

Fitting Linear models by least squares

Given TS $(x_1, y_1), \dots, (x_n, y_n)$

Suppose association looks linear
⇒ fit straight line $\ell(x) = b_0 + b_1 x$

How to find b_0, b_1 ?

Minimize resubstitution error:

$$(\hat{b}_0, \hat{b}_1) = \underset{b_0, b_1}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 \quad (*)$$

choose b_0 and b_1 that give the best predictive performance for training sample.

Differentiate (*) w.r.t b_0 and b_1

$$\hat{b}_1 = \sum_i (x_i - \bar{x})(y_i - \bar{y}) / \sum (x_i - \bar{x})^2$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

Generalize to more than 1 predictor variable

$$\ell(\underline{x}) = b_0 + b_1 x_1 + \dots + b_p x_p$$

Fitting a linear model as a linear algebra problem

Notation

$$\underline{a} = (a_1, \dots, a_p)^T$$

For \mathbf{X} $n \times p$ matrix

- \underline{x}_i is i -th row of \mathbf{X}
- \mathbf{X}^j is j -th column of \mathbf{X}

Inner product of \underline{a} and \underline{b} $\langle \underline{a}, \underline{b} \rangle = \sum_i a_i b_i$

$\|\underline{a}\|^2 = \langle \underline{a}, \underline{a} \rangle = \sum a_i^2$ squared norm of \underline{a}

$$\langle \underline{a}, \underline{b} \rangle = \|\underline{a}\| \|\underline{b}\| \cdot \cos \angle(\underline{a}, \underline{b})$$

\underline{a} and \underline{b} are called orthogonal $\underline{a} \perp \underline{b}$
if $\langle \underline{a}, \underline{b} \rangle = 0$

First consider simplest regression problem:
Fitting a line through the origin

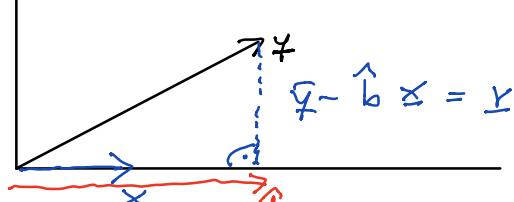
Find b such that $\sum_i (y_i - b x_i)^2 = \min !$

Differentiate wrt b :

$$\sum_i (y_i - b x_i) x_i = 0 \Leftrightarrow \underbrace{\langle \underline{y} - b \underline{x}, \underline{x} \rangle}_{\text{\underline{x} residual vector}} = 0$$

$$\langle \underline{y}, \underline{x} \rangle - b \langle \underline{x}, \underline{x} \rangle = 0$$

$$\hat{b} = \frac{\langle \underline{y}, \underline{x} \rangle}{\|\underline{x}\|^2} \quad \underline{y} = (y_1, \dots, y_n) \\ \underline{x} = (x_1, \dots, x_n)$$



$$\text{Pythagoras } \|\underline{y}\|^2 = \|\underline{b} \underline{x}\|^2 + \|\underline{\varepsilon}\|^2$$

$$\|\underline{\varepsilon}\|^2 = \|\underline{y}\|^2 - \|\underline{b} \underline{x}\|^2$$

$$= \|\underline{y}\|^2 - \frac{\langle \underline{y}, \underline{x} \rangle^2}{\|\underline{x}\|^4} \|\underline{x}\|^2$$

$$= \|\underline{y}\|^2 - \frac{\langle \underline{x}, \underline{y} \rangle^2}{\|\underline{x}\|^2}$$

Moving on: multiple predictor variables

\underline{X} : design matrix

\underline{X} has n rows

p columns $p = \text{number of predictor variables} + 1$

TS: $(\underline{x}_1, y_1), \dots, (\underline{x}_n, y_n)$

$$\underline{X}: n \left[\begin{array}{c} \underline{x}_1 \\ \vdots \\ \underline{x}_2 \\ \vdots \\ \vdots \\ \underline{x}_n \end{array} \right]$$

$$\|\underline{\varepsilon}\|^2 = \|\underline{y} - \underline{X} \underline{b}\|^2 \text{ want to minimize}$$

$$\hat{\underline{b}} = \underset{\underline{b}}{\operatorname{argmin}} \|\underline{y} - \underline{X} \underline{b}\|^2$$

clearly $(\underline{y} - \underline{X} \hat{\underline{b}})$ has to be orthogonal to $\underline{X}^1, \dots, \underline{X}^p$

Suppose \underline{X}^i was not orthogonal to

$\hat{y} - X\hat{b}$ Then

we could take \hat{y} and regress it on X^\top and reduce the residual sum of squares

