

CSE 6740 A/ISyE 6740: Computational Data Analysis: Introductory lecture

Nisha Chandramoorthy

August 28, 2023

Quick notes

- ▶ Review session tonight 7:30-9:30pm on Zoom. See Piazza/Canvas for details.
- ▶ Homework 1 will be out soon. Due on 9/13/2023.

Last time

- ▶ Empirical risk minimization, finite hypothesis classes

Last time

- ▶ Empirical risk minimization, finite hypothesis classes
- ▶ Overfitting, inductive bias, Intro to PAC learning

Last time

- ▶ Empirical risk minimization, finite hypothesis classes
- ▶ Overfitting, inductive bias, Intro to PAC learning
- ▶ Let ERM rule

$$h_S := \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_S(h),$$

where empirical error,

$$\hat{R}_S(h) := \frac{1}{m} \sum_{z \in S} \ell(z, h).$$

Let the realizability assumption be satisfied \implies ERM rule h_S has zero empirical error. Then, with probability at least $1 - \delta$, the generalization error,

$$R(h_S) := \underbrace{E_{z \in \mathcal{D}}}_{\sim} \ell(z, h_S) \leq \frac{1}{m} \log \frac{|\mathcal{H}|}{\delta}.$$

Linear models

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = w \cdot \Phi(x) + b, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

Linear models

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = w \cdot \Phi(x) + b, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

► $\Phi : \mathbb{X} \rightarrow \mathbb{R}^d$ is a set of *features*.

Linear models

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = w \cdot \Phi(x) + b, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

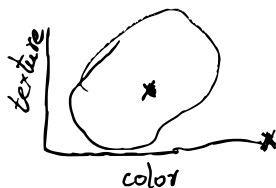
- ▶ $\Phi : \mathbb{X} \rightarrow \mathbb{R}^d$ is a set of *features*.
- ▶ Linear regression seeks ERM solution for square loss

Linear models

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = w \cdot \Phi(x) + b, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

- ▶ $\Phi : \mathbb{X} \rightarrow \mathbb{R}^d$ is a set of features.
- ▶ Linear regression seeks ERM solution for square loss

$$\rightarrow \operatorname{argmin}_{w,b} \frac{1}{m} \sum_{i=1}^m (w^\top \Phi(x_i) + b - y_i)^2$$



$$\operatorname{Supp}(\mathcal{D}) \subseteq \mathbb{X}$$

4/9

$$x \in \underline{\mathbb{X}} \equiv \mathbb{R}^d$$

$$\operatorname{argmin}_{\substack{w, b \\ \uparrow \\ \mathbb{R}^d \quad \mathbb{R}}} \frac{1}{m} \sum_{i=1}^m \underbrace{\ell(\underset{(x_i, y_i) \in S}{z_i}, h(z_i, w, b))}_{\ell(z, h(z, w, b)) = |h(z, w, b) - y|^2}$$

$$\ell(z, h(z, w, b)) = |h(z, w, b) - y|^2$$

Linear models

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = w \cdot \Phi(x) + b, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

- ▶ $\Phi : \mathbb{X} \rightarrow \mathbb{R}^d$ is a set of *features*.
- ▶ Linear regression seeks ERM solution for square loss
- ▶

$$\operatorname{argmin}_{w,b} \frac{1}{m} (w^\top \Phi(x_i) + b - y_i)^2$$

convex

- ▶ Equivalently, where X is $m \times (d+1)$ matrix with rows $X_i = (\Phi(x_i)^\top, 1)$, $W = [w_1, \dots, w_d, b]^\top$, $Y = [y_1, \dots, y_m]^\top$,

$$f(x) = x^2$$

$f \circ h$ is convex

Linear models ✓

$$\omega \cdot \Phi(x)$$

$$\mathcal{H} = \{h(\cdot, w, b) : h(x, w, b) = \underline{w \cdot \Phi(x) + b}, w \in \mathbb{R}^d, b \in \mathbb{R}\} \quad (1)$$

- ▶ $\Phi : \mathbb{X} \rightarrow \mathbb{R}^d$ is a set of *features*.
- ▶ Linear regression seeks ERM solution for square loss

$$\rightarrow \operatorname{argmin}_{w,b} \frac{1}{m} (w^\top \Phi(x_i) + b - y_i)^2 \quad \checkmark$$

