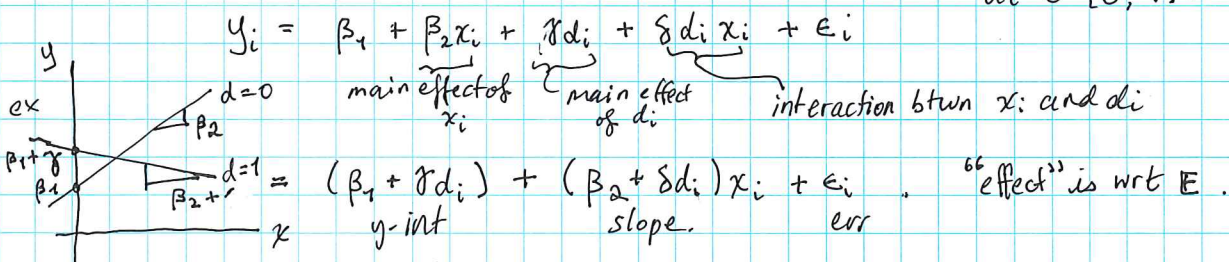
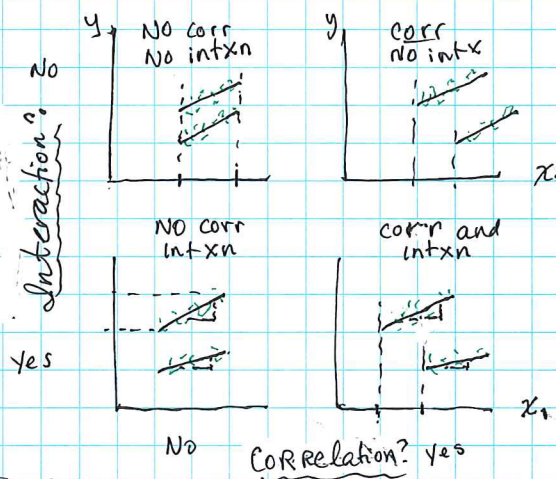


- Recall: Interaction dynamics (ex. btwn x_i, d_i), $x_i \in \mathbb{R}$, $d_i \in \{0, 1\}$.



ex. Correlation vs. interaction



- Interaction \rightarrow change slope
- Correlation \rightarrow X-shift

Recall regression model.

$$Y = X\beta + \epsilon, \text{ with } \epsilon \sim N(0, \sigma^2 I)$$

$n \times 1$ $n \times p$ $p \times 1$ $n \times 1$ $\epsilon \sim$ norm distr

$$\hat{\beta} = (X^T X)^{-1} X^T Y \text{ and}$$

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} \quad (\equiv \text{Cov of } \hat{\beta})$$

$$\text{From this, } \widehat{\text{Var}}(\hat{\beta}_j) = \hat{\sigma}^2 [(X^T X)^{-1}]_{jj}$$

where R_j^2 is the R^2 from regressing x_j on the other variables.

- We can rewrite this as

$$\widehat{\text{Var}}(\hat{\beta}_j) = \frac{\hat{\sigma}^2}{(n-1)\widehat{\text{Var}}(x_j)} \cdot \frac{1}{1-R_j^2}$$

Behavior: If ... $n \uparrow$, then $\widehat{\text{Var}}(\hat{\beta}_j) \downarrow$
 $\sigma^2 \downarrow$, then $\widehat{\text{Var}}(\hat{\beta}_j) \downarrow$
 $\text{Var}(x_j) \uparrow$, " " " \downarrow
 $R_j^2 \uparrow$, " " " \uparrow

take home message: it helps to have a good spread in X variables

So $\widehat{\text{se}}(\hat{\beta}_j) = \hat{\sigma} \sqrt{[(X^T X)^{-1}]_{jj}}$ (see R code output)

Then, we have $\frac{\hat{\beta}_j - \beta_j}{\widehat{\text{se}}(\hat{\beta}_j)} \sim t_{n-p} \rightarrow P(-t_{1-\frac{\alpha}{2}, n-p} < \frac{\hat{\beta}_j - \beta_j}{\widehat{\text{se}}(\hat{\beta}_j)} < t_{1-\frac{\alpha}{2}, n-p}) = 1-\alpha$

Rewrite this as...

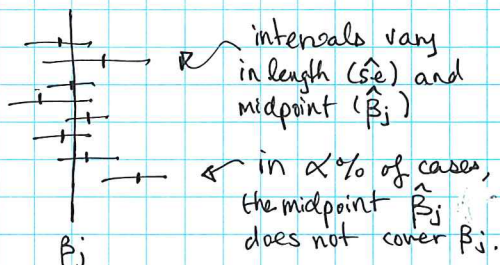
$$P(-\widehat{\text{se}} t_{1-\frac{\alpha}{2}, n-p} < \hat{\beta}_j - \beta_j < \widehat{\text{se}} t_{1-\frac{\alpha}{2}, n-p}) = 1-\alpha$$

$$P(\hat{\beta}_j - \widehat{\text{se}} t_{1-\frac{\alpha}{2}, n-p} < \beta_j < \hat{\beta}_j + \widehat{\text{se}} t_{1-\frac{\alpha}{2}, n-p}) = 1-\alpha$$

center data

Thus, CI for β_j : $(1-\alpha) \cdot 100\%$ CI for β_j : $\hat{\beta}_j \pm \widehat{\text{se}}(\hat{\beta}_j) t_{1-\frac{\alpha}{2}, n-p}$ is

ex.

