# Simple Linear Regression - Part Deux

## Just the Essentials of Linear Algebra

**Vector:** A set of *ordered* numbers, written, by default, as a column. Example: $v = \begin{bmatrix} 2 \\ -5 \\ 1 \\ 3 \end{bmatrix}$

**Transpose of a Vector:** The transpose of vector $v$, written $v^T$, is defined as the corresponding ordered row of numbers. Example: The transpose of the vector $v$ above is $v^T = \begin{bmatrix} 2 & -5 & 1 & 3 \end{bmatrix}$

**Scalar Multiplication:** Scalar (real number) multiplication distributes over components. Ex: $4v = 4\begin{bmatrix} 2 \\ -5 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 8 \\ -20 \\ 4 \\ 12 \end{bmatrix}$

**Vector Addition:** Addition is defined componentwise. Example: For vectors $v = \begin{bmatrix} 2 \\ -5 \\ 1 \\ 3 \end{bmatrix}$ and $u = \begin{bmatrix} -1 \\ 3 \\ -2 \\ 4 \end{bmatrix}$,

$$v + u = \begin{bmatrix} 2 \\ -5 \\ 1 \\ 3 \end{bmatrix} + \begin{bmatrix} -1 \\ 3 \\ -2 \\ 4 \end{bmatrix} = \begin{bmatrix} 2-1 \\ -5+3 \\ 1-2 \\ 3+4 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ -1 \\ 7 \end{bmatrix} = u + v$$

**Vector Multiplication:** The inner product (or dot product) of vectors $v$ and $u$ is defined as

$$v^T u = \begin{bmatrix} 2 & -5 & 1 & 3 \end{bmatrix}\begin{bmatrix} -1 \\ 3 \\ -2 \\ 4 \end{bmatrix} = (2)(-1)+(-5)(3)+(1)(-2)+(3)(4) = -2-15-2+12 = -7$$

*Note 1:* The inner product of two vectors is a number!

*Note 2:* $v^T u = u^T v$

**Orthogonal:** Vectors $v$ and $u$ are *orthogonal* if $v^T u = 0$. The geometrical interpretation is that the vectors are perpendicular. Orthogonality of vectors plays a big role in linear models.

**Length of a Vector:** The length of vector $v$ is given by $\|v\| = \sqrt{v^T v}$. This extends the Pythagorean theorem to vectors with an arbitrary number of components. Example: For $v^T = \begin{bmatrix} 2 & -5 & 1 & 3 \end{bmatrix}$, $\|v\| = \sqrt{v^T v} = \sqrt{39}$.

**Matrices:** $A$ is an $m$ by $n$ matrix if it is an array of numbers with $m$ rows and $n$ columns. Example: $A = \begin{bmatrix} 2 & 1 \\ -1 & 3 \\ 1 & 3 \\ 4 & 2 \end{bmatrix}$

is a $4 \times 2$ matrix.

**Scalar Multiplication:** Scalar (real number) multiplication distributes over components, as for vectors.

**Matrix Addition:** Matrix addition is componentwise, as for vectors.

*Note:* A vector is a special case of a matrix with a single column (or single row in the case of its transpose).

**Matrix Multiplication:** If $A$ is $m \times n$ and $B$ is $n \times p$, the product $AB$ is defined as the matrix of all inner products of row vectors of $A$ with column vectors of $B$. Example: If $A_{4\times2} = \begin{bmatrix} 2 & 1 \\ -1 & 3 \\ 1 & 3 \\ 4 & 2 \end{bmatrix}$ and $B_{2\times3} = \begin{bmatrix} 1 & 7 & 3 \\ -1 & 2 & -1 \end{bmatrix}$, then

$$(AB)_{4\times3} = \begin{bmatrix} 1 & 16 & 5 \\ -4 & -1 & -6 \\ -2 & 13 & 0 \\ 2 & 32 & 10 \end{bmatrix}.$$

