

# CSCI567 Machine Learning (Fall 2024)

Prof. Dani Yogatama

University of Southern California

September 6, 2024

# Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and preventing overfitting
- 4 A Detour of Numerical Optimization Methods

# Regression

## Predicting a continuous outcome variable using past observations

- Predicting future temperature (last lecture)
- Predicting the amount of rainfall
- Predicting the demand of a product
- Predicting the sale price of a house
- ...

# Regression

## Predicting a continuous outcome variable using past observations

- Predicting future temperature (last lecture)
- Predicting the amount of rainfall
- Predicting the demand of a product
- Predicting the sale price of a house
- ...

## Key difference from classification

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

# Regression

## Predicting a continuous outcome variable using past observations

- Predicting future temperature (last lecture)
- Predicting the amount of rainfall
- Predicting the demand of a product
- Predicting the sale price of a house
- ...

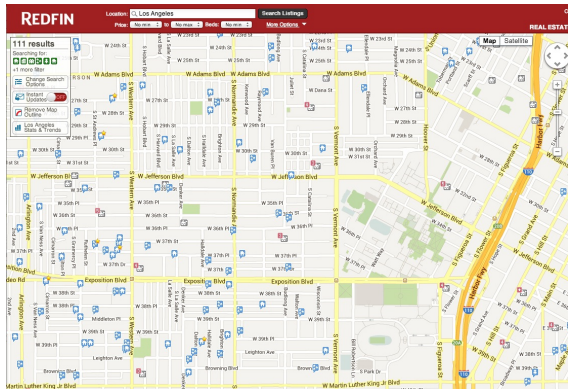
## Key difference from classification

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

**Linear Regression:** regression with linear models

# Ex: Predicting the sale price of a house

## Retrieve historical sales records (training data)



# Features used to predict

**3620 South BUDLONG**  
Los Angeles, CA 90007  
Status: Closed

**\$1,510,000**  
Last Sold Price


**14** Beds


**6** Baths

**4,418** Sq. Ft.  
\$347 / Sq. Ft.

Built: 1956 Lot Size: 9,849 Sq. Ft. Sold On: Jul 26, 2013

Overview Property Details Tour Insights Property History Public Records Activity Schools



1 of 12 

Five unit apartment complex within 2 blocks of USC campus, Gate #6. Great for students (most student leases have parents as guarantors). Most USC students live off campus, so housing units like this are always fully leased. Situated on a quiet, corner lot, and across from an elementary school, this complex was recently renovated, and has in-unit laundry hook ups, wall-unit AC, and 12 parking spaces. It is within a DPS (Department of Public Safety) and Campus Cruiser controlled area. This is a great income generating property, not to be missed!

Property Type **Multi-Family**

Community **Downtown Los Angeles**

MLS# **22176741**

Style **Two Level, Low Rise**

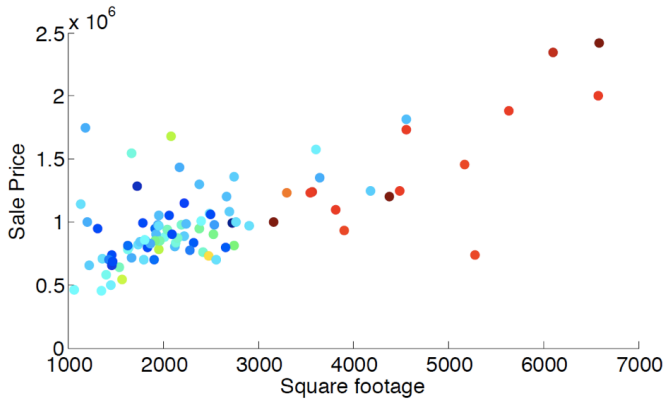
County **Los Angeles**

## Property Details for 3620 South BUDLONG, Los Angeles, CA 90007

Details provided by i-Tech MLS and may not match the public record. [Learn More](#)

Interior Features		
<b>Kitchen Information</b> <ul style="list-style-type: none"> <li>Remodeled</li> <li>Oven, Range</li> </ul>	<b>Laundry Information</b> <ul style="list-style-type: none"> <li>Inside Laundry</li> </ul>	<b>Heating &amp; Cooling</b> <ul style="list-style-type: none"> <li>Wall Cooling Unit(s)</li> </ul>
Multi-Unit Information		
<b>Community Features</b> <ul style="list-style-type: none"> <li>Units In Complex (Total): 5</li> </ul>	<b>Unit 2 Information</b> <ul style="list-style-type: none"> <li># of Beds: 3</li> <li># of Baths: 1</li> <li>Unfurnished</li> <li>Monthly Rent: \$2,250</li> </ul>	<ul style="list-style-type: none"> <li>Monthly Rent: \$2,300</li> </ul>
<b>Multi-Family Information</b> <ul style="list-style-type: none"> <li># Leased: 5</li> <li># of Buildings: 1</li> <li>Owner Pays Water</li> <li>Tenant Pays Electricity, Tenant Pays Gas</li> </ul>	<b>Unit 3 Information</b> <ul style="list-style-type: none"> <li>Unfurnished</li> </ul>	<b>Unit 5 Information</b> <ul style="list-style-type: none"> <li># of Beds: 3</li> <li># of Baths: 2</li> <li>Unfurnished</li> <li>Monthly Rent: \$2,325</li> </ul>
<b>Unit 1 Information</b> <ul style="list-style-type: none"> <li># of Beds: 2</li> <li># of Baths: 1</li> <li>Unfurnished</li> <li>Monthly Rent: \$1,700</li> </ul>	<b>Unit 4 Information</b> <ul style="list-style-type: none"> <li># of Beds: 3</li> <li># of Baths: 1</li> <li>Unfurnished</li> </ul>	<b>Unit 6 Information</b> <ul style="list-style-type: none"> <li># of Beds: 3</li> <li># of Baths: 1</li> <li>Monthly Rent: \$2,250</li> </ul>
Property / Lot Details		
<b>Property Features</b> <ul style="list-style-type: none"> <li>Automatic Gate, Card/Code Access</li> </ul>	<ul style="list-style-type: none"> <li>Automatic Gate, Lawn, Sidewalks</li> <li>Corner Lot, Near Public Transit</li> </ul>	<ul style="list-style-type: none"> <li>Tax Parcel Number: 5040017019</li> </ul>
<b>Lot Information</b> <ul style="list-style-type: none"> <li>Lot Size (Sq. Ft.): 9,849</li> <li>Lot Size (Acre): 0.2215</li> <li>Lot Size Source: Public Records</li> </ul>	<b>Property Information</b> <ul style="list-style-type: none"> <li>Updated/Remodeled</li> <li>Square Footage Source: Public Records</li> </ul>	
Parking / Garage, Exterior Features, Utilities & Financing		
<b>Parking Information</b> <ul style="list-style-type: none"> <li># of Parking Spaces (Total): 12</li> <li>Parking Space</li> <li>Gated</li> </ul>	<b>Utility Information</b> <ul style="list-style-type: none"> <li>Green Certification Rating: 0.00</li> <li>Green Location: Transportation, Walkability</li> <li>Green Walk Score: 0</li> <li>Green Year Certified: 0</li> </ul>	<b>Financial Information</b> <ul style="list-style-type: none"> <li>Capitalization Rate (%): 6.25</li> <li>Actual Annual Gross Rent: \$126,331</li> <li>Gross Rent Multiplier: 11.29</li> </ul>
<b>Building Information</b> <ul style="list-style-type: none"> <li>Total Floor: 2</li> </ul>		
Location Details, Misc. Information & Listing Information		
<b>Location Information</b> <ul style="list-style-type: none"> <li>Cross Streets: W 36th Pl</li> </ul>	<b>Expense Information</b> <ul style="list-style-type: none"> <li>Operating: \$37,664</li> </ul>	<b>Listing Information</b> <ul style="list-style-type: none"> <li>Listing Terms: Cash, Cash To Existing Loan</li> <li>Buyer Financing: Cash</li> </ul>

