

BIOS 835: Linear Regression and Bayesian Decision Theory

Alexander McLain

August 29, 2023

Outline

Introduction

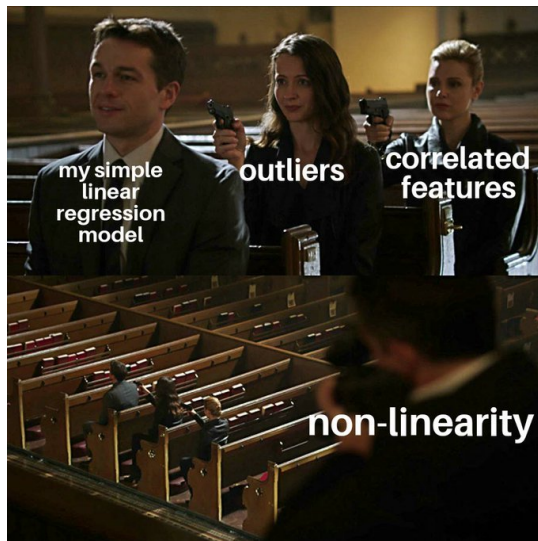
Least squares solution

Inference in linear regression

Multiple linear regression

Multivariate Regression

Bayesian Decision Theory



Linear Regression

- ▶ Linear regression models were developed before computers were around, but are still popular today.
- ▶ They are simple and (commonly) adequate and interpretable representations of the real world.
- ▶ They commonly work as well as difficult non-linear models for prediction, especially when there is limited data and low signal-to-noise ratio.
- ▶ Linear models do have assumptions, but many properties are robust to these assumptions (plus we can transform outcomes).

Linear Regression model

- ▶ A linear regression model is defined as

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon \\ &= \beta_0 + \sum_{j=1}^p \beta_j X_j + \epsilon \\ &= f(\mathbf{X}) + \epsilon \end{aligned}$$

- ▶ Y is our outcome of interest
- ▶ (X_1, X_2, \dots, X_p) are our covariates of interest
- ▶ β_0 : intercept;
- ▶ $\beta_j, j = 1, \dots, p$: regression coefficients.
- ▶ $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2$.

Linear Regression model

- We can also, define it in matrix form with the following

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \dots & X_{1p} \\ 1 & X_{21} & \dots & X_{2p} \\ \vdots & \vdots & & \\ 1 & X_{n1} & \dots & X_{np} \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

as the covariate data, observations, error term, and parameters in a linear regression model.

- The model can be expressed succinctly as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

Least squares and RSS

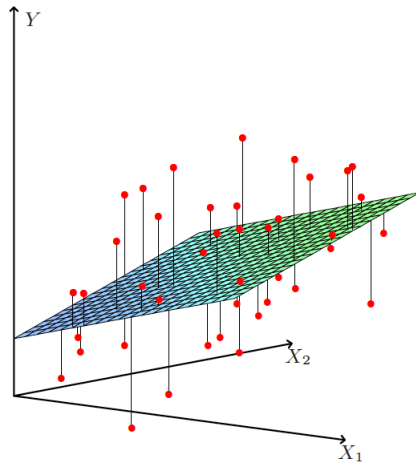
- ▶ The residual sum of squares for β is defined as

$$\begin{aligned}RSS(\beta) &= \sum_{i=1}^n (y_i - \mathbf{X}_i' \beta)^2 \\&= \sum_{i=1}^n \{y_i - f(\mathbf{X}_i)\}^2 \\&= (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2\end{aligned}$$

where $\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{ip})'$

- ▶ Using least squares criteria the estimate of β , denoted by $\hat{\beta}$, is the value that minimizes the RSS .

Least squares and RSS



The Least Squares Solution

- To minimize $RSS(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$ we differentiate with respect to β to get

$$\begin{aligned}\frac{\partial RSS(\beta)}{\partial \beta} &= -2\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) \\ \frac{\partial^2 RSS(\beta)}{\partial \beta \partial \beta'} &= 2\mathbf{X}'\mathbf{X}.\end{aligned}$$

- If \mathbf{X} has full column rank, setting the first derivative to zero yields,

$$\begin{aligned}0 &= \mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) \quad \text{and} \\ \hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.\end{aligned}$$

Predicted values

- ▶ Predicted values at a point \mathbf{x}_0 are given by $\hat{f}(\mathbf{x}_0) = \mathbf{x}_0' \hat{\beta}$.
- ▶ The predicted values for the training data are given by

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta} = \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}.$$

- ▶ $H = \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$ is referred to as the “Hat Matrix”
- ▶ Geometrically, $\hat{\beta}$ is chosen such that $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to the subspace of \mathbf{X} .

