CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Jan 22, 2021









Outline



2 Review of last lecture

3 Linear regression

Logistics

- HW 0 is due today.
- HW 1 will be released today.
- Will be releasing the schedule of lectures.

Outline







Multi-class classification

Training data (set)

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \cdots, (\boldsymbol{x}_N, y_N)\}$
- Each $x_n \in \mathbb{R}^{\mathsf{D}}$ is called a feature vector.
- Each $y_n \in [C] = \{1, 2, \cdots, C\}$ is called a label/class/category.
- They are used to learn $f : \mathbb{R}^{D} \to [C]$ for future prediction.

Special case: binary classification

- Number of classes: C = 2
- \bullet Conventional labels: $\{0,1\}$ or $\{-1,+1\}$

K-NNC: predict the majority label within the K-nearest neighbor set

Datasets

Training data

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \cdots, (\boldsymbol{x}_N, y_N)\}$
- They are used to learn $f(\cdot)$

Test data

- M samples/instances: $\mathcal{D}^{\text{TEST}} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \cdots, (\boldsymbol{x}_{\mathsf{M}}, y_{\mathsf{M}})\}$
- They are used to evaluate how well $f(\cdot)$ will do.

Development/Validation data

- L samples/instances: $\mathcal{D}^{\text{DEV}} = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \cdots, (\boldsymbol{x}_{\mathsf{L}}, y_{\mathsf{L}})\}$
- They are used to optimize hyper-parameter(s).

These three sets should *not* overlap!

S-fold Cross-validation

What if we do not have a development set?

- Split the training data into S equal parts.
- Use each part *in turn* as a development dataset and use the others as a training dataset.
- Choose the hyper-parameter leading to best *average* performance.

S = 5: 5-fold cross validation



Special case: S = N, called leave-one-out.

High level picture

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.

High level picture

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?

Outline



Review of last lecture

3 Linear regression

- Motivation
- Setup and Algorithm
- Discussions

Regression

Predicting a continuous outcome variable using past observations

- Predicting future temperature (lecture 1)
- 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 (lecture 1)
- 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.

Motivation

Regression

Predicting a continuous outcome variable using past observations

- Predicting future temperature (lecture 1)
- 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



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		
Kitchen Information • Remodeled • Oven, Range	Laundry Information Inside Laundry 	Heating & Cooling + Wall Cooling Unit(s)
Multi-Unit Information		
Community Features Utitis in Concern (Front) 5 Multi-Fundy Information Lassing Information = e of Buildings: 1 - Onane Page Weber - Tream Page Medin - Tream Page Medin - e of Basis: 2 = e of Basis: 2 - Utifumstand - Worthy Rev 51,700	List 2 Information • # of Basis: 3 • # of Basis: 3 • Automated • Nontrity Resc 52,250 List 3 Information • Ust Information • List Information • # of Basis: 3 • # of Basis: 1 • ListInstand	Monthly Rest: \$2,260 Units Enformation # of Bolds: 3 # of Bolds: 3 Unit/unitwo. Unit/unitwo. Unit/unitwo. # of Bolds: 3 # of Bolds: 4 # o
Property / Lot Details		
Property Features • Automatic Gate, Card/Code Access Lot Information • Lot Size (Sci. Pc.): 9,849 • Lot Size (Acres): 0.2215 • Lot Size Source: Public Records	Automatic Gate, Lawn, Sidewalks Corner Lot, Near Public Transit Property Information Lopated Remodeled Square Footage Source: Public Records	Tax Parcel Number: 5040017019
Parking / Garage, Exterior Features, Utilities & I	Tinancing	
Parking Information • # of Parking Spaces (Total): 12 • Parking Space • Gated Building Information • Total Floors 2	Utility Information • Green Certification Rating: 0.00 • Green Location: Transportation, Walkability • Green Walk Score: 0 • Green Year Certified: 0	Financial Information • Copitalization Rate (%): 8.25 • Actual Annual Gross Rent: \$128,331 • Gross Rent Multiplier: 11.29
Location Details, Misc. Information & Listing Inf	ormation	
Location Information Cross Streets: W 35th Pl	Expense Information • Operating: \$37,664	Listing Information • Listing Terms: Cash, Cash To Existing Listing Figure Figure Cash

Correlation between square footage and sale price



Possibly linear relationship

Sale price \approx price_per_sqft \times square_footage + fixed_expense



Possibly linear relationship



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 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: | prediction sale price |

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: | prediction sale price |
 - or *squared* error: (prediction sale price)² (most common)

