

## 3.0 Linear Regression with Matrices

- Answer Questions
- Assumptions
- Maximum Likelihood Estimators
- A Second Proof of  $\hat{\beta}$
- Residuals
- Prediction and Confidence Intervals
- Bases and Projections
- More on the Hat Matrix

## 3.1 Assumptions

Multiple linear regression assumes that:

- The vectors of explanatory variables  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$ , are not random.
- Given  $\mathbf{x}_i$ , the response  $Y_i$  is normally distributed.
- $\mathbb{E}[Y_i] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ .
- The variance of  $Y_i$ , given  $\mathbf{x}_i$ , is  $\sigma^2$  (usually unknown, but the same for all  $i$ ).
- The  $Y_1, \dots, Y_n$  are independent given the  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

The statistical model that satisfies these assumptions is called the **General Linear Model**.

In matrix notation, we write

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

where  $\mathbf{y}$  is  $n \times 1$ ,  $\mathbf{X}$  is  $n \times (p + 1)$ ,  $\boldsymbol{\beta}$  is  $(p + 1) \times 1$  and  $\boldsymbol{\epsilon}$  is  $n \times 1$ .

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

## 3.2 MLEs

In the general linear model, the distribution of the observations  $Y_i$  has density

$$f(y_i) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left[ -\frac{1}{2\sigma^2} (y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2 \right].$$

To see this, realize that the independence of the  $Y_i$  values implies their joint density is just the product of each density; i.e.,

$f(\mathbf{y}) = \prod_{i=1}^n f(y_i)$ , where each term in the product is normal with mean  $\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}$  and common variance  $\sigma^2$ .

We want to find the maximum likelihood estimates of  $\beta_0, \dots, \beta_p$ .

Inspection of the joint density shows it is sufficient to find the estimates that minimize the quadratic form

$$Q = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2.$$

The values that minimize  $Q$  are sometimes called the **Least Squares Estimates**, which here, but not in general, agree with the maximum likelihood estimates.

To solve, find the partial derivatives  $\partial Q / \partial \beta_j$  for  $j = 0, \dots, p$ . Set these partials to zero, giving  $p + 1$  equations in  $p + 1$  unknowns. Solving these gives the estimates  $\hat{\beta}_0, \dots, \hat{\beta}_p$ .

The **Design Matrix** is the  $n \times (p + 1)$  matrix  $\mathbf{X}$  whose  $i$ th row is  $(1, x_{i1}, \dots, x_{ip})$  for  $i = 1, \dots, n$ . The name comes from the fact that in some multiple regressions, the experimenter gets to select the values of the explanatory variables by design. When this is possible, one can learn more efficiently about the relationship between the response variable and the explanatory variables.

For now, assume that the set of  $p + 1$  equations has a unique solution. This implies the design matrix is full rank; i.e., has rank  $p + 1$ .

The rank of an  $n \times m$  matrix is

- the number of linearly independent rows, or
- the number of linearly independent columns, or
- the dimensions of the largest square submatrix that is invertible (nonsingular).

Some facts:

- $\text{Rank}(\mathbf{AB}) \leq \min\{\text{Rank}(\mathbf{A}), \text{Rank}(\mathbf{B})\}$ .
- If  $\mathbf{A}$  is  $n \times n$  and has determinant equal to 0, then  $\text{Rank}(\mathbf{A}) < n$ .
- If  $\mathbf{A}$  is  $m \times n$ , then  $\mathbf{A}^\top$  is the  $n \times m$  **transpose** that turns the rows into columns.

Let  $\mathbf{X}$  denote the design matrix,  $\boldsymbol{\beta}$  the true but unknown coefficients,  $\mathbf{y}$  the response vector, and  $\hat{\boldsymbol{\beta}}$  the maximum likelihood estimates of the coefficients. We shall show that

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

Define the inverse  $\mathbf{A}^{-1}$  as the square matrix such that  $\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$ , the identity matrix that is symmetric with ones on the diagonal and zeroes elsewhere.

