# orthogonalityAndOptimization

September 9, 2019

We saw that, given a model, the parameters that best explain the data can be found by minimizing the sum squares error **(SSE)**. SSE was written as a function of our parameters and minimized by taking the derivative with respect to each parameter.

Our end result was the following equation
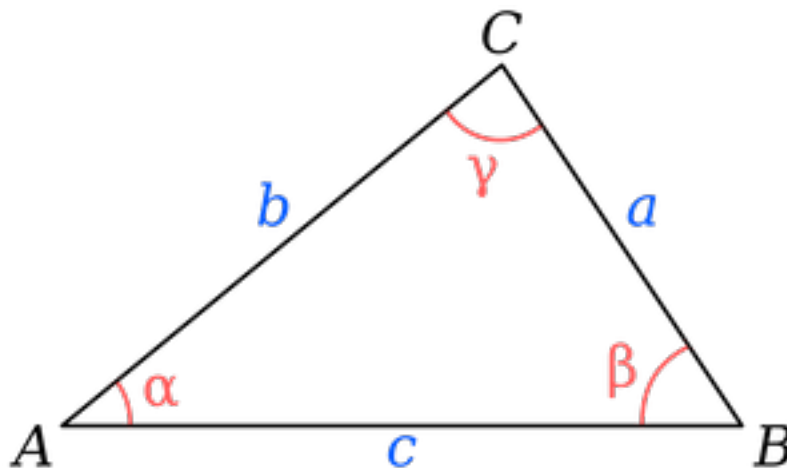
$$\beta_{\text{optimal}} = (X'X)^{-1}X'y.$$

This was a calculus approach to finding optima. We can also understand optima, and arrive at this same equation, by using linear algebra.

## 0.1 orthogonality

Two vectors, $x$ and $y$ are perpendicular to one another, or **orthogonal**, if their inner product equals 0

$$x'y = 0$$

We can use the law of cosines to see why this is the case. The law of cosines says, give a triangle with edge-lengths $a$, $b$, and $c$, and the angle $\gamma$ made by the vector $CA$ and $CB$, the edge-lengths are related like

$$C^2 = A^2 + B^2 - 2AB\cos(\gamma)$$

where capital letters denote the **length** of the triangle's sides.

If we consider $a$ and $b$ vectors, then

$$c = a - b$$

and the length of $c$ is

$$c^2 = c'c = [a_1 - b_1 a_2 - b_2] \begin{bmatrix} a_1 - b_1 \\ a_2 - b_2 \end{bmatrix} \tag{1}$$
$$= (a_1 - b_1)^2 + (a_2 - b_2)^2 \tag{2}$$
$$= (a_1^2 + a_2^2) + (b_1^2 + b_2^2) - 2(a_1 b_1 + a_2 b_2) \tag{3}$$

The above can be rewritten as the inner product of three vectors

$$c^2 = (a_1^2 + a_2^2) + (b_1^2 + b_2^2) - 2(a_1 b_1 + a_2 b_2) \tag{4}$$
$$= a'a + b'b - 2a'b \tag{5}$$

then we can relate this vector equation to our original cosine law.

$$a'a + b'b - 2a'b = A^2 + B^2 - 2AB\cos(\gamma)$$

We see that $a'a$ corresponds to the the length $A$ squared and $b'b$ corresponds to the length $B$ squared.

We define a vector's length

$$||v|| = (v'v)^{1/2},$$

and note that the length of a vector is always positive, and can only be zero if the vector has entries all zero.

The last term then relates the inner product between a and b to their lengths and the cosine of the angle they make

$$-2a'b = -2AB\cos(\gamma) \tag{6}$$
$$a'b = AB\cos(\gamma) \tag{7}$$
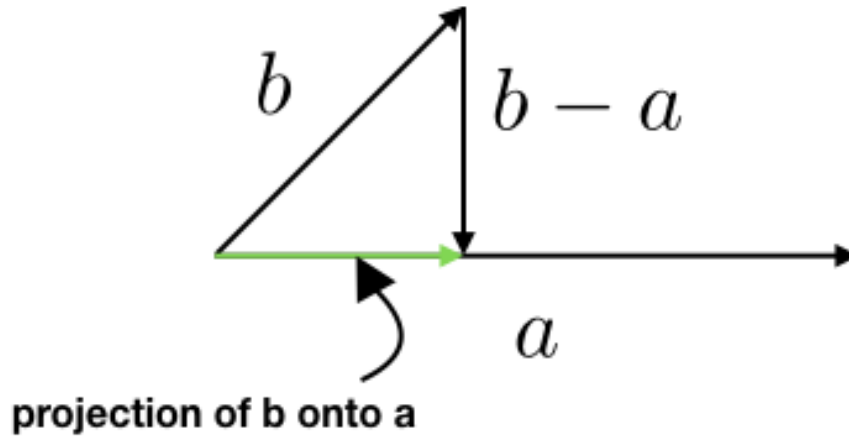$$a'b = ||a||||b||\cos(\gamma) \tag{8}$$

