

# Lecture 6

## Linear Regression

Luigi Freda

ALCOR Lab  
DIAG  
University of Rome "La Sapienza"

December 6, 2016

## 1 Intro

- Linear Regression

## 2 MLE - Least Squares

- Basic Idea
- MLE Derivation
- Geometric Interpretation

## 3 Convexity

- Idea
- Convex Set
- Convex Function
- Convexity and Optimization

## 4 Ridge Regression

- MLE Issues
- Basic Idea

## 5 Regularization Effects of Big Data

- Idea and Learning Curves

# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

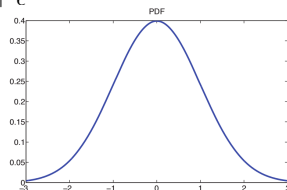
# Linear Regression

## Model Specification

### Linear Regression

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \epsilon = \sum_{j=1}^D w_j x_j + \epsilon$$

- $y \in \mathbb{R}$  and  $\mathbf{x} \in \mathbb{R}^D$
- $\mathbf{w} \in \mathbb{R}^D$  is the **weight vector**
- $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is the **residual error**



This entails

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(\mu(\mathbf{x}), \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$$

- $\mu(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = [w_0, \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}]^T$  where  $\mathbf{x} = [1, \tilde{\mathbf{x}}]^T$
- $\theta = (\mathbf{w}, \sigma^2)$  are the **model parameters**

# Linear Regression

## Polynomial Model Specification

### Polynomial Regression

if we replace  $\mathbf{x}$  by a non-linear function  $\phi(\mathbf{x})$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$$

we now have

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}), \sigma^2)$$

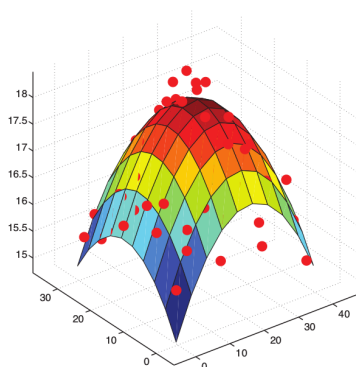
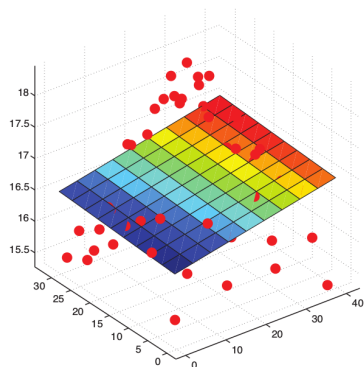
- $\mu(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$  (basis function expansion)
- if  $x \in \mathbb{R}$  we can use  $\phi(x) = [1, x, x^2, \dots, x^d]$  which is the vector of **polynomial basis functions**
- in general if  $\mathbf{x} \in \mathbb{R}^D$  we can use a **multivariate polynomial expansion**  
 $\mathbf{w}^T \phi(\mathbf{x}) = \sum w_{i_1 i_2 \dots i_D} \prod_{j=1}^D x_j^{i_j}$  up to a certain degree  $d$
- $\theta = (\mathbf{w}, \sigma^2)$  are the **model parameters**

N.B.: note that the model is still linear in the parameters  $\mathbf{w}$

# Linear Regression

## 2D Data

vertical axis is the temperature, horizontal axes are location within a room



- *left*: fitted function is a plane  $\hat{f}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$
- *right*: fitted function is quadratic  $\hat{f}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2$

# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - **Basic Idea**
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

- the MLE of a parameter  $\theta$  is computed as

$$\hat{\theta}_{MLE} \triangleq \underset{\theta}{\operatorname{argmax}} p(\mathcal{D}|\theta)$$

- given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  of iid samples, the **log-likelihood** can be computed as

$$l(\theta) \triangleq \log p(\mathcal{D}|\theta) = \sum_{i=1}^N \log p(y_i|\mathbf{x}_i, \theta)$$