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: | prediction sale price |
 - or *squared* error: (prediction sale price)² (most common)

Goal: pick the model (unknown parameters) that minimizes the average/total prediction error,

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: | prediction sale price |
 - or *squared* error: (prediction sale price)² (most common)

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

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: | prediction sale price |
 - or *squared* error: (prediction 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 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: | prediction sale price |
 - or *squared* error: (prediction 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			$0 + 67^2 + 228^2 + 740^2 + \cdots$

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

Setup and Algorithm

Formal setup for linear regression

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) **Output**: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) **Training data**: $D = \{(x_n, y_n), n = 1, 2, ..., N\}$

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

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

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

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

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

Linear model: $f : \mathbb{R}^{D} \to \mathbb{R}$, with $f(\boldsymbol{x}) = w_{0} + \sum_{d=1}^{D} w_{d}x_{d} = w_{0} + \boldsymbol{w}^{T}\boldsymbol{x}$ (superscript T stands for transpose), i.e. a *hyper-plane* parametrized by • $\boldsymbol{w} = [w_{1} \ w_{2} \ \cdots \ w_{D}]^{T}$ (weights, weight vector, parameter vector, etc) • bias w_{0}

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

Linear model: $f : \mathbb{R}^{D} \to \mathbb{R}$, with $f(\boldsymbol{x}) = w_{0} + \sum_{d=1}^{D} w_{d}x_{d} = w_{0} + \boldsymbol{w}^{T}\boldsymbol{x}$ (superscript T stands for transpose), i.e. a *hyper-plane* parametrized by • $\boldsymbol{w} = [w_{1} \ w_{2} \ \cdots \ w_{D}]^{T}$ (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: $\tilde{x} = [1 \ x_1 \ x_2 \ \dots \ x_D]^T$

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

Linear model: $f : \mathbb{R}^{D} \to \mathbb{R}$, with $f(\boldsymbol{x}) = w_{0} + \sum_{d=1}^{D} w_{d}x_{d} = w_{0} + \boldsymbol{w}^{T}\boldsymbol{x}$ (superscript T stands for transpose), i.e. a *hyper-plane* parametrized by • $\boldsymbol{w} = [w_{1} \ w_{2} \ \cdots \ w_{D}]^{T}$ (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: $\tilde{x} = [1 \ x_1 \ x_2 \ \dots \ x_{\mathsf{D}}]^{\mathrm{T}}$
- let $\tilde{w} = [w_0 \ w_1 \ w_2 \ \cdots \ w_D]^T$, a concise representation of all D+1 parameters

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

Linear model: $f : \mathbb{R}^{D} \to \mathbb{R}$, with $f(\boldsymbol{x}) = w_{0} + \sum_{d=1}^{D} w_{d}x_{d} = w_{0} + \boldsymbol{w}^{T}\boldsymbol{x}$ (superscript T stands for transpose), i.e. a *hyper-plane* parametrized by • $\boldsymbol{w} = [w_{1} \ w_{2} \ \cdots \ w_{D}]^{T}$ (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: $\tilde{x} = [1 \ x_1 \ x_2 \ \dots \ x_D]^{\mathrm{T}}$
- let $\tilde{w} = [w_0 \ w_1 \ w_2 \ \cdots \ w_D]^T$, a concise representation of all D+1 parameters
- the model becomes simply $f(\boldsymbol{x}) = \tilde{\boldsymbol{w}}^{\mathbf{T}} \tilde{\boldsymbol{x}}$

Input: $x \in \mathbb{R}^{D}$ (features, covariates, context, predictors, etc) Output: $y \in \mathbb{R}$ (responses, targets, outcomes, etc) Training data: $\mathcal{D} = \{(x_n, y_n), n = 1, 2, ..., N\}$

Linear model: $f : \mathbb{R}^{D} \to \mathbb{R}$, with $f(\boldsymbol{x}) = w_{0} + \sum_{d=1}^{D} w_{d}x_{d} = w_{0} + \boldsymbol{w}^{T}\boldsymbol{x}$ (superscript T stands for transpose), i.e. a *hyper-plane* parametrized by • $\boldsymbol{w} = [w_{1} \ w_{2} \ \cdots \ w_{D}]^{T}$ (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: $\tilde{\boldsymbol{x}} = [1 \ x_1 \ x_2 \ \dots \ x_{\mathsf{D}}]^{\mathrm{T}}$
- let $\tilde{w} = [w_0 \ w_1 \ w_2 \ \cdots \ w_D]^T$, a concise representation of all D+1 parameters
- the model becomes simply $f(\boldsymbol{x}) = \tilde{\boldsymbol{w}}^{\mathbf{T}} \tilde{\boldsymbol{x}}$
- sometimes just use ${m w},{m x},{\sf D}$ for ${m ilde w},{m ilde x},{\sf D}+1!$

