

SDS 387 Linear Models

Fall 2024

Lecture 16 - Thu, Oct 24, 2024

Instructor: Prof. Ale Rinaldo

- Announcement: I will post HW4 soon (and also post solutions to HW3)

- Last time: Ordinary Least Squares (OLS) estimator

- Setting:

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

vector of response variables

$$\mathbb{D} \quad n \times d$$

design matrix whose i^{th} row contains the i^{th} observation for the covariates (a point in \mathbb{R}^d)

Notation used in
Bach's book

- The notation Φ is non-standard in statistics but reflect the common practice in ML of turning a vector of covariates, say x_i , into a vector of features $\Phi_i \in \mathbb{R}^d$, $\Phi_i = \varphi(x_i)$.

- I will assume throughout that the first column of Φ is a vector of 1's.

This means that we always fit an intercept to our linear model

- The OLS estimator is obtained as the minimizer of the empirical risk:

$$\beta \rightarrow \hat{R}(\beta) = \frac{1}{2n} \|Y - \Phi\beta\|_2^2$$

Assuming that Φ has full-column rank ($\text{rank}(\Phi) = d$) then the solution exists and is unique and given by

$$\hat{\beta} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

$$= \sum_1^{-1} \frac{\Phi^T Y}{n} \quad \text{where } \sum_1 = \frac{\Phi^T \Phi}{n}$$

Importantly $\hat{\beta}$ satisfies the normal equations

$$\Phi^T \Phi \hat{\beta} = \Phi^T Y$$

Geometric interpretation: From the formula for $\hat{\beta}$ the vector of fitted values

$$\hat{y} = \Phi \hat{\beta} = \underbrace{\Phi (\Phi^T \Phi)^{-1} \Phi^T}_{H} y$$

is a linear function of y

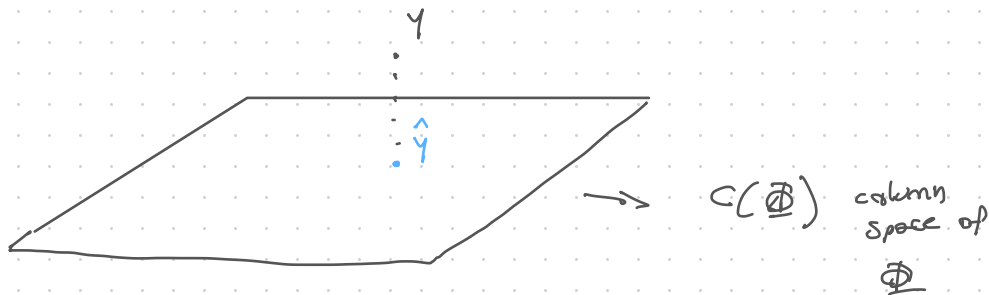
H = "hat matrix"

= prediction of the responses given by the model or the estimated value of the regression function

Note: H is a projection matrix orthogonal

$$H^2 = H \quad \text{and} \quad H^T = H$$

H projects y onto the d -dimensional linear subspace of \mathbb{R}^n spanned by the columns of Φ



\hat{y} is the point in $C(\Phi)$ that is closest to y in \mathbb{R}^n

From this we can see that the residuals

$$e = y - \hat{y} = (I - H) y \in \mathbb{R}^n$$

are the orthogonal projection of y

onto the orthogonal complement of $C(\Phi)$

Remark $I-H$ is also an orthogonal projection

and

$$\langle e, \hat{y} \rangle = \langle (I-H)y, Hy \rangle = 0$$

$$\hookrightarrow y = \underbrace{\hat{y}}_{\text{orthogonal}} + e$$

$$\hookrightarrow \|y\|^2 = \|\hat{y}\|^2 + \|e\|^2$$

\downarrow
= energy or variability explained by model

Numerical considerations

How hard is to compute $\hat{\beta}$ numerically?

To compute $\hat{\beta}$ one has to invert $\underbrace{\Phi^T \Phi}_{d \times d}$
typically requires order $\mathcal{O}(d^3)$ computations.

