

# Linear Methods for Regression

Zed

March 14, 2017

## 1 Ordinary Least Squares

We write the linear regression model

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j = X^\top \beta$$

where  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^\top$ .  $X = (1, X_1, \dots, X_p)^\top$  is a  $p+1$  column vector, with the inputs  $X_j$  being quantitative, factor variables ( $X_j = \mathbb{1}_{\{G=\mathcal{G}_j\}}$ ), transformation of quantitative (say  $\sin X_j$ ,  $\log X_j$ ), basis expansions ( $X_2 = X_1^2, X_3 = X_1^3, \dots$ ) or cross terms ( $X_3 = X_2 X_1$ ). We have a quick review of the familiar OLS estimator before proceeding to new concepts and models.

### 1.1 Algebraic Properties

**Def. Least Squares Estimator:** We choose squared error as loss function, and solve

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \beta)^2 = \underset{\beta}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)$$

by the familiar method of moments, and get  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ ;

the prediction for *training set* is  $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ , which is, geometrically, an orthogonal projection of  $\mathbf{y}$  onto the column space of  $\mathbf{X}$ , i.e.  $\mathcal{C}(\mathbf{X}) = \operatorname{span}\{\operatorname{Cols}(\mathbf{X})\}$ . A few recap and highlights:

- (*Orthogonal Projection*)  $\hat{\mathbf{y}}$  is within  $\mathcal{C}(\mathbf{X})$ , since  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ , a linear combination of the columns of  $\mathbf{X}$ . The residual  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to the subspace  $\mathcal{C}(\mathbf{X})$ , since  $\mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) = 0$ .
- (*Orthogonal Complement*) Our sample  $\mathbf{y} \in \mathbb{R}^N$ , which can always be decomposed as  $\mathbb{R}^N = V \oplus V^\perp$ , where  $V$  is a subspace,  $V^\perp$  is the orthogonal complement of  $V$ . We already have the column space  $\mathcal{C}(\mathbf{X})$ , and we can show that  $\mathcal{C}(\mathbf{X})^\perp = \mathcal{N}(\mathbf{X}^\top)$ , the null space of  $\mathbf{X}^\top$ , which has dimension  $N - p - 1$ .  
*Proof.* Suppose  $\mathbf{z} \in \mathcal{C}(\mathbf{X})^\perp$ , then  $\mathbf{z}^\top \mathbf{X}\beta = 0$  for all linear combination parameter  $\beta \neq 0$ . Hence the only way is  $\mathbf{z}^\top \mathbf{X} = \mathbf{0}$ , i.e.  $\mathbf{X}^\top \mathbf{z} = \mathbf{0}$ .  $\square$
- (*Hat Matrix*) The matrix  $\mathbf{H}_\mathbf{X} := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is called the “hat” matrix, which maps a vector to its orthogonal projection on  $\mathcal{C}(\mathbf{X})$ . (symmetric, idempotent, and maps columns of  $\mathbf{X}$  to itself.) A curious object is the trace of this matrix:

$$\operatorname{tr}(\mathbf{H}_\mathbf{X}) = \operatorname{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \operatorname{tr}(\mathbf{I}_{p+1}) = p+1$$

- (*Residual*) We are also interested in the error of the estimator *within the training set*, i.e. define  $\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}}$  as the residual term. It follows immediately that the residual sum of square  $RSS = \hat{\mathbf{u}}^\top \hat{\mathbf{u}}$ . And apply the hat matrix we see  $\hat{\mathbf{u}} = (\mathbf{I}_N - \mathbf{H}_X)\mathbf{y}$ . The object in between is also symmetric, idempotent, due to these property of  $\mathbf{H}_X$ ; consider

$$(\mathbf{I} - \mathbf{H}_X)(\mathbf{I} - \mathbf{H}_X) = \mathbf{I} - 2\mathbf{H}_X + \mathbf{H}_X$$

- (*When  $\mathbf{X}^\top \mathbf{X}$  is Singular*) When columns of  $\mathbf{X}$  are linearly dependent,  $\mathbf{X}^\top \mathbf{X}$  becomes singular, and  $\hat{\beta}$  is not uniquely defined. But  $\hat{\mathbf{y}}$  is still the orthogonal projection onto  $\mathcal{C}(\mathbf{X})$ , just with more than one way to do the projection.