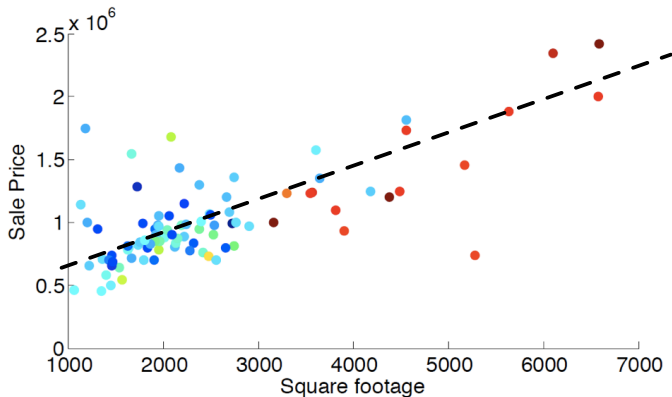
# Correlation between square footage and sale price





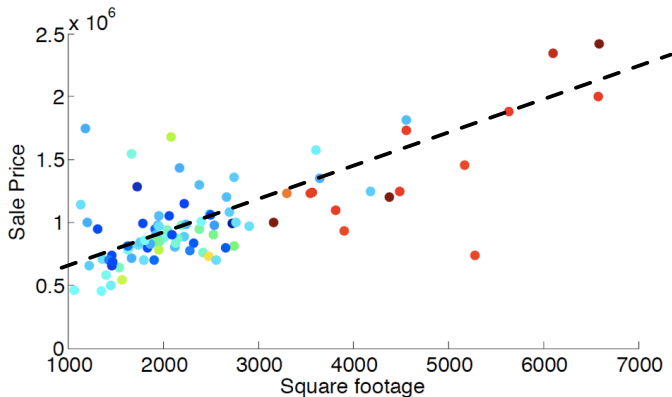
## Possibly linear relationship

Sale price  $\approx$  `price_per_sqft`  $\times$  square\_footage + `fixed_expense`



## Possibly linear relationship

Sale price  $\approx$  **price\_per\_sqft**  $\times$  square\_footage + **fixed\_expense**  
(*slope*) (*intercept*)



# 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.

# 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:  $| \text{prediction} - \text{sale price} |$

# 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**)

# 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,**

# 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?***

# 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*



# 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 ✓

## Example

Predicted price = **price\_per\_sqft** × square\_footage + **fixed\_expense**

one model: price\_per\_sqft = 0.3K, fixed\_expense = 210K

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

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

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

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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose),

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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

**NOTE:** for notation convenience, very often we

- append 1 to each  $x$  as the first feature:  $\tilde{\mathbf{x}} = [1 \ x_1 \ x_2 \ \dots \ x_D]^T$

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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

**NOTE:** for notation convenience, very often we

- append 1 to each  $x$  as the first feature:  $\tilde{\mathbf{x}} = [1 \ x_1 \ x_2 \ \dots \ x_D]^T$
- let  $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ , a concise representation of all  $D + 1$  parameters



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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

**NOTE:** for notation convenience, very often we

- append 1 to each  $x$  as the first feature:  $\tilde{\mathbf{x}} = [1 \ x_1 \ x_2 \ \dots \ x_D]^T$
- let  $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ , a concise representation of all  $D + 1$  parameters
- the model becomes simply  $f(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$

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

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$   
(superscript  $T$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

**NOTE:** for notation convenience, very often we

- append 1 to each  $x$  as the first feature:  $\tilde{\mathbf{x}} = [1 \ x_1 \ x_2 \ \dots \ x_D]^T$
- let  $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$ , a concise representation of all  $D + 1$  parameters
- the model becomes simply  $f(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$
- sometimes just use  $\mathbf{w}, \mathbf{x}, D$  for  $\tilde{\mathbf{w}}, \tilde{\mathbf{x}}, D + 1!$

# Goal

Minimize total squared error

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

# Goal

Minimize total squared error

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

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

# Goal

Minimize total squared error

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

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

- find  $\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}} \in \mathbb{R}^{D+1}}{\text{argmin}} \text{RSS}(\tilde{\mathbf{w}})$ , i.e. **least squares solution** (more generally called **empirical risk minimizer**)

# Goal

Minimize total squared error

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

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

- find  $\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}} \in \mathbb{R}^{D+1}}{\text{argmin}} \text{RSS}(\tilde{\mathbf{w}})$ , i.e. **least squares solution** (more generally called **empirical risk minimizer**)
- *reduce machine learning to optimization*

# Goal

Minimize total squared error

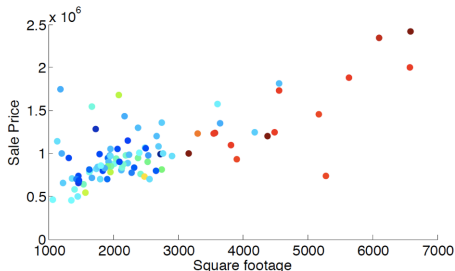
- **Residual Sum of Squares** (RSS), a function of  $\tilde{\mathbf{w}}$

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

- find  $\tilde{\mathbf{w}}^* = \underset{\tilde{\mathbf{w}} \in \mathbb{R}^{D+1}}{\text{argmin}} \text{RSS}(\tilde{\mathbf{w}})$ , i.e. **least squares solution** (more generally called **empirical risk minimizer**)
- *reduce machine learning to optimization*
- in principle can apply any optimization algorithm, but linear regression admits a *closed-form solution*

## Warm-up: $D = 0$

Only one parameter  $w_0$ : constant prediction  $f(x) = w_0$



$f$  is a horizontal line, where should it be?