if the inner product $a'b$ is zero

$$0 = ||a||||b|| \cos(\gamma) \tag{9}$$

and assuming $a$ and $b$ are not zero vectors, it must be the case that $\cos(\gamma) = 0$ and this happens when $\gamma = \frac{\pi}{2}$, a perpendicular (orthogonal) angle.

## 0.2 projection

A vector $b$ is a **orthogonal** projection onto $a$ if the inner product between $b - a$ and $a$ is 0.



**projection of b onto a**

We can derive a formula for this "green" vector. The goal is to find the number $\omega$, in the same direction as $a$, so that $b - a$ and $a$ are orthogonal.

$$(b - \omega a)'(\omega a) = \omega b'a - \omega^2 a'a = 0 \tag{10}$$
$$b'a - \omega a'a = 0 \tag{11}$$
$$\omega = \frac{b'a}{a'a} \tag{12}$$

This value $\omega = b'a \big/ a'a$ is the distance along $a$ we need to travel until $a$ and $b - a$ are orthogonal to one another.

## 0.3 orthogonal projection as minimizer

What does orthogonality and minima have to do with each other?

Suppose we want to find the vector $p \in S$ such that $p$ is closer to $y \in B$ than any other vector $v \in S$, and we assume $S \subset B$, that $y$ and any vector $v$ cannot lie in th same space.

The distance between $y$ and any vector $v$ is

$$||y - v|| = [(y - v)'(y - v)]^{1/2}, \tag{13}$$

and if a vector $p$ is closest in distance $||.||$ than it will be closes in squared distance too $||.||^2$.

So we're searching for a vector $p$ so that

$$||y - v||^2 = (y - v)'(y - v) \tag{14}$$

is as small as possible.

First we introduce this smallest vector $p$ without changing the above equation

$$||y - p + p - v||^2 = \{[(y - p) + (p - v)]'[(y - p) + (p - v)]\} \tag{15}$$
$$= (y - p)'(y - p) + (p - v)'(p - v) + 2(p - v)'(y - p) \tag{16}$$
$$= ||y - p||^2 + ||p - v||^2 + 2(p - v)'(y - p) \tag{17}$$

the first two terms here cannot be changed much, but lets look at the third term. $p$ and $v$ are both vectors in $S$ and so their subtraction is a vector in $S$. $y$ is in $B$ and $p$ is in $S$. if we suppose $p$ is the vector such that the difference $y - p$ is orthogonal to **every** possible vector in $S$ then the third term would equal $0$.

The vector $p$ is smallest if and only if the difference between $y$ and $p$ is orthogonal to every vector in $S$.

### 0.3.1 (aside) span

We can represent any vector in a space $S$ through a basis. A basis is a set of independent vectors such that every vector in $S$ is the weighted sum of basis vectors.

Suppose $a$ is in some space $V$. Then a basis is a set of vectors $v_1, v_2, \cdots, v_n$ such that

$$a = \sum_{i=1}^{N} \alpha_i v_i$$

for every vector $a \in V$.

Returning back to our vector $p$, this vector is the one so that $y - p$ is orthogonal to every vector in $S$, or

$$(y - p)' \left( \sum_{i=1}^{N} \alpha_i v_i \right) = \sum_{i=1}^{N} \alpha_i (y - p)' v_i = 0$$

of $y - p$ must be orthogonal to every basis vector.

## 0.4 reframing our problem in linear algebra

We can use material on orthogonal projections to help us understand the optimal $\beta$.

