

Ex: Predicting the sale price of a house

Retrieve historical sales records (training data)



Regression

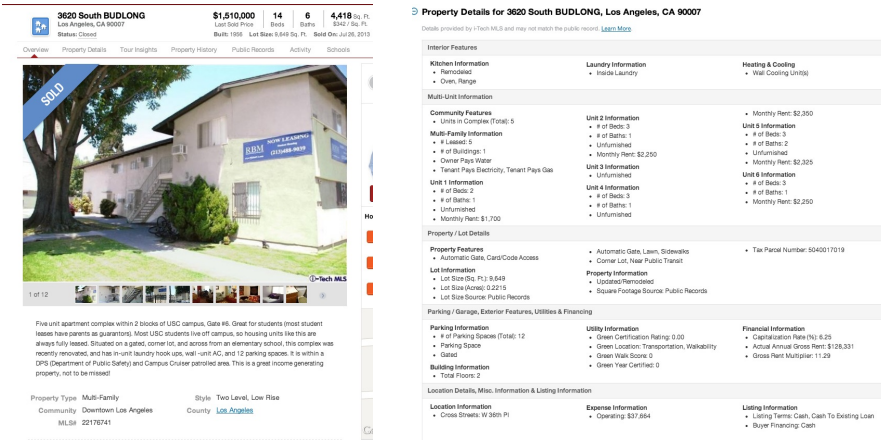
Regression is a problem of finding a mapping function f from training data $(x \in \mathbb{R}^D)$ to a *continuous* output variable $(y \in \mathbb{R})$.

- The output variable is a continuous quantity; pricing optimization, sales forecasting, rating forecasting are some examples.
- Regression predictions can be evaluated using the *mean squared error*.

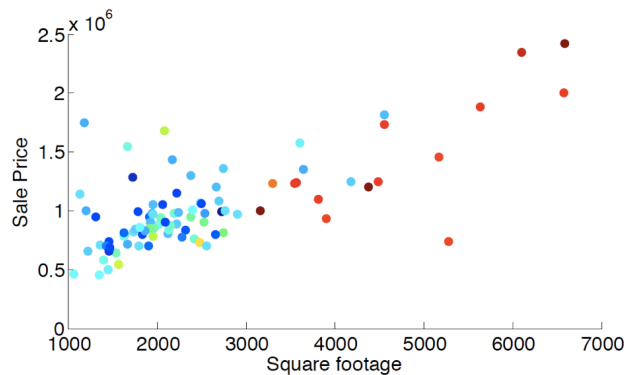
In some cases, a classification problem can be converted to a regression problem. Some algorithms do this by predicting a probability for each class.

Linear Regression: regression with linear models.

Features used to predict



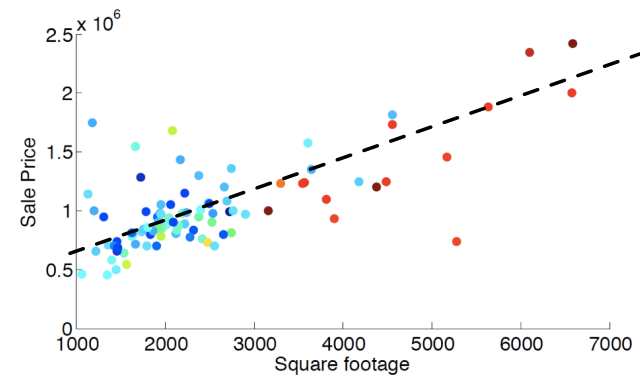
Correlation between square footage and sale price



In linear regression, the goal is to predict y from x using a linear function.

Possibly linear relationship

$$\text{Sale price} \approx \underset{\text{(slope)}}{\text{price_per_sqft}} \times \text{square_footage} + \underset{\text{(intercept)}}{\text{fixed_expense}}$$



How to learn the unknown parameters?

How to measure error for one prediction?

