Regression Notes from slides

$$\Sigma (y_i - x_i) = \Sigma y_i - \Sigma x_i$$
$$\Sigma (x_i - a) = \Sigma x_i - na$$
$$\Sigma (x_i - \bar{x})^2 = \Sigma x_i^2 - n\bar{x}^2$$
$$\bar{x} = \Sigma x_i / n$$
$$\sigma_{\bar{x}}^2 = \left(\frac{\sigma}{\sqrt{n}}\right)^2 \qquad \bar{x} \sim (\mu, \sigma/\sqrt{n})$$

t-test Statistic $\qquad t = \dfrac{\bar{x} - \mu}{S/\sqrt{n}} \qquad n-1 \ d.f.$

CI : Sample mean $\pm$ t-crit $\frac{S}{\sqrt{n}}$

Covariance : $E\left[\left(x - E(x)\right)\left(y - E(y)\right)\right]$

Sample Cov : $\dfrac{1}{N-1} \Sigma (x_i - \bar{x})(y_i - \bar{y})$

Correlation $= \dfrac{Cov(x,y)}{Sd(x) Sd(y)}$

$r = \dfrac{\Sigma (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\Sigma (x_i - \bar{x})^2}\sqrt{\Sigma (y_i - \bar{y})^2}}$   Measures direction but not slope

Does NOT imply causation

$B_1$ is the slope of the regression line and indicates the change in the mean of the distribution of $y$ for one unit increase in $x$

$B_0$ is the intercept of the regression line. If it is sensible to think of a value of $x=0$ for a particular application, then $B_0$ gives the mean distribution of $y$ at $x=0$

$$y_i = B_0 + B_1 x_i + \varepsilon_i \qquad i = 1 \ldots n$$

## Assumptions:

1. $x_i$ is the $i^{th}$ value of the predictor variable, which is a known constant for all $i$.
2. The observations $y_i$ (or $\varepsilon_i$) are independent.
3. At any given $x_i$, $y_i$ (or $\varepsilon_i$) is normally distributed.
4. The observations $y_i$ (or $\varepsilon_i$) have constant standard deviation.
5. The means of $y_i$ can be found by a straight line given as:

$$E[y_i] = \beta_0 + \beta_1 x_i$$

where $\beta_0$ and $\beta_1$ are unknown parameters, such that:

$\beta_1$ is the slope of the regression line and indicates the change in the mean distribution of $y$ per one unit increase in $x$.

$\beta_0$ is the intercept of the regression model. If it is sensible to think of a value of $x=0$ for a particular application, then $\beta_0$ gives the mean of the distribution of $y$ at $x=0$. So it is not always possible to have a physical explanation of this parameter.

## Least Squares:

Diff wrt to $\beta_0$ and $\beta_1$ $\quad \sum(y_i - \beta_0 - \beta_1 x_i)^2$

$$b_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$S_{xx} = \sum(x_i - \bar{x})^2 = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

$$S_{yy} = \sum(y_i - \bar{y})^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

$$S_{xy} = \sum(x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n}$$

## Estimating $\sigma^2$

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$y_i - \hat{y} = e_i$$

$$SSE = \sum(y_i - \hat{y})^2 = \sum(y_i - \beta_0 - \beta_1 x_i)^2 = \sum e_i^2 \quad \text{AKA RSE}$$

$$\hat{\sigma}^2 = MSE = \frac{SSE}{N-2} = \frac{\sum(y_i - \hat{y})^2}{n-2}$$

Regression Notes From Side

$$E(MSE) = \sigma^2$$

Inference about the $\beta_i$:

$E(b_i) = \beta_i$

$Var(b_i) = \dfrac{\sigma^2}{\Sigma(x-\bar{x})^2} = \dfrac{\sigma^2}{S_{xx}}$

$b_i \sim N\left(\beta_i, \dfrac{\sigma^2}{S_{xx}}\right)$

$\hat{Var}(b_i) = S_{b_i}^2 = \dfrac{MSE}{S_{xx}}$   $Se(b_i) = \sqrt{\dfrac{MSE}{S_{xx}}}$

T-test $t_{calc} = \dfrac{Estimator - value\ from\ H_0}{Se(estimator)}$   $(n-2)$ d.F.