Rank Deficiencies

- ▶ Recall, we said “if \mathbf{X} has full column rank”
- ▶ When might this not be the case?

“Assumptions” standard regression models

Linear, Independent, Normality, and Equal variance (LINE)

- ▶ L: a linear relationship exists between Y and the X 's.
- ▶ I: the data are independent.
- ▶ N: the residuals are normally distributed.
- ▶ E: the variance of the residuals is equal for all X 's.

Properties of $\hat{\mathbf{y}}$, β and σ^2

- ▶ The variance of $\hat{\beta}$ is

$$\text{Var}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1}\sigma^2$$

which is estimated using

$$\hat{\sigma}^2 = \frac{1}{n - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

- ▶ Neither of these results requires any additional assumptions.

Properties of $\hat{\mathbf{y}}$, β and σ^2

Once we impose the LINE assumptions we get the following properties

- ▶ $\hat{\beta} \sim MVN(\beta, (X'X)^{-1}\sigma^2)$, hypothesis tests and confidence intervals for $\hat{\beta}_j$ can be based on

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma} v_j} \sim t_{n-p-1}$$

where v_j is the j th diagonal element of $(X'X)^{-1}$

- ▶ The predicted value at \mathbf{x}_0 is $y \sim N\{f(\mathbf{x}_0), \sigma^2(1 + \mathbf{x}_0(X'X)^{-1}\mathbf{x}_0)\}$

Partial F-test for multiple variables

- ▶ Suppose we'd like to test $H_0 : \beta_j = \beta_{j+1} = \dots = \beta_{j+k-1} = 0$ v.s $H_1 : \beta_l \neq 0$ for some $l = j, \dots, j+k-1$ given the other $(p-k)$ variables.
- ▶ This can be tested via

$$F_0 = \frac{(RSS_0 - RSS_1)/(p_1 - p_0)}{RSS_1/(N - p_1 - 1)}$$

where RSS_1 and $p_1 + 1$ are the RSS and number of parameters for the bigger model (similarly for RSS_0 and p_0).

- ▶ Under the Gaussian assumption and the null hypotheses $F_0 \sim F_{p_1-p_0, N-p_1-1}$
- ▶ Reject if

$$F_0 > F_{p_1-p_0, N-p_1-1}(1 - \alpha)$$

Gauss-Markov Theorem

- ▶ The Gauss-Markov Theorem is one of the most famous results in statistics.
- ▶ It states that the least squares estimate $\hat{\beta}$ has the smallest variance among all linear unbiased estimates.
- ▶ Specifically, suppose we wish to estimate $\theta = a'\beta$. The least squares estimate is

$$\hat{\theta} = a'\hat{\beta} = a'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

which is unbiased when the linear model is correct.

- ▶ The Gauss-Markov Theorem shows that another estimator $\tilde{\theta} = \mathbf{c}'\mathbf{y}$, which is unbiased for $a'\beta$, will have

$$\text{Var}(a'\hat{\beta}) \leq \text{Var}(\mathbf{c}'\mathbf{y})$$

Gauss-Markov Theorem

- ▶ What does the Gauss-Markov Theorem not say?
- ▶ Consider

$$MSE(\hat{\theta}) = E\{(\hat{\theta} - \theta)^2\}$$

Simple (simple) LR and notation

- ▶ Consider the following (very) simple LR model with 1 covariate and no intercept

$$Y = X\beta + \epsilon$$

- ▶ The LS estimate and residual are

$$\hat{\beta} = \frac{\sum_1^n x_i y_i}{\sum_1^n x_i^2}$$

$$r_i = y_i - x_i \hat{\beta}$$

- ▶ It is convenient (and more general) to use the following notation

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i = \mathbf{x}' \mathbf{y}, \quad \text{so} \quad \hat{\beta} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \frac{\mathbf{x}' \mathbf{y}}{\mathbf{x}' \mathbf{x}}$$

MLR for orthogonal vectors

- ▶ Suppose we have $1 = x_0, x_1, \dots, x_p$ inputs are all orthogonal
 - ▶ $\langle \mathbf{x}_j, \mathbf{x}_k \rangle = 0$ for all $j \neq k$

- ▶ Then

$$\hat{\beta}_j = \frac{\langle \mathbf{x}_j, \mathbf{y} \rangle}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle}$$

- ▶ What are the impacts of this?
- ▶ Will this ever happen?