- The classification error (0-1 loss, i.e. *right* or *wrong*) is *inappropriate* for continuous outcomes.
- We can look at
 - ▶ *absolute* error: $|\text{prediction} - \text{sale price}|$
 - ▶ or *squared* error: $(\text{prediction} - \text{sale price})^2$ (**most common**)

Goal: pick the model (unknown parameters) that minimizes the average/total prediction error, but *on what set*?

- test set, ideal but we *cannot use test set while training*
- training set? (minimize the training error)

Example

$$\text{Predicted price} = \text{price_per_sqft} \times \text{square_footage} + \text{fixed_expense}$$

one model: $\text{price_per_sqft} = 0.3\text{K}$, $\text{fixed_expense} = 210\text{K}$

sqft	sale price (K)	prediction (K)	squared error
2000	810	810	0
2100	907	840	67^2
1100	312	540	228^2
5500	2,600	1,860	740^2
...
Total			$0 + 67^2 + 228^2 + 740^2 + \dots$

Adjust **price_per_sqft** and **fixed_expense** such that the total squared error is minimized.

Formal setup for linear regression

Input: $\mathbf{x} \in \mathbb{R}^D$ (features, covariates, context, predictors, etc)

Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc)

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

Here x_{nd} represents the d th dimension of the n th sample \mathbf{x}_n

Linear model: $f: \mathbb{R}^D \rightarrow \mathbb{R}$, with $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d$

Linear regression has been around for more than 200 years...

September 10, 2019 13 / 54

Goal

Minimize total squared error

- **Residual Sum of Squares** (RSS), a function of \mathbf{w}

$$\text{RSS}(\mathbf{w}) = \sum_{n=1}^N (f(\mathbf{x}_n) - y_n)^2 = \sum_{n=1}^N (\mathbf{x}_n^T \mathbf{w} - y_n)^2$$

- find $\mathbf{w}^* = \underset{\mathbf{w} \in \mathbb{R}^{D+1}}{\text{argmin}} \text{RSS}(\mathbf{w})$
- minimize the Euclidean distance, or find the **least squares solution**
- *reduce machine learning to optimization*

September 10, 2019 15 / 54

Notation Convenience

NOTE: for notation convenience, we will

- append 1 to each \mathbf{x}_n as the first feature: $\mathbf{x} = [1, x_1, x_2, \dots, x_D]^T$
- append w_0 to weights: $\mathbf{w} = [w_0, w_1, w_2, \dots, w_D]^T$

The model becomes

$$f: \mathbb{R}^{D+1} \rightarrow \mathbb{R}$$

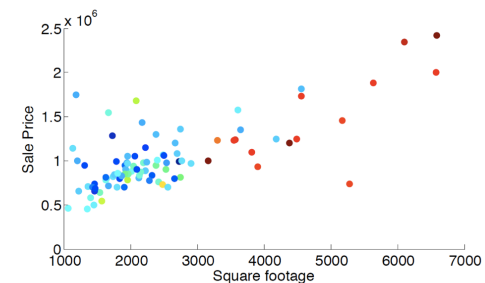
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$$

So please pay attention to notations!

September 10, 2019 14 / 54

Warm-up: $D = 0$

Only one parameter w_0 : constant prediction $f(\mathbf{x}) = w_0$



f is a horizontal line, where should it be?

Use calculus to find the value of w_0 that minimizes the RSS

September 10, 2019 16 / 54

Warm-up: $D = 0$

Optimization objective becomes

$$\begin{aligned}\text{RSS}(w_0) &= \sum_{n=1}^N (w_0 - y_n)^2 \\ \frac{\partial \text{RSS}(w_0)}{\partial w_0} &= 2 \sum_{n=1}^N (w_0 - y_n) = 0 \\ N w_0 - \sum_{n=1}^N y_n &= 0\end{aligned}$$

It follows that $w_0 = \frac{1}{N} \sum_n y_n$, i.e. the **average**

Exercise: what if we use absolute error instead of squared error?

