# Lecture 4: Linear Regression (Part 3)

CIS 4190/5190 Spring 2023

#### Administrivia

- Class roster is now stable, and add/drop deadline has passed.
- HW1 due **tonight 8 p.m.**, and HW2 will be posted tonight/tomorrow morning, on linear regression.
  - Class late policy reminder: 0.5% points for every late hour, up to a max of 48 hours.
- TA office hours:
  - Any changes will be posted to the TA office hours thread on EdSTEM, at least 48 hours ahead of time.
  - We are moving the 3 p.m. OH on Monday and Wed to 3.30 p.m. to avoid clashing. More news on that soon.
- Recitations tomorrow at 5 p.m. on Python, Numpy, Pandas, Scikit-Learn. See EdSTEM post.
- No quiz for week 1. We'll fix the webpage.

#### Last Lecture Summary

- The Train/Test Split Protocol for Measuring Underfitting / Overfitting
- Bias and variance as functions of a model class
  - Tuning them by selecting hypothesis spaces / feature maps
  - Tuning them by modifying the loss function
    - $-L_{\text{new}}(\beta; Z) = L(\beta; Z) + \lambda \cdot R(\beta)$
- Today:
  - Selecting hyperparameters like  $\lambda$
  - Finally unveil the mystery about how to find  $\hat{\beta}(Z) = \arg \min_{\beta} L(\beta; Z)$



# Cross-Validation for Model Selection

#### Hyperparameter Tuning, or "Model Selection"

- $\lambda$  is a hyperparameter that must be tuned (satisfies  $\lambda \ge 0$ )
- Naïve strategy: Try a few different candidates  $\lambda_t$  and choose the one that minimizes the test loss

#### Test Data Contamination

- Suppose you have tried 100 different hyperparameter values, that all haver the same value of generalization MSE, if evaluated on the *full* data distribution.
- But the test dataset is only a finite sample of this distribution, so test MSE is a noisy estimate of true generalization MSE. For example



Note how, in selecting based on test MSE, you have "overfit" your hyperparameter choice to your test set!

Hyperparameter values

#### Hyperparameter Tuning, or "Model Selection"

- $\lambda$  is a hyperparameter that must be tuned (satisfies  $\lambda \ge 0$ )
- Naïve strategy: Try a few different candidates  $\lambda_t$  and choose the one that minimizes the test loss
- **Problem:** We may overfit the test set!
  - Major problem if we have more hyperparameters
- Solution: A new subset of data just for selecting hyperparameters

#### Train/Val/Test Split Protocol for Model Selection

- Goal: Choose best hyperparameter  $\lambda$ 
  - Can also compare different model families, feature maps, etc.
- Solution: Optimize  $\lambda$  on a held-out validation data
  - Rule of thumb: 60/20/20 split (usually shuffle before splitting)



#### Basic Cross Validation Algorithm: "Holdout"

• Step 1: Split Z into  $Z_{\text{train}}$ ,  $Z_{\text{val}}$ , and  $Z_{\text{test}}$ 

Training data $Z_{\text{train}}$	Val data $Z_{ m val}$	Test data $Z_{\text{test}}$
----------------------------------	-----------------------	-----------------------------

- Step 2: For  $t \in \{1, ..., h\}$  hyperparameter choices:
  - Step 2a: Run linear regression with  $Z_{\text{train}}$  and  $\lambda_t$  to obtain  $\hat{\beta}(Z_{\text{train}}, \lambda_t)$
  - Step 2b: Evaluate validation loss  $L_{val}^t = L(\hat{\beta}(Z_{train}, \lambda_t); Z_{val})$
- Step 3: Use best  $\lambda_t$ 
  - Choose  $t' = \arg \min_t L_{val}^t$  with lowest validation loss
  - Re-run linear regression with  $Z_{\text{train}}$  and  $\lambda_{t'}$  to obtain  $\hat{\beta}(Z_{\text{train}}, \lambda_{t'})$

#### **Cross Validation Hygiene**



- The moment that test data is used for hyperparameter selection or to iterate on ML design choices, it should be treated as "contaminated".
- Remember: Performance on contaminated test data is an overly *optimistic* estimate of the "true" test performance.