Goal

Minimize total squared error

$$\sum_{n} (f(\boldsymbol{x}_n) - y_n)^2 = \sum_{n} (\tilde{\boldsymbol{x}}_n^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_n)^2$$
Minimize total squared error

• Residual Sum of Squares (RSS), a function of $ilde{w}$

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (f(\boldsymbol{x}_{n}) - y_{n})^{2} = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Minimize total squared error

• Residual Sum of Squares (RSS), a function of \tilde{w}

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (f(\boldsymbol{x}_n) - y_n)^2 = \sum_{n} (\tilde{\boldsymbol{x}}_n^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_n)^2$$

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

Minimize total squared error

• Residual Sum of Squares (RSS), a function of \tilde{w}

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (f(\boldsymbol{x}_{n}) - y_{n})^{2} = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

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

Minimize total squared error

• Residual Sum of Squares (RSS), a function of \tilde{w}

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (f(\boldsymbol{x}_{n}) - y_{n})^{2} = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

- find $\tilde{w}^* = \underset{\tilde{w} \in \mathbb{R}^{D+1}}{\operatorname{argmin}} \operatorname{RSS}(\tilde{w})$, i.e. least (mean) 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*

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



f is a horizontal line, where should it be?

$$\operatorname{RSS}(w_0) = \sum_n (w_0 - y_n)^2$$

(it's a *quadratic*
$$aw_0^2 + bw_0 + c$$
)

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

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

Optimization objective becomes

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

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

Optimization objective becomes

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

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?

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (w_0 + w_1 x_n - y_n)^2$$

Optimization objective becomes

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (w_0 + w_1 x_n - y_n)^2$$

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

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

Optimization objective becomes

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (w_0 + w_1 x_n - y_n)^2$$

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

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

$$\Rightarrow \begin{array}{ll} Nw_0 + w_1 \sum_n x_n &= \sum_n y_n \\ w_0 \sum_n x_n + w_1 \sum_n x_n^2 &= \sum_n y_n x_n \end{array} \quad (a \text{ linear system})$$

Optimization objective becomes

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (w_0 + w_1 x_n - y_n)^2$$

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

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

$$\Rightarrow \begin{array}{l} Nw_0 + w_1 \sum_n x_n &= \sum_n y_n \\ w_0 \sum_n x_n + w_1 \sum_n x_n^2 &= \sum_n y_n x_n \end{array} \quad (a \text{ linear system}) \\ \Rightarrow \left(\begin{array}{c} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{array} \right) \left(\begin{array}{c} w_0 \\ w_1 \end{array} \right) = \left(\begin{array}{c} \sum_n y_n \\ \sum_n x_n y_n \end{array} \right) \end{array}$$

$$\Rightarrow \left(\begin{array}{c} w_0^* \\ w_1^* \end{array}\right) = \left(\begin{array}{cc} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{array}\right)^{-1} \left(\begin{array}{c} \sum_n y_n \\ \sum_n x_n y_n \end{array}\right)$$

(assuming the matrix is invertible)

$$\Rightarrow \left(\begin{array}{c} w_0^* \\ w_1^* \end{array}\right) = \left(\begin{array}{cc} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{array}\right)^{-1} \left(\begin{array}{c} \sum_n y_n \\ \sum_n x_n y_n \end{array}\right)$$

(assuming the matrix is invertible)

Are stationary points minimizers?

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

(assuming the matrix is invertible)

Are stationary points minimizers?

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

$$\Rightarrow \left(\begin{array}{c} w_0^* \\ w_1^* \end{array}\right) = \left(\begin{array}{cc} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{array}\right)^{-1} \left(\begin{array}{c} \sum_n y_n \\ \sum_n x_n y_n \end{array}\right)$$

(assuming the matrix is invertible)

Are stationary points minimizers?