September 10, 2019 17 / 54

Warm-up: $D = 1$

Optimization objective becomes

$$\text{RSS}(\mathbf{w}) = \sum_n (w_0 + w_1 x_n - y_n)^2$$

General approach: find **stationary points**, i.e., points with **zero gradient**

$$\begin{aligned}\begin{cases} \frac{\partial \text{RSS}(\mathbf{w})}{\partial w_0} = 0 \\ \frac{\partial \text{RSS}(\mathbf{w})}{\partial w_1} = 0 \end{cases} &\Rightarrow \begin{cases} \sum_n (w_0 + w_1 x_n - y_n) = 0 \\ \sum_n (w_0 + w_1 x_n - y_n) x_n = 0 \end{cases} \\ \Rightarrow \begin{cases} N w_0 + w_1 \sum_n x_n = \sum_n y_n \\ w_0 \sum_n x_n + w_1 \sum_n x_n^2 = \sum_n y_n x_n \end{cases} &\quad \text{(a linear system)} \\ \Rightarrow \begin{pmatrix} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} &= \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}\end{aligned}$$

September 10, 2019 18 / 54

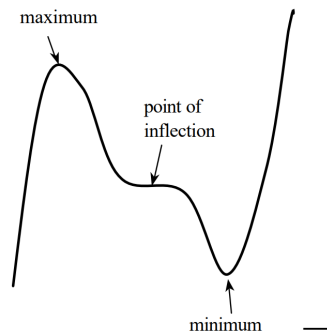
Least square solution for $D = 1$

Assuming the matrix is invertible:

$$\Rightarrow \begin{pmatrix} w_0^* \\ w_1^* \end{pmatrix} = \begin{pmatrix} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}$$

Are stationary points minimizers?

- not true in general
- yes for **convex** objectives



September 10, 2019 19 / 54

General least square solution

Objective

$$\text{RSS}(\mathbf{w}) = \sum_n (\mathbf{x}_n^T \mathbf{w} - y_n)^2$$

Again, find stationary points (**multivariate calculus**)

$$\begin{aligned}\frac{1}{2} \nabla \text{RSS}(\mathbf{w}) &= \sum_n \mathbf{x}_n (\mathbf{x}_n^T \mathbf{w} - y_n) = \left(\sum_n \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} - \sum_n \mathbf{x}_n y_n \\ &= (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}\end{aligned}$$

where

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

September 10, 2019 20 / 54

General least square solution

$$(\mathbf{X}^T \mathbf{X})\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0} \Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

assuming $\mathbf{X}^T \mathbf{X}$ (called a **covariance matrix**) is invertible for now.

Again by convexity \mathbf{w}^* is the minimizer of RSS.

Verify the solution when $D = 1$:

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \cdots & \cdots \\ 1 & x_N \end{pmatrix} = \begin{pmatrix} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}$$

when $D = 0$: $(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{N}$, $\mathbf{X}^T \mathbf{y} = \sum_n y_n$

September 10, 2019 21 / 54

Multivariate Calculus

RSS is given by

$$\text{RSS}(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$

$$\begin{aligned} \nabla \text{RSS}(\mathbf{w}) &= \nabla (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) \\ &= 0 - 2\mathbf{X}^T \mathbf{y} + \nabla (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) \\ &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} \quad (\text{prove it !}) \end{aligned}$$

It follows

$$(\mathbf{X}^T \mathbf{X})\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0} \Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

September 10, 2019 23 / 54

Another approach

RSS is the Euclidean norm squared:

$$\begin{aligned} \text{RSS}(\mathbf{w}) &= \sum_n (\mathbf{w}^T \mathbf{x}_n - y_n)^2 = \|\mathbf{X} \mathbf{w} - \mathbf{y}\|_2^2 \\ &= (\mathbf{X} \mathbf{w} - \mathbf{y})^T (\mathbf{X} \mathbf{w} - \mathbf{y}) \\ &= (\mathbf{w}^T \mathbf{X}^T - \mathbf{y}^T) (\mathbf{X} \mathbf{w} - \mathbf{y}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \end{aligned}$$

Note: $\mathbf{y}^T \mathbf{X} \mathbf{w} = (\mathbf{w}^T \mathbf{X}^T \mathbf{y})^T$

September 10, 2019 22 / 54

Multivariate Calculus

Is it a minimizer?