- maximizing the log-likelihood is equivalent to minimize the **Negative Log-Likelihood (NLL)**

$$\text{NLL}(\theta) \triangleq - \sum_{i=1}^N \log p(y_i|\mathbf{x}_i, \theta)$$



- inserting the linear regression model into the log-likelihood returns

$$\begin{aligned}l(\boldsymbol{\theta}) &= \sum_{i=1}^N \log \left[ \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left( -\frac{1}{2\sigma^2} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \right] = \\ &= -\frac{\text{RSS}(\mathbf{w})}{2\sigma^2} - \frac{N}{2} \log(2\pi\sigma^2)\end{aligned}$$

where RSS stands for **Residual Sum of Squares**

$$\text{RSS}(\mathbf{w}) \triangleq \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

aka **Sum of Squared Errors** (SSE), i.e SSE = RSS

- the Mean Squared Error (MSE) is  $\text{MSE} \triangleq \text{SSE}/N$

- the log-likelihood is

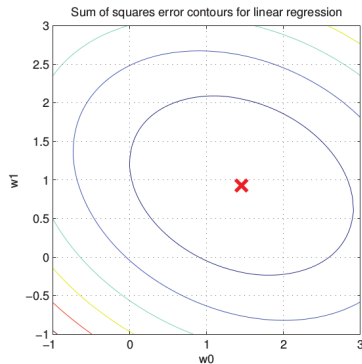
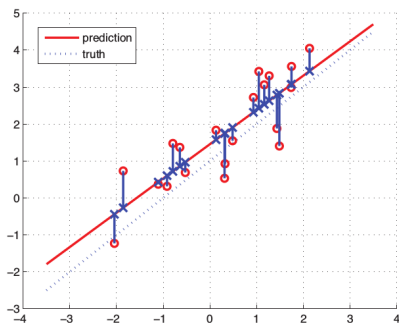
$$l(\theta) = -\frac{\text{RSS}(\mathbf{w})}{2\sigma^2} - \frac{N}{2} \log(2\pi\sigma^2)$$

- if we define the **residual errors**  $\epsilon_i \triangleq y_i - \mathbf{w}^T \mathbf{x}_i$ , one has

$$\text{RSS}(\mathbf{w}) = \|\epsilon\|_2^2 = \sum_{i=1}^N \epsilon_i^2$$

- the MLE for  $\mathbf{w}$  minimizes the RSS, and for this reason the method is called **least squares**

# MLE - Least Squares



- *left*: red circles are training points; blue crosses are approximations
- in least squares we try to minimize the sum of squared distances from each training point to its approximation (i.e. the sum of the lengths of the little vertical blue lines)
- *right*: contours of the RSS error; the surface is a **bowl** with a **unique minimum**

# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - **MLE Derivation**
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

# Derivation of the MLE

- the NLL can be rewritten as

$$\text{NLL}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \frac{1}{2} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \mathbf{w}^T (\mathbf{X}^T \mathbf{y}) + \frac{1}{2} \mathbf{y}^T \mathbf{y}$$

where  $\mathbf{y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N$ ,  $\mathbf{X} \in \mathbb{R}^{N \times D}$  is the design matrix

- in order to minimize the NLL we have to compute its gradient

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

and equating it to zero we get the **normal equation**

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

- the corresponding solution is called the **Ordinary Least Squares (OLS)**

$$\hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- considering that  $\mathbf{X}^T \mathbf{X} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$  and  $\mathbf{X}^T \mathbf{y} = \sum_i \mathbf{x}_i y_i$ , one has

$$\hat{\mathbf{w}}_{OLS} = \left( \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \sum_i \mathbf{x}_i y_i$$

# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - **Geometric Interpretation**
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

