

To minimize $\sum_{i=1}^n e_i^2$ wrt (β_0, β_1) , solve the normal equations

RSS
residual sum
of squares

$$\frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_0} = 0, \quad \frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_1} = 0,$$

$$E[Y] = \beta_0 + \beta_1 X$$

resulting in the least squares estimates

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = r S_y / S_x,$$

where r is **sample correlation**, and S_x and S_y are **standard deviations** of x and y samples, respectively.

Recall that:

$$S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2,$$

$$r = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{S_x S_y}$$

The fitted regression line

Fitted regression line: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$. Plugging-in $\hat{\beta}_0$ and $\hat{\beta}_1$,

$$\begin{aligned}\hat{Y} &= (\bar{Y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x = \bar{Y} + \hat{\beta}_1 (x - \bar{x}) \\ &= \bar{Y} + r \frac{s_y}{s_x} (x - \bar{x}) \Rightarrow \frac{(\hat{Y} - \bar{Y})}{s_y} = r \frac{(x - \bar{x})}{s_x}\end{aligned}$$

$$\frac{(x - \bar{x})}{s_x}$$

"slope of x "

implying that

\Rightarrow $\hat{Y} \approx \bar{Y}$

$$\frac{\hat{Y} - \bar{Y}}{s_y} = r \frac{(x - \bar{x})}{s_x}$$

"2-some of \hat{Y} "

$$|r| \leq 1$$

- If x is 1 SD away from its mean \bar{x} , \hat{Y} is r SD away from its mean \bar{Y} . Since $|r| \leq 1$, this means \hat{Y} is **closer** to \bar{Y} (in units of SD) than x is to \bar{x} — **regression toward mean**.
- The fitted line passes through the points (\bar{x}, \bar{y}) .
- The sign of slope $\hat{\beta}_1$ is same as the sign of r .
- The sum of residuals, $\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - \hat{y}_i) = 0$
- The average of fitted values, $(1/n) \sum_{i=1}^n \hat{Y}_i = \left(\frac{1}{n}\right) \sum_{i=1}^n y_i = \bar{Y}$

Ex: Let's get the fitted line for the house price data and add it to the scatterplot.

```
x <- house$size      - prediction  
y <- house$price     - rmse
```

```
# Get the fitted regression line  
> (house.reg <- lm (y ~ x))  
Call:  
lm(formula = y ~ x)
```

Coefficients:

| (Intercept) | x |
|-------------------|-----------------|
| 5.432 | 56.083 ✓ |
| $\hat{\beta}_0$ ✗ | $\hat{\beta}_1$ |

```
# Does R do what we expect it to do?
```

```
> c(mean(x), sd(x), mean(y), sd(y), cor(x,y))
```

```
[1] 1.8829655 0.6316624 111.0344483 40.4431900  
0.8759374
```

```
>  
> cor(x,y)*sd(y)/sd(x)  
[1] 56.08328
```