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

We will use a second derivative (the **Hessian** matrix)

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

A symmetric matrix M is said to be a **positive semi-definite** (PSD) if $\mathbf{u}^T M \mathbf{u} \geq 0$ for any vector \mathbf{u} .

Note: $\mathbf{u}^T (\mathbf{X}^T \mathbf{X}) \mathbf{u} = (\mathbf{X} \mathbf{u})^T \mathbf{X} \mathbf{u} = \|\mathbf{X} \mathbf{u}\|_2^2 \geq 0$ and is 0 if $\mathbf{u} = \mathbf{0}$.

The Hessian matrix of a convex function is positive semi-definite.

September 10, 2019 24 / 54

Computational complexity

Bottleneck of computing

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

is to invert the matrix $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$

- naively need $O(D^3)$ time
- there are many faster approaches (such as conjugate gradient)

September 10, 2019 25 / 54

How about the following?

$D = 1, N = 2$

sqft	sale price
1000	500K
1000	600K

Any line passing **the average** is a minimizer of RSS.

$D = 2, N = 3?$

sqft	#bedroom	sale price
1000	2	500K
1500	3	700K
2000	4	800K

Again *infinitely many minimizers*. How to resolve this issue?

September 10, 2019 27 / 54

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

Why would that happen?

One situation: $N < D + 1$, i.e. not enough data to estimate all parameters.

Example: $D = N = 1$

sqft	sale price
1000	500K

Any line passing through this single point is a minimizer of RSS.

September 10, 2019 26 / 54

Eigendecomposition

The decomposition of a square matrix A into matrices composed of its *eigenvectors* and *eigenvalues* is called eigendecomposition.

$$A = U \Lambda U^{-1}$$

where Λ is a diagonal matrix of eigenvalues of A , and each column of U is an eigenvector of A .

If A is symmetric $U^T U = \mathbf{I}$, then

$$A = U \Lambda U^{-1} = U \Lambda U^T = U^T \Lambda U$$

and its inverse

$$A^{-1} = U^T \Lambda^{-1} U$$

September 10, 2019 28 / 54

How to resolve this issue?

Eigendecomposition:

$$\mathbf{X}^T \mathbf{X} = \mathbf{U}^T \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_D & 0 \\ 0 & \cdots & 0 & \lambda_{D+1} \end{bmatrix} \mathbf{U}$$

where $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_{D+1} \geq 0$ are **eigenvalues**.

Inverse:

$$(\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{U}^T \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_D} & 0 \\ 0 & \cdots & 0 & \frac{1}{\lambda_{D+1}} \end{bmatrix} \mathbf{U}$$

Solution

The solution becomes

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

- not a minimizer of the original RSS

λ is a **hyper-parameter**, can be tuned by cross-validation.

How do we predict?

$$f(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x}$$

How to solve this problem?

Non-invertible \Rightarrow some eigenvalues are 0.

One natural fix: add something positive

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \mathbf{U}^T \begin{bmatrix} \lambda_1 + \lambda & 0 & \cdots & 0 \\ 0 & \lambda_2 + \lambda & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_D + \lambda & 0 \\ 0 & \cdots & 0 & \lambda_{D+1} + \lambda \end{bmatrix} \mathbf{U}$$

where $\lambda > 0$ and \mathbf{I} is the identity matrix. Now it is invertible:

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} = \mathbf{U}^T \begin{bmatrix} \frac{1}{\lambda_1 + \lambda} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2 + \lambda} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_D + \lambda} & 0 \\ 0 & \cdots & 0 & \frac{1}{\lambda_{D+1} + \lambda} \end{bmatrix} \mathbf{U}$$

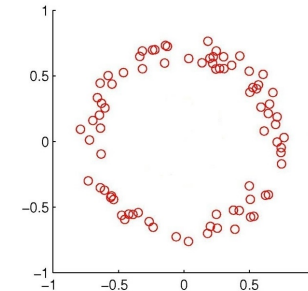
Comparison to NNC

Parametric versus non-parametric

- **Parametric methods:** the size of the model does **not grow** with the size of the training set N .
 - e.g. linear regression, Naive Bayes
- **Non-parametric methods:** the size of the model **grows** with the size of the training set.
 - NNC, Decision Trees