Gradient descent (Bach 5.2.1)

Starting from an initial point $\beta_0 \in \mathbb{R}^d$ consider the sequence of updates

$$\beta_t = \beta_{t-1} - \gamma \nabla_{\mathbb{R}^d} \ell(\beta_{t-1}) \quad (4)$$

$\hookrightarrow \gamma > 0$ stepsize

where recall that

$$\hat{R}(\beta) = \frac{1}{2n} \|Y - \Phi\beta\|^2 \quad \text{and}$$

$$\nabla \hat{R}(\beta) = \frac{1}{n} \Phi^T \Phi \beta - \frac{\Phi^T Y}{n}$$

$$\text{Hessian} \swarrow \quad \text{Hess } \hat{R}(\beta) = \frac{1}{n} \Phi^T \Phi = \sum_{i=1}^n$$

Also recall that any minimizer, say $\tilde{\beta}$, of \hat{R} satisfies

$$\sum_{i=1}^n \tilde{\beta} = \frac{\Phi^T Y}{n}$$

A solution exists always and is unique if $\Phi^T \Phi$ is invertible. If not, there exist infinitely many solutions

First off, notice that, if solution β^* exists, then

$$\hat{R}(\beta) - \hat{R}(\beta^*) = \frac{1}{2n} \|Y - \Phi\beta\|^2 - \frac{1}{2n} \|Y - \Phi\beta^*\|^2$$

$$= \frac{1}{2n} \|Y - \Phi\beta^* + \Phi(\beta^* - \beta)\|^2 - \frac{1}{2n} \|Y - \Phi\beta^*\|^2$$

$$= \frac{1}{2n} \|Y - \Phi\beta^*\|^2 + \frac{1}{2n} \|\Phi(\beta^* - \beta)\|^2 + \frac{2}{2n} \langle Y - \Phi\beta^*, \Phi(\beta^* - \beta) \rangle$$

$$- \frac{1}{2n} \|Y - \Phi\beta^*\|^2$$

$$= 0 \quad \text{because} \\ Y - \Phi\beta^* \in C(\Phi)^\perp \quad \textcircled{5}$$

$$\begin{aligned}
 &= \frac{1}{2n} \|\Phi(\beta^* - \beta)\|^2 = \frac{(\beta - \beta^*)^T \sum_1^n (\beta - \beta^*)}{2} \\
 &= \frac{\|\beta - \beta^*\|^2}{2} \sum_1^n
 \end{aligned}$$

Next, let's look at the gradient iterates:

$$\begin{aligned}
 \beta_t &= \beta_{t-1} - \gamma \nabla \hat{R}(\beta_{t-1}) = \beta_{t-1} - \gamma \left[\frac{1}{n} \Phi^T (\Phi \beta_{t-1} - y) \right] \\
 &= \beta_{t-1} - \gamma \sum_1^n (\beta_{t-1} - \beta^*)
 \end{aligned}$$

↓

because $\sum_1^n \beta^* = \frac{\Phi^T y}{n}$

$$\beta_t - \beta^* = (\mathbf{I} - \gamma \sum_1^n) (\beta_{t-1} - \beta^*)$$

This implies that

$$\|\beta_t - \beta^*\|^2 = (\beta_0 - \beta^*)^T (\mathbf{I} - \gamma \sum_1^n)^{2t} (\beta_0 - \beta^*)$$

and

$$\hat{R}(\beta_t) - \hat{R}(\beta^*) = \frac{(\beta_t - \beta^*)^T \sum_1^n (\beta_t - \beta^*)}{2}$$

$$= \frac{1}{2} (\beta_0 - \beta^*)^T (\mathbf{I} - \gamma \sum_1^n) \underbrace{\sum_1^n (\mathbf{I} - \gamma \sum_1^n)^t}_{\text{commute HW!}} (\beta_0 - \beta^*)$$

(6)

$$= \frac{1}{2} (\beta_0 - \beta^*)^T (I - \gamma \hat{\Sigma})^{2t} \hat{\Sigma} (\beta_0 - \beta^*)$$