Q: What about validation data performance then?

(yes, this is also overly optimistic)

#### Alternative Cross-Validation Algorithms

- If Z is small, then splitting it can reduce performance
  - Can use  $Z_{\text{train}} \cup Z_{\text{val}}$  in Step 3
- Alternative more thorough CV strategy: "k-fold" cross-validation
  - Split Z into Z<sub>train</sub> and Z<sub>test</sub>
  - Split  $Z_{\text{train}}$  into k disjoint sets  $Z_{\text{val}}^s$ , and let  $Z_{\text{train}}^s = \bigcup_{s' \neq s} Z_{\text{val}}^s$
  - Use  $\lambda'$  that works best on average across  $s \in \{1, ..., k\}$  with  $Z_{\text{train}}$
  - Chooses better  $\lambda'$  than above strategy

#### Example: k = 3-Fold Cross Validation





**Compute vs. accuracy tradeoff:** As  $k \rightarrow N$ , model selection becomes more accurate, but algorithm becomes more computationally expensive

#### Note: What Exactly Are "Hyperparameters"?

- Cross-Validation is a general, systematic trial-and-error procedure for selecting hyperparameters.
- Other hyperparameters too, not just the regularization  $\lambda$ .
- "Hyperparameters" are ML system properties / design choices that are not directly set in the optimization problem.

 $\hat{\beta}(Z) = \arg\min_{\beta} L(\beta; Z)$ 

- Examples of other hyperparameters you could set with cross-validation:
  - choice of feature maps in linear regression.
  - data selection and other preprocessing procedures (coming up soon).
  - Inear regression versus another ML algorithm, altogether.



#### Minimizing the MSE Loss

• Recall that linear regression minimizes the loss

$$L(\boldsymbol{\beta}; \boldsymbol{Z}) = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{y}_i - \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}_i)^2$$

- Closed-form solution: Compute a matrix expression derived using calculus
- Iterative Optimization-based solution: Search over candidate *β*



$$\begin{bmatrix} f_{\beta}(x_1) \\ \vdots \\ f_{\beta}(x_n) \end{bmatrix} = \begin{bmatrix} \beta^{\mathsf{T}} x_1 \\ \vdots \\ \beta^{\mathsf{T}} x_n \end{bmatrix}$$







$$\begin{bmatrix} f_{\beta}(x_{1}) \\ \vdots \\ f_{\beta}(x_{n}) \end{bmatrix} = \begin{bmatrix} \beta^{\mathsf{T}} x_{1} \\ \vdots \\ \beta^{\mathsf{T}} x_{n} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{d} \beta_{j} x_{1,j} \\ \vdots \\ \sum_{j=1}^{d} \beta_{j} x_{n,j} \end{bmatrix} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{bmatrix} \begin{bmatrix} \beta_{1} \\ \vdots \\ \beta_{d} \end{bmatrix} = X\beta$$

$$\begin{bmatrix} f_{\beta}(x_{1}) \\ \vdots \\ f_{\beta}(x_{n}) \end{bmatrix} = \begin{bmatrix} \beta^{\mathsf{T}} x_{1} \\ \vdots \\ \beta^{\mathsf{T}} x_{n} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{d} \beta_{j} x_{1,j} \\ \vdots \\ \sum_{j=1}^{d} \beta_{j} x_{n,j} \end{bmatrix} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{bmatrix} \begin{bmatrix} \beta_{1} \\ \vdots \\ \beta_{d} \end{bmatrix} = X\beta$$

 $\mathbf{S}$ 

 $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ 

$$\begin{bmatrix} f_{\beta}(x_{1}) \\ \vdots \\ f_{\beta}(x_{n}) \end{bmatrix} = \begin{bmatrix} \beta^{\mathsf{T}} x_{1} \\ \vdots \\ \beta^{\mathsf{T}} x_{n} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{d} \beta_{j} x_{1,j} \\ \vdots \\ \sum_{j=1}^{d} \beta_{j} x_{n,j} \end{bmatrix} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{bmatrix} \begin{bmatrix} \beta_{1} \\ \vdots \\ \beta_{d} \end{bmatrix} = X\beta$$