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and Preventing Overfitting

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



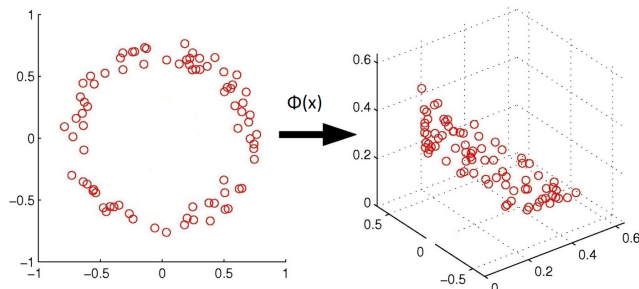
Solution: nonlinearly transformed features

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).



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:

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

Example

Polynomial basis functions for $D = 1$

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

General case:

$$\phi(x) = \left[\prod_{i=1}^D x_i^{a_i} \right], \text{ s.t. } \sum_{i=1}^D a_i \leq M$$

Learning a linear model in the new space
= learning an *M-degree polynomial model* in the original space

September 10, 2019 37 / 54

Why nonlinear?

Can I use a fancy **linear** feature map? For example,

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

No, it basically *does nothing* since

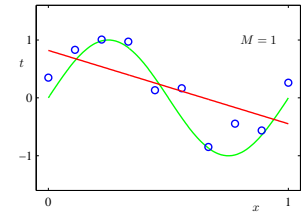
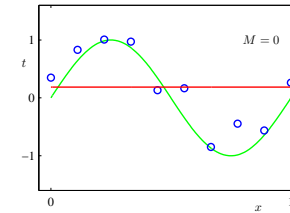
$$\min_{w \in \mathbb{R}^M} \sum_n (w^T \mathbf{A}x_n - y_n)^2 = \min_{w' \in \text{Im}(\mathbf{A}^T) \subset \mathbb{R}^D} \sum_n (w'^T x_n - y_n)^2$$

We will see more nonlinear mappings soon.

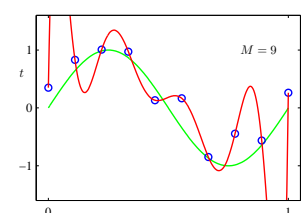
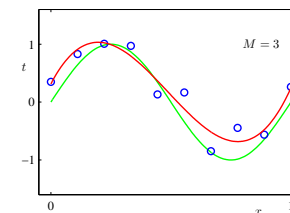
September 10, 2019 39 / 54

Example

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



M=9: overfitting



September 10, 2019 38 / 54

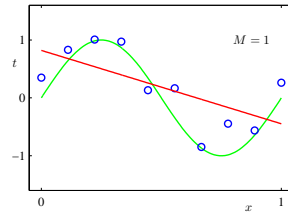
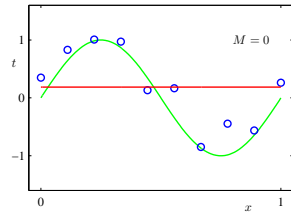
Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and Preventing Overfitting

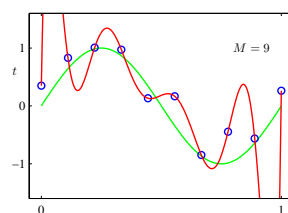
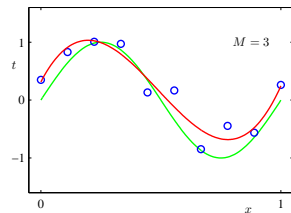
September 10, 2019 40 / 54

Should we use a very complicated mapping?