**Proof:** Write  $Q = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ . Then

$$\frac{\partial Q}{\partial \beta_k} = -2 \sum_{i=1}^n x_{ik} y_i + 2 \sum_{j=0}^p \beta_j \left( \sum_{i=1}^n x_{ik} x_{ij} \right).$$

Set the right-hand sides of the  $p + 1$  equations to zero, and then rewrite the equations as  $\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}$ . Because we assumed that  $\mathbf{X}^\top \mathbf{X}$  had rank  $p + 1$ , then from linear algebra we know it is invertible, and thus

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad \blacksquare$$

Let  $S^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \cdots - \hat{\beta}_p x_{ip})^2$ . Then one can show that the maximum likelihood estimate of  $\sigma^2$  is  $\hat{\sigma}^2 = S^2/n$ , but this is biased. An unbiased estimate is  $\hat{\sigma}^2 = S^2/(n - p - 1)$ .



**Theorem:**  $\mathbb{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ .

**Proof:**  $\mathbb{E}[\hat{\boldsymbol{\beta}}] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}[\mathbf{y}] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\beta}) = \boldsymbol{\beta}$ .

**Theorem:**  $\text{Var}[\hat{\boldsymbol{\beta}}] = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2$ .

**Proof:** Recall, from linear algebra, that if  $\mathbf{A} = \mathbf{BC}$ , then  $\mathbf{A}^\top = \mathbf{C}^\top \mathbf{B}^\top$ .

$$\begin{aligned} \text{Var}[\hat{\boldsymbol{\beta}}] &= \text{Var}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}] \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}[\mathbf{y}] [(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top]^\top \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2. \end{aligned}$$

This uses the linear algebra fact that  $\mathbf{X}^\top \mathbf{X}$  is symmetric, so its inverse is symmetric, so the transpose of the inverse is itself.

## 3.3 A Proof of $\hat{\beta}$

Some facts:

- If  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , then  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}$ .
- If  $\mathbf{A}$  is symmetric and  $\mathbf{y} = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ , then  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = 2\mathbf{x}^\top \mathbf{A}$ .
- If  $\mathbf{A}$  is symmetric, the derivative of the transpose is the transpose of the derivative.
- $(\mathbf{A}\mathbf{B})^\top = \mathbf{B}^\top \mathbf{A}^\top$ .

The rules for differentiating vectors and matrices look very much like the usual calculus for univariate quantities.

**Theorem:** If  $\mathbf{X}$  has rank  $p + 1$ , then  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

**Proof:** Recall that  $Q = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ . Thus  $Q = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}$ . Taking the derivative of  $Q$  with respect to  $\boldsymbol{\beta}$  and using the facts shows that the derivative of the first term is zero, the second term goes to  $\mathbf{y}^\top \mathbf{X}$ , the third term goes to  $(\mathbf{X}^\top \mathbf{y})^\top = \mathbf{y}^\top \mathbf{X}$ , and the last term has derivative  $\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}$ .

Set the derivative to zero and solve. One gets  $0 = -2\mathbf{y}^\top \mathbf{X} + 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}$  so  $\boldsymbol{\beta}^\top (\mathbf{X}^\top \mathbf{X}) = \mathbf{y}^\top \mathbf{X}$ . Take the transpose of both sides to get  $\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}$ . Since  $\mathbf{X}$  has rank  $p + 1$ , then  $\mathbf{X}^\top \mathbf{X}$  is invertible, and so  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ . ■

## 3.4 Residuals

What else can we do with matrix algebra in regression?

- Derive properties of regression diagnostics, such as residuals.
- Calculate prediction intervals for outcomes at a new value  $\mathbf{x}$ .
- Interpret the MLEs geometrically.