**Transpose of a Matrix:** The transpose of $m \times n$ matrix $A$ is the $n \times m$ matrix $A^T$ obtained by swapping rows for columns in $A$. Example: If $A_{4\times2} = \begin{bmatrix} 2 & 1 \\ -1 & 3 \\ 1 & 3 \\ 4 & 2 \end{bmatrix}$, then $A^T_{2\times4} = \begin{bmatrix} 2 & -1 & 1 & 4 \\ 1 & 3 & 3 & 2 \end{bmatrix}$

**Transpose of a Product:** $(AB)^T = B^T A^T$

**Symmetric Matrix:** The square matrix $A_{n\times n}$ is *symmetric* if $a_{ij} = a_{ji}$ for all $1 \leq i \leq n$, $1 \leq j \leq n$ where $a_{ij}$ is the entry in the $i^{th}$ row, $j^{th}$ column of $A_{n\times n}$. Example: In $A_{3\times3} = \begin{bmatrix} 1 & -3 & 1 \\ -3 & 7 & -2 \\ 1 & -2 & 5 \end{bmatrix}$: $a_{12} = a_{21} = -3$; $a_{13} = a_{31} = 1$; $a_{23} = a_{32} = -2$.

(**Note:** If $A_{n\times n}$ is symmetric, then $A^T = A$)

**Theorem:** for any $m \times n$ matrix $A$, $A^T A$ is an $n \times n$ symmetric matrix, and $AA^T$ is an $m \times m$ symmetric matrix.

**Linear Independence:** The set of $n$ vectors $v_1, v_2, v_3, \cdots, v_n$ form a linearly independent set if the linear combination $c_1 v_1 + c_2 v_2 + c_3 v_3 + \cdots + c_n v_n = 0$ implies that all of the constant coefficients $c_i$ equal zero.

**rank($A$):** The rank of matrix $A$, rank($A$), is the number of linearly independent columns (or rows) in the matrix. A matrix with all columns, or all rows, independent is said to be "full rank." Full rank matrices play an important role in linear models, especially regression.

**Theorem about Rank:** $\text{rank}(A^T A) = \text{rank}(AA^T) = \text{rank}(A)$. This theorem will be used in regression.

**Inverse of a Matrix:** The inverse of the $n \times n$ square matrix $A$ is the $n \times n$ square matrix $A^{-1}$ such that $AA^{-1} = A^{-1}A =$

$I$, where $I$ is the $n \times n$ "Identity" matrix, $I_{n \times n} = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix}_{n \times n}$

**Theorem:** The $n \times n$ square matrix $A$ has an inverse if and only if it is full rank, i.e., rank($A$) = $n$.

**Theorem:** The invertible $2 \times 2$ matrix $A_{2 \times 2} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ has $2 \times 2$ inverse $A_{2 \times 2}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$. Example: If

$A = \begin{bmatrix} 2 & 4 \\ 3 & 1 \end{bmatrix}$, then $A^{-1} = \frac{1}{(2)(1) - (4)(3)} \begin{bmatrix} 1 & -4 \\ -3 & 2 \end{bmatrix} = \frac{1}{-10} \begin{bmatrix} 1 & -4 \\ -3 & 2 \end{bmatrix} = -\frac{1}{10} \begin{bmatrix} 1 & -4 \\ -3 & 2 \end{bmatrix}$. You should carry out

the products $AA^{-1}$ and $A^{-1}A$ to confirm that they equal $I_{2 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$.


## Solving a System of $n$ Linear Equations in $n$ Unknowns

The system of $n$ linear equations in $n$ unknowns,
$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$$
$$a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2$$
$$\vdots$$
$$a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n$$
, can be represented by the

matrix equation $Ax = b$, where $A_{n \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}_{n \times n}$, $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, and $b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$ are the matrix of

coefficients, vector of unknowns, and vector of constants, respectively. The system has a unique solution *if and only if* $A$ is invertible, i.e., $A^{-1}$ exists. In that case, the solution is given by $x = A^{-1}b$.