 $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = Y$ 

 $\mathbf{N}$ 

#### Summary: $Y \approx X\beta$

Note: n equations, d variables

#### $Y \approx X\beta$



#### Vectorizing Mean Squared Error

#### Vectorizing Mean Squared Error

 $L(\beta; \mathbf{Z})$ 

#### Vectorizing Mean Squared Error

$$L(\boldsymbol{\beta}; \boldsymbol{Z}) = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{y}_i - \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}_i)^2$$

# Vectorizing Mean Squared Error $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \begin{bmatrix} f_{\beta}(x_1) \\ \vdots \\ f_{\beta}(x_n) \end{bmatrix}$



#### Intuition on Vectorized Linear Regression

• Rewriting the vectorized loss:

$$n \cdot L(\beta; Z) = \|Y - X\beta\|_2^2 = \|Y\|_2^2 - 2Y^{\mathsf{T}}X\beta + \|X\beta\|_2^2$$
$$= \|Y\|_2^2 - 2Y^{\mathsf{T}}X\beta + \beta^{\mathsf{T}}(X^{\mathsf{T}}X)\beta$$

- Side note: Quadratic function of  $\beta$  with leading "coefficient"  $X^{\top}X$ 
  - In one dimension, "width" of parabola  $ax^2 + bx + c$  is  $a^{-1}$
  - In multiple dimensions, "width" along direction  $v_i$  is  $\lambda_i^{-1}$ , where  $v_i$  is an eigenvector of  $X^{\top}X$  with eigenvalue  $\lambda_i$ 
    - Large width (small  $\lambda_i$ ) along a direction  $v_i$  implies that parameter values along that direction affect the loss value less.
    - This will be interesting to us later in class ("PCA")

#### Intuition on Vectorized Linear Regression



Directions/magnitudes are given by eigenvectors/eigenvalues of  $X^{\top}X$ 

• Recall that linear regression minimizes the loss:

$$L(\boldsymbol{\beta}; \boldsymbol{Z}) = \frac{1}{n} \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2$$

• Minimum solution has gradient equal to zero:

$$\nabla_{\beta} L(\hat{\beta}(Z); Z) = 0$$



• Recall that linear regression minimizes the loss

$$L(\boldsymbol{\beta}; \boldsymbol{Z}) = \frac{1}{n} \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2$$

• Minimum solution has gradient equal to zero:

$$\nabla_{\beta} L(\hat{\beta}; \mathbf{Z}) = 0$$



• The gradient is

 $\nabla_{\boldsymbol{\beta}} L(\boldsymbol{\beta}; \mathbf{Z})$ 

• The gradient is

$$\nabla_{\beta} L(\beta; \mathbf{Z}) = \nabla_{\beta} \frac{1}{n} \|\mathbf{Y} - \mathbf{X}\beta\|_{2}^{2}$$

• The gradient is

$$\nabla_{\beta} L(\beta; Z) = \nabla_{\beta} \frac{1}{n} ||Y - X\beta||_{2}^{2} = \nabla_{\beta} \frac{1}{n} (Y - X\beta)^{\mathsf{T}} (Y - X\beta)$$
$$= \frac{2}{n} [\nabla_{\beta} (Y - X\beta)^{\mathsf{T}}] (Y - X\beta)$$
$$= -\frac{2}{n} X^{\mathsf{T}} (Y - X\beta)$$
$$= \left[-\frac{2}{n} X^{\mathsf{T}} Y + \frac{2}{n} X^{\mathsf{T}} X\beta\right]$$

• The gradient is

$$\nabla_{\beta}L(\beta; Z) = \nabla_{\beta} \frac{1}{n} \|Y - X\beta\|_2^2 = -\frac{2}{n} X^{\mathsf{T}}Y + \frac{2}{n} X^{\mathsf{T}}X\beta$$