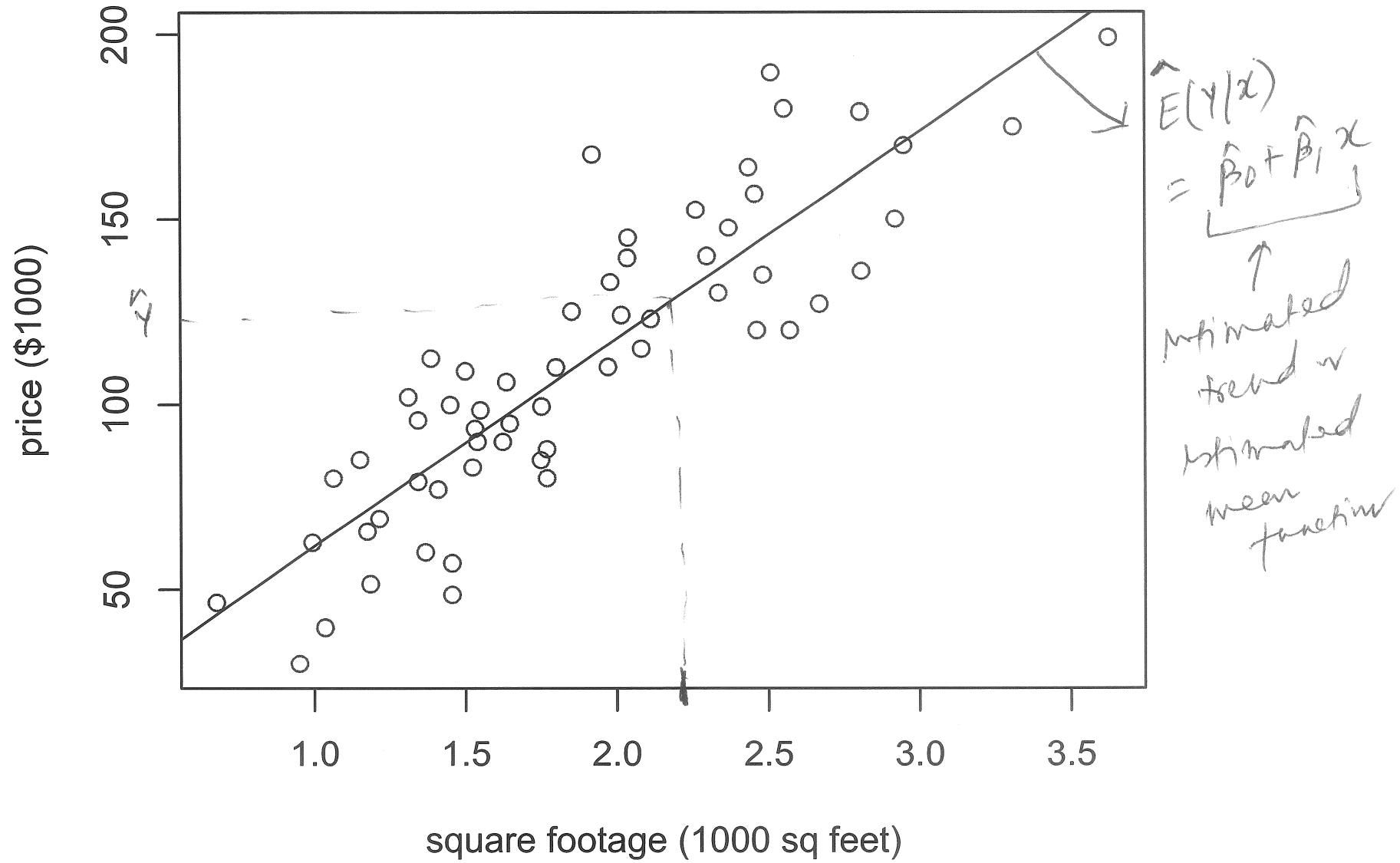
```
>  
> mean(y)-(cor(x,y)*sd(y)/sd(x))*mean(x)  
[1] 5.431568  
>
```

```
# Add the line to the plot  
plot(x, y, xlab="square footage (1000 sq feet)",  
ylab="price ($1000)")  
abline(house.reg)
```

$$\hat{\beta}_1 = r \frac{s_y}{s_x}$$

$$r_{\beta_0} = \bar{y} - \hat{\beta}_1 \bar{x}$$

Fitted regression for house price data



The estimated regression coefficients are:

$$\hat{\beta}_0 = 5.432, \hat{\beta}_1 = 56.083$$

Q: How do we interpret these coefficients? What is the predicted price of a house that is 3200 square feet?

$$\begin{aligned}\hat{y} &= \hat{\beta}_0 + \hat{\beta}_1 x \\ &= (5.432) + (56.083)(3.2) \\ &= ? \text{ (in \$1,000)}\end{aligned}$$

↑
multiply this by 1000 to get
the price in the original scale.

$$\begin{aligned}E[y|x] &= \beta_0 + \beta_1 x \\ E[y|x=x+1] &= \beta_0 + \beta_1(x+1)\end{aligned}$$

$$\begin{aligned}E[y|x=x+1] - E[y|x=x] \\ &= \beta_1\end{aligned}$$

Issue: How well does the fitted regression line describe the data?