## Warm-up: $D = 0$

### Optimization objective becomes

$$\text{RSS}(w_0) = \sum_n (w_0 - y_n)^2 \quad (\text{it's a } \textit{quadratic} \text{ } aw_0^2 + bw_0 + c)$$

## Warm-up: $D = 0$

### Optimization objective becomes

$$\begin{aligned} \text{RSS}(w_0) &= \sum_n (w_0 - y_n)^2 && \text{(it's a *quadratic* } aw_0^2 + bw_0 + c) \\ &= Nw_0^2 - 2 \left( \sum_n y_n \right) w_0 + \text{cnt.} \end{aligned}$$

## Warm-up: $D = 0$

### Optimization objective becomes

$$\begin{aligned}\text{RSS}(w_0) &= \sum_n (w_0 - y_n)^2 && \text{(it's a *quadratic* } aw_0^2 + bw_0 + c) \\ &= Nw_0^2 - 2 \left( \sum_n y_n \right) w_0 + \text{cnt.} \\ &= N \left( w_0 - \frac{1}{N} \sum_n y_n \right)^2 + \text{cnt.}\end{aligned}$$

## Warm-up: $D = 0$

### Optimization objective becomes

$$\begin{aligned}\text{RSS}(w_0) &= \sum_n (w_0 - y_n)^2 && \text{(it's a *quadratic* } aw_0^2 + bw_0 + c) \\ &= Nw_0^2 - 2 \left( \sum_n y_n \right) w_0 + \text{cnt.} \\ &= N \left( w_0 - \frac{1}{N} \sum_n y_n \right)^2 + \text{cnt.}\end{aligned}$$

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

## Warm-up: $D = 0$

### Optimization objective becomes

$$\begin{aligned}\text{RSS}(w_0) &= \sum_n (w_0 - y_n)^2 && \text{(it's a *quadratic* } aw_0^2 + bw_0 + c) \\ &= Nw_0^2 - 2 \left( \sum_n y_n \right) w_0 + \text{cnt.} \\ &= N \left( w_0 - \frac{1}{N} \sum_n y_n \right)^2 + \text{cnt.}\end{aligned}$$

It is clear 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?*

## Warm-up: $D = 1$

### Optimization objective becomes

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

## Warm-up: $D = 1$

### Optimization objective becomes

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

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

$$\left\{ \begin{array}{l} \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_0} = 0 \\ \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_1} = 0 \end{array} \right. \Rightarrow \begin{array}{l} \sum_n (w_0 + w_1 x_n - y_n) = 0 \\ \sum_n (w_0 + w_1 x_n - y_n) x_n = 0 \end{array}$$

## Warm-up: $D = 1$

### Optimization objective becomes

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

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

$$\begin{cases} \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_0} = 0 \\ \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_1} = 0 \end{cases} \Rightarrow \begin{cases} \sum_n (w_0 + w_1 x_n - y_n) = 0 \\ \sum_n (w_0 + w_1 x_n - y_n) x_n = 0 \end{cases}$$

$$\Rightarrow \begin{cases} N w_0 + w_1 \sum_n x_n = \sum_n y_n \\ w_0 \sum_n x_n + w_1 \sum_n x_n^2 = \sum_n y_n x_n \end{cases} \quad (\text{a linear system})$$



# Warm-up: $D = 1$

## Optimization objective becomes

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

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

$$\begin{cases} \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_0} = 0 \\ \frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial w_1} = 0 \end{cases} \Rightarrow \begin{cases} \sum_n (w_0 + w_1 x_n - y_n) = 0 \\ \sum_n (w_0 + w_1 x_n - y_n) x_n = 0 \end{cases}$$

$$\Rightarrow \begin{cases} N w_0 + w_1 \sum_n x_n = \sum_n y_n \\ w_0 \sum_n x_n + w_1 \sum_n x_n^2 = \sum_n y_n x_n \end{cases} \quad (\text{a linear system})$$

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

# Least square solution for $D = 1$

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

(assuming the matrix is invertible)

# Least square solution for $D = 1$

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

(assuming the matrix is invertible)

*Are stationary points minimizers?*

## Least square solution for $D = 1$

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

(assuming the matrix is invertible)

*Are stationary points minimizers?*

- yes for **convex** objectives (RSS is convex in  $\tilde{w}$ )

## Least square solution for $D = 1$

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

(assuming the matrix is invertible)

*Are stationary points minimizers?*

- yes for **convex** objectives (RSS is convex in  $\tilde{w}$ )
- not true in general

# General least square solution

## Objective

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

# General least square solution

## Objective

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

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

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

# General least square solution

## Objective

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

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

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



# General least square solution

## Objective

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

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

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

where

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

# General least square solution

## Objective

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

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

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

where

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

## General least square solution

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

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

## General least square solution

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

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

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

## General least square solution

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

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

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

**Verify the solution when  $D = 1$ :**

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

## General least square solution

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

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

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

**Verify the solution when  $D = 1$ :**

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

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

## Another approach

RSS is a **quadratic**, so let's complete the square:

$$\text{RSS}(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2 = \|\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}\|_2^2$$

## Another approach

RSS is a **quadratic**, so let's complete the square:

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



## Another approach

RSS is a **quadratic**, so let's complete the square:

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

## Another approach

RSS is a **quadratic**, so let's complete the square:

$$\begin{aligned}\text{RSS}(\tilde{\mathbf{w}}) &= \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2 = \|\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}\|_2^2 \\ &= (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}) \\ &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \mathbf{y} + \text{cnt.} \\ &= \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right) + \text{cnt.}\end{aligned}$$

## Another approach

RSS is a **quadratic**, so let's complete the square:

$$\begin{aligned}
 \text{RSS}(\tilde{\mathbf{w}}) &= \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2 = \|\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}\|_2^2 \\
 &= (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}) \\
 &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \mathbf{y} + \text{cnt.} \\
 &= \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right) + \text{cnt.}
 \end{aligned}$$

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

## Another approach

RSS is a **quadratic**, so let's complete the square:

$$\begin{aligned}
 \text{RSS}(\tilde{\mathbf{w}}) &= \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2 = \|\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}\|_2^2 \\
 &= (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}) \\
 &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \mathbf{y}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \mathbf{y} + \text{cnt.} \\
 &= \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \left( \tilde{\mathbf{w}} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y} \right) + \text{cnt.}
 \end{aligned}$$

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

So  $\tilde{\mathbf{w}}^* = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$  is the minimizer.

# Computational complexity

**Bottleneck** of computing

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

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

- naively need  $O(D^3)$  time

# Computational complexity

**Bottleneck** of computing

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

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

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

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

What does that imply?

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

**What does that imply?**

Recall  $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \mathbf{w}^* = \tilde{\mathbf{X}}^T \mathbf{y}$ .



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

## What does that imply?

Recall  $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \mathbf{w}^* = \tilde{\mathbf{X}}^T \mathbf{y}$ . If  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  not invertible, this equation has

- no solution

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

## What does that imply?

Recall  $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \mathbf{w}^* = \tilde{\mathbf{X}}^T \mathbf{y}$ . If  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  not invertible, this equation has

- no solution
- or infinitely many solutions

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

## What does that imply?

Recall  $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \mathbf{w}^* = \tilde{\mathbf{X}}^T \mathbf{y}$ . If  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  not invertible, this equation has

- no solution ( $\Rightarrow$  RSS has no minimizer?  $\times$ )
- or infinitely many solutions ( $\Rightarrow$  infinitely many minimizers  $\checkmark$ )

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

Why would that happen?

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

## Why would that happen?

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

# What if $\tilde{X}^T \tilde{X}$ 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

# What if $\tilde{X}^T \tilde{X}$ 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 this single point is a minimizer of RSS.