The  $i$ th residual  $\hat{\epsilon}_i$  is the difference between the actual value of  $y_i$  in the training sample and the predicted value  $\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$ . In matrix notation, this is

$$\begin{aligned}\hat{\boldsymbol{\epsilon}} &= \mathbf{y} - \hat{\mathbf{y}} \\ &= \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y}.\end{aligned}$$

**Theorem:**  $\mathbb{E}[\hat{\boldsymbol{\epsilon}}] = \mathbf{0}$ .

**Proof:**  $\hat{\boldsymbol{\epsilon}} = (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y}$  so

$$\begin{aligned} \mathbb{E}[\hat{\boldsymbol{\epsilon}}] &= \mathbb{E}[(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{y}] \\ &= (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbb{E}[\mathbf{y}] \\ &= (\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{X} \boldsymbol{\beta} \\ &= \mathbf{X} \boldsymbol{\beta} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top * \mathbf{X} \boldsymbol{\beta} \\ &= \mathbf{X} \boldsymbol{\beta} - \mathbf{X} \boldsymbol{\beta} \\ &= \mathbf{0}. \quad \blacksquare \end{aligned}$$

Of course, this is completely unsurprising.

Consider a random vector  $\mathbf{y} \in \mathbb{R}^n$  with mean  $\boldsymbol{\mu}$ . Its covariance matrix,  $\text{Cov}(\mathbf{y})$ , is  $\boldsymbol{\Sigma}$  where

$$\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})^\top].$$

This definition does not require multivariate normality.

Here  $\boldsymbol{\Sigma}$  is  $n \times n$  and symmetric; i.e.,  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^\top$ . Its diagonal entries are the variances of the  $i$  random variable  $y_i$ . The  $(i, j)$ th entry is the covariance between the random variables  $y_i$  and  $y_j$ .

**Theorem:**  $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{y} \mathbf{y}^\top] - \boldsymbol{\mu} \boldsymbol{\mu}^\top$ .

The proof is an easy exercise. The logic is exactly like the corresponding univariate case, with  $\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$ .

If  $\mathbf{A}$  is not random, it is pretty clear that  $\mathbb{E}[\mathbf{A}\mathbf{y}] = \mathbf{A}\mathbb{E}[\mathbf{y}] = \mathbf{A}\boldsymbol{\mu}$ . For the covariance of a matrix product, it is a little more interesting.

**Theorem:**  $\text{Cov}(\mathbf{A}\mathbf{y}) = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top$ .

**Proof:** Use the definition.

$$\begin{aligned}\text{Cov}(\mathbf{A}\mathbf{y}) &= \mathbb{E}[\mathbf{A}\mathbf{y}\mathbf{y}^\top\mathbf{A}^\top] - \mathbb{E}[\mathbf{A}\mathbf{y}](\mathbb{E}[\mathbf{A}\mathbf{y}])^\top \\ &= \mathbf{A}\mathbb{E}[\mathbf{y}\mathbf{y}^\top]\mathbf{A}^\top - \mathbf{A}\mathbb{E}[\mathbf{y}]\mathbb{E}[\mathbf{y}]^\top\mathbf{A}^\top \\ &= \mathbf{A}(\mathbb{E}[\mathbf{y}\mathbf{y}^\top] - \mathbb{E}[\mathbf{y}]\mathbb{E}[\mathbf{y}]^\top)\mathbf{A}^\top \\ &= \mathbf{A}\text{Cov}(\mathbf{y})\mathbf{A}^\top. \quad \blacksquare\end{aligned}$$

Note that in the second line we used the fact that  $(\mathbf{A}\mathbf{B})^\top = \mathbf{B}^\top\mathbf{A}^\top$ .

Some facts: The inverse of a symmetric matrix is symmetric. And  $\mathbf{X}^\top \mathbf{X}$  is symmetric, since its transpose equals itself.

Let  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ . This called the **hat matrix** since  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ .