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

Objective

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Objective

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Again, find stationary points (multivariate calculus)

$$abla ext{RSS}(ilde{m{w}}) = 2\sum_n ilde{m{x}}_n (ilde{m{x}}_n^{ ext{T}} ilde{m{w}} - y_n)$$

Objective

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Again, find stationary points (multivariate calculus)

$$\nabla \text{RSS}(\tilde{\boldsymbol{w}}) = 2\sum_{n} \tilde{\boldsymbol{x}}_{n} (\tilde{\boldsymbol{x}}_{n}^{\text{T}} \tilde{\boldsymbol{w}} - y_{n}) \propto \left(\sum_{n} \tilde{\boldsymbol{x}}_{n} \tilde{\boldsymbol{x}}_{n}^{\text{T}}\right) \tilde{\boldsymbol{w}} - \sum_{n} \tilde{\boldsymbol{x}}_{n} y_{n}$$

Objective

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Again, find stationary points (multivariate calculus)

$$\nabla \text{RSS}(\tilde{\boldsymbol{w}}) = 2 \sum_{n} \tilde{\boldsymbol{x}}_{n} (\tilde{\boldsymbol{x}}_{n}^{\text{T}} \tilde{\boldsymbol{w}} - y_{n}) \propto \left(\sum_{n} \tilde{\boldsymbol{x}}_{n} \tilde{\boldsymbol{x}}_{n}^{\text{T}} \right) \tilde{\boldsymbol{w}} - \sum_{n} \tilde{\boldsymbol{x}}_{n} y_{n}$$
$$= (\tilde{\boldsymbol{X}}^{\text{T}} \tilde{\boldsymbol{X}}) \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{X}}^{\text{T}} \boldsymbol{y}$$

where

$$\tilde{\boldsymbol{X}} = \begin{pmatrix} \tilde{\boldsymbol{x}}_1^{\mathrm{T}} \\ \tilde{\boldsymbol{x}}_2^{\mathrm{T}} \\ \vdots \\ \tilde{\boldsymbol{x}}_{\mathsf{N}}^{\mathrm{T}} \end{pmatrix} \in \mathbb{R}^{\mathsf{N} \times (D+1)}, \quad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{\mathsf{N}} \end{pmatrix} \in \mathbb{R}^{\mathsf{N}}$$

Objective

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2}$$

Again, find stationary points (multivariate calculus)

$$\nabla \text{RSS}(\tilde{\boldsymbol{w}}) = 2 \sum_{n} \tilde{\boldsymbol{x}}_{n} (\tilde{\boldsymbol{x}}_{n}^{\text{T}} \tilde{\boldsymbol{w}} - y_{n}) \propto \left(\sum_{n} \tilde{\boldsymbol{x}}_{n} \tilde{\boldsymbol{x}}_{n}^{\text{T}} \right) \tilde{\boldsymbol{w}} - \sum_{n} \tilde{\boldsymbol{x}}_{n} y_{n}$$
$$= (\tilde{\boldsymbol{X}}^{\text{T}} \tilde{\boldsymbol{X}}) \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{X}}^{\text{T}} \boldsymbol{y} = \boldsymbol{0}$$

where

$$\tilde{\boldsymbol{X}} = \begin{pmatrix} \tilde{\boldsymbol{x}}_1^{\mathrm{T}} \\ \tilde{\boldsymbol{x}}_2^{\mathrm{T}} \\ \vdots \\ \tilde{\boldsymbol{x}}_{\mathsf{N}}^{\mathrm{T}} \end{pmatrix} \in \mathbb{R}^{\mathsf{N} \times (D+1)}, \quad \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{\mathsf{N}} \end{pmatrix} \in \mathbb{R}^{\mathsf{N}}$$

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

assuming $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ is invertible for now.

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

assuming $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ is invertible for now.

Again by convexity $ilde{w}^*$ is the minimizer of RSS.

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

assuming $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ is invertible for now.

Again by convexity $ilde{w}^*$ is the minimizer of RSS.

Verify the solution when D = 1:

$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_{\mathsf{N}} \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \cdots & \cdots \\ 1 & x_{\mathsf{N}} \end{pmatrix} = \begin{pmatrix} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}$$

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

assuming $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ is invertible for now.

Again by convexity $ilde{w}^*$ is the minimizer of RSS.