## 1.2 Statistical Properties

**(Linear Assumptions)** To discuss statistical properties of  $\hat{\beta}$ , we assume that the linear model is the true model for the mean, i.e. the conditional expectation of  $Y$  is  $X\beta$ , and that the deviation of  $Y$  from the mean is additive, distributed as  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . That is

$$Y = \mathbb{E}[Y|X] + \epsilon = X\beta + \epsilon$$

We further assume that the inputs  $\mathbf{X}$  in the training set are fixed (non-random).

Under these assumptions, a few other highlights on statistical properties of OLS estimator:

- (*Expectation of  $\hat{\beta}$* )  $\mathbb{E}(\hat{\beta}) = \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \epsilon)] = \beta$ , i.e. it is an unbiased estimator.
- (*Variance of  $\hat{\beta}$* )  $\text{Var}(\hat{\beta}) = \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \epsilon^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ . That is, the estimator  $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$
- (*Residual Revisited*) With the assumption of the real model of  $\mathbf{y}$ , we can further write  $\hat{\mathbf{u}} = (\mathbf{I} - \mathbf{H}_X)\mathbf{y} = (\mathbf{I} - \mathbf{H}_X)(\mathbf{X}\beta + \epsilon) = (\mathbf{I} - \mathbf{H}_X)\epsilon$ . It is easy to see that  $\mathbb{E}[\hat{\mathbf{u}}] = \mathbb{E}[\mathbf{X}(\beta - \hat{\beta}) + \epsilon] = 0$ . And therefore

$$\text{Var}[\hat{\mathbf{u}}] = \mathbb{E}[\hat{\mathbf{u}}\hat{\mathbf{u}}^\top] = \mathbb{E}[(\mathbf{I} - \mathbf{H}_X)\epsilon\epsilon^\top(\mathbf{I} - \mathbf{H}_X)] = \sigma^2(\mathbf{I} - \mathbf{H}_X)$$

So, although the errors  $\epsilon$  are i.i.d., residuals  $\hat{\mathbf{u}}$  are correlated.

- (*Individual Residual Term*) Pick any individual residual  $\hat{u}_i$ ,  $\text{Var}[\hat{u}_i] = \sigma^2(1 - h_i)$ , where  $h_i$  is the  $i$ -th diagonal entry of  $\mathbf{H}_X$ . Furthermore  $\text{Cov}[\hat{u}_i, \hat{u}_j] = \sigma^2 h_{ij}$ ,  $i \neq j$ ,  $h_{ij}$  is the row  $i$ , column  $j$  entry in  $\mathbf{H}_X$ .

An unbiased estimator of residual variance (square of residual standard error:  $RSE^2$ ) is

$$\hat{\sigma}^2 = \frac{RSS}{N - p - 1} = \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{N - p - 1}$$

*Prop.*  $\mathbb{E}(\hat{\sigma}^2) = \sigma^2$ . We present two proofs.

*Proof (1).*

$$\mathbb{E}[\hat{\mathbf{u}}^\top \hat{\mathbf{u}}] = \mathbb{E}\left[\sum_{i=1}^N \hat{u}_i^2\right] = \sum_{i=1}^N \text{Var}[\hat{u}_i] = \sum_{i=1}^N \sigma^2(1 - h_i) \quad (1)$$

By the trace formula we have discussed,  $\sum h_i = \text{tr}(\mathbf{H}_X) = p+1$ . Hence  $(2) = \sigma^2(N-p-1)$ . We conclude that

$$(N - p - 1)\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}[\epsilon^\top (\mathbf{I} - \mathbf{H}_X) \epsilon] = (N - p - 1)\sigma^2 \quad \square.$$

Before the second proof, we present a lemma.

