


# Towards Regression

## Review LS 7 Comparing Means Hypothesis Tests

- Perform hypothesis testing and estimation to generalize
- Important to know and sample from target popl'n
  - Often the available popl'n  $\neq$  target popl'n
- Null hypotheses are crisp assertion, e.g.
  - Null hypothesis;  $H(0): \theta = a$  ( $a$  is a real constant) vs
  - Select apriori alternative one of  $H(a): \theta \neq a$  or  $H(a): \theta < a$  or  $H(a): \theta > a$
  - Use  $P\_value$  determined from computed value of  $Z^*$  or  $t^*$  and compared against distribution of transformed sampling distribution under  $H(0)$
- Other concerns
  - Data types: metric or ordinal
  - Test types: unpaired (two groups), paired (diff)
  - Variance: pooled vs not pooled
  - Assumption: normality, iid, skewness, metric  $\rightarrow$  parametric test, otherwise non-parametric

## Preview of LS 8 Regression

- Start the last of the course's major topics
- Set up the approach and machinery for creating a formal structure to model a functional relationship between elements or variables
- Starting out focusing on simple linear regression (SLR) → one predictor, one response variable and move on to >1 predictor and one response variable multiple linear regression
- Use an estimation method called Ordinary Least Squares (OLS) which is built upon assumptions of MOM estimators BP
- Recall for the  the following result  $E[(Y-g(X))^2 | X]$  called mean square error was the minimized when  $g(x) = E[Y | X=x]$
- In the world of statistics we want  $BLP = \sum (Y - E[Y | X])^2 = \sum (\epsilon^2)$  which under the assumptions  $\sum(\epsilon)=0$ ,  $cov(\epsilon, X) = 0$  and  $Y$  (and  $\epsilon$ )  $\sim$  IID →  $E[Y | X] = \beta(0) + \beta(1)*X(i)$
- The model correct on average; error term is unpredictable by the X's, referred to as exogeneity, otherwise OLS attributes (incorrectly) some of the variance the error explains to the predictors.

# Moment Conditions

- Moment Condition Assumptions  $\sum(\epsilon)=0$ ,  $\text{cov}(\epsilon, X) = 0$
- Let  $Y(i) = \beta(0) + \beta(1)*X(i) + \epsilon(i)$   
$$\bar{Y} = 1/n * \sum(\beta(0) + \beta(1)*X(i) + \epsilon(i))$$
$$= \beta(0) + \beta(1)*\sum(X(i))/n + \sum(\epsilon(i))$$
- $\text{Cov}(X, Y) = \text{Cov}(X, \beta(0) + \beta(1)*X + \epsilon)$   
$$= \text{Cov}(X, \beta(0)) + \text{Cov}(X, \beta(1)*X) + \text{Cov}(X, \epsilon)$$
$$\text{Cov}(X, Y) = \beta(1) * \text{Var}(X)$$

## OLS Normal Equations

- $y(i) = \beta(0) + \beta(1)*x(i) + \epsilon(i)$  or  $\epsilon(i) = y(i) - \beta(0) - \beta(1)*x(i)$
- Let  $Q = \sum \epsilon(i)^2$ ; want to find  $\beta(0), \beta(1)$  which minimize  $Q$  or  $\min \sum (y(i) - \beta(0) - \beta(1)*x(i))^2$
- For  $\beta(0)$ :  $\partial Q / \partial \beta(0) = 2 * \sum (y(i) - \beta(0) - \beta(1)*x(i)) * (-1)$
- $= -2n * (\bar{y} - \beta(0) - \beta(1)*\bar{x})$
- Min:  $0 = \bar{y} - \hat{\beta}(0) - \hat{\beta}(1)*\bar{x} \rightarrow \hat{\beta}(0) = \bar{y} - \hat{\beta}(1)*\bar{x}$  (note hat)
- For  $\beta(1)$ :  $\partial Q / \partial \beta(1) = -2 * \sum (y(i) - \beta(0) - \beta(1)*x(i)) * x(i)$
- Min:  $0 = -2 \sum x(i)(y(i) - \hat{\beta}(0) - \hat{\beta}(1)*x(i))$
- $= \sum x(i)y(i) - \hat{\beta}(0) \sum x(i) - \hat{\beta}(1) * \sum x(i)^2 \rightarrow$
- $\sum x(i)y(i) - \hat{\beta}(1) * \sum x(i)^2 - (\bar{y} - \hat{\beta}(1)*\bar{x}) * \sum x(i)$
- Then substituting in  $\hat{\beta}(0)$  & rearranging (left as exercise for you)
- $\rightarrow \hat{\beta}(1) = (\sum (x(i) - \bar{x})(y(i) - \bar{y})) / \sum (x(i) - \bar{x})^2$