# Geometric Interpretation

## Least Squares

- let's consider the optimization problem we solve

$$\hat{\mathbf{w}}_{OLS} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

- this optimization problem is equivalent to solve the **normal equation**

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

- these equations have an elegant geometric interpretation
- the columns of  $\mathbf{X} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D] \in \mathbb{R}^{N \times D}$  define a linear subspace in  $\mathbb{R}^N$
- N.B.: the columns  $\tilde{\mathbf{x}}_j \in \mathbb{R}^N$  of  $\mathbf{X}$  are different from the rows  $\mathbf{x}_i \in \mathbb{R}^D$  of  $\mathbf{X}$
- the vector  $\mathbf{y}$  lives in  $\mathbb{R}^N$
- the least squares problem above is equivalent to

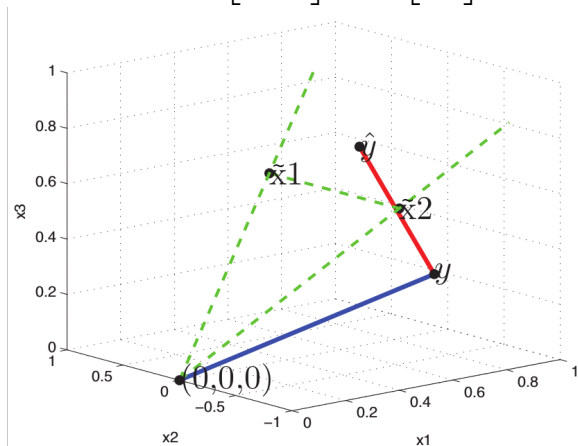
$$\underset{\hat{\mathbf{y}} \in \operatorname{span}(\{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D\})}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2 \quad (\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} = w_1\tilde{\mathbf{x}}_1 + \dots + w_D\tilde{\mathbf{x}}_D)$$

# Geometric Interpretation

## Least Squares

- let's assume  $N = 3$  and  $D = 2$

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 1 & -2 \\ 1 & 2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 8.89 \\ 0.61 \\ 1.77 \end{bmatrix}$$



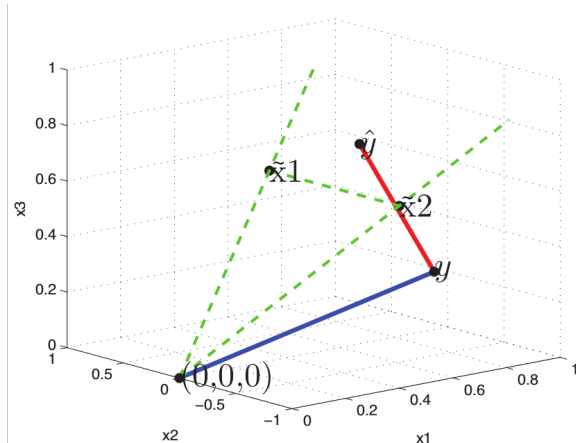


# Geometric Interpretation

## Least Squares

- we seek a vector  $\hat{\mathbf{y}} \in \mathbb{R}^N$  that lies in the linear subspace defined by the columns of  $\mathbf{X}$  and is close as possible to  $\mathbf{y} \in \mathbb{R}^N$

$$\underset{\hat{\mathbf{y}} \in \text{span}\{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D\}}{\text{argmin}} \quad \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2 \quad (\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} = w_1\tilde{\mathbf{x}}_1 + \dots + w_D\tilde{\mathbf{x}}_D)$$

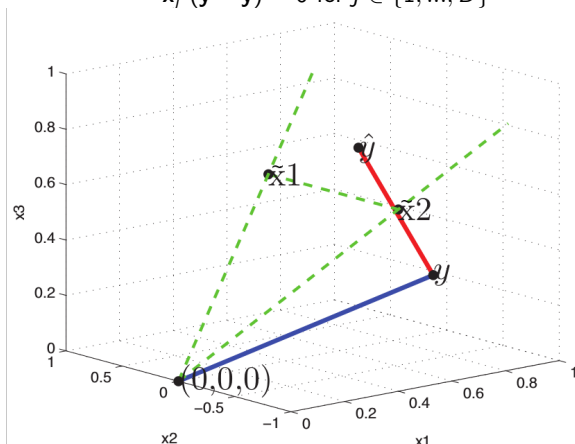


# Geometric Interpretation

## Least Squares

- if we want to minimize the residual  $\operatorname{argmin}_{\hat{\mathbf{y}} \in \operatorname{span}(\{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D\})} \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$  then we also want the residual vector  $\mathbf{y} - \hat{\mathbf{y}}$  to be orthogonal to the linear space (hyperplane) defined by the columns of  $\mathbf{X}$

$$\tilde{\mathbf{x}}_j^T (\hat{\mathbf{y}} - \mathbf{y}) = 0 \text{ for } j \in \{1, \dots, D\}$$