Verify the solution when D = 1:

$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_{\mathsf{N}} \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \cdots & \cdots \\ 1 & x_{\mathsf{N}} \end{pmatrix} = \begin{pmatrix} N & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}$$

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

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2}$$

$$\begin{aligned} &\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2} \\ &= \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right) \end{aligned}$$

$$\begin{aligned} &\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2} \\ &= \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right) \\ &= \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y} + \operatorname{cnt.} \end{aligned}$$

$$\begin{aligned} &\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2} \\ &= \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right) \\ &= \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y} + \operatorname{cnt.} \\ &= \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}}\right) \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right) + \operatorname{cnt.} \end{aligned}$$

$$\begin{aligned} &\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2} \\ &= \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right) \\ &= \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y} + \operatorname{cnt.} \\ &= \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}}\right) \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right) + \operatorname{cnt.} \end{aligned}$$

Note:
$$\boldsymbol{u}^{\mathrm{T}}\left(\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}}\right)\boldsymbol{u} = \left(\tilde{\boldsymbol{X}}\boldsymbol{u}\right)^{\mathrm{T}}\tilde{\boldsymbol{X}}\boldsymbol{u} = \|\tilde{\boldsymbol{X}}\boldsymbol{u}\|_{2}^{2} \geq 0$$
 and is 0 if $\boldsymbol{u} = 0$.

$$\operatorname{RSS}(\tilde{\boldsymbol{w}}) = \sum_{n} (\tilde{\boldsymbol{x}}_{n}^{\mathrm{T}} \tilde{\boldsymbol{w}} - y_{n})^{2} = \|\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\|_{2}^{2}$$
$$= \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}\right)$$
$$= \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \boldsymbol{y}^{\mathrm{T}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{w}} - \tilde{\boldsymbol{w}}^{\mathrm{T}} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y} + \operatorname{cnt.}$$
$$= \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right)^{\mathrm{T}} \left(\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}}\right) \left(\tilde{\boldsymbol{w}} - (\tilde{\boldsymbol{X}}^{\mathrm{T}} \tilde{\boldsymbol{X}})^{-1} \tilde{\boldsymbol{X}}^{\mathrm{T}} \boldsymbol{y}\right) + \operatorname{cnt.}$$

Note:
$$\boldsymbol{u}^{\mathrm{T}}\left(\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}}\right)\boldsymbol{u} = \left(\tilde{\boldsymbol{X}}\boldsymbol{u}\right)^{\mathrm{T}}\tilde{\boldsymbol{X}}\boldsymbol{u} = \|\tilde{\boldsymbol{X}}\boldsymbol{u}\|_{2}^{2} \geq 0$$
 and is 0 if $\boldsymbol{u} = 0$.
So $\tilde{\boldsymbol{w}}^{*} = (\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}})^{-1}\tilde{\boldsymbol{X}}^{\mathrm{T}}\boldsymbol{y}$ is the minimizer.

Computational complexity

Bottleneck of computing

$$ilde{oldsymbol{w}}^* = \left(ilde{oldsymbol{X}}^{ ext{T}} ilde{oldsymbol{X}}
ight)^{-1} ilde{oldsymbol{X}}^{ ext{T}} oldsymbol{y}$$

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

- aka pseudo-inverse¹ denoted by $(\cdot)^{\dagger}$, i.e. $\tilde{X}^{\dagger} = \left(\tilde{X}^{\mathrm{T}}\tilde{X}\right)^{-1}\tilde{X}^{\mathrm{T}}$
- naively need $O(\mathsf{D}^3)$ time

see https://en.wikipedia.org/wiki/Moore-Penrose_inverse

Computational complexity

Bottleneck of computing

$$ilde{oldsymbol{w}}^* = \left(ilde{oldsymbol{X}}^{ ext{T}} ilde{oldsymbol{X}}
ight)^{-1} ilde{oldsymbol{X}}^{ ext{T}} oldsymbol{y}$$

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

- aka pseudo-inverse¹ denoted by $(\cdot)^{\dagger}$, i.e. $ilde{m{X}}^{\dagger} = \left(ilde{m{X}}^{\mathrm{T}} ilde{m{X}}\right)^{-1} ilde{m{X}}^{\mathrm{T}}$
- naively need $O(\mathsf{D}^3)$ time
- there are many faster approaches

see https://en.wikipedia.org/wiki/Moore-Penrose_inverse

What if $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ is not invertible

What does that imply?

See https://sites.math.washington.edu/~burke/crs/308/LeastSquares.pdf
What does that imply?

$$\mathsf{Recall}\,\left(ilde{oldsymbol{X}}^{\mathrm{T}} ilde{oldsymbol{X}}
ight)oldsymbol{w}^{*}= ilde{oldsymbol{X}}^{\mathrm{T}}oldsymbol{y}.$$

What does that imply?

Recall $\left(\tilde{X}^{\mathrm{T}} \tilde{X} \right) w^* = \tilde{X}^{\mathrm{T}} y$. If $\tilde{X}^{\mathrm{T}} \tilde{X}$ not invertible, this equation aka Normal Equations has

• infinitely many solutions

What does that imply?