In matrix form

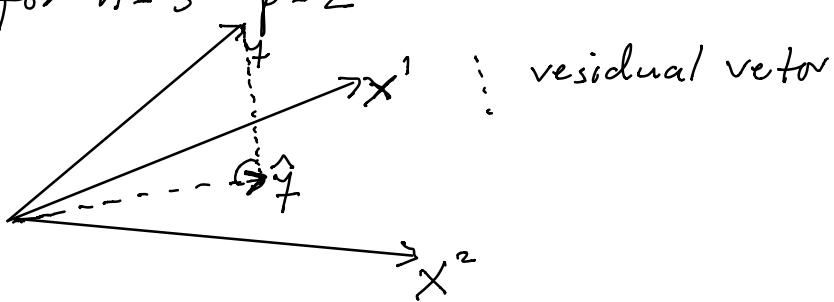
$$\begin{aligned} \cdot X^\top (\hat{y} - X\hat{b}) &= 0 \rightarrow \text{residual vector } \perp \text{ to} \\ X^\top \hat{y} - X^\top X \hat{b} &= 0 \quad \text{or columns of } X. \end{aligned}$$

$$\hat{b} = (X^\top X)^{-1} X^\top \hat{y} \quad \text{"normal equations"}$$

Define $\hat{y} = X\hat{b}$ vector of predicted values

$y = \hat{y} + \hat{\epsilon}$ $\perp [X^1 \dots X^p]$ space spanned by $X^1 \dots X^p$

Picture for $n=3$, $p=2$



$\Rightarrow \hat{y}$ = projection of y on $[X^1 \dots X^p]$

$$\begin{aligned} \hat{y} &= X\hat{b} = \underbrace{X(X^\top X)^{-1} X^\top}_H \hat{y} \\ &= H \cdot \text{hat matrix} \end{aligned}$$

H is a projection: H is symmetric

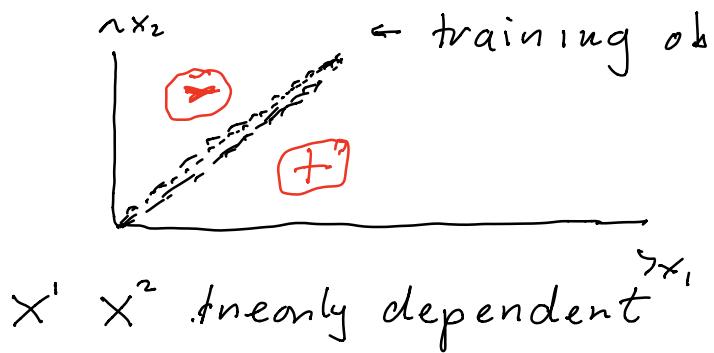
$H^\top = H$ and idempotent $H^2 = H$

Verify idempotence:

$$H^2 = X(X^T X)^{-1} \underbrace{X^T}_{\leftarrow} \underbrace{X(X^T X)^{-1}}_{\leftarrow} X^T$$
$$= X(X^T X)^{-1} X^T = H$$

Note: If $\text{rank } X < p \Rightarrow$ problem

why: $\text{rank}(X^T X) = \text{rank}(X) \Rightarrow (X^T X)^{-1}$
does not exist $\Rightarrow \hat{b}$ is not unique.



Coefficient vector \hat{b} not unique but
 \hat{y} = projection of y on $X^1 \dots X^p$ always unique.

If $\text{rank}(X) < p$ then \hat{y} is unique but
there are infinitely many ways to represent
 \hat{y} as a linear combination of $X^1 \dots X^p$

Handling categorical predictors

X is predictor variable with values a_1, \dots, a_k
categorical.

Convert X into $(K-1)$ indicator variables

$$z^1 \dots z^{K-1}$$

$$z_i^j = 1 \text{ if } x_i = a_j$$

↓
value of z^j for observation i

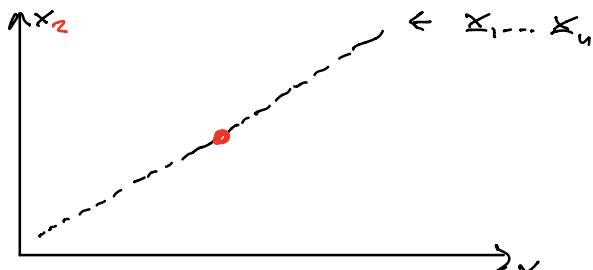
Why $K-1$? If we had K indicator variables, then $\sum(z^1 + z^2 \dots + z^K) = 1_n$

But design matrix already has column n of $1 \Rightarrow X$ would be rank deficient

Interpreting regression coefficients

$$y(\underline{x}) = \langle \hat{b}, \underline{x} \rangle = \hat{b}_1 x_1 + \dots + \hat{b}_p x_p$$

Interpretation of b_i = rate of change in $y(\underline{x})$ if x_i is changed but the values of all other predictors are held fixed.



No information in data about what would happen to y if I changed x_i but x_i was held fixed.