## 4.6 Potential Issues

Categorical predictors take on finitely many unordered values, sometimes called factors.

Ex) If $x_i$ = male/female has only two levels.

Need a convention to numerically code $x$, using a dummy variable, or indicator variable [or in ML, one-hot encoding]. E.g.,

$$x_i = \begin{cases} 1 & i\text{th person female} \\ 0 & \text{''} \quad \text{''} \quad \text{male} \end{cases}$$

Then

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & i\text{th person female} \\ \beta_0 + \varepsilon_i & \text{''} \quad \text{''} \quad \text{male} \end{cases}$$

where

$$\beta_0 = \text{Avg response for males}$$

$$\beta_0 + \beta_1 = \text{''} \quad \text{''} \quad \text{''} \quad \text{females}$$

$$\beta_1 = \text{Change in avg resp. for females vs. males}$$

Note we could have coded

$$x_i = \begin{cases} 1 & i^{th} \text{ person female} \\ -1 & \text{''} \quad \text{''} \quad \text{male} \end{cases}$$

$$\Rightarrow \quad Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & i^{th} \text{ person female} \\ \beta_0 - \beta_1 + \varepsilon_i & \text{''} \quad \text{''} \quad \text{male} \end{cases}$$

so   $\beta_0 =$ Average overall response

$\beta_1 =$ difference from the average response for m/f.

Important: predictions will not change with different encodings, but interpretations will!

More than 2 levels requires multiple dummy variables

Ex   $Y =$ life expectancy in country

$X \in \{$ Africa, OECD, other $\}$ (3 levels)

Could use

$$X_{i1} = \begin{cases} 1 & \text{ith country in OECD} \\ 0 & \text{,, ,, not ,, ,,} \end{cases}$$

$$x_{i2} = \begin{cases} 1 & i\text{th country in other} \\ 0 & \text{"} \quad \text{"} \quad \text{not "} \quad \text{"} \end{cases}$$

$$\Rightarrow y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & i\text{th country OECD} \\ \beta_0 + \beta_2 + \varepsilon_i & \text{"} \quad \text{"} \quad \text{other} \\ \beta_0 + \varepsilon_i & \text{"} \quad \text{"} \quad \text{Africa} \end{cases}$$

$\Rightarrow \quad \beta_0 = $ Avg life exp for countries in Africa

$\beta_1 = $ Change in l.e. for a country in OECD over Africa

$\beta_2 = \quad$ " $\quad$ " $\quad$ " $\quad$ " $\quad$ " in other $\quad$ " $\quad$ "

$\boxed{\text{Bad idea!}}$

$$x_i = \begin{cases} 0 & i\text{th country in Africa} \\ 1 & \text{"} \quad \text{"} \quad \text{"} \quad \text{other} \\ 2 & \text{"} \quad \text{"} \quad \text{"} \quad \text{OECD} \end{cases}$$

$y = \beta_0 + \beta_1 x + \varepsilon$

$= \begin{cases} \beta_0 + \varepsilon & \text{Afr} \\ \beta_0 + \beta_1 + \varepsilon & \text{other} \\ \beta_0 + 2\beta_1 + \varepsilon & \text{OECD} \end{cases}$

# Beyond additivity & linearity

Additivity: Effect of a predictor on $Y$ is independent of values of other predictors. To overcome, use interactions:

Ex

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \varepsilon = \begin{cases} \beta_0 + \beta_1 x_1 + \beta_2 + \beta_3 x_1 + \varepsilon & m \\ \quad = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x_1 + \varepsilon \\ \beta_0 + \beta_1 x_1 + \varepsilon & f \end{cases}$$

$Y$ = salary    $x_1$ = yrs since degree    $x_2 = \begin{cases} 1 & \text{male} \\ 0 & \text{female} \end{cases}$

$\beta_1 = $ Avg raise/yr for females

$(\beta_1 + \beta_3) = $ "   "   "   "  males

$\beta_0 = $ Starting salary for females

$(\beta_0 + \beta_2) = $ "   "   "   males

**Linearity:** Y depends linearly on X. Polynomial regression is an easy way to overcome:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$$

or

$$\log Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\Rightarrow \quad Y = e^{\beta_0} e^{\beta_1 X} e^{\varepsilon}$$

1) Y depends exponentially on X

2) errors are multiplicative

# Degrees of freedom

$X \quad n \times (p+1)$

The model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \qquad \text{with } p \text{ features}$$

uses $(p+1)$ degrees of freedom. If $H$ is the hat matrix, then

$$\text{tr } H = \sum_{i=1}^{n} H_{ii} = \text{tr}\left( X (X^\top X)^{-1} X^\top \right)$$

$$= \text{tr}\left( X^\top X (X^\top X)^{-1} \right) = \text{tr}\left( I_{p+1} \right) = p+1$$

For any linear predictor $\hat{Y} = My$, define $\text{tr}(M) =$ degrees of freedom