## How about the following?

$$D = 1, N = 2$$

sqft	sale price
1000	500K
1000	600K



## 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.

## 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

## 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?

**Intuition:** what does inverting  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  do?

**eigendecomposition:**  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{U}^T$

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

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

## How to resolve this issue?

**Intuition:** what does inverting  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  do?

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

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

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

*i.e. just invert the eigenvalues*

## How to solve this problem?

Non-invertible  $\Rightarrow$  some eigenvalues are 0.

## How to solve this problem?

Non-invertible  $\Rightarrow$  some eigenvalues are 0.

**One natural fix: add something positive**

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

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

## How to solve this problem?

Non-invertible  $\Rightarrow$  some eigenvalues are 0.

**One natural fix: add something positive**

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

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

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



## Fix the problem

The solution becomes

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

## Fix the problem

The solution becomes

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

- not a minimizer of the original RSS

## Fix the problem

The solution becomes

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

- not a minimizer of the original RSS
- more than an arbitrary hack (as we will see soon)

## Fix the problem

The solution becomes

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

- not a minimizer of the original RSS
- more than an arbitrary hack (as we will see soon)

$\lambda$  is a *hyper-parameter*, can be tuned by cross-validation.

# Comparison to NNC

## Non-parametric versus Parametric

- **Non-parametric methods:** the size of the model *grows* with the size of the training set.
  - e.g. NNC, the training set itself needs to be kept in order to predict. Thus, the size of the model is the size of the training set.

# Comparison to NNC

## Non-parametric versus Parametric

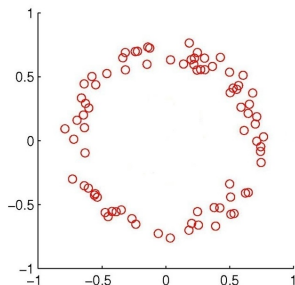
- **Non-parametric methods:** the size of the model *grows* with the size of the training set.
  - e.g. NNC, the training set itself needs to be kept in order to predict. Thus, the size of the model is the size of the training set.
- **Parametric methods:** the size of the model does *not grow* with the size of the training set  $N$ .
  - e.g. linear regression,  $D + 1$  parameters, independent of  $N$ .

# Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis**
- 3 Overfitting and preventing overfitting
- 4 A Detour of Numerical Optimization Methods

# What if linear model is not a good fit?

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





## Solution: nonlinearly transformed features

### 1. Use a nonlinear mapping

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

to transform the data to a more complicated feature space

## Solution: nonlinearly transformed features

### 1. Use a nonlinear mapping

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

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).

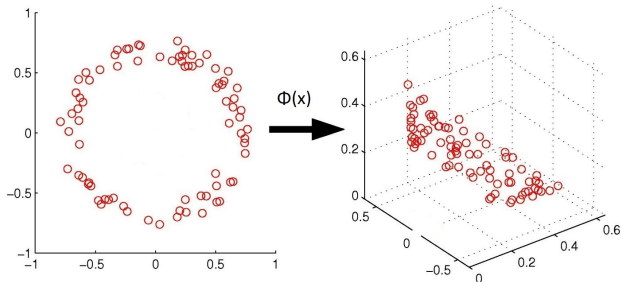
# Solution: nonlinearly transformed features

## 1. Use a nonlinear mapping

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

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

# Regression with nonlinear basis

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

**Objective:**

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

# Regression with nonlinear basis

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

**Objective:**

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

**Similar least square solution:**

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad \text{where} \quad \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

# Example

## Polynomial basis functions for $D = 1$

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

# Example

## Polynomial basis functions for $D = 1$

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

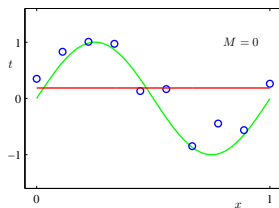
Learning a linear model in the new space

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



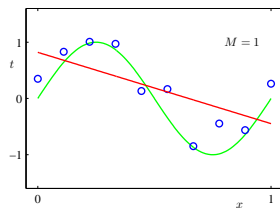
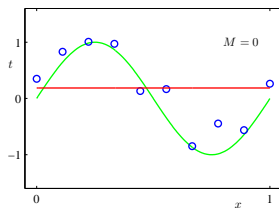
# Example

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



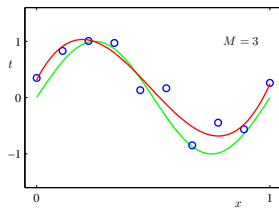
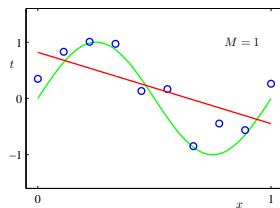
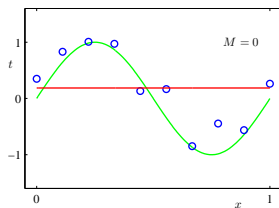
# Example

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



# Example

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



## Why nonlinear?

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

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

## Why nonlinear?

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

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

No, it basically *does nothing* since

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

## Why nonlinear?

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

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

No, it basically *does nothing* since

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

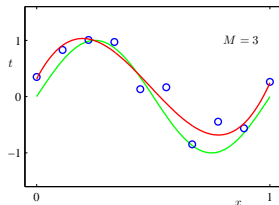
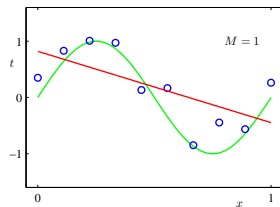
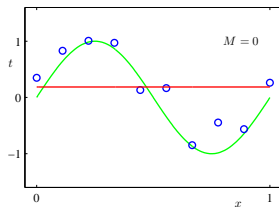
We will see more nonlinear mappings soon.

# Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and preventing overfitting**
- 4 A Detour of Numerical Optimization Methods

# Should we use a very complicated mapping?

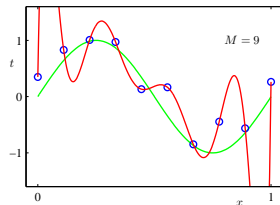
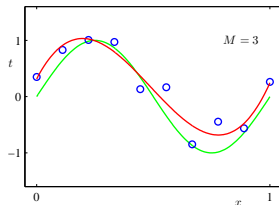
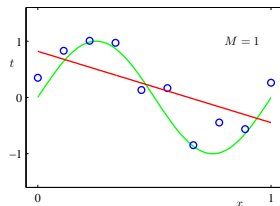
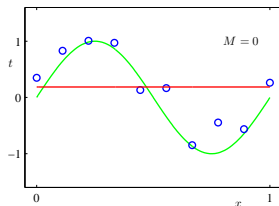
Ex: fitting a noisy sine function with a polynomial:





# Should we use a very complicated mapping?

Ex: fitting a noisy sine function with a polynomial:



