

Linear Regression (continued)

Professor Ameet Talwalkar

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Linear regression
- 4 Nonlinear basis functions

Announcements

- HW2 will be returned in section on Friday
- HW3 due in class next Monday
- Midterm is next Wednesday (will review in more detail next class)

Outline

- 1 Administration
- 2 Review of last lecture
 - Perceptron
 - Linear regression introduction
- 3 Linear regression
- 4 Nonlinear basis functions

Perceptron Main idea

Consider a linear model for binary classification

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

We use this model to distinguish between two classes $\{-1, +1\}$.

One goal

$$\varepsilon = \sum_n \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)]$$

i.e., to minimize errors on the training dataset.

Hard, but easy if we have only one training example

How can we change w such that

$$y_n = \text{sign}(w^T x_n)$$

Two cases

- If $y_n = \text{sign}(w^T x_n)$, do nothing.
- If $y_n \neq \text{sign}(w^T x_n)$, $w^{\text{NEW}} \leftarrow w^{\text{OLD}} + y_n x_n$
 - ▶ Guaranteed to make progress, i.e., to get us closer to $y(w^T x) > 0$

Hard, but easy if we have only one training example

How can we change \mathbf{w} such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

Two cases

- If $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, do nothing.
- If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, $\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w}^{\text{OLD}} + y_n \mathbf{x}_n$
 - ▶ Guaranteed to make progress, i.e., to get us closer to $y(\mathbf{w}^T \mathbf{x}) > 0$

What does update do?

$$y_n [(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

Hard, but easy if we have only one training example

How can we change \mathbf{w} such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

Two cases

- If $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, do nothing.
- If $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$, $\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w}^{\text{OLD}} + y_n \mathbf{x}_n$
 - ▶ Guaranteed to make progress, i.e., to get us closer to $y(\mathbf{w}^T \mathbf{x}) > 0$

What does update do?

$$y_n [(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

We are adding a positive number, so it's possible that $y_n(\mathbf{w}^{\text{NEW}T} \mathbf{x}_n) > 0$

Perceptron algorithm

Iteratively solving one case at a time

- REPEAT
- Pick a data point \mathbf{x}_n (can be a fixed order of the training instances)
- Make a prediction $y = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ using the *current* \mathbf{w}
- If $y = y_n$, do nothing. Else,

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

- UNTIL converged.

Properties

Perceptron algorithm

Iteratively solving one case at a time

- REPEAT
- Pick a data point \mathbf{x}_n (can be a fixed order of the training instances)
- Make a prediction $y = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ using the *current* \mathbf{w}
- If $y = y_n$, do nothing. Else,

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

- UNTIL converged.

Properties

- This is an online algorithm.
- If the training data is linearly separable, the algorithm stops in a finite number of steps (we proved this).
- The parameter vector is always a linear combination of training instances (requires initialization of $\mathbf{w}_0 = 0$).

Regression

Predicting a continuous outcome variable

- Predicting shoe size from height, weight and gender
- Predicting song year from audio features

Regression

Predicting a continuous outcome variable

- Predicting shoe size from height, weight and gender
- Predicting song year from audio features

Key difference from classification

Regression

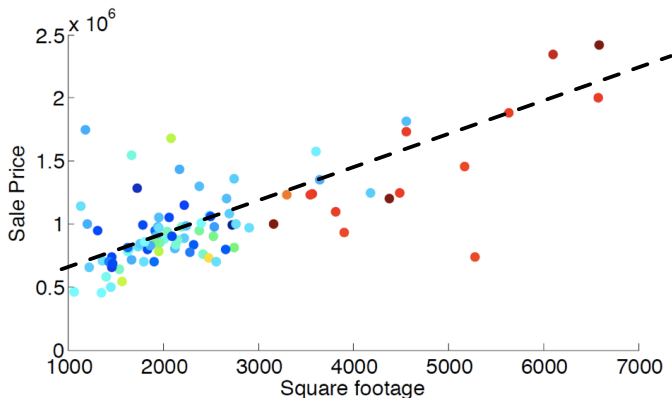
Predicting a continuous outcome variable

- Predicting shoe size from height, weight and gender
- Predicting song year from audio features