- ▶ Equivalently, where X is $m \times (d+1)$ matrix with rows $X_i = (\Phi(x_i)^\top, 1)$, $W = [w_1, \dots, w_d, b]^\top$, $Y = [y_1, \dots, y_m]^\top$,

$$\rightarrow \operatorname{argmin}_W \frac{1}{m} \|XW - Y\|^2$$

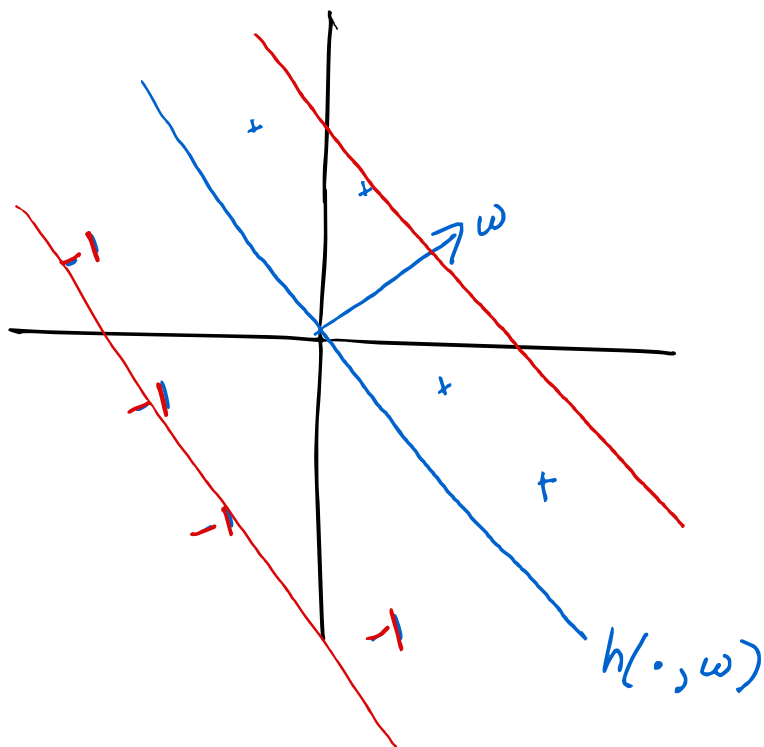
- ▶ Features may be defined by kernels

$$X = \begin{matrix} m \times (d+1) \end{matrix} \begin{bmatrix} \Phi(x_1)^\top & \textcircled{1} \\ \vdots & \\ \Phi(x_m)^\top & 1 \end{bmatrix}$$

$$XW \downarrow$$

$$\begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix}$$

$$\|X_{\omega} - Y\|^2$$



$$x^T \omega$$

$$\mathcal{H}_{d+1}^{\text{lin}}$$

$$\mathcal{H}_d^{\text{aff}}$$

Least squares solutions: the optimization way

- ▶ Convex, differentiable function of W – composition of convex, differentiable functions

Least squares solutions: the optimization way

- ▶ Convex, differentiable function of W – composition of convex, differentiable functions
- ▶ Global minimum is the extremal point where derivative vanishes

Least squares solutions: the optimization way

- ▶ Convex, differentiable function of W – composition of convex, differentiable functions
- ▶ Global minimum is the extremal point where derivative vanishes

- ▶ $\nabla \frac{1}{m} \|XW - Y\|^2 = \frac{2}{m} X^T (XW - Y)$

Matrix calculus

Least squares solutions: the optimization way

- ▶ Convex, differentiable function of W – composition of convex, differentiable functions
- ▶ Global minimum is the extremal point where derivative vanishes
- ▶ $\nabla \frac{1}{m} \|XW - Y\|^2 = \frac{2}{m} X^T (XW - Y)$
- ▶ $X^T XW = X^T Y$. Can also get this by differentiating before writing in matrix form

$$X^T (X \overset{W}{w} - Y) = 0$$

w is
the minimizer

$$X^T X w^{\text{OLS}} = X^T Y$$

Ordinary least squares

$$w^{\text{OLS}} = (X^T X)^{-1} X^T Y$$

Least squares solutions: the optimization way

- ▶ Convex, differentiable function of W – composition of convex, differentiable functions
- ▶ Global minimum is the extremal point where derivative vanishes
- ▶ $\nabla \frac{1}{m} \|XW - Y\|^2 = \frac{2}{m} X^T (XW - Y)$
- ▶ $X^T XW = X^T Y$. Can also get this by differentiating before writing in matrix form
- ▶ When is $X^T X = \sum_{i=1}^m \Phi(x_i) \Phi(x_i)^T$ invertible? When the training features span \mathbb{R}^d .