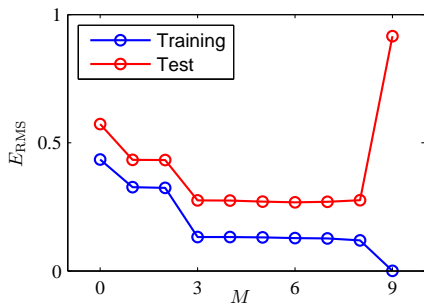
# Underfitting and Overfitting

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

- large training error
- large test error

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

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



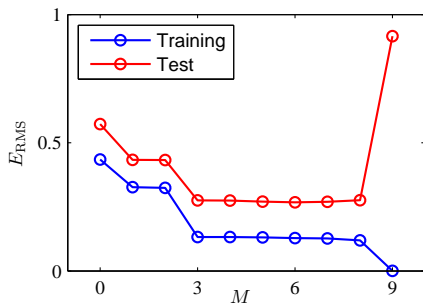
# Underfitting and Overfitting

$M \leq 2$  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*

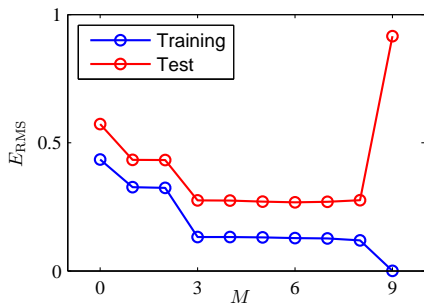
# Underfitting and Overfitting

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

- large training error
- large test error

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

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

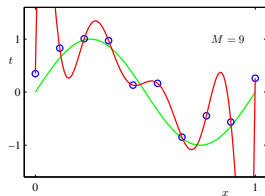


*More complicated models  $\Rightarrow$  larger gap between training and test error*

How to prevent overfitting?

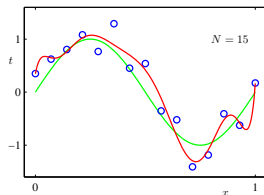
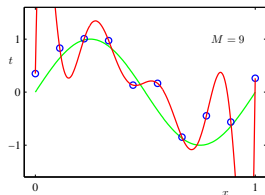
# Method 1: use more training data

The more, the merrier



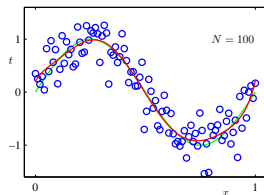
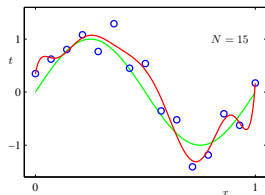
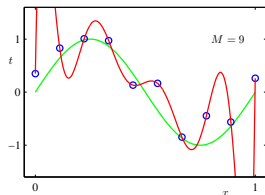
# Method 1: use more training data

The more, the merrier



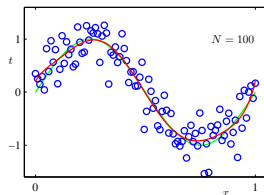
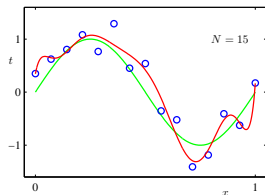
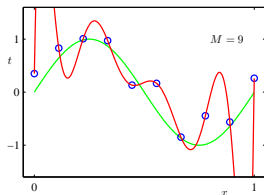
# Method 1: use more training data

The more, the merrier



# Method 1: use more training data

## The more, the merrier



*More data  $\Rightarrow$  smaller gap between training and test error*



## Method 2: control the model complexity

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

- use cross-validation to pick hyper-parameter  $M$

## Method 2: control the model complexity

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

- use cross-validation to pick hyper-parameter  $M$

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

# Magnitude of weights

Least square solution for the polynomial example:

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

## Magnitude of weights

Least square solution for the polynomial example:

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

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

## How to make $w$ small?

**Regularized linear regression:** new objective

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

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

## How to make $w$ small?

**Regularized linear regression:** new objective

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

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

- $R : \mathbb{R}^D \rightarrow \mathbb{R}^+$  is the *regularizer*
  - measure how complex the model  $w$  is, penalize complex models
  - common choices:  $\|w\|_2^2$ ,  $\|w\|_1$ , etc.

## How to make $w$ small?

**Regularized linear regression:** new objective

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

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

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

# The effect of $\lambda$

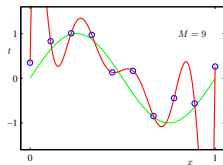
when we increase regularization coefficient  $\lambda$

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0$	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



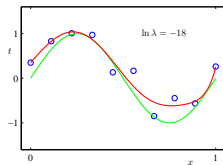
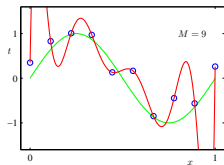
# The trade-off

when we increase regularization coefficient  $\lambda$



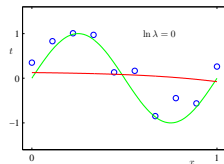
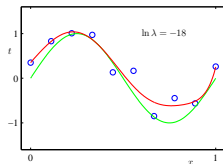
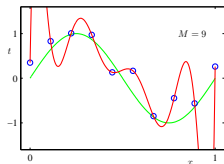
# The trade-off

when we increase regularization coefficient  $\lambda$



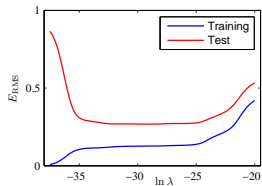
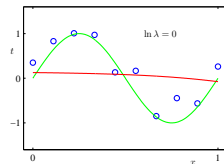
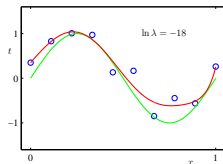
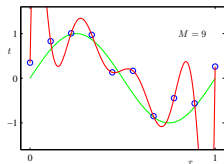
# The trade-off

when we increase regularization coefficient  $\lambda$



# The trade-off

when we increase regularization coefficient  $\lambda$



## How to solve the new objective?

**Simple for**  $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ :

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

## How to solve the new objective?

**Simple for**  $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ :

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

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

## How to solve the new objective?

**Simple for**  $R(\mathbf{w}) = \|\mathbf{w}\|_2^2$ :

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

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

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

## How to solve the new objective?

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

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

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

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

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



## How to solve the new objective?

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

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

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

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

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

*Note the same form as in the fix when  $\mathbf{X}^T\mathbf{X}$  is not invertible!*

## How to solve the new objective?

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

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

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

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

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

*Note the same form as in the fix when  $\mathbf{X}^T\mathbf{X}$  is not invertible!*

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

## Equivalent form

Regularization is also sometimes formulated as

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

where  $\beta$  is some hyper-parameter.

## Equivalent form

Regularization is also sometimes formulated as

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

where  $\beta$  is some hyper-parameter.

Finding the solution becomes a *constrained optimization problem*.

## Equivalent form

Regularization is also sometimes formulated as

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

where  $\beta$  is some hyper-parameter.

Finding the solution becomes a *constrained optimization problem*.

Choosing either  $\lambda$  or  $\beta$  can be done by cross-validation.

# Summary

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

# Summary

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

*Important to understand the derivation than remembering the formula*

# Summary

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

*Important to understand the derivation than remembering the formula*

**Overfitting:** small training error but large test error



# Summary

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

*Important to understand the derivation than remembering the formula*

**Overfitting:** small training error but large test error

**Preventing Overfitting:** more data + regularization

## Recall the question

**Typical steps** of developing a machine learning system:

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

How to do the *red part* exactly?