Key difference from classification

- We can measure 'closeness' of prediction and labels
 - ▶ Predicting shoe size: better to be off by one size than by 5 sizes
 - ▶ Predicting song year: better to be off by one year than by 20 years
- As opposed to 0-1 classification error, we will focus on squared difference, i.e., $(\hat{y} - y)^2$

1D example: predicting the sale price of a house



Sale price \approx price_per_sqft \times square_footage + fixed_expense

Minimize squared errors

Our model

Sale price = price_per_sqft \times square_footage + fixed_expense + unexplainable_stuff

Training data

sqft	sale price	prediction	error	squared error
2000	810K	720K	90K	8100
2100	907K	800K	107K	107^2
1100	312K	350K	38K	38^2
5500	2,600K	2,600K	0	0
...	...			
Total				$8100 + 107^2 + 38^2 + 0 + \dots$

Minimize squared errors

Our model

Sale price = price_per_sqft \times square_footage + fixed_expense + unexplainable_stuff

Training data

sqft	sale price	prediction	error	squared error
2000	810K	720K	90K	8100
2100	907K	800K	107K	107^2
1100	312K	350K	38K	38^2
5500	2,600K	2,600K	0	0
...	...			
Total				$8100 + 107^2 + 38^2 + 0 + \dots$

Aim

Adjust price_per_sqft and fixed_expense such that the sum of the squared error is minimized — i.e., unexplainable_stuff is minimized.

Linear regression

Setup

- Input: $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- Output: $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- Model: $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$

Linear regression

Setup

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

Linear regression

Setup

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

Linear regression

Setup

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

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

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

Linear regression

Setup

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

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

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

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

Probabilistic interpretation

- Noisy observation model

$$Y = w_0 + w_1X + \eta$$

where $\eta \sim N(0, \sigma^2)$ is a Gaussian random variable

Probabilistic interpretation

- Noisy observation model

$$Y = w_0 + w_1 X + \eta$$

where $\eta \sim N(0, \sigma^2)$ is a Gaussian random variable

- Likelihood of one training sample (x_n, y_n)

$$p(y_n|x_n) = N(w_0 + w_1 x_n, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[y_n - (w_0 + w_1 x_n)]^2}{2\sigma^2}}$$

Maximum likelihood estimation

- Maximize over w_0 and w_1

$$\max \log P(\mathcal{D}) \Leftrightarrow \min \sum_n [y_n - (w_0 + w_1 x_n)]^2 \leftarrow \text{That is RSS}(\tilde{\mathbf{w}})!$$

Maximum likelihood estimation

- Maximize over w_0 and w_1

$$\max \log P(\mathcal{D}) \Leftrightarrow \min \sum_n [y_n - (w_0 + w_1 x_n)]^2 \leftarrow \text{That is RSS}(\tilde{\mathbf{w}})!$$

- Maximize over $s = \sigma^2$

$$\begin{aligned} \frac{\partial \log P(\mathcal{D})}{\partial s} &= -\frac{1}{2} \left\{ -\frac{1}{s^2} \sum_n [y_n - (w_0 + w_1 x_n)]^2 + \mathbf{N} \frac{1}{s} \right\} = 0 \\ \rightarrow \sigma^{*2} = s^* &= \frac{1}{\mathbf{N}} \sum_n [y_n - (w_0 + w_1 x_n)]^2 \end{aligned}$$

How does this probabilistic interpretation help us?

How does this probabilistic interpretation help us?

- It gives a solid footing to our intuition: minimizing $\text{RSS}(\tilde{\mathbf{w}})$ is a sensible thing based on reasonable modeling assumptions
- Estimating σ^* tells us how much noise there could be in our predictions. For example, it allows us to place confidence intervals around our predictions.

