

## Factor Predictor

A factor predictor is a predictor that takes a discrete set values on "nominal" scale (non numerical)

e.g. Drug A, B, C

Femail, Male

Consider the case where the predictor takes  $M$  possible values.

$$x = \begin{cases} 1 \\ 2 \\ \vdots \\ M \end{cases}$$

We define  $M$  "dummy" or indicator variables

$$x_m = \begin{cases} 1 & x=m \\ 0 & \text{otherwise} \end{cases}$$

e.g.

|        | $x^{(1)}$ | $x^{(2)}$ | $x^{(3)}$ |
|--------|-----------|-----------|-----------|
| Drug A | 1         | 0         | 0         |
| B      | 0         | 1         | 0         |
| C      | 0         | 0         | 1         |

redundant

Remove last column and using only  $M-1$  variables.

$$x^{(l)} = \begin{cases} 1 & x=l \\ 0 & \text{otherwise} \end{cases} \quad \text{for } l=0 \dots M-1$$

when  $l=M \Rightarrow$  baseline level

We may write  $X_1$  - categorial with  $M$  levels.

$$E[Y | X_1] = \beta_0 + \sum_{l=1}^{M-1} \beta_l x_1^{(l)}$$

$$(1) X=1 \quad (X^{(1)}, X^{(2)} \dots X^{(M-1)}) = (1 \ 0 \ 0 \ \dots \ 0)$$

$$E[Y|X=1] = \beta_0 + \beta_1$$

$$(2) X=2 \quad (X^{(1)}, X^{(2)} \dots X^{(M-1)}) = (0 \ 1 \ 0 \ \dots \ 0)$$

$$E[Y|X=2] = \beta_0 + \beta_2$$

$$\vdots$$

$$(M-1) \quad X=M-1 \quad (X^{(1)}, X^{(2)} \dots X^{(M-1)}) = (0 \ 0 \ \dots \ 0 \ 1)$$

$$E[Y|X=M-1] = \beta_0 + \beta_{M-1}$$

$$(M) \quad X=M \quad (X^{(1)}, X^{(2)} \dots X^{(M)}) = (0 \ 0 \ 0 \ \dots \ 0)$$

$$E[Y|X=M] = \beta_0$$

Therefore,

$$\beta_1 = E[Y|X=1] - E[Y|X=M]$$

$$\beta_2 = E[Y|X=2] - E[Y|X=M]$$

$\vdots$

$$\beta_{M-1} = E[Y|X=3] - E[Y|X=M]$$

$$\beta_0 = E[Y|X=M] \quad \text{baseline level effect}$$

$$\Rightarrow \beta'_l = \begin{cases} \beta_0 & l=M \\ \beta_l + \beta_0 & l=1 \dots M-1 \end{cases}$$

$$\Rightarrow E[Y|X] = \sum_{l=1}^M \beta'_l X^{(l)} \quad \text{b.c. } X = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\text{e.g. } X=M \quad (X^{(1)} \dots X^{(M)}) = (0 \ 0 \ \dots \ 0 \ 1)$$

$$E[Y|X=M] = \beta'_M = \beta_0$$

$$X=1 \quad (X^{(1)} \dots X^{(M)}) = (1 \ 0 \ \dots \ 0 \ 0)$$

$$E[Y|X=1] = \beta'_1 = \beta_1 + \beta_0$$

In this case,  $\beta'_m$  can be interpreted directly as the effect of level  $m$ , for  $m=1 \dots M$

## Interaction

For two distinct predictors, we may consider the joint effect of the predictors simultaneously.

### continuous predictors

We just include product term e.g.  $x_1 \cdot x_2$

$$\Rightarrow E[Y|X] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 \cdot x_2$$

modifications of the effect of outcome of  $x_1$  on the presence of  $x_2$  (or vice versa)

interaction

$$\begin{aligned} E[Y|X] &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 \cdot x_2 \\ &= \beta_0 + (\beta_1 + \beta_{12} x_2) x_1 + \beta_2 x_2 \\ &= \beta_0 + \beta_1 x_1 + (\beta_2 + \beta_{12} x_1) x_2 \end{aligned}$$

$\left[ \begin{array}{c} \beta_0 \\ \beta_1 + \beta_{12} x_2 \\ \beta_2 \end{array} \right]$        $\left[ \begin{array}{c} \beta_0 \\ \beta_1 \\ \beta_2 + \beta_{12} x_1 \end{array} \right]$

### interaction and correlation

correlation  $x_1 \perp\!\!\!\perp x_2 \quad \text{cov}(x_1, x_2)$

interaction  $D(Y, x_1) | x_2$

### factor predictor

For the factor predictors with  $M_1$  and  $M_2$  levels, an interaction term consider all possible combinations of the level of each factor.