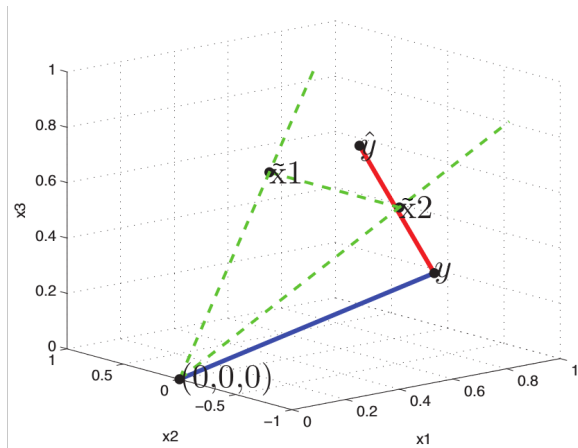


# Geometric Interpretation

## Least Squares

- hence, solving  $\operatorname{argmin}_{\hat{\mathbf{y}} \in \operatorname{span}(\{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_D\})} \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$  is equivalent to solve the normal equation

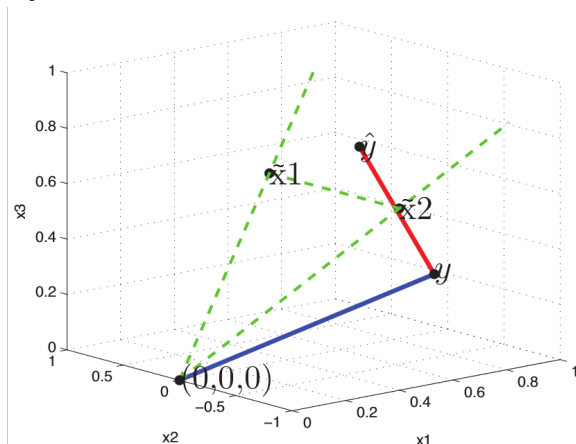
$$\tilde{\mathbf{x}}_j^T (\hat{\mathbf{y}} - \mathbf{y}) = 0 \text{ for } j \in \{1, \dots, D\} \implies \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y}) = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$$



# Geometric Interpretation

## Least Squares

- solving  $\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$  returns  $\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$
- we can get the **orthogonal projection**  $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$
- the projection matrix  $\mathbf{P} \triangleq \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  is called the **hat matrix**, since it "puts the hat" on  $\mathbf{y}$



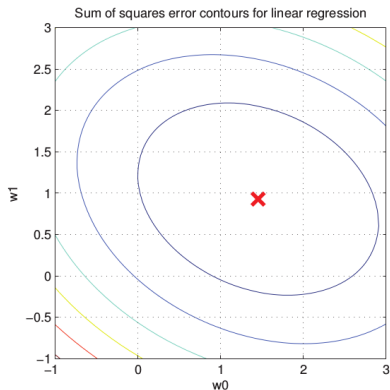
# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

# Convexity

## Intro

- when solving least squares, we noted that the NLL had a bowl shape with a **unique minimum**
- the technical term for functions like this is **convex**
- convex functions play a very important role in machine learning



# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - **Convex Set**
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

# Convexity

## Convex Set Definition

- a set  $\mathcal{S}$  is **convex** if for any  $\theta, \theta' \in \mathcal{S}$  we have

$$\lambda\theta + (1 - \lambda)\theta' \in \mathcal{S} \quad \forall \lambda \in [0, 1]$$

that is, if we draw a line from  $\theta$  to  $\theta'$ , all points on the line lie inside the set  $\mathcal{S}$