• Setting 
$$\nabla_{\beta} L(\hat{\beta}; Z) = 0$$
, we have  $X^{\top} X \hat{\beta} = X^{\top} Y$   
Compare this to:  $Y \approx X\beta$ 

• Assuming  $X^{\top}X$  is invertible, we have

$$\hat{\beta}(Z) = (X^{\top}X)^{-1}X^{\top}Y$$
This is called the  
"pseudoinverse" of X

• Setting 
$$\nabla_{\beta} L(\hat{\beta}; Z) = 0$$
, we have  $X^{\top} X \hat{\beta} = X^{\top} Y$ 

• Assuming  $X^{\top}X$  is invertible, we have

 $\hat{\beta}(Z) = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}Y$ 

## Note on Invertibility

- Closed-form solution only **unique** if  $X^{\top}X$  (size dxd) is invertible
  - Otherwise, multiple solutions exist to  $X^{\top}X\hat{\beta} = X^{\top}Y$
  - Intuition: Underconstrained system of linear equations
- Example:

$$\begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

• In this case, any  $\hat{\beta}_2 = 2 - \hat{\beta}_1$  is a solution

## When Can this Happen?

- Case 1
  - Fewer data examples than feature dimension (i.e., n < d)
    - Remember: we are solving something like  $Y \approx X\beta$
  - **Solution:** Remove features so  $d \le n$
  - Solution: Collect more data until  $d \le n$
- Case 2: Some feature is a linear combination of the others
  - Special case (duplicated feature): For some j and j',  $x_{i,j} = x_{i,j'}$  for all i
  - Solution: Remove linearly dependent features
  - Solution: Use L<sub>2</sub> regularization (we will soon see why)

## Shortcomings of Closed-Form Solution

- Computing  $\hat{\beta}(Z) = (X^{\top}X)^{-1}X^{\top}Y$  can be challenging
- Computing  $(X^{\top}X)^{-1}$  is  $O(d^3)$ 
  - $d = 10^4$  features  $\rightarrow O(10^{12})$
  - Even storing  $X^{\top}X$  requires a lot of memory
- Numerical accuracy issues due to "ill-conditioning"
  - $X^{\top}X$  is "barely" invertible
  - Then,  $(X^{\top}X)^{-1}$  has large variance along some dimension
  - Regularization helps (more on this later)



#### **Iterative Optimization Algorithms**

• Recall that linear regression minimizes the loss

$$L(\boldsymbol{\beta}; \boldsymbol{Z}) = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{y}_i - \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}_i)^2$$

- Iteratively optimize  $\beta$ 
  - Initialize  $\beta_1 \leftarrow \text{Init}(...)$
  - For some number of iterations T, update  $\beta_t \leftarrow \text{Step}(...)$
  - Return  $\beta_T$

#### **Iterative Optimization Algorithms**

- **Global search**: Try random values of  $\beta$  and choose the best
  - I.e.,  $\beta_t$  independent of  $\beta_{t-1}$
  - Very unstructured, can take a long time (especially in high dimension d)!
- Local search: Start from some initial  $\beta$  and make local changes
  - I.e.,  $\beta_t$  is computed based on  $\beta_{t-1}$
  - What is a "local change", and how do we find good one?

• Gradient descent: Update  $\beta$  based on gradient  $\nabla_{\beta} L(\beta; Z)$  of  $L(\beta; Z)$ :

$$\beta_{t+1} \leftarrow \beta_t - \alpha \cdot \nabla_\beta L(\beta_t; \mathbf{Z})$$

- Intuition: The gradient is the direction along which  $L(\beta; Z)$  changes most quickly as a function of  $\beta$
- $\alpha \in \mathbb{R}$  is a hyperparameter called the **learning rate** 
  - More on this later

- Choose initial value for  $\beta$
- Until we reach a minimum:
  - Choose a new value for  $\beta$  to reduce  $L(\beta; \mathbb{Z})$



Figure by Andrew Ng

- Choose initial value for  $\beta$
- Until we reach a minimum:
  - Choose a new value for  $\beta$  to reduce  $L(\beta; \mathbb{Z})$