$$X = \begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_m) \end{bmatrix}$$

$$\Phi(x) = x \\ \in \mathbb{R}^d$$

$$x_i \sim \mathcal{D}$$

$$X X^T = \sum_{i=1}^m \underbrace{x_i x_i^T}$$

Least squares solutions: the linear algebraic way

- ▶ Case 1: $d + 1 = m$, X is invertible. $W = X^{-1}Y$ ✓

$$XW = Y$$

$$m = d$$

Least squares solutions: the linear algebraic way

- ▶ Case 1: $d + 1 = m$, X is invertible. $W = X^{-1} Y$
- ▶ Case 2: $d + 1 > m$, underdetermined/overparameterized. If X has full row rank, then, min norm solution

$$W = X^{\top} (X X^{\top})^{-1} Y$$

Least squares solutions: the linear algebraic way

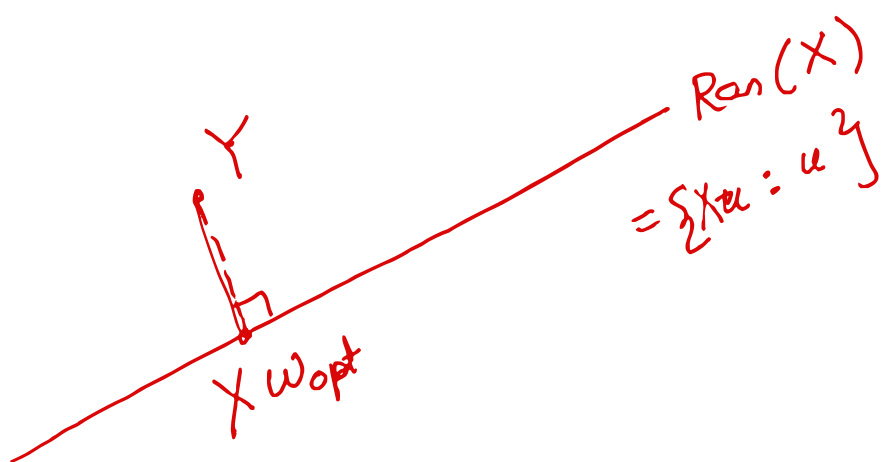
- ▶ Case 1: $d + 1 = m$, X is invertible. $W = X^{-1} Y$
- ▶ Case 2: $d + 1 > m$, underdetermined/overparameterized. If X has full row rank, then, min norm solution

$$W = X^{\top} (X X^{\top})^{-1} Y$$

- ▶ Case 3: $m > d + 1$, overdetermined. If X has full col rank, then, many solutions. Min norm solution

$$\underline{W = (X^{\top} X)^{-1} X^{\top} Y}$$

$$\operatorname{argmin}_W \|XW - Y\|^2$$



$$\langle Y - Xw_{\text{opt}}, Xu \rangle = 0 \quad \forall u$$

$$\langle X^T(Y - Xw_{\text{opt}}), u \rangle = 0 \quad \forall u$$

$$X^T Y = X^T X w_{\text{opt}}$$

$$w_{\text{opt}} = (X^T X)^{-1} X^T Y$$

Recall

When is $X^T X$ invertible?

$$u^T X^T X u = \|Xu\|^2 \geq 0$$

SPSD

$$m < d$$

$$\operatorname{rank}(X^T X)$$

$$X^T X = \sum_{i=1}^m x_i x_i^T$$

if x_i are linearly independent

Least squares solutions: the linear algebraic way

- ▶ Case 1: $d + 1 = m$, X is invertible. $W = X^{-1} Y$
- ▶ Case 2: $d + 1 > m$, underdetermined/overparameterized. If X has full row rank, then, min norm solution

$$W = X^T (X X^T)^{-1} Y$$

- ▶ Case 3: $m > d + 1$, overdetermined. If X has full col rank, then, many solutions. Min norm solution

$$W = (\underbrace{X^T X})^{-1} X^T Y$$

- ▶ Can solve normal equations above directly, or use iterative methods for linear systems. Cost $\mathcal{O}(\underline{d^3})$



6/9

$$\underbrace{X^T X} W = \underbrace{X^T Y}$$

X
 $m \times d$

$$W = (X^T X)^{\dagger} X^T Y$$

d : input dim/
feature dim

Gauss Markov theorem

- ▶ Take noisy $y_i = x_i^\top W + \epsilon_i$, with $E\epsilon_i = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$; x_i is non-random.

Gauss Markov theorem

$$\omega^{\text{OLS}} = \underline{(X^T X)^{-1} X^T y}$$

- ▶ Take noisy $y_i = x_i^T W + \epsilon_i$, with $E\epsilon_i = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$; x_i is non-random.
- ▶ Statement: the OLS estimator is the best linear unbiased estimator (blue). It has the lowest variance.