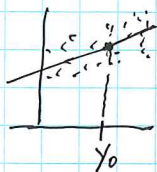


point estimate. estimated std dev. quantile

a two-sided confidence interval which covers the true β_j with prob $1-\alpha$; here $t_{1-\frac{\alpha}{2}, n-p}$ denotes the $1-\frac{\alpha}{2}$ quantile of a t_{n-p} distribution.

- Getting the CI for Y_0 . Now consider a new point $X_0^T = (X_{01}, \dots, X_{0p})$. We are interested in $E(Y_0)$ or Y_0 . What is the difference?

$E(Y_0)$ will be the value on the regression line. Y_0 is a point and has more variability (since we have E_0); won't necessarily be on regression line.



$$Y_0 = X_0^T \beta + \varepsilon \quad \text{and} \quad E(Y_0) = X_0^T \beta$$

$$\hat{Y}_0 = X_0^T \hat{\beta} + \varepsilon \quad = E(\hat{Y}_0)$$

$$Y = \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$

$$\text{lm}(X_j \sim X_1 + \dots + X_{j-1} + X_{j+1} + \dots + X_p) \quad (\text{R code})$$

R_j^2 represents the proportion of variance in X_j that can be explained by a linear regression on the other predictors

- Note that both quantities have expectation $E(Y_0) = X_0^T \beta$, so as a point estimate/best guess, we use $X_0^T \hat{\beta}$.
- Since $E(X_0^T \hat{\beta}) = X_0^T \beta$, this is an unbiased estimate.

- Note: $\text{Var}(X_0^T \hat{\beta}) = X_0^T \text{Var}(\hat{\beta}) X_0 = \sigma^2 X_0^T (X^T X)^{-1} X_0$.

- Hence, $\frac{X_0^T \hat{\beta} - E(Y_0)}{\sigma \sqrt{X_0^T (X^T X)^{-1} X_0}} \sim t_{n-p}$ So the $(1-\alpha)$, 100% CI for $E(Y_0)$ is:

divide by Var to normalize $\frac{X_0^T \hat{\beta} - Y_0}{\text{Var}(X_0^T \hat{\beta} - Y_0)} \sim t_{n-p}$ $X_0^T \hat{\beta} \pm \sqrt{\hat{\sigma}^2 X_0^T (X^T X)^{-1} X_0} t_{1-\frac{\alpha}{2}, n-p}$

- Now for Y_0 : $\text{Var}(X_0^T \hat{\beta} - Y_0) = \text{Var}(X_0^T \hat{\beta} - (X_0^T \beta + \varepsilon)) = \text{Var}(X_0^T \hat{\beta} - \varepsilon) =$

accounts for widening of CI as X data spreads further from center

$$= \text{Var}(X_0^T \hat{\beta}) + \text{Var}(\varepsilon) - 2 \text{Cov}(X_0^T \hat{\beta}, \varepsilon)$$

$$= [\sigma^2 X_0^T (X^T X)^{-1} X_0] + [\sigma^2] - 0$$

$$= \sigma^2 (1 + X_0^T (X^T X)^{-1} X_0) = \text{Var}(X_0^T \hat{\beta} - Y_0)$$

So, ① $(1-\alpha)$ CI for Y_0 is $X_0^T \hat{\beta} \pm \sigma \sqrt{1 + X_0^T (X^T X)^{-1} X_0} t_{1-\frac{\alpha}{2}, n-p}$

② $\frac{X_0^T \hat{\beta} - Y_0}{\sigma \sqrt{1 + X_0^T (X^T X)^{-1} X_0}}$ converges to zero as sample size increases

(aka "prediction interval")

CI of Y_0 is strictly larger than that of $E(Y_0)$

* Bias-Variance Tradeoff (ISLR 2.2.1-2)

- Let $Y_i = f(x_i) + \varepsilon_i$. Consider an estimator $\hat{f}(x_i)$. Then the mean squared error (MSE) is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

(More precisely, the training MSE as the points x_i are the data used to train the estimator.

How does \hat{f} perform on new data? (ex. tomorrow's stock market values).

- (now) from training data $(x_1, y_1), \dots, (x_n, y_n)$ used to obtain estimator \hat{f} , we would like for $y_i \approx f(x_i)$, i.e. small training MSE. Now, is $\hat{f}(x_0) \approx y_0$ for a new pair (x_0, y_0) ? Test MSE at $x_0 = (y_0 - \hat{f}(x_0))^2$.

(from script, $\text{MSE}(x) = E[(\hat{f}(x) - f(x))^2] = (E(\hat{f}(x)) - f(x))^2 + \text{Var}(\hat{f}(x))$).