Our **design** matrix $X$ times $\beta$ can be thought of as a basis.

$$X\beta = \begin{bmatrix} x_{1,1}\beta_1 + x_{1,2}\beta_2 + \cdots + \beta_n x_{1,n} \\ x_{2,1}\beta_1 + x_{2,2}\beta_2 + \cdots + \beta_n x_{2,n} \\ x_{3,1}\beta_1 + x_{3,2}\beta_2 + \cdots + \beta_n x_{3,n} \\ \vdots \\ x_{m,1}\beta_1 + x_{m,2}\beta_2 + \cdots + \beta_n x_{m,n} \end{bmatrix} = \beta_1 \begin{bmatrix} x_{1,1} \\ x_{2,1} \\ \vdots \\ x_{m,1} \end{bmatrix} + \beta_2 \begin{bmatrix} x_{1,2} \\ x_{2,2} \\ \vdots \\ x_{m,2} \end{bmatrix} + \cdots + \beta_n \begin{bmatrix} x_{1,n} \\ x_{2,n} \\ \vdots \\ x_{m,n} \end{bmatrix} = \sum_{i=1}^{N} \beta_i x_{;,i}$$

The $y$ observations can also be considered a single $m$-dimensional vector.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}.$$

Instead of asking for the *beta* that minimizes the **SSE**, let's instead ask for the vector that is a member of the space spanned by the columns of $X$ and closest to the vector $y$. This could be an alternative expression for "good fit to the data".

This best vector's (denoted $b$ for best) difference from $y$ must be orthogonal to all vectors in the space, or equivalently all vectors in the basis.

$$(y - b)' x_{;1} = 0 \tag{18}$$
$$(y - b)' x_{;2} = 0 \tag{19}$$
$$(y - b)' x_{;3} = 0 \tag{20}$$
$$\vdots = 0 \tag{21}$$
$$(y - b)' x_{;1} = 0 \tag{22}$$

or

$$y' x_{;1} - b' x_{;1} = 0 \tag{23}$$
$$y' x_{;2} - b' x_{;2} = 0 \tag{24}$$
$$y' x_{;3} - b' x_{;3} = 0 \tag{25}$$
$$\vdots = 0 \tag{26}$$
$$y' x_{;n} - b' x_{;n} = 0 \tag{27}$$
$$\tag{28}$$

rearranging terms

$$y'x_{;1} = b'x_{;1} \tag{29}$$
$$y'x_{;2} = b'x_{;2} \tag{30}$$
$$y'x_{;3} = b'x_{;3} \tag{31}$$
$$\vdots = 0 \tag{32}$$
$$y'x_{;n} = b'x_{;n} \tag{33}$$
$$\tag{34}$$

We can rewrite both sides of the above equation as a matrix times a vector.

$$X'y = X'b \tag{35}$$

We can take this equation further by remembering $b$ must be a member of the space created by the columns of $X$. That is $b$ is a weighted sum of the columns of $X$, for weights (let's say) $\beta$.

$$b = \sum_{i=1}^{N} \beta_i x_{;i} \tag{36}$$

$$= \beta_1 \begin{bmatrix} x_{1,1} \\ x_{2,1} \\ \vdots \\ x_{m,1} \end{bmatrix} + \beta_2 \begin{bmatrix} x_{1,2} \\ x_{2,2} \\ \vdots \\ x_{m,2} \end{bmatrix} + \cdots + \beta_n \begin{bmatrix} x_{1,n} \\ x_{2,n} \\ \vdots \\ x_{m,n} \end{bmatrix} = X\beta \tag{37}$$

and the above equation now becomes

$$X'y = X'b \tag{38}$$
$$X'y = X'X\beta \tag{39}$$

This is **exactly** the same equation as before. Minimizing the sum squares of error is the same as finding the vector $b$, constrained to be a weighted sum of the columns of $X$, closest to the vector $y$.

$$\beta = (X'X)^{-1}X'y \tag{40}$$

## 0.5 hat matrix

Now that we know how to compute optimal weights ($\beta$) for our vector $b$, we see the vector closest to $y$ is

$$b = X\beta, \tag{41}$$

but this vector is just the functional form we specified for our model, minus the error. The vector $b$ is used to make predictions about $y$ given data $X$, so that

$$\hat{y} = Xb \tag{42}$$

$$\hat{y} = X\left[(X'X)^{-1}X'y\right] \tag{43}$$

$$\hat{y} = \left[X(X'X)^{-1}X'\right]y \tag{44}$$

$$\tag{45}$$

Considered a function, the matrix

$$H = \left[X(X'X)^{-1}X'\right]$$

is called the **hat matrix** because it transforms $y$ into the vector $\hat{y}$, it places the "hat" on $y$.