

The fitted regression line $E[Y|X=x] = \beta_0 + \beta_1 x$

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

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

Recall:

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= r \frac{s_y}{s_x}\end{aligned}$$

implying that

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

$$|r| \leq 1.$$

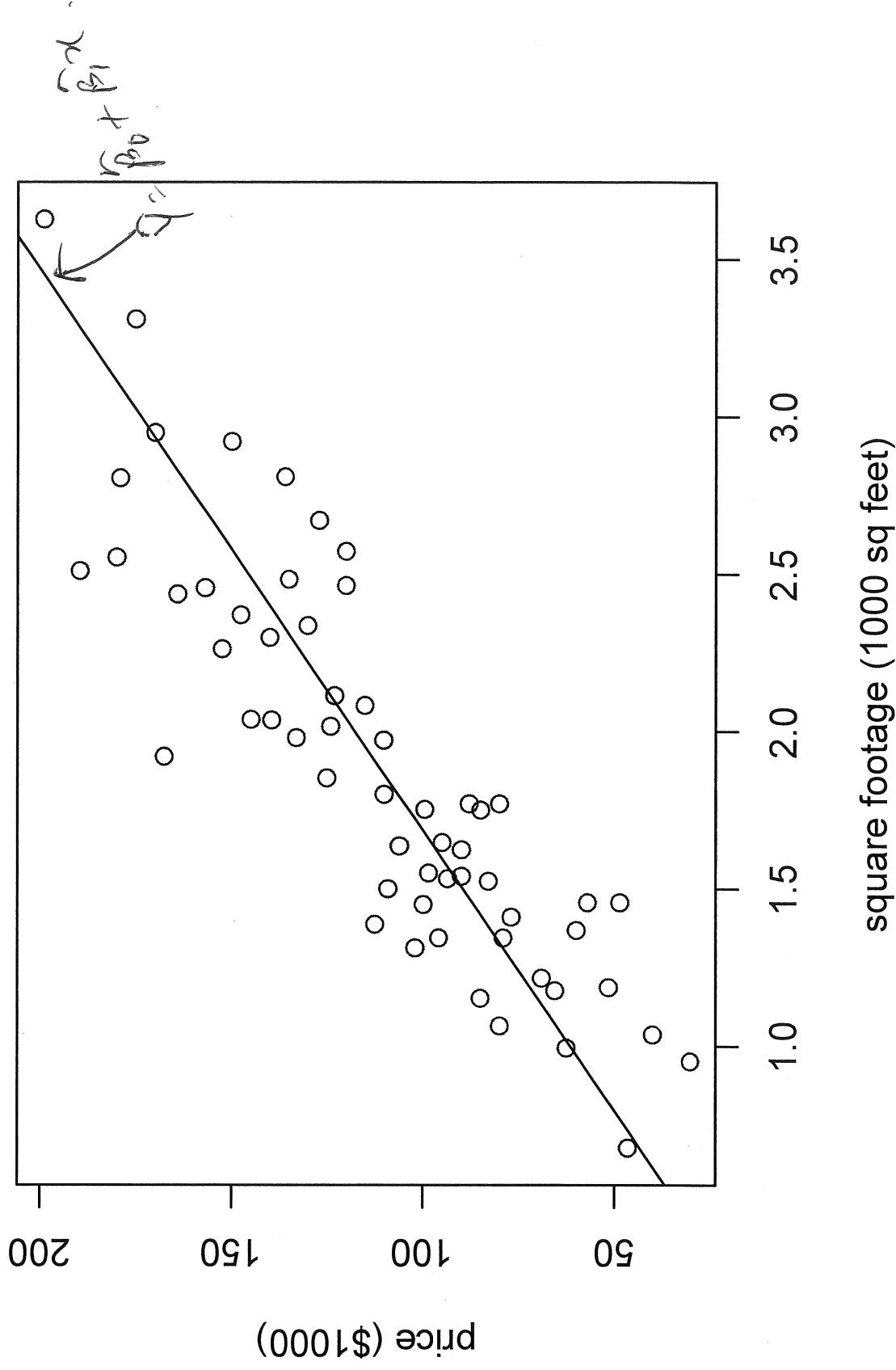
- If x is 1 SD away from its mean \bar{x} , \hat{Y} is r SD away from its mean \bar{Y} . Since $|r| \leq 1$, this means \hat{Y} is closer to \bar{Y} (in units of SD) than x is to \bar{x} — **regression toward mean**.
- The fitted line passes through the points (\bar{x}, \bar{y}) .
- The sign of slope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