# General idea to derive 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.  $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^M\}$

## General idea to derive 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**  $L(y', y)$

## General idea to derive 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)**:

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

## General idea to derive ML algorithms

- 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.  $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^M\}$
- Define **error/loss**  $L(y', y)$
- 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**:

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

## General idea to derive ML algorithms

- 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.  $\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) \mid \mathbf{w} \in \mathbb{R}^M\}$
- Define **error/loss**  $L(y', y)$
- 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*

# Outline

- 1 Linear regression
- 2 Linear regression with nonlinear basis
- 3 Overfitting and preventing overfitting
- 4 A Detour of Numerical Optimization Methods**



# Numerical optimization

## Problem setup

- Given: a function  $F(\mathbf{w})$
- Goal: minimize  $F(\mathbf{w})$  (approximately)

# First-order optimization methods

Two simple yet extremely popular methods

- **Gradient Descent (GD)**: simple and fundamental
- **Stochastic Gradient Descent (SGD)**: faster, effective for large-scale problems

# First-order optimization methods

Two simple yet extremely popular methods

- **Gradient Descent (GD)**: simple and fundamental
- **Stochastic Gradient Descent (SGD)**: faster, effective for large-scale problems

Gradient is sometimes referred to as *first-order* information of a function. Therefore, these methods are called *first-order methods*.

# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where  $\eta > 0$  is called step size or learning rate

# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where  $\eta > 0$  is called step size or learning rate

- in theory  $\eta$  should be set in terms of some parameters of  $F$

# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where  $\eta > 0$  is called step size or learning rate

- in theory  $\eta$  should be set in terms of some parameters of  $F$
- in practice we just try several small values

# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where  $\eta > 0$  is called step size or learning rate

- in theory  $\eta$  should be set in terms of some parameters of  $F$
- in practice we just try several small values
- might need to be **changing** over iterations (think  $F(w) = |w|$ )



# Gradient Descent (GD)

**GD**: keep moving in the *negative gradient direction*

Start from some  $\mathbf{w}^{(0)}$ . For  $t = 0, 1, 2, \dots$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)})$$

where  $\eta > 0$  is called step size or learning rate

- in theory  $\eta$  should be set in terms of some parameters of  $F$
- in practice we just try several small values
- might need to be **changing** over iterations (think  $F(w) = |w|$ )
- adaptive and automatic step size tuning is an active research area

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ .

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ . Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \qquad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ . Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize  $w_1^{(0)}$  and  $w_2^{(0)}$  (to be 0 or *randomly*),  $t = 0$

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ . Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize  $w_1^{(0)}$  and  $w_2^{(0)}$  (to be 0 or *randomly*),  $t = 0$
- do

$$w_1^{(t+1)} \leftarrow w_1^{(t)} - \eta \left[ 2(w_1^{(t)2} - w_2^{(t)})w_1^{(t)} + w_1^{(t)} - 1 \right]$$

$$w_2^{(t+1)} \leftarrow w_2^{(t)} - \eta \left[ -(w_1^{(t)2} - w_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

## An example

Example:  $F(\mathbf{w}) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$ . Gradient is

$$\frac{\partial F}{\partial w_1} = 2(w_1^2 - w_2)w_1 + w_1 - 1 \quad \frac{\partial F}{\partial w_2} = -(w_1^2 - w_2)$$

GD:

- Initialize  $w_1^{(0)}$  and  $w_2^{(0)}$  (to be 0 or *randomly*),  $t = 0$
- do

$$w_1^{(t+1)} \leftarrow w_1^{(t)} - \eta \left[ 2(w_1^{(t)2} - w_2^{(t)})w_1^{(t)} + w_1^{(t)} - 1 \right]$$

$$w_2^{(t+1)} \leftarrow w_2^{(t)} - \eta \left[ -(w_1^{(t)2} - w_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

- until  $F(\mathbf{w}^{(t)})$  **does not change much** or  $t$  **reaches a fixed number**

# Why GD?

Intuition: by first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

## Why GD?

Intuition: by first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

GD ensures

$$F(\mathbf{w}^{(t+1)}) \approx F(\mathbf{w}^{(t)}) - \eta \|\nabla F(\mathbf{w}^{(t)})\|_2^2 \leq F(\mathbf{w}^{(t)})$$



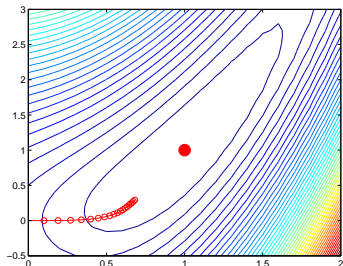
# Why GD?

Intuition: by first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

GD ensures

$$F(\mathbf{w}^{(t+1)}) \approx F(\mathbf{w}^{(t)}) - \eta \|\nabla F(\mathbf{w}^{(t)})\|_2^2 \leq F(\mathbf{w}^{(t)})$$



reasonable  $\eta$  decreases function value

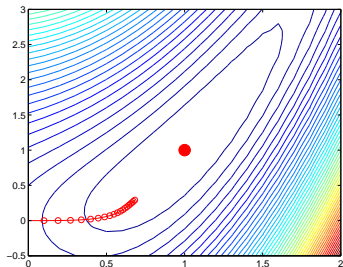
# Why GD?

Intuition: by first-order **Taylor approximation**

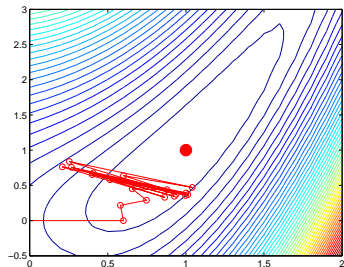
$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

GD ensures

$$F(\mathbf{w}^{(t+1)}) \approx F(\mathbf{w}^{(t)}) - \eta \|\nabla F(\mathbf{w}^{(t)})\|_2^2 \leq F(\mathbf{w}^{(t)})$$



reasonable  $\eta$  decreases function value



but large  $\eta$  is unstable

# Stochastic Gradient Descent (SGD)

GD: keep moving in the negative gradient direction

**SGD**: keep moving in some *noisy* negative gradient direction

# Stochastic Gradient Descent (SGD)

GD: keep moving in the negative gradient direction

**SGD**: keep moving in some *noisy* negative gradient direction

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where  $\tilde{\nabla} F(\mathbf{w}^{(t)})$  is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} \left[ \tilde{\nabla} F(\mathbf{w}^{(t)}) \right] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

# Stochastic Gradient Descent (SGD)

GD: keep moving in the negative gradient direction

**SGD**: keep moving in some *noisy* negative gradient direction

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \tilde{\nabla} F(\mathbf{w}^{(t)})$$

where  $\tilde{\nabla} F(\mathbf{w}^{(t)})$  is a random variable (called **stochastic gradient**) s.t.

$$\mathbb{E} \left[ \tilde{\nabla} F(\mathbf{w}^{(t)}) \right] = \nabla F(\mathbf{w}^{(t)}) \quad (\text{unbiasedness})$$

Key point: it could be *much faster to obtain a stochastic gradient!*  
(examples coming soon)

# Convergence guarantees — convex objectives

**Many** for both GD and SGD on **convex objectives**.

# Convergence guarantees — convex objectives

**Many** for both GD and SGD on **convex objectives**.

They tell you how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \epsilon$$

# Convergence guarantees — convex objectives

**Many** for both GD and SGD on **convex objectives**.

They tell you how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \epsilon$$

- usually SGD needs more iterations



# Convergence guarantees — convex objectives

**Many** for both GD and SGD on **convex objectives**.

They tell you how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$F(\mathbf{w}^{(t)}) - F(\mathbf{w}^*) \leq \epsilon$$