Simple Gram-Schmidt process

- ▶ Now, suppose we have the SLR model

$$Y = \beta_0 + X\beta_1 + \epsilon$$

How can we make 1 and X orthogonal?

- ▶ By orthogonalizing we get the coefficient as

$$\hat{\beta}_1 = \frac{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle}{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{x} - \bar{x}\mathbf{1} \rangle}$$

which is the coefficient for the model $Y = (x - \bar{x})\beta_1 + \epsilon$

Simple Gram-Schmidt process

► So a way to calculate β_1 is to do the following:

1. regress \mathbf{x} on $\mathbf{1}$, and get the residual $\mathbf{z} = \mathbf{x} - \bar{x}\mathbf{1}$
2. regress \mathbf{y} on the residual \mathbf{z} to give the coefficient $\hat{\beta}_1$

where regress \mathbf{y} on \mathbf{z} means to estimate the “no intercept model” we initially started with.

Gram-Schmidt algorithm

We can generalize this for p non-orthogonal inputs $\mathbf{1} = x_0, x_1, \dots, x_p$

1. Set $\mathbf{z}_0 = \mathbf{x}_0 = \mathbf{1}$
2. For $j = 1, 2, \dots, p$:
regress \mathbf{x}_j on $\mathbf{z}_0, \dots, \mathbf{z}_{j-1}$, and get the residual $\mathbf{z}_j = \mathbf{x}_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} \mathbf{z}_k$ where $\hat{\gamma}_{kj} = \langle \mathbf{z}_k, \mathbf{x}_j \rangle / \langle \mathbf{z}_k, \mathbf{z}_k \rangle$.
3. regress \mathbf{y} on the residual \mathbf{z}_p to give the coefficient

$$\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle}$$

Implications

- ▶ Another way to express the variability in $\hat{\beta}_p$ is

$$\text{Var}(\hat{\beta}_p) = \frac{\sigma^2}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle} = \frac{\sigma^2}{\|\mathbf{z}_p\|^2}$$

- ▶ What will \mathbf{z}_p look like if \mathbf{x}_p is accurately predicted from $x_0, x_1, \dots, x_{(p-1)}$?
- ▶ What about $\|\mathbf{z}_p\|^2$?
- ▶ How do we get back to the original \mathbf{x} coefficients?
 - ▶ Let $\mathbf{x}_{0,1,\dots,(j-1),(j+1),\dots,p}$ denote the residual after regression \mathbf{x}_j on $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{(j-1)}, \mathbf{x}_{(j+1)}, \dots, \mathbf{x}_p$
 - ▶ The j th multiple regression coefficient is the univariate regression coefficient of \mathbf{y} on $\mathbf{x}_{0,1,\dots,(j-1),(j+1),\dots,p}$.

Matrix GS with QR decomp

- ▶ In matrix form, we can represent the second step of the GS algorithm with

$$\mathbf{X} = \mathbf{Z}\mathbf{\Gamma}$$

where \mathbf{Z} has columns \mathbf{z}_j , and $\mathbf{\Gamma}$ is the upper triangular matrix with entries $\hat{\gamma}_{kj}$.

- ▶ Let's add in \mathbf{D} a diagonal matrix with entries $D_{jj} = \|\mathbf{z}_j\|$ to get

$$\begin{aligned}\mathbf{X} &= \mathbf{Z}\mathbf{D}^{-1}\mathbf{D}\mathbf{\Gamma} \\ &= \mathbf{Q}\mathbf{R}\end{aligned}$$

The \mathbf{QR} decomposition of \mathbf{X}

Matrix GS with QR decomp

- Using the **QR** decomposition of **X** we get the following form for $\hat{\beta}$ and \hat{y}

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \{(\mathbf{QR})'\mathbf{QR}\}^{-1}(\mathbf{QR})'\mathbf{y} = \mathbf{R}^{-1}\mathbf{Q}'\mathbf{y}, \\ \hat{y} &= \mathbf{Q}\mathbf{Q}'\mathbf{y}\end{aligned}$$

This form of $\hat{\beta}$ is a lot easier to solve because **R** is upper triangular.