**Lemma. (Distribution of Quadratic Form)**

- If an  $n$ -vector  $\mathbf{x}$  is distributed as  $\mathcal{N}(\mathbf{0}, \Sigma)$ , then the quadratic form  $\mathbf{x}^\top \Sigma^{-1} \mathbf{x} \sim \chi^2(n)$ .
- If an  $n$ -vector  $\mathbf{x}$  is standard multivariate normal:  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\mathbf{H}_Z$  is a projection matrix onto the column space of  $\mathbf{Z}$ , which has dimension  $r$  (i.e. consider  $\mathbf{Z}$  is a  $n \times r$  matrix, and  $\mathbf{Z}$  and  $\mathbf{H}_Z$  both have rank  $r$ ); then the quadratic form  $\mathbf{x}^\top \mathbf{H}_Z \mathbf{x} \sim \chi^2(r)$ .

*Proof of lemma. (First Part)* Since  $\Sigma$  is symmetric positive definite, we have *Cholesky decomposition*  $\Sigma = \mathbf{Q}\mathbf{Q}^\top$ , where  $\mathbf{Q}$  is  $n \times n$  lower triangular.

$$\mathbf{x}^\top \Sigma^{-1} \mathbf{x} = \mathbf{x}^\top \mathbf{Q}^{-\top} \mathbf{Q}^{-1} \mathbf{x} = (\mathbf{Q}^{-1} \mathbf{x})^\top (\mathbf{Q}^{-1} \mathbf{x}) = \mathbf{z}^\top \mathbf{z}$$

in which we let  $\mathbf{z} := \mathbf{Q}^{-1} \mathbf{x}$ . It is clear that  $\mathbb{E}[\mathbf{z}] = \mathbf{Q}^{-1} \mathbb{E}[\mathbf{x}] = \mathbf{0}$ . And

$$\text{Var}[\mathbf{z}] = \mathbb{E}[\mathbf{Q}^{-1} \mathbf{x} (\mathbf{Q}^{-1} \mathbf{x})^\top] = \mathbf{Q}^{-1} \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{Q}^{-\top} = \mathbf{Q}^{-1} \text{Var}[\mathbf{x}] \mathbf{Q}^{-\top} = \mathbf{Q}^{-1} \Sigma \mathbf{Q}^{-\top} = \mathbf{I}$$

which indicates that  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  is an  $n$ -variate standard normal. It follows that  $\mathbf{z}^\top \mathbf{z} \sim \chi^2(n)$ .  $\square$

*(Second Part)*

$$\mathbf{x}^\top \mathbf{H}_Z \mathbf{x} = \mathbf{x}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{x} = \mathbf{y}^\top \mathbf{\Omega}^{-1} \mathbf{y}$$

in which we let  $\mathbf{y} := \mathbf{Z}^\top \mathbf{x}$  (an  $r \times 1$  vector), and  $\mathbf{\Omega} := \mathbf{Z}^\top \mathbf{Z}$  (an  $r \times r$  matrix). This is exactly the form in part 1. And the linear transform of  $n$ -variate normal:  $\mathbf{Z}^\top \mathbf{x}$  is distributed as  $r$ -variate normal  $\mathcal{N}(\mathbf{0}, \mathbf{Z}^\top \mathbf{Z})$ . By the result of part 1  $\Rightarrow \mathbf{x}^\top \mathbf{H}_Z \mathbf{x} \sim \chi^2(r)$ .  $\square$