- usually SGD needs more iterations
- but then again each iteration takes less time

## Convergence guarantees — nonconvex objectives

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$\|\nabla F(\mathbf{w}^{(t)})\| \leq \epsilon$$

## Convergence guarantees — nonconvex objectives

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$\|\nabla F(\mathbf{w}^{(t)})\| \leq \epsilon$$

- that is, how close  $\mathbf{w}^{(t)}$  is as an **approximate stationary point**

## Convergence guarantees — nonconvex objectives

Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$\|\nabla F(\mathbf{w}^{(t)})\| \leq \epsilon$$

- that is, how close  $\mathbf{w}^{(t)}$  is as an **approximate stationary point**
- for convex objectives, stationary point  $\Rightarrow$  global minimizer

## Convergence guarantees — nonconvex objectives

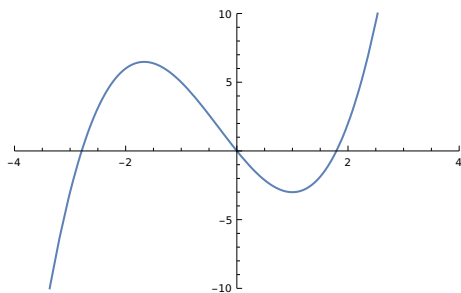
Even for *nonconvex objectives*, some guarantees exist: e.g. how many iterations  $t$  (in terms of  $\epsilon$ ) needed to achieve

$$\|\nabla F(\mathbf{w}^{(t)})\| \leq \epsilon$$

- that is, how close  $\mathbf{w}^{(t)}$  is as an **approximate stationary point**
- for convex objectives, stationary point  $\Rightarrow$  global minimizer
- for nonconvex objectives, *what does it mean?*

# Convergence guarantees — nonconvex objectives

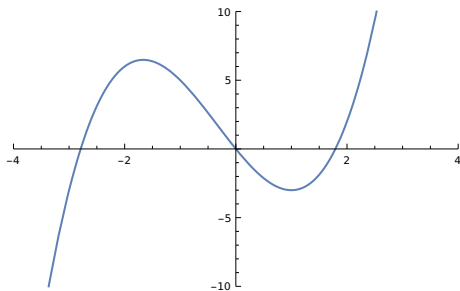
A stationary point can be a **local minimizer**



$$f(w) = w^3 + w^2 - 5w$$

## Convergence guarantees — nonconvex objectives

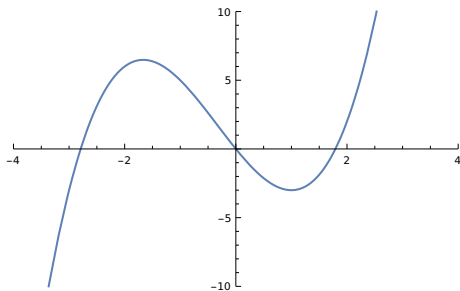
A stationary point can be a **local minimizer** or even a **local/global maximizer**



$$f(w) = w^3 + w^2 - 5w$$

## Convergence guarantees — nonconvex objectives

A stationary point can be a **local minimizer** or even a **local/global maximizer** (but the latter is not an issue for GD/SGD).



$$f(w) = w^3 + w^2 - 5w$$



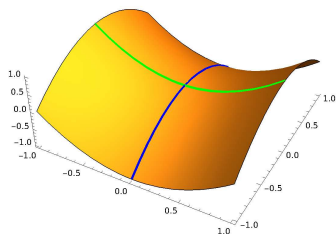
## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

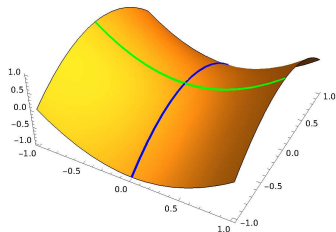
- $f(\mathbf{w}) = w_1^2 - w_2^2$



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

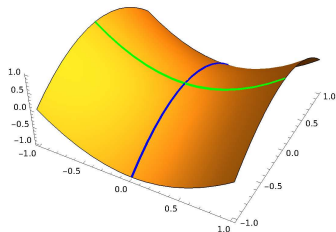
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

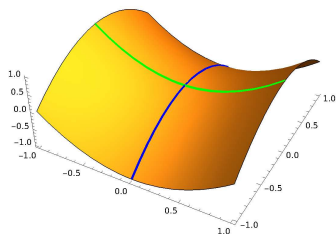
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

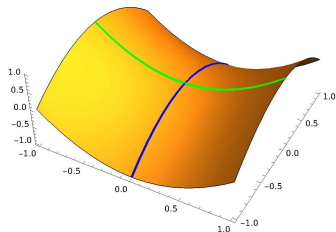
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for **blue direction** ( $w_1 = 0$ )



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!*

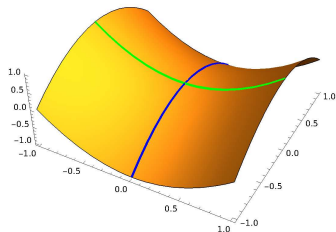
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for **blue direction** ( $w_1 = 0$ )
- local min for **green direction** ( $w_2 = 0$ )



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!* This is called a **saddle point**.

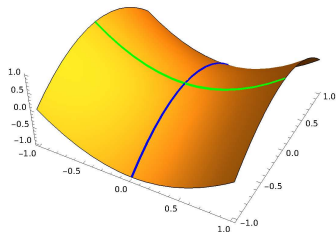
- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for **blue direction** ( $w_1 = 0$ )
- local min for **green direction** ( $w_2 = 0$ )



## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!* This is called a **saddle point**.

- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for **blue direction** ( $w_1 = 0$ )
- local min for **green direction** ( $w_2 = 0$ )
- but GD gets stuck at  $(0, 0)$  only if initialized along the **green direction**

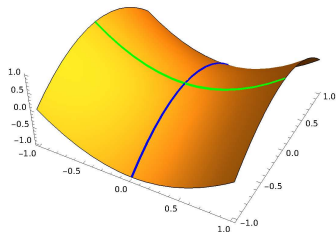




## Convergence guarantees — nonconvex objectives

A stationary point can also be *neither a local minimizer nor a local maximizer!* This is called a **saddle point**.