Ex: fitting a sine function with a polynomial:



Training error is zero:



September 10, 2019 41 / 54

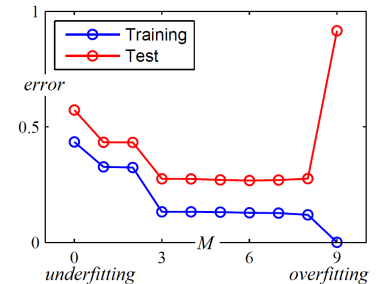
Underfitting and Overfitting

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

- large training error
- large test error

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

- small training error
- large test error



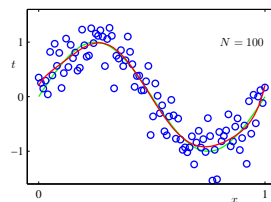
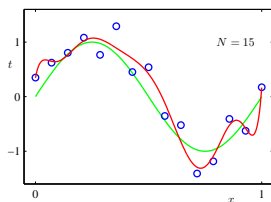
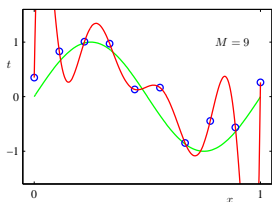
More complicated models \Rightarrow larger gap between training and test error

How to prevent overfitting?

September 10, 2019 42 / 54

Method 1: use more training data

The more, the merrier. We increase N - the number of training points.



More data \Rightarrow smaller gap between training and test error

September 10, 2019 43 / 54

Method 2: control the model complexity

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

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

Use cross-validation to pick hyperparameter M

Are there still other ways to control complexity?

September 10, 2019 44 / 54

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

September 10, 2019 45 / 54

The effect of λ

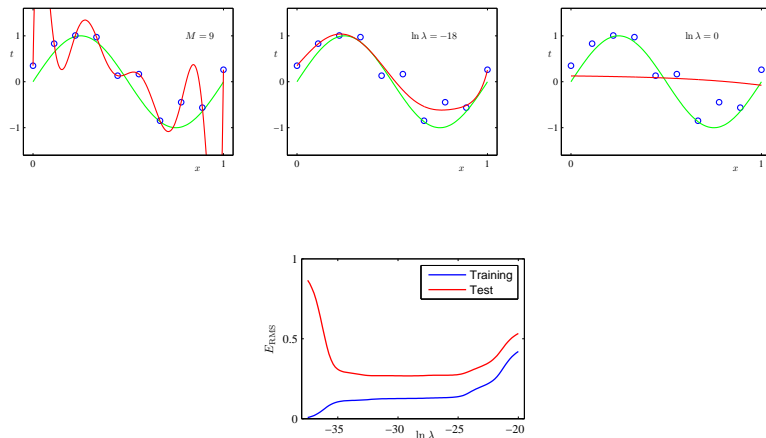
when we increase regularization coefficient λ

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

September 10, 2019 46 / 54

The trade-off

When we increase regularization coefficient λ , overfitting decreases:



September 10, 2019 47 / 54

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 a hyperparameter. To tune it,

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

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

September 10, 2019 48 / 54

The root of overfitting

Dealing with over and underfitting is really about dealing with bias and variance.

Mathematically, the expected prediction error can be decomposed into bias and variance components.

Simpler models have a smaller variance but a larger bias.

Complex models have a larger variance but a smaller bias.

Thus, we balance bias and variance by choosing λ .

Regularization reduces variance (because they lead to simpler models) but then increase the bias.

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 hyperparameter.

Finding the solution becomes a *constrained optimization problem*.

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

How to solve the new objective?

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

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

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

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

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

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.

Summary

Linear regression summarized:

$$w^* = (X^T X)^{-1} X^T y$$

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

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

It is important to understand the derivation.

Overfitting: small training error but large test error.

Preventing Overfitting: more data and/or regularization.

Typical steps

Typical steps of developing a machine learning system:

- Collect data, split into training, development, and test sets.
- *Train an ML model* with training data to learn from.
- Evaluate it using the test data and report performance.
- Use the model to predict future/make decisions.

How to do the *red part* exactly?

General idea to provide ML algorithms

1. Pick a set of **models** \mathcal{F}
 - e.g. $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \mid \mathbf{w} \in \mathbb{R}^D\}$
 - e.g. $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^M\}$

2. Define **error/loss** $L(y', y)$

3. Find **empirical risk minimizer (ERM)**:

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

or **regularized empirical risk minimizer**:

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

ML becomes optimization