## Representing the Simple Linear Regression Model as a Matrix Equation

For a sample of $n$ observations on the bivariate distribution of the variables $X$ and $Y$, the simple linear regression model $Y = \beta_0 + \beta_1 X + \varepsilon$ leads to the system of $n$ equations

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 X_1 + \varepsilon_1 \\ \beta_0 + \beta_1 X_2 + \varepsilon_2 \\ \vdots \\ \beta_0 + \beta_1 X_n + \varepsilon_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 X_1 \\ \beta_0 + \beta_1 X_2 \\ \vdots \\ \beta_0 + \beta_1 X_n \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix},$$

which can be written $\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$. The shorthand for this is $Y = X\beta + \varepsilon$, where the $n \times 2$ matrix

$X$ is called the "design" matrix because the values of the variable $X$ are often fixed by the researcher in the design of the experiment used to investigate the relationship between the variables $X$ and $Y$.

## Fitting the Best Least Squares Regression Line to the $n$ Observations

Ideally, we would solve the matrix equation $Y = X\beta + \varepsilon$ for the vector of regression coefficients $\beta$, i.e., the true intercept $\beta_0$ and true slope $\beta_1$ in the model $Y = \beta_0 + \beta_1 X + \varepsilon$. In practice, however, this is never possible because the equation $Y = X\beta + \varepsilon$ *has no knowable solution*! (Why?) When faced with an equation that we cannot solve, we do what mathematicians usually do: we find a related equation that we *can* solve. First, since we don't know the errors in our observations, we forget about the vector of the errors $\varepsilon$. (Here, it is important that we not confuse the errors $\varepsilon_i$, which we never know, with the residuals $e_i$ associated with our estimated regression line.)

Unable to determine the parameters $\beta_0$ and $\beta_1$, we look to estimate them, i.e., solve for $\hat{\beta}_0$ and $\hat{\beta}_1$ in the matrix equation $Y = X\hat{\beta}$, but the system involves $n$ equations in the two unknowns $\hat{\beta}_0$ and $\hat{\beta}_1$. If you remember your algebra, overdetermined systems of linear equations rarely have solutions. What we need is a *related* system of linear equations that we *can* solve. (Of course, we'll have to show that the solution to the new system has relevance to our problem of estimating $\beta_0$ and $\beta_1$.) Finally (drum roll please), the system we'll actually solve is,

(1.1) $$X^T X \hat{\beta} = X^T Y$$

Next, we show what the system above looks like, and, as we go, we'll see why it's solvable and why it's relevant.

(1.2) $$X^T X = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ X_1 & X_2 & \cdots & X_n \end{bmatrix} \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ 1 & \vdots \\ 1 & X_n \end{bmatrix} = \begin{bmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{bmatrix}$$

$$(1.3) \qquad \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \begin{bmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} n\hat{\beta}_0 + \hat{\beta}_1 \sum X_i \\ \hat{\beta}_0 \sum X_i + \hat{\beta}_1 \sum X_i^2 \end{bmatrix}$$

$$(1.4) \qquad \mathbf{X}^T\mathbf{Y} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ X_1 & X_2 & \cdots & X_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \sum Y_i \\ \sum X_i Y_i \end{bmatrix}$$

Now we're in a position to write out the equations for the system $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$,

$$(1.5) \qquad n\hat{\beta}_0 + \hat{\beta}_1 \sum X_i = \sum Y_i$$

$$\hat{\beta}_0 \sum X_i + \hat{\beta}_1 \sum X_i^2 = \sum X_i Y_i$$

If these equations look familiar, they are the equations derived in the previous notes by applying the Least Squares criterion for the best fitting regression line. Thus, the solution to the matrix equation (1.1) is the least squares solution to the problem of fitting a line to data! (Although we could have established this directly using theorems from linear algebra and vector spaces, the calculus argument made in Part I of the regression notes is simpler and probably more convincing.) Finally, the system (1.6) of two equations in the unknowns $\hat{\beta}_0$ and $\hat{\beta}_1$ has a solution!