*Proof (2).*

$$\begin{aligned} (N - p - 1) \hat{\sigma}^2 &= \hat{\mathbf{u}}^\top \hat{\mathbf{u}} = \mathbf{y}^\top (\mathbf{I} - \mathbf{H}_X)^\top (\mathbf{I} - \mathbf{H}_X) \mathbf{y} \\ &= \boldsymbol{\epsilon}^\top (\mathbf{I} - \mathbf{H}_X) \boldsymbol{\epsilon} = \boldsymbol{\epsilon}^\top \mathbf{H}_Z \boldsymbol{\epsilon} \end{aligned} \quad (2)$$

in which we let  $\mathbf{H}_Z := \mathbf{I} - \mathbf{H}_X$ . By previous result, this is also symmetric, idempotent, and projects any vector to the null space of  $\mathbf{X}^\top$ , the orthogonal complement of  $\mathcal{C}(\mathbf{X})$ . We can always compose a matrix  $\mathbf{Z}$  whose columns are the general solutions of  $\mathbf{X}^\top \mathbf{z} = 0$ . Clearly it has  $N - p - 1$  columns, since the orthogonal complement has dimension  $N - p - 1$ . Hence  $\mathbf{H}_Z$  has  $(N - p - 1)$  rank. Moreover,  $\boldsymbol{\epsilon}^\top \mathbf{H}_Z \boldsymbol{\epsilon} = \boldsymbol{\epsilon}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \boldsymbol{\epsilon}$ , and  $\mathbf{Z}^\top \mathbf{Z}$  is of  $(N - p - 1) \times (N - p - 1)$ . By lemma, and multiply a normalization factor  $\Rightarrow \mathbf{Z}^\top \boldsymbol{\epsilon} / \sigma \sim \mathcal{N}(\mathbf{0}, (\mathbf{Z}^\top \mathbf{Z}))$ ,  $\frac{1}{\sigma^2} \boldsymbol{\epsilon}^\top \mathbf{H}_Z \boldsymbol{\epsilon} \sim \chi^2(N - p - 1)$ . So:

$$\mathbb{E}[\boldsymbol{\epsilon}^\top \mathbf{H}_Z \boldsymbol{\epsilon}] = \sigma^2 (N - p - 1) \quad \square$$

Proof (2) gives us a stronger result:

*Prop. (Distribution of Sample Estimator of Variance)* The residual sum of square is Chi squared distributed with degree of freedom  $(N - p - 1)$ .

$$(N - p - 1) \hat{\sigma}^2 = RSS \sim \sigma^2 \chi^2(N - p - 1)$$

In addition,  $\hat{\beta}$  and  $\hat{\sigma}$  are independent.

### 1.3 Hypothesis Tests

**(t Statistic)** The  $t(n)$  distribution is defined as  $t(n) \sim \frac{\mathcal{N}(0,1)}{\sqrt{\chi^2(n)/n}}$ . To test hypothesis that a particular coefficient  $\beta_j = 0$ , we formulate the statistic

$$t_j = \frac{\hat{\beta}_j / \text{se}(\hat{\beta}_j)}{\sqrt{(N - p - 1) \hat{\sigma}^2 / (N - p - 1) \sigma^2}} = \frac{\hat{\beta}_j}{\hat{\sigma} \cdot \text{se}(\hat{\beta}_j) / \sigma} = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}}$$

where  $\hat{\sigma} = \sqrt{RSS/(N-p-1)}$ ,  $\sqrt{v_j}$  is the  $j$ -th diagonal element of  $(\mathbf{X}^\top \mathbf{X})^{-1}$ . And we know that  $\hat{\beta}_j/\text{se}(\hat{\beta}_j) \sim \mathcal{N}(\beta_j/\text{se}(\hat{\beta}_j), 1)$  and that  $\sqrt{(N-p-1)\hat{\sigma}^2/(N-p-1)\sigma^2} \sim \sqrt{\chi_{N-p-1}^2/(N-p-1)}$ . Under the null hypothesis  $\beta_j = 0$ ,  $\hat{\beta}_j/\text{se}(\hat{\beta}_j) \sim \mathcal{N}(0, 1)$ . We have  $t_j \sim t(N-p-1)$ . If we know  $\sigma$  before hand, we just use it instead of  $\hat{\sigma}$ . And  $t_j$  reduces to  $\hat{\beta}_j/\text{se}(\hat{\beta}_j) \sim \mathcal{N}(0, 1)$ . Where  $\text{se}(\hat{\beta}_j) = \sigma\sqrt{v_j}$ .

