

# ضروریات آمار برای علم داده

گمشدگی داده‌ها و برازش رگرسیون

## گمشدگی داده‌ها

$$D = \{A_1, A_2, \dots, A_r\}$$

$$A_j = \{A_j^{obs}, A_j^{mis}\}$$

$$D = \{D^{obs}, D^{mis}\}$$

$$R_{ij} = \begin{cases} 0 & \text{if } v_{ij} \text{ is missing} \\ 1 & \text{if } v_{ij} \text{ is observed} \end{cases}$$



## انواع گمشدگی داده‌ها

$$\Pr(R|D^{\text{mis}}, D^{\text{obs}}) = \Pr(R)$$

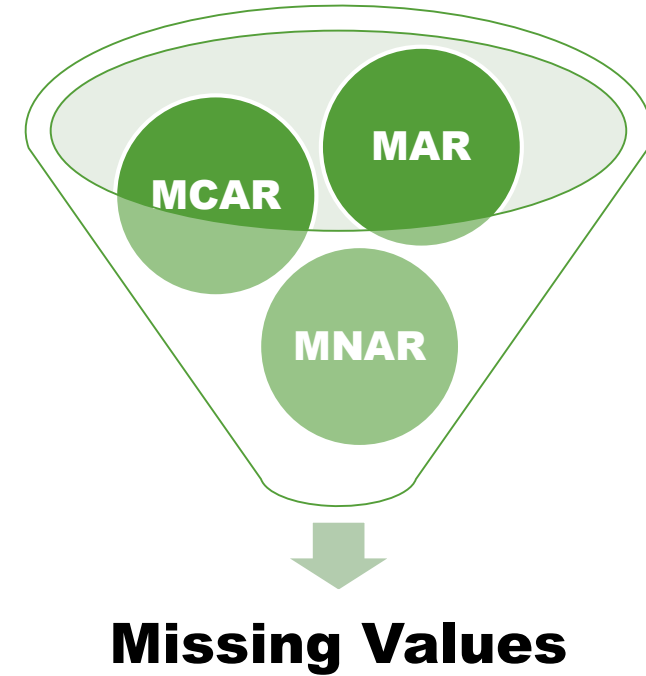
گمشدگی کاملاً تصادفی

$$\Pr(R|D^{\text{mis}}, D^{\text{obs}}) = \Pr(R|D^{\text{obs}})$$

گمشدگی تصادفی

$$\Pr(R|D^{\text{mis}}, D^{\text{obs}})$$

گمشدگی غیرتصادفی

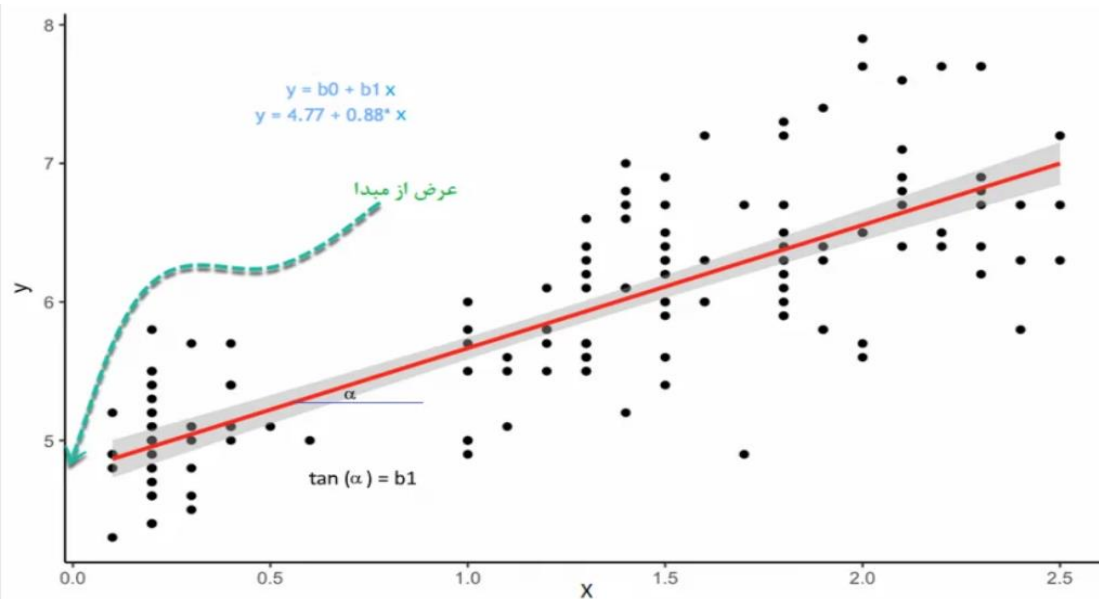


## روش‌های مواجهه با داده‌های ناکامل



- تحلیل موارد کامل
- تحلیل موارد موجود وزن دار شده
- روش‌های مبنی بر درست‌نمایی
- روش جان‌هی
- میانگین
- میانگین + عامل تصادفی
- رگرسیون