Let's first look at convergence to the minimizer. Assume that $\hat{\Sigma}$ is invertible. The eigenvalues of

$$(I - \gamma \hat{\Sigma})^{2t} \text{ are of the form } (1 - \gamma \lambda)^{2t}$$

where λ is an eigenvalue of $\hat{\Sigma}$. So all the

eigenvalues of $(I - \gamma \hat{\Sigma})^{2t}$ are less than

$$\max_{\lambda} |1 - \gamma \lambda|^{2t}$$

$$\lambda_{\min}(\hat{\Sigma}) \leq \lambda \leq \lambda_{\max}(\hat{\Sigma})$$

Now let's choose γ to be $\frac{1}{\lambda_{\max}(\hat{\Sigma})}$. Then

the expression above is equal to

$$\left(1 - \frac{\lambda_{\min}(\hat{\Sigma})}{\lambda_{\max}(\hat{\Sigma})}\right)^{2t} = \left(1 - \frac{1}{\kappa}\right)^{2t}$$

$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$
 \downarrow
condition number of $\hat{\Sigma}$

Putting everything together:

$$\|\beta_{2t} - \beta^*\|^2 \leq \left(1 - \frac{1}{\kappa}\right)^{2t} \|\beta_0 - \beta^*\|^2$$

$$\leq e^{-2t/\kappa} \|\beta_0 - \beta^*\|^2$$

exponential / geometric / linear convergence

If $\kappa = \infty$ (i.e. $\text{dmax}(\hat{\Sigma}) = 0$) this will simply say that $\|\beta_t - \beta^*\|^2 \leq \|\beta_0 - \beta^*\|^2$ all t .

Let's look at the convergence of the objective function \hat{R}

$$\hat{R}(\beta_t) - \hat{R}(\beta^*) = \frac{1}{2} (\beta_0 - \beta^*)^T (\mathbb{I} - \gamma \hat{\Sigma})^{2t} \hat{\Sigma} (\beta_0 - \beta^*)$$

$$= \frac{1}{2} \text{tr} \left((\mathbb{I} - \gamma \hat{\Sigma})^{2t} \hat{\Sigma} (\beta_0 - \beta^*) (\beta_0 - \beta^*)^T \right)$$

$x^T A x = \text{tr}(A x x^T)$

because

$\text{tr}(AB) \leq \|A\|_{\text{op}} \text{tr}(B)$
 \downarrow
 pd



$$\leq \frac{1}{2} \|\mathbb{I} - \gamma \hat{\Sigma}\|_{\text{op}}^{2t} \text{tr} \left(\hat{\Sigma} (\beta_0 - \beta^*) (\beta_0 - \beta^*)^T \right)$$

$$= \frac{1}{2} \max_{\lambda \text{ eig. of } \hat{\Sigma}} |1 - \gamma \lambda|^{2t} \underbrace{\left[(\beta_0 - \beta^*)^T \hat{\Sigma} (\beta_0 - \beta^*) \right]}_{\hat{R}(\beta_0) - \hat{R}(\beta^*)}$$

if $\gamma = \frac{1}{\text{dmax}(\hat{\Sigma})}$

$$\leq \exp\left\{-\frac{2t}{\kappa}\right\}$$

$$\leq \frac{e^{-2t/\kappa}}{2} \left[\hat{R}(\beta_0) - \hat{R}(\beta^*) \right]$$

- What if $k = \infty$ (i.e. $\hat{\Sigma}^t$ is not invertible)?

$$\hat{R}(\beta_t) - \hat{R}(\beta^*) \leq \frac{1}{2} \|(\mathbb{I} - \gamma \hat{\Sigma}^t)^{2t} \hat{\Sigma}^t\|_{\text{op}} \|\beta_0 - \beta^*\|^2$$

setting $\delta = \frac{1}{\lambda_{\max}(\hat{\Sigma}^t)}$

$$\leq \max_{\lambda} \delta (1 - \gamma \lambda)^{2t} \|\beta_0 - \beta^*\|^2$$

∴

$$\leq \frac{\delta \max_{\lambda} x}{8t} \|\beta_0 - \beta^*\|^2$$



convergence is polynomial in t