**(F Statistic)** The  $\mathcal{F}(n_1, n_2)$  distribution is defined as  $\mathcal{F}(n_1, n_2) \sim \frac{\chi^2(n_1)/n_1}{\chi^2(n_2)/n_2}$ . To test hypothesis that  $k$  coefficients  $\beta_{[1]} = \dots = \beta_{[k]} = 0$  simultaneously, we formulate the statistic

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

Where the bigger model 1 has  $p_1 + 1$  parameters, the smaller model 0 (corresponds to null hypothesis  $H_0$ ) has  $p_0 + 1$  parameters,  $p_1 - p_0 = k$ . We have  $F \sim \mathcal{F}(p_1 - p_0, N - p_1 - 1)$  under the null hypothesis.

**(Confidence Interval)** We can isolate  $\beta_j$  to form a  $1 - 2\alpha$  confidence interval

$$\beta_j \in (\hat{\beta}_j - z_{(1-\alpha)}\sqrt{v_j}\hat{\sigma}, \hat{\beta}_j + z_{(1-\alpha)}\sqrt{v_j}\hat{\sigma})$$

*Proof.* We know that  $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$ , a multivariate normal. So isolating  $\hat{\beta}_j$ , we have  $\hat{\beta}_j \sim \mathcal{N}(\beta_j, \sigma^2 v_j)$ , where, as before,  $v_j$  is the  $j$ -th diagonal element of the covariance matrix of  $\hat{\beta}$ .  $\text{se}(\hat{\beta}_j) = \sigma\sqrt{v_j}$ . And hence  $\frac{\hat{\beta}_j - \beta_j}{\sigma\sqrt{v_j}} \sim \mathcal{N}(0, 1)$ .

$$1 - 2\alpha = \mathbb{P}\left(\left|\frac{\hat{\beta}_j - \beta_j}{\sigma\sqrt{v_j}}\right| > z_{(1-\alpha)}\right) = \mathbb{P}\left(\hat{\beta}_j - z_{(1-\alpha)}\sqrt{v_j}\sigma < \beta_j < \hat{\beta}_j + z_{(1-\alpha)}\sqrt{v_j}\sigma\right)$$

And substitute  $\sigma$  with the estimate  $\hat{\sigma}$ , yields the result.  $\square$

**(Confidence Region)** We also obtain a confidence set for the entire parameter vector  $\beta$ ,

$$\beta \in C_\beta = \{(\hat{\beta} - \beta)^\top \mathbf{X}^\top \mathbf{X}(\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi_{p+1, (1-\alpha)}^2\}$$

*Proof.* We know  $\hat{\beta} - \beta \sim \mathcal{N}(\mathbf{0}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$ , by *lemma* (Dist of quadratic form) part 1,  $(\hat{\beta} - \beta)^\top \frac{1}{\sigma^2}(\mathbf{X}^\top \mathbf{X})(\hat{\beta} - \beta) \sim \chi^2(p+1)$ . Hence

$$1 - \alpha = \mathbb{P}\left((\hat{\beta} - \beta)^\top \frac{1}{\sigma^2}(\mathbf{X}^\top \mathbf{X})(\hat{\beta} - \beta) \leq \chi_{p+1, (1-\alpha)}^2\right) = \mathbb{P}\left((\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X})(\hat{\beta} - \beta) \leq \sigma^2 \chi_{p+1, (1-\alpha)}^2\right)$$

And substitute  $\sigma$  with the estimate  $\hat{\sigma}$ , yields the result.  $\square$

## 1.4 Gauss Markov Theorem

*Thm. (Gauss-Markov)* the least squares estimator has smallest variance among all *linear unbiased* estimates.

*Proof.* Let  $\tilde{\beta}$  be an unbiased linear estimator other than  $\hat{\beta}$ , which is the ols estimator. By linearity:  $\tilde{\beta} = \mathbf{A}\mathbf{y}$ , where  $\mathbf{A}$  is some (non-random) matrix. Hence we may decompose  $\tilde{\beta} = ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{C})\mathbf{y} = \hat{\beta} + \mathbf{C}\mathbf{y}$ , where we let  $\mathbf{C} := \mathbf{A} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .

By unbiasedness:  $\beta = \mathbb{E}[\tilde{\beta}] = \mathbb{E}[\mathbf{A}\mathbf{y}] = \mathbb{E}[\mathbf{A}(\mathbf{X}\beta + \epsilon)] = \mathbf{A}\mathbf{X}\beta + \mathbf{A}\mathbb{E}[\epsilon]$ . Since the last

term has mean  $\mathbf{0}$ , this requires  $\mathbf{A}\mathbf{X} = \mathbf{I} \Rightarrow \mathbf{C}\mathbf{X} = \mathbf{O}$ . Hence  $\mathbf{C}\mathbf{y} = \mathbf{C}(\mathbf{X}\beta + \epsilon) = \mathbf{C}\epsilon$ . Therefore

$$\begin{aligned}\text{Cov}[\hat{\beta}, \mathbf{C}\mathbf{y}] &= \text{Cov}[\hat{\beta}, \mathbf{C}\epsilon] = \mathbb{E}[(\hat{\beta} - \mathbb{E}\hat{\beta})(\mathbf{C}\epsilon - \mathbf{C}\mathbb{E}\epsilon)^\top] = \mathbb{E}[(\hat{\beta} - \beta)\epsilon^\top \mathbf{C}^\top] \\ &= \mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \epsilon^\top \mathbf{C}^\top] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{C}\mathbf{X})^\top = \mathbf{O}\end{aligned}\quad (3)$$

So:

$$\text{Var}[\tilde{\beta}] = \text{Var}[\hat{\beta} + \mathbf{C}\mathbf{y}] = \text{Var}[\hat{\beta} + \mathbf{C}\epsilon] = \text{Var}[\hat{\beta}] + \sigma^2 \mathbf{C}\mathbf{C}^\top \quad \square$$

## 1.5 Algorithm for Multiple Regression

For the univariate regression (with no intercept), we calculate ols estimator as:

$$\hat{\beta}_1 = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y} = \frac{\langle \mathbf{x}, \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{y} \rangle}$$

And the residual  $\mathbf{r} = \mathbf{y} - \mathbf{x}\hat{\beta}$ . Suppose  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0$ , i.e.  $\mathbf{X}$  is an orthogonal matrix, then  $\hat{\beta}_j = \langle \mathbf{x}_j, \mathbf{y} \rangle / \langle \mathbf{x}_j, \mathbf{x}_j \rangle$ , just write down  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  and use the fact that  $\mathbf{X}$  is orthogonal we can easily get the result. This implies that when the inputs are orthogonal, they have no effect on each other's parameter estimates in the model.

For non-orthogonal  $\mathbf{X}$ , we perform the *Gram-Schmidt* orthogonalization procedure:

*Algo.* (*Gram-Schmidt*) Suppose  $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p)$ .

1. Let  $\mathbf{z}_0 \leftarrow \mathbf{x}_0 \leftarrow \mathbf{1}$ .
2. For  $j = 1:p$ : Regress  $\mathbf{x}_j$  on  $\mathbf{z}_0, \dots, \mathbf{z}_{j-1}$  respectively to produce coefficients  $\hat{\gamma}_{ij} \leftarrow \langle \mathbf{z}_i, \mathbf{x}_j \rangle / \langle \mathbf{z}_i, \mathbf{z}_i \rangle$ ,  $i = 0, 1, \dots, j-1$ ;  $\hat{\gamma}_{jj} \leftarrow 1$ .
3. Calculate residual  $\mathbf{z}_j \leftarrow \mathbf{x}_j - \sum_{i=0}^{j-1} \hat{\gamma}_{ij} \mathbf{z}_i$
4. Regress  $\mathbf{y}$  on the residual  $\mathbf{z}_j$  to produce  $\hat{\beta}_j \leftarrow \langle \mathbf{z}_j, \mathbf{y} \rangle / \langle \mathbf{z}_j, \mathbf{z}_j \rangle$

*Prop.*  $\mathbf{Z} = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_p)$  is orthogonal.