**Theorem:**  $\text{Cov}(\mathbf{r}) = \sigma^2(\mathbf{I} - \mathbf{H})$ .

**Proof:** Since  $\mathbf{r} = (\mathbf{I} - \mathbf{H})\mathbf{y}$ , then

$$\begin{aligned}\text{Cov}(\mathbf{r}) &= (\mathbf{I} - \mathbf{H})\text{Cov}(\mathbf{y})(\mathbf{I} - \mathbf{H})^\top \\ &= (\mathbf{I} - \mathbf{H})\sigma^2\mathbf{I}(\mathbf{I} - \mathbf{H})^\top \\ &= \sigma^2(\mathbf{I} - \mathbf{H})(\mathbf{I}^\top - \mathbf{H}^\top) \\ &= \sigma^2(\mathbf{I}\mathbf{I}^\top - \mathbf{I}\mathbf{H}^\top - \mathbf{H}\mathbf{I}^\top + \mathbf{H}\mathbf{H}^\top).\end{aligned}$$

Note that  $\mathbf{I}^\top = \mathbf{I}$  and that  $\mathbf{H}\mathbf{H}^\top = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top [\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top]^\top$  which simplifies to  $\mathbf{H}$ . ■



One implication of the theorem is that the variance of the  $i$ th residual is  $(1 - h_{ii})\sigma^2$ , which would be estimated by  $(1 - h_{ii})\hat{\sigma}^2$ , where  $\hat{\sigma}^2$  is the unbiased estimator

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

The **standardized residual** is  $(y_i - \hat{y}_i) / \sqrt{(1 - h_{ii})\hat{\sigma}^2}$ . This follows a  $t$  distribution with  $n - p - 1$  degrees of freedom.

For example, one could use the standardized residual to decide whether an observation might be an outlier.

## 3.5 Confidence and Prediction Intervals

Suppose we want to predict the response  $Y$  for a new vector of covariates  $\mathbf{x}_{\text{new}}$ . The point estimate is easy:  $\mathbf{x}_{\text{new}}^\top \hat{\boldsymbol{\beta}}$ .

The variance depends upon whether we want to estimate the average response for an observation with the value  $\mathbf{x}_{\text{new}}$  (i.e., the location of response on the estimated hyperflat) or the value of an individual with  $\mathbf{x}_{\text{new}}$ .

This is the distinction we saw previously for simple linear regression between setting a confidence interval on the response  $\hat{\beta}_0 + \hat{\beta}_1 x$  and a prediction interval on an individual that has explanatory variable  $x$ .

The variance for the prediction interval should be larger, since the individual is unlikely to have the average value.

For the average response at  $\mathbf{x}_{\text{new}}$  we can calculate

$$\begin{aligned}\text{Var}(\mathbf{x}_{\text{new}}^\top \hat{\beta}) &= \mathbf{x}_{\text{new}}^\top \text{Var}(\hat{\beta}) (\mathbf{x}_{\text{new}}^\top)^\top \\ &= \mathbf{x}_{\text{new}}^\top \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{x}_{\text{new}}^\top)^\top \\ &= \sigma^2 \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}.\end{aligned}$$

For  $\sigma^2$ , we use the usual estimate

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

The two-sided  $(1 - \alpha)100\%$  confidence interval is just

$$U, L = \mathbf{x}_{\text{new}}^\top \hat{\beta} \pm \hat{\sigma} \sqrt{\mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}} \times t_{n-p-1, \alpha/2}.$$

For predictions on individuals, there are two sources that contribute to the variance:

- the variance due to estimating the average response
- the variance of an individual's value around the average response.

The prediction interval must combine both:

$$\begin{aligned}\text{Var}(\hat{y}_{\text{new,ind}}) &= \text{Var}(\hat{y}_{\text{new}} + \epsilon_{\text{new}}) \\ &= \text{Var}(\hat{y}_{\text{new}}) + \text{Var}(\epsilon_{\text{new}}) \\ &= \sigma^2 \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}} + \sigma^2 \\ &= \sigma^2 (\mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}} + 1).\end{aligned}$$

This uses the fact that  $\hat{y}_{\text{new}}$  and  $\epsilon_{\text{new}}$  are independent.

The  $\hat{y}_{\text{new}}$  and  $\epsilon_{\text{new}}$  are independent since

- $\epsilon_{\text{new}}$  was not used in fitting the model **and**
- the general linear model assumes that the errors ( $\epsilon$ s) are independent.