i.e. There are  $M_1 \times M_2$  possible combinations

$$E[Y|X] = \frac{\sum_{m=1}^{M_1} \beta_m x^{(m)} + \sum_{m'=1}^{M_2} \beta_{m'} x^{(m')}}{x_1 x_2} + \sum_{m=1}^{M_1} \sum_{m'=1}^{M_2} \beta_{mm'} x^{(m)} x^{(m')}$$

If the intercept is used, then we could just use  $(M_1-1)(M_2-1)$  terms to represent the interaction.

$$E[Y|X] = \beta_0 + \sum_{m=1}^{M_1-1} \beta_m x^{(m)} + \sum_{m'=1}^{M_2-1} \beta_{m'} x^{(m')} + \sum_{m=1}^{M_1-1} \sum_{m'=1}^{M_2-1} \beta_{mm'} x^{(m)} x^{(m')}$$

e.g.  $X_1$  - categorical A, B, C

$X_2$  - categorical Female, Male

| $X_1$ | $X_1^{(1)}$ | $X_1^{(2)}$ | $X_1^{(3)}$ | $X_2$  | $X_2^{(1)}$ | $X_2^{(2)}$ |
|-------|-------------|-------------|-------------|--------|-------------|-------------|
| A     | 1           | 0           | 0           | Female | 1           | 0           |
| B     | 0           | 1           | 0           | Male   | 0           | 1           |
| C     | 0           | 0           | 1           |        |             |             |

Model 1 (with main effect only)

$$E[Y|X] = \frac{\beta_1 X_1^{(1)} + \beta_2 X_1^{(2)} + \beta_3 X_1^{(3)}}{X_1} + \frac{\beta_4 X_2^{(1)} + \beta_5 X_2^{(2)}}{X_2}$$

Model 2 (with main effect only)

$$E[Y|X] = \beta_0' + \beta_1' X_1^{(1)} + \beta_2' X_1^{(2)} + \beta_4' X_2^{(1)}$$

| $X_1$ | $X_2$  | $X_1^{(1)}$ | $X_1^{(2)}$ | $X_2^{(1)}$ | Model                            |
|-------|--------|-------------|-------------|-------------|----------------------------------|
| A     | Female | 1           | 0           | 1           | $\beta_0' + \beta_1' + \beta_4'$ |
| B     | Female | 0           | 1           | 1           | $\beta_0' + \beta_2' + \beta_4'$ |
| C     | Female | 0           | 0           | 1           | $\beta_0' + \beta_4'$            |
| A     | Male   | 1           | 0           | 0           | $\beta_0' + \beta_1'$            |
| B     | Male   | 0           | 1           | 0           | $\beta_0' + \beta_2'$            |
| C     | Male   | 0           | 0           | 0           | $\beta_0'$                       |

$$\beta_2' = E[Y | X_1=B, X_2=\text{Male}] - E[Y | X_1=C, X_2=\text{Male}]$$

the effect contrast between.

## Model 2 (with interaction term)

$$E[Y|X] = \beta_0 + \underbrace{\beta_1' X_1^{(1)} + \beta_2' X_1^{(2)}}_{X_1} + \frac{\beta_3' X_2^{(1)}}{X_2} + \frac{\beta_{12} X_1^{(1)} X_2^{(1)}}{A} + \frac{\beta_{23} X_1^{(2)} X_2^{(2)}}{B}$$

not allowed (?)

When there is  $\beta_0$ , total number of interactions between  $X_1$  ( $M_1$  levels) and  $X_2$  ( $M_2$  levels) is

$$(M_1 - 1)(M_2 - 1) = 2 \times 1 = 2$$

| $X_1$ | $X_2$  | $X_1^{(1)}$ | $X_1^{(2)}$ | $X_2^{(1)}$ | $E[Y X_1, X_2]$                                |
|-------|--------|-------------|-------------|-------------|--|
| A     | Female | 1           | 0           | 1           | $\beta_0' + \beta_1' + \beta_3' + \beta_{12}'$ |
| B     | Female | 0           | 1           | 1           | $\beta_0' + \beta_2' + \beta_3' + \beta_{23}'$ |
| C     | Female | 0           | 0           | 1           | $\beta_0' + \beta_3'$                          |
| A     | Male   | 1           | 0           | 0           | $\beta_0' + \beta_1'$                          |
| B     | Male   | 0           | 1           | 0           | $\beta_0' + \beta_2'$                          |
| C     | Male   | 0           | 0           | 0           | $\beta_0'$                                     |

## Model selection

criterion: generalization error

We estimate the model by minimizing least-squares

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \frac{1}{n} (Y - X\beta)^T (Y - X\beta)$$

$$\text{let } \hat{m}(x) = X\hat{\beta}$$

To measure how well  $\hat{m}(x)$  can predict, we define generalization error

Imagine there are new data points  $(Y^0, X^0)$   
 generalization error  $G = E[(Y^0 - \hat{m}(X^0))^2]$   $X^0 = (1, X_1^0 \dots X_K^0)$