Approach 1: Consider r^2 .

- High r^2 (and hence $|r|$) \implies points are tightly clustered around the line \implies predicted Y s are close to observed Y s \implies residuals are small \implies fit is good

Approach 2: Consider the variability in Y s explained by regression. To understand this, let's think about why the house prices are different. This is because the houses may have

- different square-footage \rightarrow in the model
 - different locations
 - different years of sale
 - other known/unknown reasons
- not in the model.*

Analysis of Variance (ANOVA)

- Total variability in Y s:

$$SS_{TOT} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (n - 1)S_y^2 \quad \text{total SS}$$

- A part of SS_{TOT} is explained by the fitted regression:

$$SS_{REG} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \quad \text{SS due to regression}$$

- The rest is error variability:

$$SS_{ERR} = SS_{TOT} - SS_{REG} = \sum_{i=1}^n e_i^2 \quad \text{error SS}$$

- ANOVA Identity: $SS_{TOT} = SS_{REG} + SS_{ERR}$.

This suggests proportion of total variation explained,

Inferior prediction:
 $R^2 = 1$ [because
 $SS_{ERR} = 0$]

$$R^2 = \frac{SS_{REG}}{SS_{TOT}}, \quad 0 \leq R^2 \leq 1$$

In general,
 $R^2 \neq S^2$,
but if there
is just one predictor (simple reg.),

as a measure of goodness of fit of the fitted regression.

- Also called **coefficient of determination**

- Between 0 and 1, with high values suggesting a good fit.

Simple linear regression ($E(Y|x) = \beta_0 + \beta_1 x$)

- $SS_{TOT} = (n - 1)S_y^2$
- $SS_{REG} = r^2(n - 1)S_y^2$
- $SS_{ERR} = (1 - r^2)(n - 1)S_y^2$
- $\underline{R^2 = r^2}$ — a reasonable measure from Approach 1 also.

Ex: For house price data: $r^2 = 0.88^2 \approx 0.77$

$\Rightarrow 77\%$ of the variability in house price is being explained by size of house.

Formulate the problem of checking if a predictor is important or not as a testing problem.

Alternative form for a regression model

Regression model: Models mean response — $E(Y|X = x)$ — as a function of x

Alternative form:
$$Y = \boxed{E(Y|X = x) + \epsilon}$$

- $E(Y|x)$ is modeled as before
- $\epsilon = Y - E(Y|X = x) = \text{error}$ — a catchall for everything that causes the observed response to differ from its mean — e.g., random variability, effect of missing predictors, etc.
- $E(\epsilon) = 0$, $\text{var}(\epsilon) = \sigma^2 = \text{var}[Y] = \sigma^2 \rightarrow$ unknown that needs to be estimated.

Model for data: $Y_i = E(Y_i|X = x_i) + \epsilon_i, i = 1, \dots, n$

Regression assumptions: The errors ϵ_i have mean zero, variance σ^2 , and are independent. No additional assumptions are needed to estimate regression coefficients by least squares.

Additional assumption: Errors follow a **normal** distribution — needed for testing hypotheses and constructing confidence intervals. This means

$$\epsilon_i \sim \text{i.i.d. } N(0, \sigma^2), \quad i = 1, \dots, n$$

Simple Linear Regression with Normality

Assumed model: $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, $\epsilon_i \sim \text{i.i.d. } N(0, \sigma^2)$,
 $i = 1, \dots, n$.

Note: The values x_1, \dots, x_n of predictor X are known and fixed (i.e., non-random), and are assumed to be measured without error.

Properties:

- $E(Y_i|x_i) =$
- $\text{var}(Y_i|x_i) =$
- $Y_i|x_i \sim \text{independent } N(\beta_0 + \beta_1 x_i, \sigma^2)$