Recall $\left(\tilde{X}^{\mathrm{T}} \tilde{X} \right) w^* = \tilde{X}^{\mathrm{T}} y$. If $\tilde{X}^{\mathrm{T}} \tilde{X}$ not invertible, this equation aka Normal Equations has

• infinitely many solutions

What does that imply?

Recall $\left(\tilde{X}^{\mathrm{T}}\tilde{X}\right)w^* = \tilde{X}^{\mathrm{T}}y$. If $\tilde{X}^{\mathrm{T}}\tilde{X}$ not invertible, this equation aka *Normal Equations* has

- infinitely many solutions (\Rightarrow infinitely many minimizers)
- This is because *Normal Equations* are always *consistent*² meaning a solution *always* exists! It may not be unique though.

Why would that happen?

Why would that happen?

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

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

Why would that happen?

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

Example: D = N = 1

sqft	sale price
1000	500K

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

How to resolve this issue?

Intuition: what does inverting $ilde{X}^{\mathrm{T}} ilde{X}$ do?

eigendecomposition:
$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_{1} & 0 & \cdots & 0 \\ 0 & \lambda_{2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_{\mathrm{D}} & 0 \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}+1} \end{bmatrix} \boldsymbol{U}$$

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

How to resolve this issue?

Intuition: what does inverting $ilde{m{X}}^{\mathrm{T}} ilde{m{X}}$ do?

eigendecomposition:
$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_{1} & 0 & \cdots & 0 \\ 0 & \lambda_{2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_{\mathsf{D}} & 0 \\ 0 & \cdots & 0 & \lambda_{\mathsf{D}+1} \end{bmatrix} \boldsymbol{U}$$

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

inverse:
$$(\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}})^{-1} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \frac{1}{\lambda_{1}} & 0 & \cdots & 0\\ 0 & \frac{1}{\lambda_{2}} & \cdots & 0\\ \vdots & \vdots & \vdots & \vdots\\ 0 & \cdots & \frac{1}{\lambda_{\mathrm{D}}} & 0\\ 0 & \cdots & 0 & \frac{1}{\lambda_{\mathrm{D}+1}} \end{bmatrix} \boldsymbol{U}$$

i.e. just inverse 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

$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} + \lambda \boldsymbol{I} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_{1} + \lambda & 0 & \cdots & 0 \\ 0 & \lambda_{2} + \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_{\mathsf{D}} + \lambda & 0 \\ 0 & \cdots & 0 & \lambda_{\mathsf{D}+1} + \lambda \end{bmatrix} \boldsymbol{U}$$

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

How to solve this problem?

Non-invertible \Rightarrow some eigenvalues are 0.

One natural fix: add something positive

$$\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} + \lambda \boldsymbol{I} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_{1} + \lambda & 0 & \cdots & 0 \\ 0 & \lambda_{2} + \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_{\mathrm{D}} + \lambda & 0 \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}+1} + \lambda \end{bmatrix} \boldsymbol{U}$$

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

$$(\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}} + \lambda \boldsymbol{I})^{-1} = \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \frac{1}{\lambda_{1}+\lambda} & 0 & \cdots & 0\\ 0 & \frac{1}{\lambda_{2}+\lambda} & \cdots & 0\\ \vdots & \vdots & \vdots & \vdots\\ 0 & \cdots & \frac{1}{\lambda_{\mathsf{D}}+\lambda} & 0\\ 0 & \cdots & 0 & \frac{1}{\lambda_{\mathsf{D}+1}+\lambda} \end{bmatrix} \boldsymbol{U}$$

Fix the problem

The solution becomes

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

• not a minimizer of the original RSS

Fix the problem

The solution becomes

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

• not a minimizer of the original RSS

This in fact comes from minimizing **regularized** RSS (covered in next lecture)!

$$\min_{\boldsymbol{w}} \|\tilde{\boldsymbol{X}}\tilde{\boldsymbol{w}} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{y}\|_2^2$$

Fix the problem

The solution becomes

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

• not a minimizer of the original RSS

This in fact comes from minimizing **regularized** RSS (covered in next lecture)!

$$\min_{\boldsymbol{w}} \|\tilde{\boldsymbol{X}}\tilde{\boldsymbol{w}} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{y}\|_2^2$$

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

Comparison to NNC

Parametric versus non-parametric

- **Parametric methods**: the size of the model does *not grow* with the size of the training set N.
 - e.g. linear regression, D + 1 parameters, independent of N.
- 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.