Therefore the two-sided  $(1 - \alpha)100\%$  prediction interval is just

$$U, L = \mathbf{x}_{\text{new}}^\top \hat{\boldsymbol{\beta}} \pm \hat{\sigma} \sqrt{1 + \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}} \times t_{n-p-1, \alpha/2}.$$

For the case of simple linear regression, with one explanatory variable, it is easy to show that the confidence interval formula and the prediction interval formula agree with the formulae given in Lecture 2. It just takes some algebra.

## 3.6 Bases and Projections

A **basis** of a vector space  $\mathbb{R}^p$  is a set of vectors  $\mathbf{x}_1, \dots, \mathbf{x}_p$  such that any vector  $\mathbf{v} \in \mathbb{R}^p$  can be expressed as  $\mathbf{v} = \sum_{i=1}^p a_i \mathbf{x}_i$  but none of the  $\mathbf{x}_i$  is a linear combination of the other basis vectors.

For example, a standard choice of basis is  $(1, 0, 0, \dots, 0)$ ,  $(0, 1, 0, \dots, 0)$ ,  $(0, 0, 1, \dots, 0)$  up to  $(0, 0, 0, \dots, 1)$ . In fact, this is an **orthonormal** basis since each member is orthogonal to the others (i.e., had dot product equal to zero) and is normalized (has unit length).

A subspace of  $\mathbb{R}^p$  is  $\mathbb{R}^q$  for  $q < p$ . Any vector in  $\mathbb{R}^q$  can be written as a linear combination of  $q$  basis vectors.

A **projection matrix**  $\mathbf{P}$  is a square matrix such that  $\mathbf{P}^2 = \mathbf{P}$ . For example, the projection matrix

$$\mathbf{P} = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$$

projects  $\mathbb{R}^2$  onto the subspace  $\mathbb{R}^1$  corresponding to the line  $x_1 = x_2$ , since  $\mathbf{P}(x_1, x_2)^\top = (x_2, x_2)^\top$ .

**Note:** The line  $x_1 = x_2$  is a proper subspace of  $\mathbb{R}^2$ . It contains  $\mathbf{0}$ , and  $(1, 1)^\top$  is a basis element.

**Note:** It makes sense that  $\mathbf{P}^2 = \mathbf{P}$  since if one applies the same projection twice, nothing changes.

A projection matrix is symmetric if and only if the vector space projection is orthogonal. In an **orthogonal projection** of  $\mathbb{R}^p$  onto  $\mathbb{R}^q$  for  $q < p$ , any vector  $\mathbf{v}$  can be written as  $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2$  where  $\mathbf{v}_1 = \mathbf{P}\mathbf{v}$  and  $\mathbf{v}_2 = (\mathbf{I} - \mathbf{P})\mathbf{v}$ .

This implies that the dot product  $\mathbf{v}_1^\top \mathbf{v}_2 = 0$ .

Consider the projection matrix

$$\tilde{\mathbf{P}} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

This is symmetric and projects  $(x_1, x_2)^\top$  onto  $(x_1, 0)^\top$ . And it is an orthogonal projection since

$$(x_1, x_2)^\top = \tilde{\mathbf{P}}(x_1, x_2)^\top + (\mathbf{I} - \tilde{\mathbf{P}})(x_1, x_2)^\top = (x_1, 0)^\top + (0, x_2)^\top$$

with  $\langle (x_1, 0), (0, x_2) \rangle = 0$ .



In contrast,  $\mathbf{P}$  is not an orthogonal projection since

$$\mathbf{P}(x_1, x_2)^\top + (\mathbf{I} - \mathbf{P})(x_1, x_2)^\top = (x_2, x_2)^\top + (x_1 - x_2, 0)^\top$$

and  $\langle (x_2, x_2), (x_1 - x_2, 0) \rangle \neq 0$ .

The orthogonal projection of the hat matrix minimizes the sum of the squared vertical distances onto the subspace. Recall that in multiple linear regression we assume the explanatory variables are measured without error, and thus we want to minimize the sum of the squared vertical distances.

The hat matrix is a projection.

$$\mathbf{H}^2 = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \times \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{H}.$$

But to what subspace does it project? To the column space of  $\mathbf{X}$ .