- for instance in this figure



we have a convex set on the left and a nonconvex set on the right



# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - **Convex Function**
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

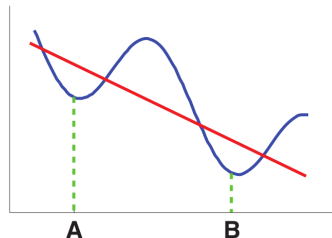
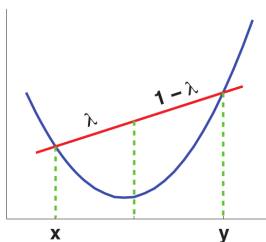
# Convexity

## Convex Function Definition

- a function  $f(\theta)$  is **convex** if it is defined on a convex set and if, for any  $\theta, \theta' \in \mathcal{S}$ , and for any  $\lambda \in [0, 1]$ , we have

$$f(\lambda\theta + (1 - \lambda)\theta') \leq \lambda f(\theta) + (1 - \lambda)f(\theta')$$

- a function  $f(\theta)$  is **strictly convex** if the inequality is strict
- a function  $f(\theta)$  is **concave** if  $-f(\theta)$  is convex
- for instance, in the following figure

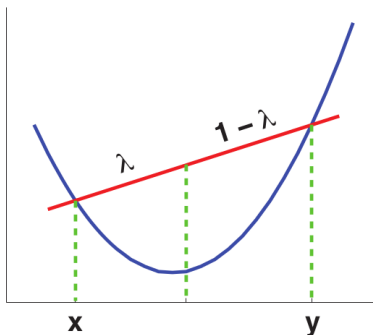


the function on the left is convex and the function on the right is neither concave nor convex

# Convexity

## Convex Function Definition

- examples of scalar convex functions include  $\theta^2$ ,  $e^\theta$  and  $\theta \log \theta$  (for  $\theta > 0$ )
- examples of concave functions include  $\log \theta$  and  $\sqrt{\theta}$

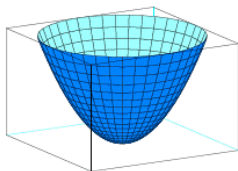


- intuitively, a (strictly) convex function has a *bowl shape*, and hence has a **unique global minimum**  $\theta^*$  corresponding to the bottom of the bowl
- in the scalar case, a convex function has its second derivative positive everywhere, i.e.  $\frac{d}{d\theta^2} f(\theta) > 0$

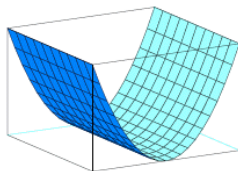
# Convexity

## Convex Function Examples

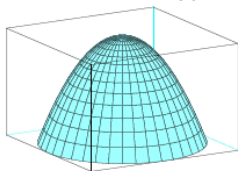
Convex  $f = x^2 + y^2$



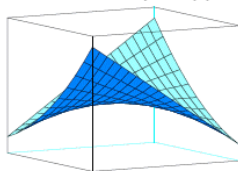
Convex (degenerate)  $f = x^2$



Concave  $f = -x^2 - y^2$



Nonconvex  $f = x^2 + 0.3y^2$



- convex(concave) functions have a unique global minimum(maximum)

# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

### Theorem 1

a twice-continuously differentiable, multivariate function  $f(\boldsymbol{\theta}) \in \mathbb{R}$  is convex **iff** its Hessian is positive definite for all  $\boldsymbol{\theta}$

- the **Hessian matrix**  $\mathbf{H} = \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^2}$  of a function  $f(\boldsymbol{\theta}) \in \mathbb{R}$  is defined as follows (element-wise)

$$\mathbf{H}_{jk} = \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_j \partial \theta_k}$$