Gauss Markov theorem

$$W^{OLS} = (X^T X)^{-1} X^T Y = E[\epsilon] = 0$$

- ▶ Take noisy $y_i = x_i^T W + \epsilon_i$, with $E\epsilon_i = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$; x_i is non-random.
- ▶ Statement: the OLS estimator is the best linear unbiased estimator (blue). It has the lowest variance.
- ▶ Proof: consider another linear estimator $W' = CY$, for some $(d+1) \times n$ matrix $C = (X^T X)^{-1} X^T + D$.

OLS is an unbiased estimator
 Let W be true value that makes
 $y_i = x_i^T W + \epsilon_i$

$$E[W^{OLS}] = E[(X^T X)^{-1} X^T Y]$$

$$= E[(X^T X)^{-1} X^T (XW + \epsilon)]$$

$$= \underbrace{W}_{\text{true value}} + \underbrace{E[\epsilon]}_0$$

$$= W$$

$$\text{Bias of } W^{OLS} = E[W^{OLS}] - W = 0$$

$$E[X + Y] = E[X] + E[Y]$$

Gauss Markov theorem

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_m \end{bmatrix} \quad \mathbb{E}[Y_i] = 0$$

- Take noisy $y_i = x_i^\top \underbrace{W}_{\text{true}} + \epsilon_i$, with $E\epsilon_i = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$; x_i is non-random.
- Statement: the OLS estimator is the best linear unbiased estimator (blue). It has the lowest variance.
- Proof: consider another linear estimator $W' = CY$, for some $(d+1) \times n$ matrix $C = \underline{(X^\top X)^{-1} X^\top} + \underline{D}$.
- For W' to be unbiased, show $DX = 0$. Then show,

$$\rightarrow \text{Var}(W') = \text{Var}(W) + \sigma^2 DD^\top$$

CY

$$\mathbb{E}[W'] = \mathbb{E}[CY] \quad \text{Var}(W'_i) \geq \text{Var}(W_i^{\text{OLS}})$$

$$\text{Var}(W) = \mathbb{E}[(W - \mathbb{E}[W])^2]_{1 \leq i \leq m}$$

$$\text{Var}(W') = \text{Var}(W^{\text{OLS}}) + \sigma^2 DD^\top$$

Gauss Markov theorem

- ▶ Take noisy $y_i = x_i^\top W + \epsilon_i$, with $E\epsilon_i = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$; x_i is non-random.
- ▶ Statement: the OLS estimator is the best linear unbiased estimator (blue). It has the lowest variance.
- ▶ Proof: consider another linear estimator $W' = CY$, for some $(d+1) \times n$ matrix $C = (X^\top X)^{-1}X^\top + D$.
- ▶ For W' to be unbiased, show $DX = 0$. Then show, $\text{Var}(W') = \text{Var}(W) + \sigma^2 DD^\top$.
- ▶ Since DD^\top is positive semi-definite, qed.

Ridge regression

- ▶ Motivation: unbiased estimation does not mean least mean-squared error

Ridge regression

Regularization

- Motivation: unbiased estimation does not mean least mean-squared error
- Let true $h_W \in \mathcal{H}$. Mean-squared error of statistical estimator of W , \hat{W} , is its variance + bias-squared

$$E[(\hat{W}_i - W_i)^2] = \text{Var}(\hat{W}_i) + (E[\hat{W}_i] - W_i)^2.$$

►

$$\operatorname{argmin}_W \frac{1}{m} \|XW - Y\|^2 + \lambda \|W\|^2.$$

Ex :

$$E((\hat{W}_i - W_i)^2) = \underbrace{\text{Var}(\hat{W}_i)}_{\downarrow} + \underbrace{(E[\hat{W}_i] - W_i)^2}_{\uparrow}$$

8/9

Ridge regression

- ▶ Motivation: unbiased estimation does not mean least mean-squared error
- ▶ Let true $h_W \in \mathcal{H}$. Mean-squared error of statistical estimator of W , \hat{W} , is its variance + bias-squared

$$E[(\hat{W}_i - W_i)^2] = \text{Var}(\hat{W}_i) + (E[\hat{W}_i] - W_i)^2.$$



$$\operatorname{argmin}_W \underbrace{\frac{1}{m} \|XW - Y\|^2}_{\text{variance}} + \underbrace{\lambda \|W\|^2}_{\text{bias-squared}}.$$

- ▶ penalizes l^2 norm of W . Still convex problem.