# رگرسیون خطی



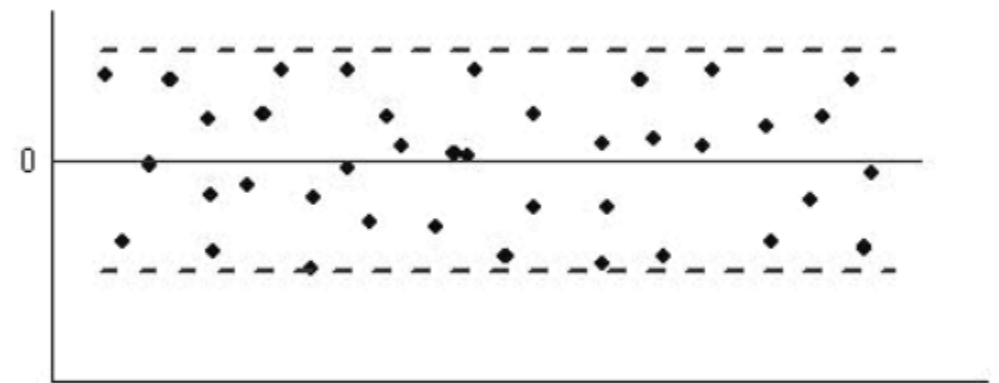
$$y = \beta_0 + \beta_1x + \epsilon$$

✓ جملات خطا از یکدیگر مستقل و دارای توزیع نرمال هستند.

✓ میانگین جمله خطا صفر است.

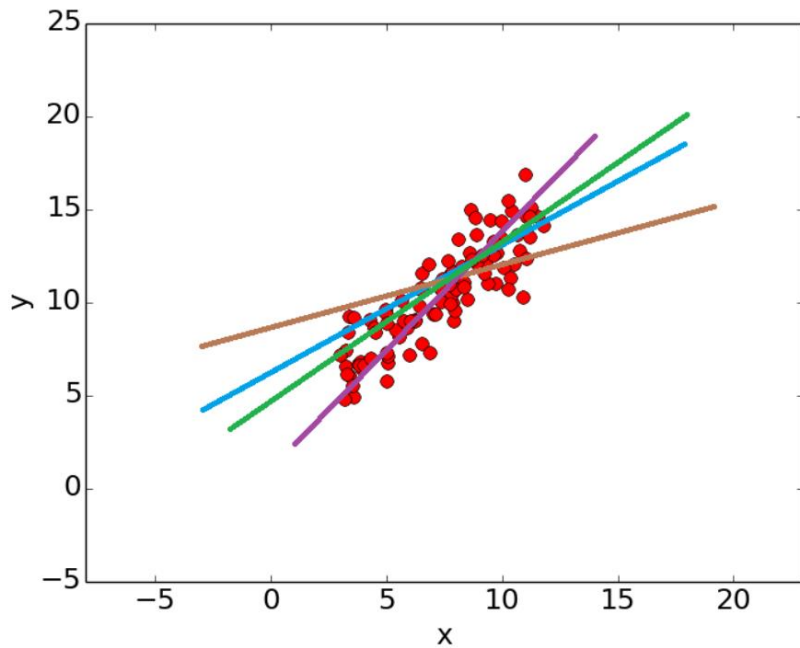
✓ واریانس هر مولفه از جمله خطا ثابت و برابر با  $\sigma^2$  است.

باقی مانده‌ها



مقدارهای پیش‌بینی شده

## برآورد پارامترهای رگرسیون خطی



$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\beta_1 = \frac{cov(x, y)}{var(x)}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

مینیم سازی مجموع مربعات خطا به منظور برآورد پارامترهای مدل

$$\sum \epsilon^2$$

$$\sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

$$\sum (-y_i + \hat{\beta}_0 + \hat{\beta}_1 x_i) = 0$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$-\sum (x_i y_i + (\bar{y} - \hat{\beta}_1 \bar{x}) \sum x_i + \hat{\beta}_1 \sum x_i^2) = 0$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

مجموعه داده ۵۰ خانه شامل قیمت (به میلیون ریال) و متراژ (متر مربع) در شهر تهران جمع‌آوری شده است.

از مجموعه داده `mtcars` که اطلاعات و ویژگی‌هایی مختلف ۳۲ خودرو در آن ثبت شده کمک می‌گیریم و مدل رگرسیونی چند گانه را برای پیش‌بینی مسافت طی شده با یک گالن سوخت ایجاد می‌کنیم. این مجموعه داده به طور خودکار در `R` بارگذاری شده است.



# رگرسیون لجستیک

$$p(x) = \hat{Y} = E(Y = 1|X = x) = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}}$$