To solve the system $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$, we note that the $2 \times 2$ matrix $\mathbf{X}^T\mathbf{X}$ has an inverse. We know this since rank($\mathbf{X}^T\mathbf{X}$) = rank($\mathbf{X}$) = 2, and full rank square matrices have inverses. (Note: rank($\mathbf{X}$) = 2 because the two columns of the design matrix $\mathbf{X}$ are linearly independent.) The solution to the matrix equation (1.1) has the form,

$$(1.7) \qquad \hat{\boldsymbol{\beta}} = \left( \mathbf{X}^T\mathbf{X} \right)^{-1} \mathbf{X}^T\mathbf{Y}$$

where $\left( \mathbf{X}^T\mathbf{X} \right)^{-1} = \dfrac{1}{n\sum X_i^2 - \left( \sum X_i \right)^2} \begin{bmatrix} \sum X_i^2 & -\sum X_i \\ -\sum X_i & n \end{bmatrix} = \dfrac{1}{n\sum \left( X_i - \bar{X} \right)^2} \begin{bmatrix} \sum X_i^2 & -\sum X_i \\ -\sum X_i & n \end{bmatrix}$. After a

fair amount of algebra, equation (1.6) returns the estimated intercept $\hat{\beta}_0$ and slope $\hat{\beta}_1$ given in Part I of the regression notes:

$$(1.8) \qquad \hat{\beta}_1 = \frac{\sum X_i Y_i - \frac{1}{n}\left[ \left( \sum X_i \right)\left( \sum Y_i \right) \right]}{\sum X_i^2 - \frac{1}{n}\left( \sum X_i \right)^2} = \frac{\sum \left( X_i - \bar{X} \right)\left( Y_i - \bar{Y} \right)}{\sum \left( X_i - \bar{X} \right)^2}$$

$$(1.9) \qquad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

## Example:

Although you will have plenty of opportunities to use software to determine the best fitting regression line for a sample of bivariate data, it might be useful to go through the process once by hand for the "toy" data set:

| $x$ | 1 | 2 | 4 | 5 |
|---|---|---|---|---|
| $y$ | 8 | 4 | 6 | 2 |

Rather than using the equations (1.7) and (1.8), start with $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$ and find (a) $\mathbf{X}^T\mathbf{X}$, (b) $\mathbf{X}^T\mathbf{Y}$, (c) $\left(\mathbf{X}^T\mathbf{X}\right)^{-1}$ (using only the definition of the inverse of a $2\times2$ matrix given in the linear algebra review), and (d) $\left(\mathbf{X}^T\mathbf{X}\right)^{-1}\mathbf{X}^T\mathbf{Y}$. Perform the analysis again using software or a calculator to confirm the answer.

## The (Least Squares) Solution to the Regression Model

The equation $\mathbf{Y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}$ is the estimate to the simple linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix}$ is the least squares estimate of the intercept and slope of the *true* regression line $y = \beta_0 + \beta_1 x$, and $\mathbf{e}$ is the $n\times1$ vector of residuals. Now, $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ can be written $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e}$, where $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is the $n\times1$ vector of predictions made for the observations by the *estimated* regression line $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$. The $n$ points $\left(X_i, \hat{Y}_i\right)$ will, of course, all lie on the estimated regression line.

## The Hat Matrix

A useful matrix that shows up again and again in regression analysis is the $n\times n$ matrix $\mathbf{H} = \mathbf{X}\left(\mathbf{X}^T\mathbf{X}\right)^{-1}\mathbf{X}$, called the "hat" matrix. To see how it gets its name, note that $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\left(\mathbf{X}^T\mathbf{X}\right)^{-1}\mathbf{X}\mathbf{Y} = \mathbf{H}\mathbf{Y}$. Thus, the matrix $\mathbf{H}$ puts a "hat" on the vector $\mathbf{Y}$.