- $f(\mathbf{w}) = w_1^2 - w_2^2$
- $\nabla f(\mathbf{w}) = (2w_1, -2w_2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- local max for **blue direction** ( $w_1 = 0$ )
- local min for **green direction** ( $w_2 = 0$ )
- but GD gets stuck at  $(0, 0)$  only if initialized along the **green direction**
- so not a real issue especially *when initialized randomly*



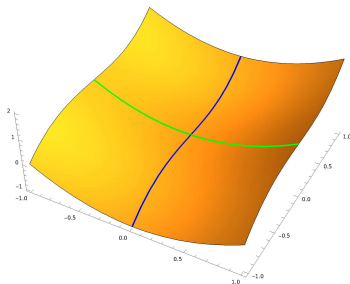
## Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

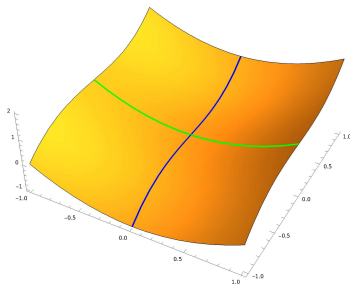
- $f(\mathbf{w}) = w_1^2 + w_2^3$



# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

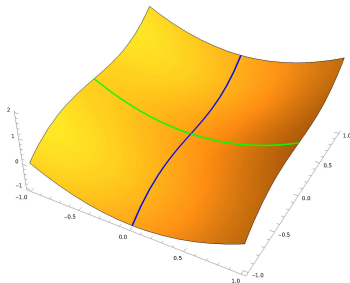
- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$



# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

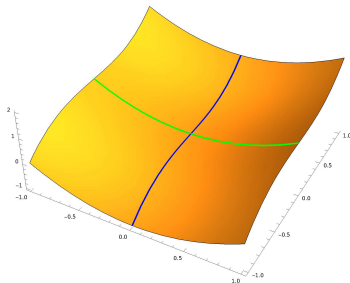
- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary



# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

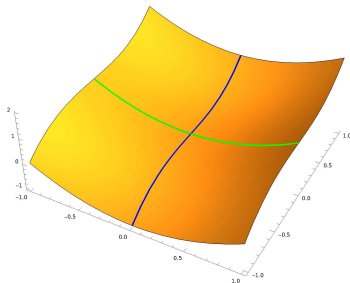
- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- not local min/max for **blue direction**  
( $w_1 = 0$ )



# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

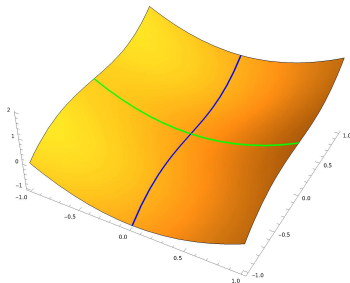
- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- not local min/max for **blue direction** ( $w_1 = 0$ )
- GD gets stuck at  $(0, 0)$  for *any initial point with  $w_2 \geq 0$  and small  $\eta$*



# Convergence guarantees — nonconvex objectives

But not all saddle points look like a “saddle” ...

- $f(\mathbf{w}) = w_1^2 + w_2^3$
- $\nabla f(\mathbf{w}) = (2w_1, 3w_2^2)$
- so  $\mathbf{w} = (0, 0)$  is stationary
- not local min/max for **blue direction** ( $w_1 = 0$ )
- GD gets stuck at  $(0, 0)$  for *any initial point with  $w_2 \geq 0$  and small  $\eta$*



Even worse, distinguishing local min and saddle point is generally *NP-hard*.



# Convergence guarantees

## Summary:

- GD/SGD converges to a stationary point

# Convergence guarantees

## Summary:

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need

# Convergence guarantees

## Summary:

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or “bad” saddle points (random initialization escapes “good” saddle points)

# Convergence guarantees

## Summary:

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or “bad” saddle points (random initialization escapes “good” saddle points)
- recent research shows that *many problems have no “bad” saddle points or even “bad” local minimizers*

# Convergence guarantees

## Summary:

- GD/SGD converges to a stationary point
- for convex objectives, this is all we need
- for nonconvex objectives, can get stuck at local minimizers or “bad” saddle points (random initialization escapes “good” saddle points)
- recent research shows that *many problems have no “bad” saddle points or even “bad” local minimizers*
- justify the practical effectiveness of GD/SGD (default method to try)

## Second-order methods

Recall the intuition of GD: we look at first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

## Second-order methods

Recall the intuition of GD: we look at first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

What if we look at *second-order* Taylor approximation?

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}_t (\mathbf{w} - \mathbf{w}^{(t)})$$

## Second-order methods

Recall the intuition of GD: we look at first-order **Taylor approximation**

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)})$$

What if we look at *second-order* Taylor approximation?

$$F(\mathbf{w}) \approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}_t (\mathbf{w} - \mathbf{w}^{(t)})$$

where  $\mathbf{H}_t = \nabla^2 F(\mathbf{w}^{(t)}) \in \mathbb{R}^{D \times D}$  is the *Hessian* of  $F$  at  $\mathbf{w}^{(t)}$ , i.e.,

$$H_{t,ij} = \left. \frac{\partial^2 F(\mathbf{w})}{\partial w_i \partial w_j} \right|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

(think “second derivative” when  $D = 1$ )



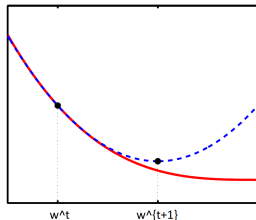
## Newton method

If we minimize the second-order approximation (via “complete the square”)

$$F(\mathbf{w})$$

$$\approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}_t (\mathbf{w} - \mathbf{w}^{(t)})$$

$$= \frac{1}{2} \left( \mathbf{w} - \mathbf{w}^{(t)} + \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \right)^T \mathbf{H}_t \left( \mathbf{w} - \mathbf{w}^{(t)} + \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \right) + \text{const.}$$



## Newton method

If we minimize the second-order approximation (via “complete the square”)

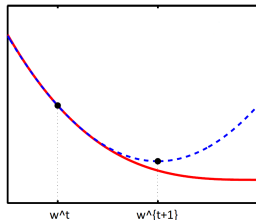
$$F(\mathbf{w})$$

$$\approx F(\mathbf{w}^{(t)}) + \nabla F(\mathbf{w}^{(t)})^T (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}_t (\mathbf{w} - \mathbf{w}^{(t)})$$

$$= \frac{1}{2} \left( \mathbf{w} - \mathbf{w}^{(t)} + \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \right)^T \mathbf{H}_t \left( \mathbf{w} - \mathbf{w}^{(t)} + \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \right) + \text{const.}$$

for convex  $F$  (so  $H_t$  is *positive semidefinite*)  
we obtain **Newton method**:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)})$$



## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures,

## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\eta$  (so **no tuning needed!**)

## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\eta$  (so **no tuning needed!**)
- converges **super fast** in terms of #iterations (for convex objectives)

# Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\eta$  (so **no tuning needed!**)
- converges **super fast** in terms of #iterations (for convex objectives)
  - e.g. how many iterations needed when applied to a quadratic?

## Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\eta$  (so **no tuning needed!**)
- converges **super fast** in terms of #iterations (for convex objectives)
  - e.g. how many iterations needed when applied to a quadratic?
- computing Hessian in each iteration is **very slow** though

# Comparing GD and Newton

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla F(\mathbf{w}^{(t)}) \quad (\text{GD})$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)}) \quad (\text{Newton})$$

Both are iterative optimization procedures, but Newton method

- has no learning rate  $\eta$  (so **no tuning needed!**)
- converges **super fast** in terms of #iterations (for convex objectives)
  - e.g. how many iterations needed when applied to a quadratic?
- computing Hessian in each iteration is *very slow* though
- does not really make sense for *nonconvex objectives* (but generally Hessian can be useful for escaping saddle points)