The **column space**  $\mathcal{C}(\mathbf{A})$  is the space spanned by the columns of  $\mathbf{A}$ . It is formed by taking all possible linear combinations of the columns of  $\mathbf{A}$ . One views the columns of  $\mathbf{A}$  as the basis elements of the space. If  $\mathbf{A}$  has rank  $n$ , then the column space has dimension  $n$ .

The **null space** of an  $m \times n$  matrix  $\mathbf{A}$  is the space of vectors in  $\mathbb{R}^n$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ . If  $\mathbf{A}$  is invertible, then the null space is just  $\mathbf{0}$ . For

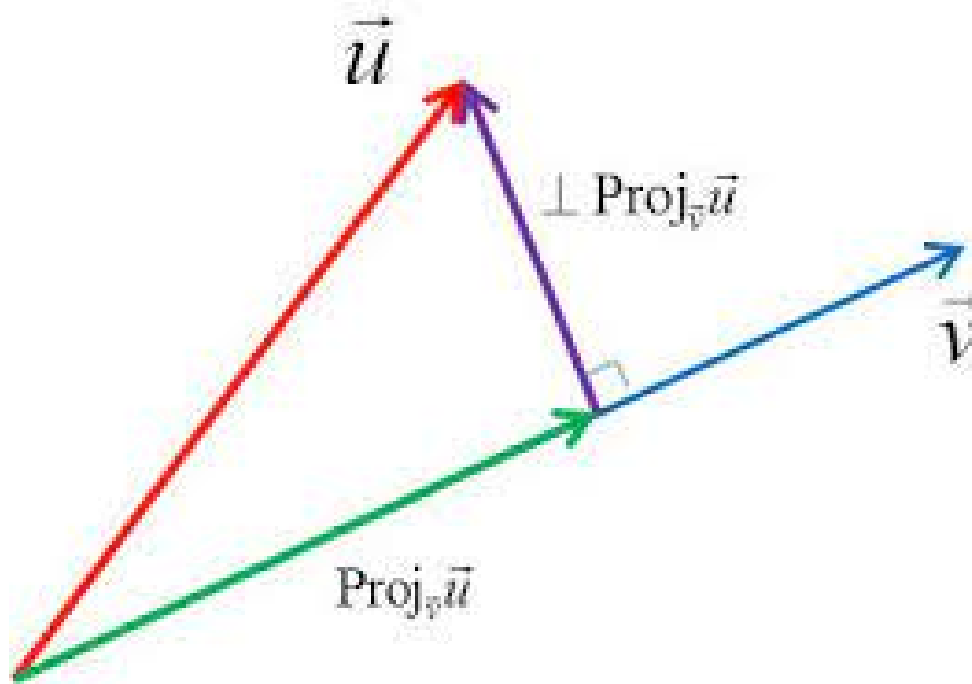
$$\mathbf{P} = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix},$$

the null space is  $\{(x_1, 0)^\top\}$ .

The **row space**  $\mathcal{R}(\mathbf{A})$  is the space formed by taking all possible linear combinations of the rows of  $\mathbf{A}$ . If  $\mathbf{A}$  has rank  $m$ , then the row space has dimension  $m$ .

Some properties of orthogonal projections  $\mathbf{H}$ :

- $\mathbf{v}_1 = \mathbf{H}\mathbf{v}$ ,  $\mathbf{v}_2 = (\mathbf{I} - \mathbf{H})\mathbf{v}$  and  $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2$ .
- $\mathbf{v}_1$  is orthogonal to  $\mathbf{v}_2$ , so their dot product is zero. We write this as  $\mathbf{v}_1 \perp \mathbf{v}_2$ .
- $\mathbf{H}$  and  $\mathbf{I} - \mathbf{H}$  are symmetric, and their squares equal themselves.
- In the hat matrix, if  $\mathbf{X}$  is  $n \times (p + 1)$  of rank  $p + 1$ , then  $\mathbf{H}$  has rank  $p + 1$ .
- In the hat matrix, the eigenvalues of  $\mathbf{H}$  consist of  $p + 1$  ones and  $n - p - 1$  zeroes.
- $\mathbf{H}\mathbf{X} = \mathbf{X}$ ,  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$ .
- For the hat matrix, the trace of  $\mathbf{H}$  equals the rank of  $\mathbf{X}$ .