- the Hessian is symmetric since  $\frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_j \partial \theta_k} = \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_k \partial \theta_j}$  (Schwartz's theorem)
- recall that a matrix  $\mathbf{H}$  is **positive definite** iff  $\mathbf{v}^T \mathbf{H} \mathbf{v} > 0$  for any  $\mathbf{v} \neq 0$
- a convex function has a bowl shape

# Convexity

## Convex Function

- a convex function can be approximated about its unique global minimum  $\theta^*$  with a bowl shaped quadratic function (paraboloid in 3D)

$$\begin{aligned}f(\theta) &\approx f(\theta^*) + \left. \frac{\partial f}{\partial \theta} \right|_{\theta^*} (\theta - \theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H}(\theta^*) (\theta - \theta^*) = \\ &= f(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H}(\theta^*) (\theta - \theta^*)\end{aligned}$$

where we used the fact that at the minimum  $\theta^*$  one has  $\left. \frac{\partial f}{\partial \theta} \right|_{\theta^*} = 0$  and we know that  $\mathbf{H}(\theta^*) > 0$

# Outline

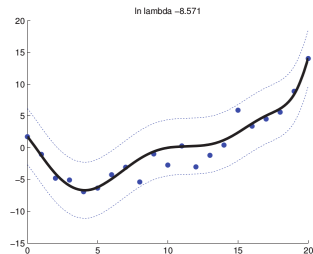
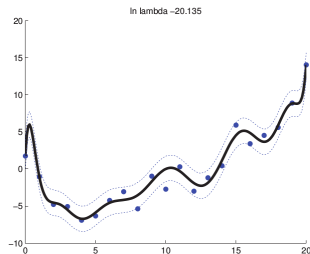
- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves



# Ridge Regression

## MLE Issues

- one problem with ML estimation is that it can result in **overfitting**
- the reason that the MLE can overfit is that it is picking the parameter values that are the best for modeling the **training data**
- but if the data is **noisy**, such parameters often result in complex functions
- as a simple example, suppose we fit a degree 14 polynomial to  $N = 21$  data points using least squares. The resulting curve is very "wiggly"
- in this case, overfitting equals interpolating noise
- if we change the data a little bit our estimate of  $\mathbf{w}$  will change a lot (unstable estimation)



# Ridge Regression

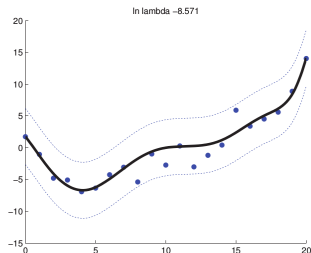
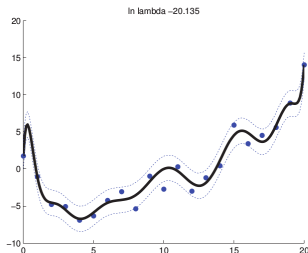
## Basic Idea

- the corresponding least squares coefficients  $\mathbf{w}$  (excluding  $w_0$ ) are as follows

6.560, -36.934, -109.255, 543.452, 1022.561, -3046.224, -3768.013,

8524.540, 6607.897, -12640.058, -5530.188, 9479.730, 1774.639, -2821.526

- there are many large positive and negative numbers
- these variations balance out exactly to make the curve “wiggle” in just the right way so that it almost perfectly interpolates the data



# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - **Basic Idea**
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

# Ridge Regression

- we can encourage the parameters to be small, thus resulting in a smoother curve, by using a zero-mean Gaussian prior

$$p(\mathbf{w}) = \prod_j \mathcal{N}(w_j | \mathbf{0}, \tau^2)$$

where the precision  $1/\tau^2$  controls the strength of the prior

- the corresponding MAP estimation problem becomes

$$\begin{aligned} \hat{\mathbf{w}}_{MAP} &= \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w} | \mathcal{D}) = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathcal{D} | \mathbf{w}) p(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \log p(\mathcal{D} | \mathbf{w}) + \log p(\mathbf{w}) = \\ &= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^N \log \mathcal{N}(y_i | w_0 + \mathbf{w}^T \mathbf{x}_i, \sigma^2) + \sum_{i=1}^D \log \mathcal{N}(w_i | \mathbf{0}, \tau^2) \end{aligned}$$