CI : $(estimator) \pm t_{critic}\ Se(Estimator)$

Important, want to test $\beta_i = 0$, if it equals zero, no relationship exists

Inference About intercept $\beta_0$

$\beta_0 \Rightarrow b_0 = \bar{y} - \bar{x}b_1$

$E(b_0) = \beta_0$

$Var(b_0) = \sigma^2\left[\dfrac{1}{n} + \dfrac{\bar{x}^2}{S_{xx}}\right]$

$b_0 \sim N\left(\beta_0, \sigma^2\left[\dfrac{1}{n} + \dfrac{\bar{x}^2}{S_{xx}}\right]\right)$

$Se(b_0) = \sqrt{\sigma\left[\dfrac{1}{n} + \dfrac{\bar{x}^2}{S_{xx}}\right]}$

Sampling Distribution of $\dfrac{b_0 - \beta_0}{Se(b_0)}$ is t-dist with $n-2$ d.F.

4

## Inference for E[y] at X = Xsome val

Make inference about mean distribution of y (e.E[y]) at x=Some x.
- Want to construct point estimate and CI for E(y|x=x).

$$E[y|x'] = \beta_0 + \beta_1 x'$$
$$\hat{y}' = b_0 + b_1 x'$$

$$Var[\hat{y}'] = \sigma^2 \left[ \frac{1}{n} + \frac{(x'-\bar{x})^2}{Sxx} \right] \leftarrow \quad se(\hat{y}') \text{ is } \sigma \text{ of error}$$

$$\hat{y}' \sim N\left( (\beta_0 + \beta_1 x'), \sigma^2 \left[ \frac{1}{n} + \frac{(x'-\bar{x})^2}{Sxx} \right] \right)$$

### Hypothesis testing for $\hat{y}$.

$$H_0 : \hat{y}' = y' \quad vs \quad H_1 : \hat{y}' \neq y'$$

Test statistic is $\dfrac{\hat{y}' - y'}{se(\hat{y}')}$

CI $\quad \hat{y}' + se(\hat{y}') t_{critical}$

## Prediction of a new observation

In first case we refer to the mean distribution of y for a particular value of X, whereas in PI, we predict an individual outcome drawn from the distribution of y for a given x

$$\hat{y}_{new} = b_0 + b_1 x' \quad E[\hat{y}_{new}] = y'_{new}$$
$$Var(\hat{y}_{new}) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x'-\bar{x})^2}{Sxx} \right]$$

Variance has 2 components:
1. Variance of the sampling distribution of fitted value $\hat{y}$.
2. Variance of distribution of y at some X

Prediction interval $\quad \hat{y}_{new} \pm t_{critical} \sqrt{MSE \left[ 1 + \frac{1}{n} + \frac{(x'-\bar{x})^2}{Sxx} \right]}$

Regression Notes Slides

$$y_i = \beta_0 + \beta_1 x_1 + \varepsilon_i \quad \text{or} \quad E(Y_i | x_i) = \beta_0 + \beta_1 x_i$$

## ANOVA

SSTO = $\sum (y_i - \bar{y})^2$ Total uncert $\sum y_i^2$

IF SSTO all respond are equal

SSE $\sum (y_i - \hat{y}_i)^2$ variability around fitted to

SSTO = SSE + SSR

SSR = $\sum (\hat{y}_i - \bar{y})$ variably assoc with the regression line

$$y_i - \bar{y} = \hat{y}_i - \bar{y} + y_i - \hat{y}_i$$

SSTO    SSE      SSR

$$\sum y = \sum \hat{y}_i$$

| Source of variation | SS | DF | MS | F |
|---|---|---|---|---|
| Regression | SSR $\sum(\hat{y}-\bar{y})$ | 1 | MSR $SSR/1$ | $F_{calc} = \dfrac{MSR}{MSE}$ |
| Error | SSE $\sum(y-\hat{y})^2$ | n-2 | MSE $SSE/n-2$ | |
| Total | $\sum(y_i-\bar{y})^2$ | n-1 | | |

$F_{calc}$ can be used to test $H_0 : \beta_1 = 0$ vs $H_1 : \beta_1 \neq 0$ follow $F_{1, n-2}$ distr.

If $F_{calc} \leq F_{(1-\alpha, 1, n-2)}$ do not reject $H_0$

$F = (t)^2$

$$\sum \left( y_i - b_0 - b_1 x_i \right)^2$$

$\dfrac{db_0}{\boxed{b_0}}$

$$-2 \sum \left( y_i - b_0 - b_1 x_i \right) = 0$$

$$\sum \left( y_i - b_0 - b_1 x_i \right) = 0$$

$$\sum y_i - \sum b_0 - b_1 \sum x_i = 0 \qquad \dfrac{\sum x_i}{n} = \bar{x}$$

$$n b_0 = \sum y_i - b_1 \sum x_i$$

$$b_0 = \dfrac{\sum y_i}{n} - b_1 \dfrac{\sum x_i}{n}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$\dfrac{db_1}{\boxed{b_1}} =$

