# Coefficient of Determination $(R^2)$ $r^2$

$R$ = Correlation Coefficient = $\sum_{i=1}^{n} \left(\frac{X_i - \bar{X}}{S_x}\right) \times \left(\frac{Y_i - \bar{Y}}{S_y}\right)$

$= \sum_{i=1}^{n} Z_{X_i} \times Z_{Y_i}$

43.8% of var of y is explained by X.

$X, Y$
we want to predict Y from X.

$Z_x, Z_y$

$(\bar{X}, \bar{Y})$

Residual is the difference between the observed value $Y_i$ and $\hat{Y}_i = \hat{b}_0 + \hat{b}_1 \cdot X_i$

$(Y_i - \hat{Y}_i)$

$R = .662$ $R^2 = .438$

**Predicting bikeshare rental from temperature**

Var of y

This is 60% the var of y

$Y_i \approx 7.5k$

$\hat{Y}_i \approx 4.5k$

$(7.5k - 4.5k) = 3.0k$ residual

Rentals (y-axis: 2000, 6000)

Temp (Celcius) (x-axis: 10, 15, 20, 25, 30)

The line of best fit $\left(\hat{b}_1 = \frac{S_y}{S_x} \cdot R, \ \hat{b}_0 = \bar{Y} - \hat{b}_1 \bar{X}\right)$ minimizes the sum of squared residuals $\left(\sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2\right) = \sum_{i=1}^{n} (Y_i - (\hat{b}_0 + \hat{b}_1 X_i))$

This is our best guess informed by $X_i$.

$R^2 = R^2$ ← Square of the correlation coefficient

$= 1 - \frac{RSS}{TSS}$

RSS ← Residual sum of squares

TSS ← Total sum of squares (for y)

This tells us what percentage of variability is still unexplained.

$\sum_{i=1}^{n} (Y_i - \bar{Y})^2$

This is our best guess of $Y_i$ for $X_i$ unknown

i.e. not using regression

If RSS = 0  $R^2 = 100\%$

If RSS = TSS  $R^2 = 0\%$

$\hat{y} = \hat{b}_0 + \hat{b}x$

Var of y {

$X_i$ explain 100% of the variability of $y_i$.

Var of y {

$\hat{y} = \hat{b_0} + \hat{b}x$

Var of y for $x = x_i$

$x_i$

$R^2$ tells us what <u>percentage</u> of the variability in $y$ is explained by $X$.