Ridge regression

- ▶ Motivation: unbiased estimation does not mean least mean-squared error
- ▶ Let true $h_W \in \mathcal{H}$. Mean-squared error of statistical estimator of W , \hat{W} , is its variance + bias-squared

$$E[(\hat{W}_i - W_i)^2] = \text{Var}(\hat{W}_i) + (E[\hat{W}_i] - W_i)^2.$$

▶

$$\text{Ridge regression}$$
$$\underline{\underset{W}{\operatorname{argmin}} \frac{1}{m} \|XW - Y\|^2 + \lambda \|W\|^2}.$$

- ▶ penalizes ℓ^2 norm of W . Still convex problem.
- ▶ to derive OLS, also can take derivative and set it to zero. Similarly here.

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$\|W\|^2 = \sum_{i=1}^d |w_i|^2$$

Ridge regression

- ▶ Motivation: unbiased estimation does not mean least mean-squared error
- ▶ Let true $h_W \in \mathcal{H}$. Mean-squared error of statistical estimator of W , \hat{W} , is its variance + bias-squared

$$E[(\hat{W}_i - W_i)^2] = \text{Var}(\hat{W}_i) + (E[\hat{W}_i] - W_i)^2.$$



$$\operatorname{argmin}_W \frac{1}{m} \|XW - Y\|^2 + \lambda \|W\|^2.$$

- ▶ penalizes l^2 norm of W . Still convex problem.
- ▶ to derive OLS, also can take derivative and set it to zero. Similarly here.
- ▶ Equivalent formulation: $\min_W \sum_{i=1}^m (w^\top \Phi(x_i) - y_i)^2$ subject to $\|w\|^2 \leq \Lambda^2$

Ridge regression optimization $d \gg m$

Primal :

$$\min_{w \in \mathbb{R}^d} (\|Xw - Y\|^2 + \lambda \|w\|^2)$$

Lagrangian



Dual:

$$\max_{\alpha \in \mathbb{R}^m} \underbrace{-\alpha^T (XX^T + \lambda I) \alpha + 2\alpha^T Y}_{\text{concave}}$$

$XX^T \qquad X^T X$

Closed form solution

$$W^{\text{Ridge}} = \underbrace{(X^T X + \lambda I)^{-1}}_{\text{always invertible}} X^T Y$$

Shrinkage

$$W^{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T Y$$

$$Y_{\text{pred}}^{\text{Ridge}} = X W^{\text{Ridge}}$$

$$= X (X^T X + \lambda I)^{-1} X^T Y$$

$$X_{m \times d} = U \Sigma V^T$$

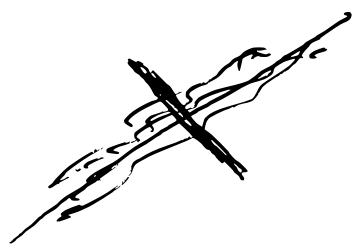
$$U = [u_1 | u_2 | \dots | u_d] \\ u_i \in \mathbb{R}^m$$

Ex:

$$Y_{\text{pred}}^{\text{Ridge}} = U \Sigma V^T (V^T \Sigma V + \lambda I)^{-1} V \Sigma V^T Y$$

$$= \sum_{i=1}^d \underbrace{u_i}_{\uparrow} \left(\frac{d_i^2}{\lambda + d_i^2} \right) \underbrace{u_i^T Y}$$

$$\Sigma = \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_d \end{bmatrix}$$



Ridge regression solution and interpretation

- ▶ Revisit convex optimization. Derive solution using KKT conditions.

Ridge regression solution and interpretation

- ▶ Revisit convex optimization. Derive solution using KKT conditions.
- ▶ Now simply use convexity and differentiability to obtain global minimum:

$$W = (X^T X + \lambda I)^{-1} X^T Y.$$

Ridge regression solution and interpretation

- ▶ Revisit convex optimization. Derive solution using KKT conditions.
- ▶ Now simply use convexity and differentiability to obtain global minimum:

$$W = (X^T X + \lambda I)^{-1} X^T Y.$$

- ▶ LASSO: with l^1 regularization.

Ridge regression solution and interpretation

- ▶ Revisit convex optimization. Derive solution using KKT conditions.
- ▶ Now simply use convexity and differentiability to obtain global minimum:

$$W = (X^T X + \lambda I)^{-1} X^T Y.$$

- ▶ LASSO: with l^1 regularization.
- ▶ Generalization bounds for bounded regression problems.

Ridge regression solution and interpretation

- ▶ Revisit convex optimization. Derive solution using KKT conditions.
- ▶ Now simply use convexity and differentiability to obtain global minimum:

$$W = (X^T X + \lambda I)^{-1} X^T Y.$$

- ▶ LASSO: with l^1 regularization.
- ▶ Generalization bounds for bounded regression problems.
- ▶ Shrinkage by l^2 regularization.