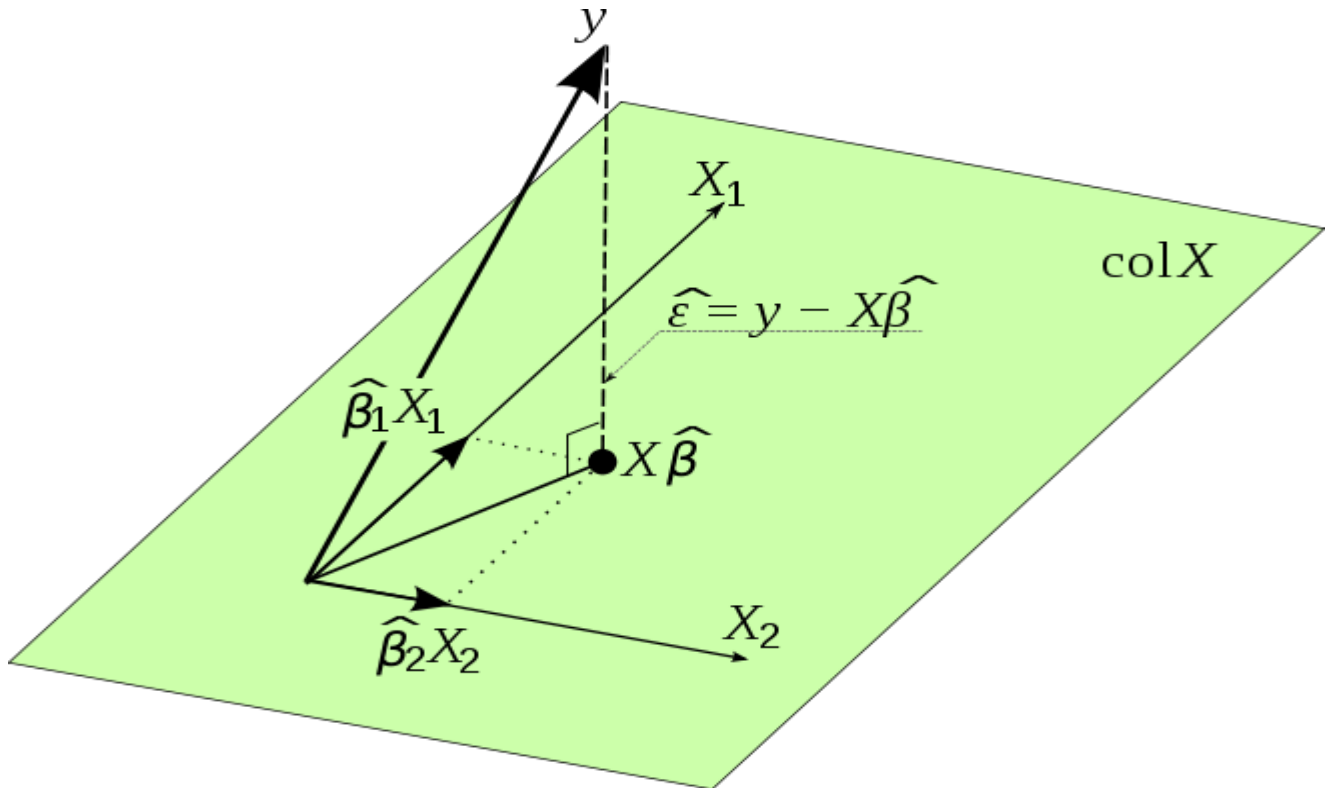The hat matrix $\mathbf{H}$ has many nice properties. These include:

- $\mathbf{H}$ is symmetric, i.e., $\mathbf{H}^T = \mathbf{H}$, as is easily proven.

- $\mathbf{H}$ is idempotent, which is a fancy way of saying that $\mathbf{H}\mathbf{H} = \mathbf{H}$. The shorthand for this is $\mathbf{H}^2 = \mathbf{H}$. This is also easily proven.

- The matrix $(\mathbf{I}\text{-}\mathbf{H})$, where $\mathbf{I}$ is the $n\times n$ identity matrix, is both symmetric and idempotent as well.

## The Error Sum of Squares, *SSE*, in Matrix Form

From the equation $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e}$, we get $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$. Using the hat matrix, this becomes $\mathbf{e} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$, where we have right-factored the vector $\mathbf{Y}$. The error sum of squares, *SSE*, is just the squared length of the vector of residuals, i.e., $SSE = \sum e_i^2 = \mathbf{e}^T\mathbf{e}$. In terms of the hat matrix this becomes $SSE = \mathbf{Y}^T(\mathbf{I} - \mathbf{H})\mathbf{Y}$, where we have used the symmetry and idempotency of the matrix $(\mathbf{I} - \mathbf{H})$ to simplify the result.