This image uses the notation  $\perp \mathbf{v}$  to denote the component of  $\mathbf{u}$  that is orthogonal to the orthogonal projection of  $\mathbf{u}$  onto the subspace containing  $\mathbf{v}$ . In this class, we shall usually denote this by  $\mathbf{v}^\perp$ .

More results from linear algebra:

**Theorem:** A vector that is orthogonal to the column space of an  $n \times m$  matrix  $\mathbf{X}$  lies in the null space of  $\mathbf{X}^\top$ .

**Proof:** Let  $\mathbf{u}$  be orthogonal to the column space of an  $n \times m$  matrix  $\mathbf{X}$ . So for any vector  $\mathbf{v}$  in the column space of  $\mathbf{X}$ ,  $\mathbf{v}^\top \mathbf{u} = 0$ . Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$  be the columns of  $\mathbf{X}$ . Each  $\mathbf{x}_i$  is in the column space of  $\mathbf{X}$ , so  $\mathbf{x}_i^\top \mathbf{u} = 0 \forall i$ . Since the rows of  $\mathbf{X}^\top$  are the columns of  $\mathbf{X}$ , then  $\mathbf{u}$  lies in the null space of  $\mathbf{X}^\top$ . ■

**Theorem:** A vector in the column space of an  $n \times m$  matrix  $\mathbf{X}$  is orthogonal to a vector in the null space of  $\mathbf{X}^\top$ .

**Proof:** Let  $\mathbf{v}$  be in the column space of  $\mathbf{X}$  and  $\mathbf{u}$  be in the null space of  $\mathbf{X}^\top$ . Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$  be the columns of  $\mathbf{X}$ . Since  $\mathbf{v}$  the columns of  $\mathbf{X}$  are a basis for the column space, then there exists  $\mathbf{a}$  such that

$$\mathbf{v} = \sum_{i=1}^m a_i \mathbf{x}_i = \mathbf{X} \mathbf{a}.$$

So the inner product  $\mathbf{u}^\top \mathbf{v} = \mathbf{u}^\top \mathbf{X} \mathbf{a} = (\mathbf{X}^\top \mathbf{u}) \mathbf{a}$  which is zero since we previously showed that  $\mathbf{u}$  is in the null space of  $\mathbf{X}^\top$ . ■

**Note:** Thus  $\hat{\boldsymbol{\beta}}$  is obtained by projecting  $\mathbf{y}$  onto the column space of  $\mathbf{X}$ .

## 3.7 More on the Hat Matrix

The hat matrix projects  $\mathbf{y} \in \mathbb{R}^n$  onto the column space of  $\mathbf{X}$ , which has dimension  $p + 1$ . That imposes  $p + 1$  linear constraints on what the response can be, accounting for the loss of  $p + 1$  degrees of freedom.

The diagonal entry of the hat matrix lies in  $[0, 1]$  since  $\mathbf{H}^2 = \mathbf{H}$ . A diagonal element of  $\mathbf{H}$  is  $h_{ii}$ . The corresponding diagonal element of  $\mathbf{H}^2$  is  $\mathbf{h}_i^\top \mathbf{h}_i = \sum h_{ij}^2 = h_{ii}^2 + \sum_{i \neq j} h_{ij}^2$ . Thus  $h_{ii}^2 \leq h_{ii}$  so  $0 \leq h_{ii} \leq 1$ .

An **eigenvalue** of  $\mathbf{H}$  satisfies  $\mathbf{H}\mathbf{x} = \lambda\mathbf{x}$ . Since  $\mathbf{H}^2\mathbf{x} = \mathbf{H}\mathbf{x} = \lambda\mathbf{H}\mathbf{x} = \lambda^2\mathbf{x}$ , it follows that  $\lambda = \lambda^2$ . An eigenvalue of a projection matrix is zero or one.