Figure by Andrew Ng

- Choose initial value for  $\beta$
- Until we reach a minimum:
  - Choose a new value for  $\beta$  to reduce  $L(\beta; \mathbb{Z})$



# Linear regression loss is convex, so no local minima

Indexing iteration now, rather than parameter vector element

- Initialize  $\beta_1^{\prime} = \vec{0}$
- Repeat until convergence:

 $\beta_{t+1} \leftarrow \beta_t - \alpha \cdot \nabla_\beta L(\beta_t; \mathbf{Z})$ 

 For linear regression, know the gradient from strategy 1



For in-place updates  $\beta \leftarrow \beta - \alpha \cdot \nabla_{\beta} L(\beta; \mathbb{Z})$ , compute all components of  $\nabla_{\beta} L(\beta; \mathbb{Z})$  before modifying  $\beta$ 

- Initialize  $\beta_1 = \vec{0}$
- Repeat until convergence:

$$\beta_{t+1} \leftarrow \beta_t - \alpha \cdot \nabla_\beta L(\beta_t; \mathbf{Z})$$

• For linear regression, know the gradient from strategy 1



- Initialize  $\beta_1 = \vec{0}$
- Repeat until  $\|\beta_t \beta_{t+1}\|_2 \le \epsilon$ :

 $\beta_{t+1} \leftarrow \beta_t - \alpha \cdot \nabla_\beta L(\beta_t; Z)$ 

• For linear regression, know the gradient from strategy 1

Hyperparameter defining convergence



#### Aside: Gradient As Sum of Sample-Wise Gradients

(Equivalent to our earlier matrix expression of gradient)

• By linearity of the gradient, we have

$$\nabla_{\beta} L(\beta; Z) = \sum_{i=1}^{n} \nabla_{\beta} (y_i - \beta^{\mathsf{T}} x_i)^2 = \sum_{i=1}^{n} 2(y_i - \beta^{\mathsf{T}} x_i) x_i$$

 $-\frac{2}{n}X^{\mathsf{T}}Y + \frac{2}{n}X^{\mathsf{T}}X\beta$ 

• The gradient term induced by a single training data sample is:

$$\nabla_{\beta}(y_i - \beta^{\mathsf{T}} x_i)^2 = 2(y_i - \beta^{\mathsf{T}} x_i)x_i$$

• I.e., the current error  $y_i - \beta^T x_i$  times the feature vector  $x_i$ 

"Large error samples induce large changes to  $\beta$ , proportional to their feature values."



















Minimizer of loss function

#### Choice of Learning Rate $\alpha$



**Problem:**  $\alpha$  too small

•  $L(\beta; Z)$  decreases slowly

**Problem:**  $\alpha$  too large •  $L(\beta; Z)$  increases!

Plot  $L(\beta_t; Z_{\text{train}})$  vs. t to diagnose these problems



#### Choice of Learning Rate $\alpha$

- α is a hyperparameter for gradient descent that we need to choose
   Can set just based on training data
- Rule of thumb
  - α too small: Loss decreases slowly
  - α too large: Loss increases!
- Try rates  $\alpha \in \{1.0, 0.1, 0.01, ...\}$  (can tune further once one works)

#### **Comparison of Strategies**

#### Closed-form solution

- No hyperparameters
- Slow if n or d are large
- Gradient descent
  - Need to tune α
  - Scales to large n and d
- For linear regression, there are better optimization algorithms, but gradient descent is very general
  - Accelerated gradient descent is an important tweak that improves performance in practice (and in theory)



#### Loss Minimization View of ML

- Two design decisions
  - Model family: What are the candidate models f? (E.g., linear functions)
  - Loss function: How to define "approximating"? (E.g., MSE loss)

#### Loss Minimization View of ML

#### Three design decisions

- Model family: What are the candidate models f? (E.g., linear functions)
- Loss function: How to define "approximating"? (E.g., MSE loss)
- Optimizer: How do we minimize the loss? (E.g., gradient descent)