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Linear regression**
 - Multivariate solution in matrix form
 - Computational and numerical optimization
 - Ridge regression
- 4 Nonlinear basis functions

LMS when x is D-dimensional

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

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

LMS when \mathbf{x} is D-dimensional

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

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

where we have redefined some variables (by augmenting)

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

LMS when x is D-dimensional

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

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

where we have redefined some variables (by augmenting)

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

which leads to

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

LMS when \mathbf{x} is D-dimensional

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

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

where we have redefined some variables (by augmenting)

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

which leads to

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

LMS when \mathbf{x} is D-dimensional

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

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

where we have redefined some variables (by augmenting)

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

which leads to

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

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$9 \times 1 + 3 \times 3 + 5 \times 2 = 28$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & \\ & \end{bmatrix}$$

$$9 \times 1 + 3 \times 3 + 5 \times 2 = 28$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & \quad \\ \quad & \quad \end{bmatrix}$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & 18 \end{bmatrix}$$

Matrix Multiplication via Inner Products

Each entry of output matrix is result of inner product of inputs matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & 18 \\ 11 & 9 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 \\ 4 \end{bmatrix} \begin{bmatrix} 3 & 5 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix} + \begin{bmatrix} 9 & -15 \\ 3 & -5 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix} + \begin{bmatrix} 9 & -15 \\ 3 & -5 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix} + \begin{bmatrix} 9 & -15 \\ 3 & -5 \end{bmatrix} + \begin{bmatrix} 10 & 15 \\ 4 & 6 \end{bmatrix}$$

Matrix Multiplication via Outer Products

Output matrix is **sum of outer products** between corresponding rows and columns of input matrices

$$\begin{bmatrix} 9 & 3 & 5 \\ 4 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -5 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 28 & 18 \\ 11 & 9 \end{bmatrix}$$

$$\begin{bmatrix} 9 & 18 \\ 4 & 8 \end{bmatrix} + \begin{bmatrix} 9 & -15 \\ 3 & -5 \end{bmatrix} + \begin{bmatrix} 10 & 15 \\ 4 & 6 \end{bmatrix}$$

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

From previous slide

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

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

From previous slide

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

Design matrix and target vector

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

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

From previous slide

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

Design matrix and target vector

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

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

From previous slide

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

Design matrix and target vector

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

Compact expression

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

Solution in matrix form

Compact expression

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

Solution in matrix form

Compact expression

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

Gradients of Linear and Quadratic Functions

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

Solution in matrix form

Compact expression

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

Gradients of Linear and Quadratic Functions

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

Normal equation

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

Solution in matrix form

Compact expression

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

Gradients of Linear and Quadratic Functions

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

Normal equation

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

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

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

Mini-Summary

- Linear regression is the linear combination of features
 $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$
- If we minimize residual sum of squares as our learning objective, we get a closed-form solution of parameters
- Probabilistic interpretation: maximum likelihood if assuming residual is Gaussian distributed

Computational complexity

Bottleneck of computing the solution?

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

Computational complexity

Bottleneck of computing the 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 $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$

How many operations do we need?

Computational complexity

Bottleneck of computing the 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 $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$

How many operations do we need?

- $O(ND^2)$ for matrix multiplication
- $O(D^3)$ (e.g., using Gauss-Jordan elimination) or $O(D^{2.373})$ (recent theoretical advances) for matrix inversion
- Impractical for very large D or N

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (e.g., randomly); set $t = 0$; choose $\eta > 0$

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

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

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

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

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

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

Alternative method: an example of using numerical optimization

(Batch) Gradient descent

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

What is the complexity of each iteration?

Why would this work?

Why would this work?

If gradient descent converges, it will converge to the same solution as using matrix inversion.

This is because $RSS(\tilde{\mathbf{w}})$ is a convex function in its parameters $\tilde{\mathbf{w}}$

Why would this work?

If gradient descent converges, it will converge to the same solution as using matrix inversion.

This is because $RSS(\tilde{\mathbf{w}})$ is a convex function in its parameters $\tilde{\mathbf{w}}$

Hessian of RSS

$$\begin{aligned}RSS(\tilde{\mathbf{w}}) &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} + \text{const} \\ \Rightarrow \frac{\partial^2 RSS(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}} \tilde{\mathbf{w}}^T} &= 2 \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\end{aligned}$$

Why would this work?

If gradient descent converges, it will converge to the same solution as using matrix inversion.