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

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

$$\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y} \quad \text{verify}$$

Fitted regression for house price data



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 $\xrightarrow{\text{only this so far.}}$
- different locations $\left. \vphantom{\begin{matrix} \text{different square-footage} \\ \text{different locations} \end{matrix}} \right\} \text{predictors}$
- different years of sale
- other known/unknown reasons

Analysis of Variance (ANOVA)

- Total variability in Y s:
$$\underline{SS_{TOT}} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (n-1)S_y^2$$
 — total SS
- A part of SS_{TOT} is explained by the fitted regression:
$$SS_{REG} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$
 — SS due to regression
- The rest is error variability:
$$SS_{ERR} = SS_{TOT} - SS_{REG} = \sum_{i=1}^n e_i^2$$
 — error SS
- **ANOVA Identity:** $SS_{TOT} = SS_{REG} + SS_{ERR}$.

This suggests *proportion of total variation explained*,

$$R^2 = \frac{SS_{REG}}{SS_{TOT}}$$

[In general, $R^2 \neq r^2$]
sample
Cov.

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

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

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

Alternative form for a regression model

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

Alternative form: $Y = E(Y|X = x) + \epsilon$ — ϵ is random, $E(Y|X = x)$ is fixed.

- $E(Y|x)$ 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.
- $E(\epsilon) = 0, \text{var}(\epsilon) = \sigma^2$

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

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.

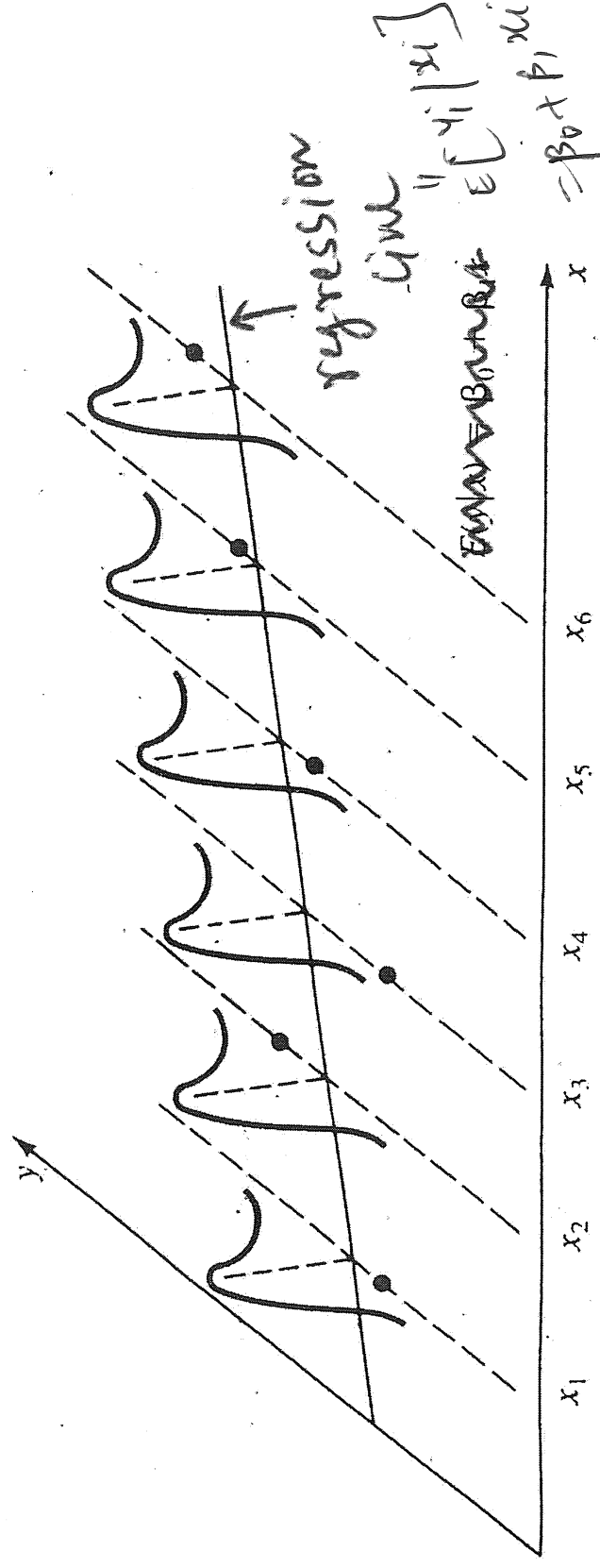
$$(\hat{\beta}_0, \hat{\beta}_1) = \text{LSE} \neq (\beta_0, \beta_1).$$