## The Geometry of Least Squares



$$(1.10) \qquad \mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e}$$

One of the attractions of linear models, and especially of the least squares solutions to them, is the wealth of geometrical interpretations that spring from them. The $n\times 1$ vectors $\mathbf{Y}$, $\hat{\mathbf{Y}}$, and $\mathbf{e}$ are vectors in the vector space $\mathbb{R}^n$, (an $n$-dimensional space whose components are real numbers, as opposed to complex numbers). From equation (1.9) above, we know that the vectors $\hat{\mathbf{Y}}$ and $\mathbf{e}$ sum to $\mathbf{Y}$, but we can show a much more surprising result that will eventually lead to powerful conclusions about the sums of squares *SSE*, *SSR*, and *SST*.

First, we have to know where each of the vectors $\mathbf{Y}$, $\hat{\mathbf{Y}}$, and $\mathbf{e}$ "live":

- The $n$-vector of observations $\mathbf{Y}$ has no restrictions placed on it and therefore can lie anywhere in $\mathbb{R}^n$.

- $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\beta}_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \beta_1 \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$ is restricted to the two-dimensional subspace of $\mathbb{R}^n$ spanned by the columns

  of the design matrix $\mathbf{X}$, called (appropriately) the **column space** of $\mathbf{X}$.

- We will shortly show that $\mathbf{e}$ lives in a subspace of $\mathbb{R}^n$ **orthogonal** to the column space of $\mathbf{X}$.

Next, we derive the critical result that the vectors $\hat{\mathbf{Y}}$ and $\mathbf{e}$ are orthogonal (perpendicular) in $\mathbb{R}^n$. (Remember: vectors $\mathbf{v}$ and $\mathbf{u}$ are *orthogonal* if and only if $\mathbf{v}^T\mathbf{u} = 0$.)

- $\hat{\mathbf{Y}}^T\mathbf{e} = (\mathbf{HY})^T (\mathbf{I}-\mathbf{H})\mathbf{Y} = \mathbf{Y}^T\mathbf{H}^T (\mathbf{I}-\mathbf{H})\mathbf{Y} = \mathbf{Y}^T (\mathbf{H}-\mathbf{H})\mathbf{Y} = 0$

(Where we made repeated use of the fact that $\mathbf{H}$ is symmetric and idempotent.) Therefore, $\mathbf{e}$ is restricted to an $n$ - 2 dimensional subspace of $\mathbb{R}^n$ (because it must be orthogonal to every vector in the column space of $\mathbf{X}$).
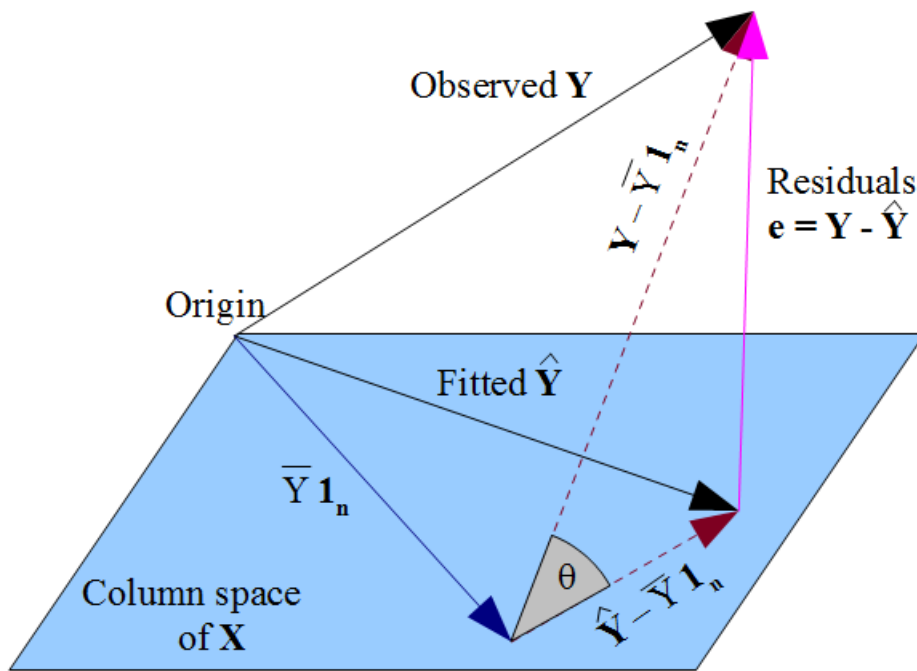
Combining the facts that the vectors $\hat{\mathbf{Y}}$ and $\mathbf{e}$ sum to $\mathbf{Y}$ *and* are orthogonal to each other, we conclude that $\hat{\mathbf{Y}}$ and $\mathbf{e}$ form the legs of a right triangle (in $\mathbb{R}^n$) with hypotenuse $\mathbf{Y}$. By the Pythagorean Theorem for right triangles, $\left\|\hat{\mathbf{Y}}\right\|^2 + \left\|\mathbf{e}\right\|^2 = \left\|\mathbf{Y}\right\|^2$ or equivalently,