This is because $RSS(\tilde{\mathbf{w}})$ is a convex function in its parameters $\tilde{\mathbf{w}}$

Hessian of RSS

$$\begin{aligned}RSS(\tilde{\mathbf{w}}) &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} + \text{const} \\ \Rightarrow \frac{\partial^2 RSS(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}} \tilde{\mathbf{w}}^T} &= 2 \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\end{aligned}$$

$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is positive semidefinite, because for any \mathbf{v}

$$\mathbf{v}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \mathbf{v} = \|\tilde{\mathbf{X}}^T \mathbf{v}\|_2^2 \geq 0$$

Stochastic gradient descent

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

- Initialize $\tilde{\mathbf{w}}$ to some $\tilde{\mathbf{w}}^{(0)}$; set $t = 0$; choose $\eta > 0$

Stochastic gradient descent

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

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

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

Stochastic gradient descent

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

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

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

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

Stochastic gradient descent

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

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

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

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

Stochastic gradient descent

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

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

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

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

How does the complexity per iteration compare with gradient descent?

Stochastic gradient descent

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

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

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

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

How does the complexity per iteration compare with gradient descent?

- $O(ND)$ for gradient descent versus $O(D)$ for SGD

Mini-summary

- Batch gradient descent computes the exact gradient.
- Stochastic gradient descent approximates the gradient with a single data point; Its expectation equals the true gradient.
- Mini-batch variant: trade-off between accuracy of estimating gradient and computational cost
- Similar ideas extend to other ML optimization problems.
 - ▶ For large-scale problems, stochastic gradient descent often works well.

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

Why might this happen?

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

Why might this happen?

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

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

Why might this happen?

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

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

Ridge regression

Intuition: what does a non-invertible $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ mean? Consider the SVD of this matrix:

Ridge regression

Intuition: what does a non-invertible $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ mean? Consider the SVD of this matrix:

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{U} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \mathbf{U}^T$$

where $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r > 0$ and $r < D$.

Ridge regression

Intuition: what does a non-invertible $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ mean? Consider the SVD of this matrix:

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{U} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \mathbf{U}^T$$

where $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r > 0$ and $r < D$.

Fix the problem by ensuring all singular values are non-zero

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I} = \mathbf{U} \text{diag}(\lambda_1 + \lambda, \lambda_2 + \lambda, \cdots, \lambda) \mathbf{U}^T$$

where $\lambda > 0$ and \mathbf{I} is the identity matrix

Regularized least square (ridge regression)

Solution

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

Regularized least square (ridge regression)

Solution

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

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

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

Regularized least square (ridge regression)

Solution

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

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

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

Benefits

- Numerically more stable, invertible matrix
- Prevent overfitting — more on this later

How to choose λ ?

λ is referred as *hyperparameter*

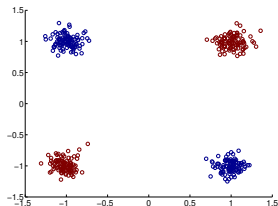
- In contrast w is the parameter vector
- Use validation set or cross-validation to find good choice of λ

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Linear regression
- 4 Nonlinear basis functions**

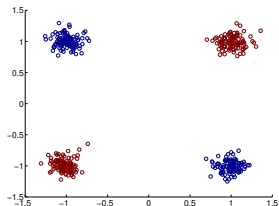
Is a linear modeling assumption always a good idea?

Example of nonlinear classification

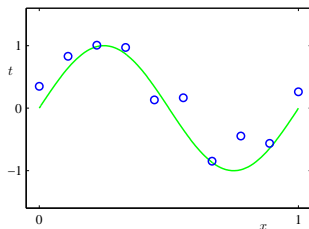


Is a linear modeling assumption always a good idea?

Example of nonlinear classification



Example of nonlinear regression



Nonlinear basis for classification

Transform the input/feature

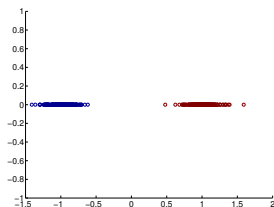
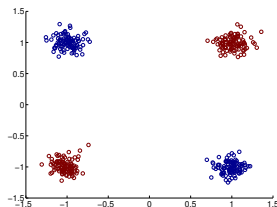
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Nonlinear basis for classification

Transform the input/feature

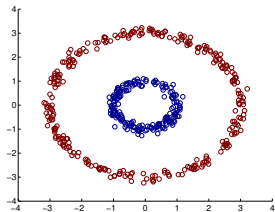
$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Transformed training data: linearly separable!