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: $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n.$

Simple linear regression model



Simple Linear Regression with Normality

Assumed model: $\underline{Y_i} = \underbrace{(\beta_0 + \beta_1 x_i)}_{\leftarrow E[\epsilon_i]} + \epsilon_i, \epsilon_i \sim \text{i.i.d. } N(0, \sigma^2),$
 $i = 1, \dots, n.$ $\uparrow \text{var}(\epsilon_i)$

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:

- $E(Y_i|x_i) = \underbrace{\beta_0 + \beta_1 x_i}_{\rightarrow \text{fixed}} + E[\epsilon_i] = \beta_0 + \beta_1 x_i$
- $\text{var}(Y_i|x_i) = \text{var}[E\epsilon_i] = \sigma^2$
- $Y_i|x_i \sim \text{independent } N(\beta_0 + \beta_1 x_i, \sigma^2) \xrightarrow{\text{NOK!}} \text{variance is constant.}$
- The least squares estimators $(\hat{\beta}_0, \hat{\beta}_1)$ of (β_0, β_1) are also maximum likelihood estimators.

• $\hat{\beta}_1 \sim N(\hat{\beta}_1, \underbrace{\sigma^2 / \{(n-1)S_x^2\}}_{\uparrow \text{var}[\hat{\beta}_1]})$

• $E[\hat{\beta}_1] = \hat{\beta}_1$ is unbiased and $\text{SE}(\hat{\beta}_1) = \frac{\sigma}{\sqrt{(n-1)S_x^2}}$

\Rightarrow

- Define: $\hat{\sigma}^2 = SS_{\text{ERR}} / (n - 2)$. Then, $E(\hat{\sigma}^2) = \sigma^2$.
 → losing 2 d.f. because estimating β_0, β_1 .
- An unbiased estimator of σ^2 is $\hat{\sigma}^2$.
- **Note:** The sample variance S_y^2 is no longer unbiased for σ^2 . This is because y_i 's are $\text{Recall: } S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$ not identically distributed anymore.
- $SS_{\text{ERR}} / \sigma^2 = (n - 2) \hat{\sigma}^2 / \sigma^2$ follows a χ^2 distribution with $(n - 2)$ degrees of freedom.

ANOVA table: A standard summary of regression fit. Here we have “simple linear regression” — i.e., two regression coefficients, β_0 and β_1 .

regression coefficients - 1 = square
coefficients → near square

Source	SS	d.f.	MS	F
Model	SS_{REG}	1	$MS_{\text{REG}} = \frac{SS_{\text{REG}}}{1}$	$\frac{MS_{\text{REG}}}{MS_{\text{ERR}}}$
Error	SS_{ERR}	$n - 2$	$MS_{\text{ERR}} = \frac{SS_{\text{ERR}}}{n-2}$	
Total	SS_{TOT}	$n - 1$		

Recall that:

- $SS_{\text{TOT}} = \sum_{i=1}^n (Y_i - \bar{Y})^2$
- $SS_{\text{REG}} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$
- $SS_{\text{ERR}} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$