

SLR, MR, Dummy variable, Quadratic term, Model selection, Logistic regression.

* SLR $y = \alpha + \beta x + \epsilon$, $\epsilon \stackrel{iid}{\sim} N(0, \sigma^2)$

- estimates: $a = \bar{y} - b\bar{x}$, $b = r \frac{S_y}{S_x}$, $s = \sqrt{MSE}$
- interpretation: y -int (a), slope (b)
- prediction: $\hat{y} = a + bx \rightarrow$ residual: $y_i - \hat{y}_i$
- correlation: $-1 \leq r \leq 1$
- coefficient of determination: $R^2 = (r)^2 = SSR/SS_T$
- inference: C.I., t-test, ANOVA [For SLR, t-test \Leftrightarrow ANOVA]

$$1 \quad b \pm t_{n-2, \frac{\alpha}{2}} \cdot S.E.b$$

$$2 \quad t_{obs} = \frac{b}{S.E.b}$$

3	Source	df	SS	MS	F	p-val
	Regression	1	SSR	MSR	MSR/MSE	$P(X > F)$
	Error	$n-2$	SSE	MSE		where $X \sim F_{1, n-2}$
	Total	$n-1$	SST			

- Understanding of SS (Page 74)
- CI v.s. PI for response interpretations $\begin{cases} \text{average response (CI)} \\ \text{new individual (PI)} \end{cases}$
- Diagnostics via residual analysis [Check $\epsilon \stackrel{iid}{\sim} N(0, \sigma^2)$]
 - Histogram of residual
 - NPP
 - Residual vs fit
 - Residual vs Order/ x (Page 79)

* Multiple regression $y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$, $\epsilon \stackrel{iid}{\sim} N(0, \sigma^2)$

y Quant x can be categorical

- estimates: $a \rightarrow \alpha$, $b_i \rightarrow \beta_i$, $\sqrt{MSE} \rightarrow \sigma$

- interpretation: similar to SLR

- prediction

- inference: ANOVA \rightarrow any good predictor?

t-test for $\beta_i \rightarrow$ Is X_i good pred AFTER all other pred variables are accounted for?

Source	df	...
Reg	p	
Error	$n-p-1$...
Total	$n-1$	

- oversaturated model
- multicollinearity
- R^2 vs. R_{adj}^2 (pros & cons)

} both bad.

- Order of testing.

* Categorical variables in regression

- Dummy variable $X = \{0, 1\}$
- Baseline group/model
- Interpretation of coefficients in the full model (α, β_i)
- Interpretation of estimates (a, b_i)
- Prediction
- Interaction (between quantitative terms and categorical terms)
- Type of lines the model allows (parallel? Line or curve?)
- How to test intercept/slope for ... is sig diff than ...

* Quadratic regression $y = \alpha + \beta_1 X + \beta_2 X^2 + \varepsilon$

- When to use it?
- Interpretation, only for β_2 [sig diff from zero? Sign?] $\cup \cap \swarrow$
- Prediction
- Comparison to SLR.

* Selecting the best model

- Backwards elimination, forward selection, best subsets regression
- Goal: Simple model that does a decent job of prediction.
- R^2_{adj} , C_p , AIC/BIC
- Data cleaning: - Unusual observations $\left\{ \begin{array}{l} \text{"R"} - \text{large standardized residuals} \\ \text{"X"} - \text{large leverage} \end{array} \right.$
- Outlier, influential outlier

* Logistic regression $P(Y=1) = \frac{e^{\alpha + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\alpha + \beta_1 X_1 + \dots + \beta_p X_p}}$ Y : Binary We predict a probability

- Interpretation: a - Don't care, b_i - (Sign and) log odds ratio interpretation
- Properties of $f(x) = \frac{e^x}{1 + e^x}$ and $\text{logit}(x) = \log\left(\frac{x}{1-x}\right)$ $[\ln(\frac{x}{1-x})]$
- Equivalent formulation: $-P(Y=1) = \frac{1}{1 + e^{-(\alpha + \beta_1 X_1 + \dots + \beta_p X_p)}}$
 $-\text{logit}(P(Y=1)) = \alpha + \beta_1 X_1 + \dots + \beta_p X_p$
- Prediction (\hat{p}) • Solve for x when $\hat{p} = 0.5$.