Training error

$$T = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{m}(X_i))^2$$

$T$  is a poor estimation of  $\alpha$ , usually  $T < \alpha$

The training data is obtained from the model

$$Y = X\beta + \varepsilon$$

From this data we obtain  $\hat{\beta}$ . Our fitted model is

$$\hat{m}(X) = X\hat{\beta}$$

The training error (in-sample error)

$$MSE = \frac{SS_{RES}}{n} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i\hat{\beta})^2$$

Now imagine we get a new data at the same  $X_i'$

$$X_i'^0 = X_i \quad Y^0 = \underset{n \times 1}{X^0 \beta} + \underset{n \times 1}{\varepsilon^0} = X\beta + \varepsilon^0$$

Testing error (out-of-sample prediction error)

$$\frac{1}{n} \sum_{i=1}^n (Y_i^0 - \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i^0 - X_i'^0 \hat{\beta})^2 = \frac{1}{n} \sum_{i=1}^n (Y_i^0 - X_i \hat{\beta})^2$$

We can see that

$$\begin{aligned} \underset{\text{testing error}}{E\left[\frac{1}{n} \sum_{i=1}^n (Y_i^0 - \hat{Y}_i)^2\right]} &= \frac{1}{n} \sum_{i=1}^n E[(Y_i^0 - \hat{Y}_i)^2] \\ &= \frac{1}{n} \cdot n \underset{\alpha}{E[(Y^0 - \hat{m}(X))^2]} \end{aligned}$$

$$\Rightarrow \underset{\text{testing error}}{E\left[\frac{1}{n} \sum_{i=1}^n (Y_i^0 - \hat{Y}_i)^2\right]} = \underset{\text{training error}}{E\left[\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2\right]} + \frac{1}{n} \alpha^2 p = \alpha$$

proof

$$E[(Y_i - \hat{Y}_i)^2] = \text{Var}(Y_i) + \text{Var}(\hat{Y}_i) - 2\text{Cov}(Y_i, \hat{Y}_i) + (E(Y_i) - E(\hat{Y}_i))^2$$

$$E[(Y_i^0 - \hat{Y}_i)^2] = \text{Var}(Y_i^0) + \text{Var}(\hat{Y}_i) - 2\text{Cov}(Y_i^0, \hat{Y}_i) + (E(Y_i^0) - E(\hat{Y}_i))^2$$

$Y_i^0$  is independent to  $Y_i$ , but has the same distribution

$$\Rightarrow E(Y_i^0) = E(Y_i) \quad \text{Var}(Y_i^0) = \text{Var}(Y_i) \quad \text{Cov}(Y_i^0, Y_i) = 0$$

$$\begin{aligned} E[(Y_i^0 - \hat{Y}_i)^2] &= \text{Var}(Y_i) + \text{Var}(\hat{Y}_i) + (E(Y_i) - E(\hat{Y}_i))^2 \\ &= E[(Y_i - \hat{Y}_i)^2] + 2\text{Cov}(Y_i, \hat{Y}_i) \end{aligned}$$

Average over data points for new observations.

$$E\left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right] = E\left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right] + \frac{1}{n} \sum_{i=1}^n \text{Cov}(y_i, \hat{y}_i)$$

For a linear model,  $\text{Cov}(y_i, \hat{y}_i) = \alpha^2 H_{ii}$

$$\Rightarrow E\left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right] = \frac{E\left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right]}{\text{training error}} + \frac{\frac{1}{n} \cdot \alpha^2 \text{tr } H}{\frac{\sum_{i=1}^n h_{ii}}{n} = P}$$

$$\Rightarrow a = \text{MSE} + \text{penalty}$$

Mallow's Cp

$$C_p = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{2}{n} \hat{\alpha}^2 P$$

Smaller Cp are preferred.

R<sup>2</sup> / adjusted R<sup>2</sup>

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{T}}} = 1 - \frac{\text{MSE}}{SST/n} \quad \text{ignore } P$$

$$\text{adj } R^2 = 1 - \frac{\text{MSE} \cdot \frac{n}{n-P}}{SST/n} \rightarrow \text{MSE} \cdot \frac{n}{n-P} \approx \text{MSE} \left(1 + \frac{P}{n}\right) \\ = \text{MSE} + \text{MSE} \cdot \frac{P}{n} \\ \xrightarrow{n \rightarrow \infty} \text{MSE} + \alpha^2 \frac{P}{n} = a$$

- Maximize R<sup>2</sup> / adj R<sup>2</sup>

AIC & BIC

$$AIC = n \cdot \ln(\text{MSE}) + 2 \cdot P = n \ln \frac{SS_{\text{Res}}}{n} + 2P$$

$$BIC = n \cdot \ln(\text{MSE}) + \ln(n) \cdot P$$

- Minimize AIC & BIC