$$-2 \sum x_i \left( y_i - b_0 - b_1 x_i \right)$$

$$-2 \sum \left( y_i x_i - b_0 x_i - b_1 x_i^2 \right) = 0 \qquad b_0 =$$

$$\sum y_i x_i - b_0 \sum x_i - b_1 \sum x_i^2 = 0$$

$$b_1 \sum x_i^2 = \sum y_i x_i - b_0 \sum x_i$$

$$b_1 \sum x_i^2 = \sum y_i x_i - \left[ \dfrac{\sum y_i}{n} - b_1 \dfrac{\sum x_i}{n} \right] \sum x_i$$

$$b_1 \sum x_i^2 = \sum y_i x_i - \dfrac{\sum y_i x_i}{n} + \dfrac{b_1 \sum x_i^2}{n}$$

$$b_1 \sum x_i^2 - b_1 \dfrac{\sum x_i^2}{n} = \sum y_i x_i - \dfrac{\sum y_i x_i}{n}$$

$$b_1 \left[ \sum x_i^2 - \dfrac{\sum x_i^2}{n} \right] = \sum y_i x_i - \dfrac{\sum y_i x_i}{n}$$

$$b_1 = \dfrac{\sum y_i x_i - \dfrac{\sum y_i x_i}{n}}{\sum x_i^2 - \dfrac{\sum x_i^2}{n}}$$

$\bullet\, y_i$

$\bullet\, \hat{y}_i$

$\bullet\, \bar{y}$

$$\Sigma(y_i - \bar{y})^2 \quad = \quad \Sigma(y_i - \hat{y})^2 + \Sigma(\hat{y} - \bar{y})^2$$

$$\Sigma[(y_i - \hat{y}_i)(\hat{y}_i - \bar{y})]^2$$

$$= \Sigma(y_i - \hat{y})^2 + \Sigma(\hat{y}_i - \bar{y})^2 + 2\Sigma(y_i - \hat{y}_i)(\hat{y} - \bar{y})$$

$$2\Sigma(y_i - \hat{y}_i)(\hat{y} - \bar{y})$$
$$\Sigma[\hat{y}(y_i - \hat{y}_i) - \bar{y}(y_i - \hat{y}_i)]$$

$$\Sigma[(\bar{y} + b_1(x_i - \bar{x}))(y_i - \hat{y}_i) - \bar{y}(y_i - \hat{y}_i)]$$
$$\Sigma \bar{y}(y_i - \hat{y}_i) + \Sigma b_1(x_i - \bar{x})(y_i - \hat{y})\; -\; \Sigma \bar{y}\, y_i - \hat{y}_i$$
$$cuy$$

Right column:
$$\hat{y}_i = b_0 + b_1 x_i$$
$$b_0 = \bar{y} - b_1 \bar{x}_i$$
$$\hat{y}_i = \bar{y} - b_1 \bar{x}_i + b_1 x_i$$
$$\hat{y}_i = \bar{y} + b_1(x_i - \bar{x})$$

$$\Sigma b_1(x_i - \bar{x})(y_i - \hat{y}_i)$$
$$b_1 \Sigma x_i(y_i - \hat{y}_i) - \Sigma \bar{x}(y_i - \hat{y})$$
$$= 0$$
$$-\Sigma \bar{x} y_i + \Sigma \bar{x} \hat{y}_i \quad \Sigma \hat{y}_i = \hat{y}_i$$
$$= 0$$

Right column:
$$\Sigma y_i = \Sigma \bar{y} + b_1(x_i - \bar{x})$$
$$\Sigma \bar{y} + \Sigma b_1(x - \bar{x})$$
$$= 0$$
$$\Sigma \hat{y} = \Sigma \bar{y}$$
$$= n\bar{y} =$$
$$\Sigma \hat{y}_i = \Sigma y_i$$

$$b_1 \Sigma x_i(y_i - \hat{y}_i)$$
$$b_1 \Sigma x_i(y_i - [\bar{y} + b_1(x_i - \bar{x})])$$
$$b_1 \Sigma x_i(y_i - \bar{y} - b_1(x_i - \bar{x}))$$
$$b_1 \Sigma x_i(y_i - \bar{y}) - b_1 \Sigma x_i(x_i - \bar{x})$$

$$\Sigma(y_i - \bar{y})(x_i - \bar{x}) - b_1 \Sigma(x_i - \bar{x})^2$$

$$b_1 = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\Sigma(x_i - \bar{x})^2}$$

$$\Sigma(y_i - \bar{y})(x_i - \bar{x}) - \Sigma(y_i - \bar{y})(x_i - \bar{x}) = 0$$