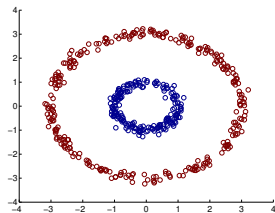


Another example

How to transform the input/feature?



Another example

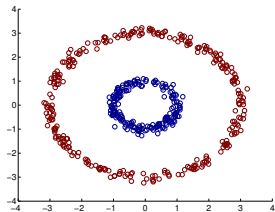


How to transform the input/feature?

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow \mathbf{z} = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

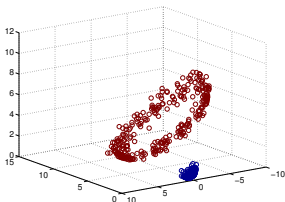
Another example

How to transform the input/feature?



$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^2 \rightarrow \mathbf{z} = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3$$

Transformed training data: linearly separable



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 $\phi(\mathbf{x})$).

- M could be greater than, less than, or equal to D

General nonlinear basis functions

We can use a nonlinear mapping

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

where M is the dimensionality of the new feature/input \mathbf{z} (or $\phi(\mathbf{x})$).

- M could be greater than, less than, or equal to D

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 $\phi(\mathbf{x})$.

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

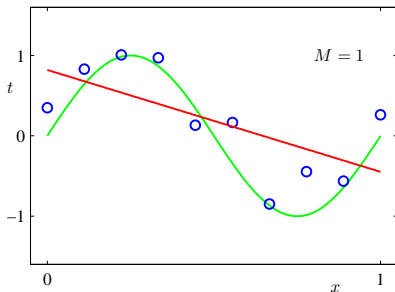
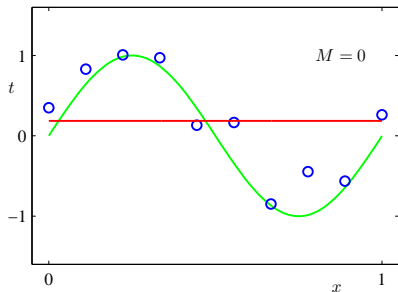
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Example with regression

Polynomial basis functions

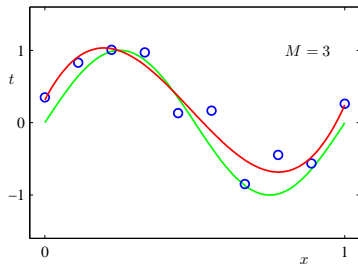
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

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



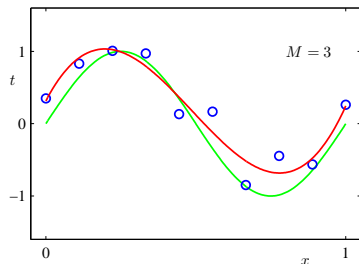
Adding high-order terms

$M=3$

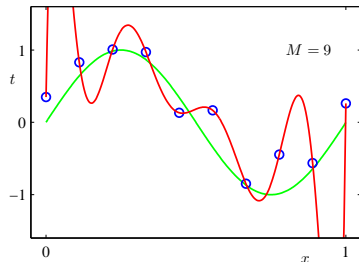


Adding high-order terms

$M=3$



$M=9$: *overfitting*



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

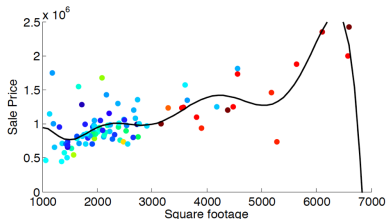
Overfitting

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 large M



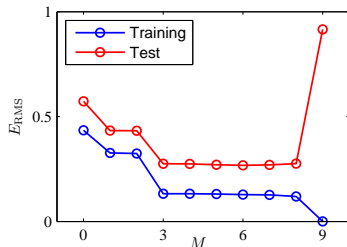
Predicted price goes to zero (and is ultimately negative) if you buy a big enough house!

This is called poor generalization/overfitting.

Detecting overfitting

Plot model complexity versus objective function

- X axis: model complexity, e.g., M
- Y axis: error, e.g., RSS, RMS (square root of RSS), 0-1 loss



As a model increases in complexity:

- Training error keeps improving
- Test error may first improve but eventually will deteriorate