(1.11)
$$\hat{\mathbf{Y}}^T\hat{\mathbf{Y}} + \mathbf{e}^T\mathbf{e} = \mathbf{Y}^T\mathbf{Y}$$

A slight modification of the argument above shows that the vectors $\hat{\mathbf{Y}} - \bar{\mathbf{Y}}$ and $\mathbf{e}$ sum to $\mathbf{Y} - \bar{\mathbf{Y}}$ *and* are orthogonal to each other, hence we conclude that $\hat{\mathbf{Y}} - \bar{\mathbf{Y}}$ and $\mathbf{e}$ form the legs of a right triangle (in $\mathbb{R}^n$) with hypotenuse $\mathbf{Y} - \bar{\mathbf{Y}}$. By the Pythagorean Theorem for right triangles,

(1.12)
$$\left(\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\right)^T \left(\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\right) + \mathbf{e}^T\mathbf{e} = \left(\mathbf{Y} - \bar{\mathbf{Y}}\right)^T \left(\mathbf{Y} - \bar{\mathbf{Y}}\right)$$

or, equivalently, $\boxed{SSR + SSE = SST}$. This last equality is the famous one involving the three sums of squares in regression. (See picture below.)

We've actually done more than just derive the equation involving the three sums of squares. It turns out that the dimensions of the subspaces the vectors "live" in also determines their "degrees of freedom," so we've also shown that *SSE* has *n* - 2 degrees of freedom because the vector of residuals **e** is restricted to an *n* - 2 dimensional subspace of $\mathbb{R}^n$ . (The term "degrees of freedom" is actually quite descriptive because the vector **e** is only "free" to assume values in this *n* - 2 dimensional subspace.)

## The Analysis of Variance (ANOVA) Table

The computer output of a regression analysis always contains a table containing the sums of squares *SSR*, *SSE*, and *SST*. The table is called an analysis of variance, or ANOVA, table for reasons we will see later in the course. The table has the form displayed below. (Note: The mean square of the residuals is the "mean squared error" *MSE*, i.e., the estimate of the variance $\sigma^2$ of the error variable $\varepsilon$ in the regression model.)

| Source | Degrees of Freedom (df) | Sums of Squares (SS) | Mean Square (MS) = SS/df |
|---|---|---|---|
| Regression Model | 1 | $SSR = \sum\left(\hat{Y}_i \text{-} \bar{Y}\right)^2 = \left(\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\right)^{\mathbf{T}}\left(\hat{\mathbf{Y}} - \bar{\mathbf{Y}}\right)$ | $MSR = \sum\left(\hat{Y}_i - \bar{Y}\right)^2$ |
| Residual Error | *n* - 2 | $SSE = \sum\left(Y_i \text{-} \hat{Y}_i\right)^2 = \mathbf{e}^{\mathbf{T}}\mathbf{e}$ | $MSE = \dfrac{\sum\left(Y_i - \hat{Y}_i\right)^2}{n-2}$ |
| Total | *n* - 1 | $SST = \sum\left(Y_i \text{-} \bar{Y}\right)^2 = \left(\mathbf{Y} - \bar{\mathbf{Y}}\right)^{\mathbf{T}}\left(\mathbf{Y} - \bar{\mathbf{Y}}\right)$ | |