*Proof.* We show by induction proof. Firstly, it is easy to see that

$$\langle \mathbf{z}_0, \mathbf{z}_1 \rangle = \langle \mathbf{z}_0, \mathbf{x}_1 - \frac{\langle \mathbf{z}_0, \mathbf{x}_1 \rangle}{\langle \mathbf{z}_0, \mathbf{z}_0 \rangle} \mathbf{z}_0 \rangle = \langle \mathbf{z}_0, \mathbf{x}_1 \rangle - \langle \mathbf{z}_0, \mathbf{x}_1 \rangle = 0$$

We assume  $\langle \mathbf{z}_0, \mathbf{z}_k \rangle = 0$  for all  $1 < k \leq j < p$ . Then for  $k = j+1$ :

$$\langle \mathbf{z}_0, \mathbf{z}_{j+1} \rangle = \langle \mathbf{z}_0, \mathbf{x}_{j+1} - \sum_{l=0}^j \frac{\langle \mathbf{z}_l, \mathbf{x}_{j+1} \rangle}{\langle \mathbf{z}_l, \mathbf{z}_l \rangle} \mathbf{z}_l \rangle = \langle \mathbf{z}_0, \mathbf{x}_{j+1} \rangle - \langle \mathbf{z}_0, \frac{\langle \mathbf{z}_0, \mathbf{x}_{j+1} \rangle}{\langle \mathbf{z}_0, \mathbf{z}_0 \rangle} \mathbf{z}_0 \rangle = 0$$

So we conclude that  $\langle \mathbf{z}_0, \mathbf{z}_j \rangle = 0$  for  $j = 1, 2, \dots, p$ . Do the same induction for  $\mathbf{z}_1$  as follows:

- Base case, using the fact (what we already known):  $\langle \mathbf{z}_0, \mathbf{z}_1 \rangle = 0$

$$\langle \mathbf{z}_1, \mathbf{z}_2 \rangle = \langle \mathbf{z}_1, \mathbf{x}_2 - \frac{\langle \mathbf{z}_0, \mathbf{x}_2 \rangle}{\langle \mathbf{z}_0, \mathbf{z}_0 \rangle} \mathbf{z}_0 - \frac{\langle \mathbf{z}_1, \mathbf{x}_2 \rangle}{\langle \mathbf{z}_1, \mathbf{z}_1 \rangle} \mathbf{z}_1 \rangle = \langle \mathbf{z}_1, \mathbf{x}_2 \rangle - \langle \mathbf{z}_1, \mathbf{x}_2 \rangle = 0$$

- The induction, assume  $\langle \mathbf{z}_1, \mathbf{z}_k \rangle = 0$  for all  $2 < k \leq j < p$ . Then for  $k = j+1$ :

$$\langle \mathbf{z}_1, \mathbf{z}_{j+1} \rangle = \langle \mathbf{z}_1, \mathbf{x}_{j+1} - \sum_{l=0}^j \frac{\langle \mathbf{z}_l, \mathbf{x}_{j+1} \rangle}{\langle \mathbf{z}_l, \mathbf{z}_l \rangle} \mathbf{z}_l \rangle = \langle \mathbf{z}_1, \mathbf{x}_{j+1} \rangle - \langle \mathbf{z}_1, \frac{\langle \mathbf{z}_1, \mathbf{x}_{j+1} \rangle}{\langle \mathbf{z}_1, \mathbf{z}_1 \rangle} \mathbf{z}_1 \rangle = 0$$

So we conclude that  $\langle \mathbf{z}_1, \mathbf{z}_j \rangle = 0$  for  $j = 2, \dots, p$ . And the induction for  $\mathbf{z}_i$ ,  $i = 2, 3, \dots, p-1$  in the same fashion, we have  $\mathbf{Z}$  is orthogonal.  $\square$