- it is simple to show that this is equivalent to minimize

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where  $\lambda \triangleq \sigma^2/\tau^2$  and  $\|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$

- the first term is the MSE/NLL as usual, and the second term with  $\lambda > 0$  is a **complexity penalty**

- minimization problem

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

- the solution is

$$\hat{\mathbf{w}}_{ridge} = (\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- this technique is known as **ridge regression**

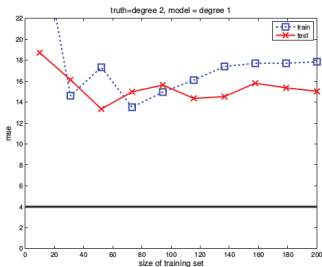
# Outline

- 1 Intro
  - Linear Regression
- 2 MLE - Least Squares
  - Basic Idea
  - MLE Derivation
  - Geometric Interpretation
- 3 Convexity
  - Idea
  - Convex Set
  - Convex Function
  - Convexity and Optimization
- 4 Ridge Regression
  - MLE Issues
  - Basic Idea
- 5 Regularization Effects of Big Data
  - Idea and Learning Curves

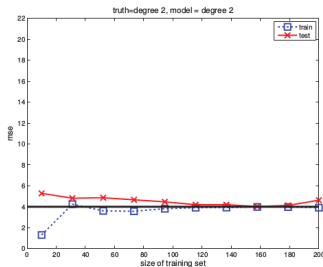
# Regularization Effects of Big Data

- **regularization** is the most common way to avoid **overfitting**
- another effective approach — which is not always available — is to use **lots of data**
- intuitively, the more training data we have, the better we will be able to learn
- let's consider a **polynomial regression** and let's plot the Mean Squared Error (MSE) incurred on the test set achieved by models of different degrees vs  $N = |\mathcal{D}|$
- a plot of error vs training set size is known as a **learning curve**

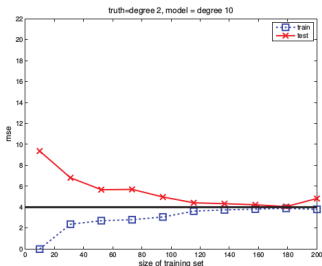
# Regularization Effects of Big Data



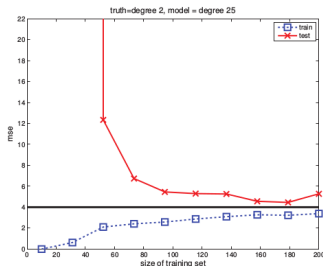
(a)



(b)



(c)

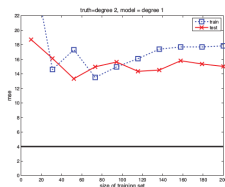


(d)

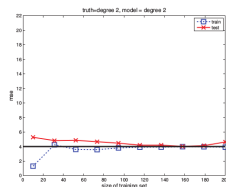


# Regularization Effects of Big Data

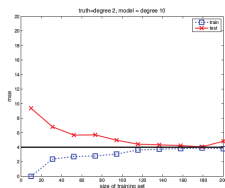
- the truth is a degree 2 polynomial
- we try fitting polynomials of degrees 1, 2, 10 and 25 to this data (respectively models  $\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_{10}$  and  $\mathcal{M}_{25}$ )



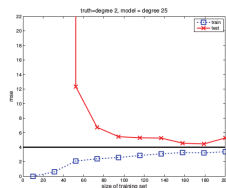
(a)



(b)



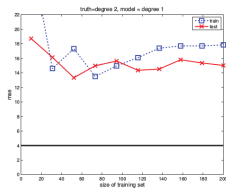
(c)



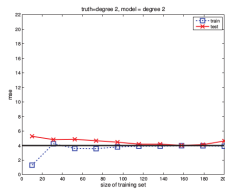
(d)