Introduction to Multivariate Regression

- Now let's consider the more classical multivariate case where we have Y_1, Y_2, \dots, Y_K and $X_0, \mathbf{X}_1, \dots, \mathbf{X}_p$ and the linear model

$$Y_k = \beta_{0k} + \sum_{j=1}^p X_j \beta_{jk} + \epsilon_k = f_k(\mathbf{X}) + \epsilon_k \quad \text{for } k = 1, 2, \dots, K$$

or in matrix form

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$$

where \mathbf{Y} is $n \times K$ are the outcomes, \mathbf{X} is $n \times (p + 1)$ are the inputs, \mathbf{B} is an $(p + 1) \times K$ matrix of β coefficients and \mathbf{E} is an $n \times K$ matrix of residuals.

RSS and solution

- ▶ The RSS for the multivariate case is given by

$$\begin{aligned}RSS(\mathbf{B}) &= \sum_{k=1}^K \sum_{i=1}^n \{y_{ik} - f_k(x_i)\}^2 \\&= \text{tr}\{(\mathbf{Y} - \mathbf{XB})'(\mathbf{Y} - \mathbf{XB})\}\end{aligned}$$

- ▶ It's rather straightforward to show that $RSS(\mathbf{B})$ is minimized with

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

Implications

- ▶ More specifically, for the k th outcome the coefficient is

$$\hat{\beta}_k = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}_k$$

- ▶ If we were to use the general linear model where $\text{Cov}(\epsilon) = \Sigma$ and

$$RSS(\mathbf{B}, \Sigma) = \text{tr}\{(\mathbf{Y} - \mathbf{XB})'\Sigma^{-1}(\mathbf{Y} - \mathbf{XB})\},$$

..... the solution is still the same.

- ▶ What does this imply about doing regression with multiple outcomes?

Introduction

- ▶ **Bayesian Decision Theory** is a decision-making framework that combines probability theory and decision theory.
- ▶ It is based on the Bayesian approach to probability, where prior beliefs are updated with new evidence in the form of data to produce posterior beliefs.
- ▶ Bayesian Decision Theory provides a formal mechanism for making decisions under uncertainty, weighing the risks of different actions against the posterior beliefs about the state of the world.

Probabilistic Setup

- ▶ Let θ be a parameter or state of nature we want to learn about (it could represent various things: the true state of a system, a model parameter, etc.)
- ▶ $p(\theta)$ is the prior probability distribution over θ which reflects our beliefs about θ before observing data.
- ▶ $p(x|\theta)$ is the likelihood function which provides the probability of observing data x given θ .
- ▶ $p(\theta|x)$ is the posterior probability distribution over θ given data x , which is computed using Bayes' theorem:

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$$

where $p(x)$ is the evidence or marginal likelihood.

Decision Making

- ▶ Suppose we have a set of possible actions A we can take.
- ▶ A loss function $L(\theta, a)$ represents the cost or loss of taking action a when the true state is θ .
- ▶ The loss function quantifies the “penalty” for making decisions; related to the norms we discussed but more general (e.g., for binary data).
- ▶ The goal in Bayesian decision theory is to minimize the expected loss. The decision rule is given by:

$$a^* = \arg \min_a \mathbb{E}[L(\theta, a)|x]$$

where the expectation is over the posterior distribution of θ given data x .

Example: Binary Hypothesis Testing

Let's consider a simple case where θ can be either H_0 or H_1 .

- ▶ Prior Probabilities: $p(H_0)$ and $p(H_1)$.
- ▶ Likelihoods: $p(x|H_0)$ and $p(x|H_1)$.
- ▶ Actions: here we'll consider $A \in \{a_0, a_1\}$, i.e., a “hard” decision.
- ▶ Loss function might be defined such that $L(H_0, a_0) = 0$ (no loss for correctly deciding H_0), $L(H_0, a_1) = 1$ (loss for wrongly deciding H_1 when truth is H_0), and similar definitions for decisions regarding H_1 .
 - ▶ if $A \in (0, 1)$ the other loss functions would be used.

Example: Binary Hypothesis Testing

- The Bayesian decision rule will decide a_1 , i.e., H_1 , if:

$$\frac{p(x|H_1)p(H_1)}{p(x)} > \frac{p(x|H_0)p(H_0)}{p(x)}$$

This simplifies to:

$$\frac{p(x|H_1)}{p(x|H_0)} > \frac{p(H_0)}{p(H_1)}$$

Here, $\frac{p(x|H_1)}{p(x|H_0)}$ is the likelihood ratio and $\frac{p(H_0)}{p(H_1)}$ is the prior odds ratio.