Another observation is that  $\mathbf{x}_j$  is a linear combination of  $\mathbf{z}_k$ , for  $k \leq j$ . Hence  $\mathbf{Z}$  is a orthogonal basis for the column space of  $\mathbf{X}$ . Let  $\mathbf{D} = \text{diag}(\|\mathbf{z}_j\|)$ , then  $\mathbf{Z}\mathbf{D}^{-1}$  gives the *orthonormal basis* of column space of  $\mathbf{X}$ . We denote  $\mathbf{Q} := \mathbf{Z}\mathbf{D}^{-1}$ , which is also an orthogonal matrix.

By writing the algo in a matrix form, we denote  $\mathbf{\Gamma} = \{\hat{\gamma}_{ij}\}$ , which is an upper triangular matrix with main diagonal entries being 1s. And hence we have

$$\mathbf{X} = \mathbf{Z}\mathbf{\Gamma} = \mathbf{Z}\mathbf{D}^{-1}\mathbf{D}\mathbf{\Gamma} =: \mathbf{Q}\mathbf{R}$$

And the ols estimator given by

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

## 2 Subset Selection

- (*Best-Subset Selection*) Look at all possible models at every given number ( $k$ ) of variables chosen. (computationally expensive, becomes infeasible for  $p$  much larger than 30-40 or so)
- (*Forward-Stepwise Selection*) Rather than search through all possible subsets, we want to seek a path through them. FSS proceeds by sequentially adds into the model the predictor that most improves the fit. This is characterized as a *greedy algorithm*, which must produce a nested sequence of models, i.e. it may not find the best model, when, for example, the best subset of size 2 does not include that of size 1 (which may happen). However, it has lower variance compared with best-subset.
- (*Backward-Stepwise Selection*) Starts with the full model, and sequentially deletes the predictors that has the least impact on the fit. Can only be used for  $N > p$ .
- (*Forward-Stagewise (FS) Selection*) Start as the forward-stepwise, with intercept  $\bar{y}$ , and centered predictors with coefficients initially set as 0. Then at each step, choose the variable that are most *correlated* with the current residual, then compute simple regression param  $\gamma$  of residual on this variable, add this to the current  $\beta_j$ , i.e.  $\beta_j \leftarrow \beta_j + \gamma$ . Continues until none are correlated with the residual. The convergence of this algorithm can be slow, but it has good performance for problems with high dimensionality.

Subset selection is a *discrete* process, we either include a variable or exclude it. As a result it often exhibits high variance. Shrinkage methods are more continuous, and do not suffer from high variability.

## 3 Shrinkage Methods

The motivation of various shrinkage methods is to overcome the *combinatorial explosion* of the number of possible subsets (when  $p$  large) by converting the discrete problem to a continuous one, which turn out to be simpler to solve.

### 3.1 Ridge Regression

The ridge regression shrinks the regression coeffs by imposing a penalty on the magnitudes of these coefficients. Denote the ridge estimator  $\hat{\beta}^{ridge}$ , it minimizes a penalized sum of squares:

$$\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad (*)$$

where  $\lambda$  is a hyperparameter (complexity parameter) that controls the amount of shrink. The solution of (\*) are not equivalent under the scaling of inputs. Hence we usually standardize the input before solving (\*). In addition, we don't penalize the magnitude  $\beta_0$

The standardization procedure is done as: calculate the centered inputs as  $x_{ij} - \bar{x}_j$ , (in the following text we assume  $\{\mathbf{X}\}_{ij}$  is this, has  $p$  columns without  $\mathbf{1}$ ), and estimate  $\hat{\beta}_0$  by  $\sum_1^N y_i / N$ .

**Def. Ridge Regression Estimator:** We minimize loss function with penalization:

$$\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^\top \beta \}$$

$$\text{We have } \partial RSS(\lambda) / \partial \beta = 2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) + 2\lambda\beta = 0 \Rightarrow \hat{\beta}^{ridge} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$