# Regularization Effects of Big Data

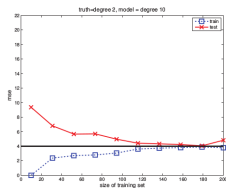
- the level of the plateau for the test error consists of two terms
  - noise floor**: an irreducible component that all models incur, due to the intrinsic variability of the generating process
  - structural error**: a component that depends on the discrepancy between the generating process (the “truth”) and the model



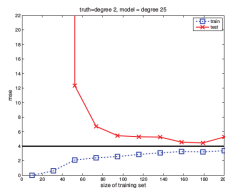
(a)



(b)



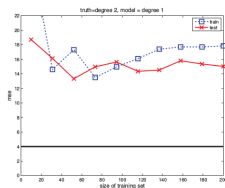
(c)



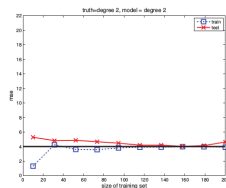
(d)

# Regularization Effects of Big Data

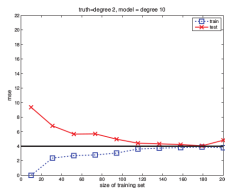
- the structural error for models  $\mathcal{M}_2$ ,  $\mathcal{M}_{10}$  and  $\mathcal{M}_{25}$  is zero, since both are able to capture the true generating process
- the structural error for  $\mathcal{M}_1$  is substantial, which is evident from the fact that the plateau occurs high above the noise floor



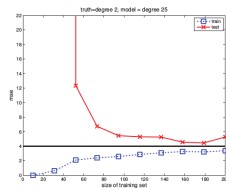
(a)



(b)



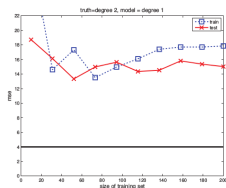
(c)



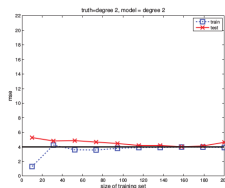
(d)

# Regularization Effects of Big Data

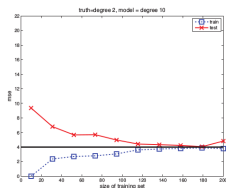
- for any model that is expressive enough to capture the truth (i.e., one with small structural error), the **test error** will go to the **noise floor** as  $N \rightarrow \infty$
- the error will typically go to zero faster for simpler models (fewer parameters to estimate)



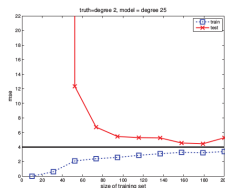
(a)



(b)



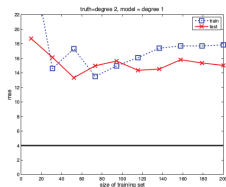
(c)



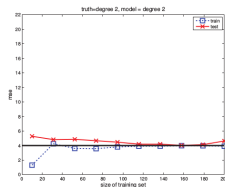
(d)

# Regularization Effects of Big Data

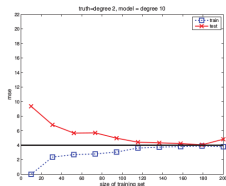
- for finite training sets, there will be some discrepancy between the parameters that we estimate and the best parameters that we could estimate given the particular model class
- this discrepancy is called **approximation error**, and goes to zero as  $N \rightarrow \infty$ , but it goes to zero faster for simpler models



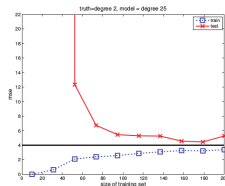
(a)



(b)



(c)



(d)

# Regularization Effects of Big Data

- in domains with lots of data, **simple methods** can work surprisingly well
- however, there are still reasons to study more **sophisticated learning methods**, because there will always be problems for which we have **little data**
- for example, even in such a data-rich domain as web search, as soon as we want to start **personalizing the results**, the amount of data available for any given user starts to look small again (relative to the complexity of the problem)

- Kevin Murphy's book