## RISK DIFFERENCE :

$X \rightarrow$ Exposed $\begin{cases} \text{Yes } E \\ \text{No } \sim E \end{cases}$

$Y \rightarrow$ Disease $\begin{cases} \text{Yes, } D \\ \text{No, } \sim D \end{cases}$

|  | D | $\sim$D |  |
|---|---|---|---|
| E | 30 | 70 | 100 |
| $\sim$E | 20 | 80 | 100 |
|  | 50 | 150 |  |

$$P(D|E) = \frac{30}{100} = 0.30$$

$$P(D|\sim E) = \frac{20}{100} = 0.20$$

On the
Additive Scale : Risk Difference (RD)
or
Attributable Risk (AR) $= P(D|E) - P(D|\sim E) = 0.30 - 0.20 = 0.10 = 10\%$.

Prob. that an exposure leads to 10%
inc. in D. (Extra Risk)

On Relative Scale :

Relative Risk $= \frac{P(D|E)}{P(D|\sim E)} = \frac{0.30}{0.20} = 1.5 = 15\%$

15% inc. on a relative scale
10% inc. on additive scale

$\text{odds (A)}$
$= \frac{\text{Probability of A happening}}{\text{Prob. of A not happening}}$

$$\text{Odds Ratio (OR)} = \frac{\text{Odds}(D|E)}{\text{Odds}(D|\sim E)} = \frac{P(D|E)/P(D|\sim E)}{P(D|\sim E)/P(\sim D|\sim E)}$$

$$= \frac{0.30/0.70}{0.20/0.80} = 1.71\%$$

Odds of D for E is 1.71 times of some one $\sim$E

and $1 - 1.71 = 0.71 \Rightarrow$ odds of D inc by 71%.

Odds Ratio for Case-control Study Design : $OR = \frac{\text{Odds}(D|E)}{\text{Odds}(D|\sim E)}$

|  | D | $\sim$D |  |
|---|---|---|---|
| E | a | b | a+b |
| $\sim$E | c | d | c+d |
|  | a+c | b+d | a+b+c+d |

$OR = \frac{P(D|E)/P(D|\sim E)}{P(D|\sim E)/P(\sim D|\sim E)} = \frac{\left(\frac{a}{a+b}\right)\left(\frac{b}{a+b}\right)}{\left(\frac{c}{c+d}\right)\left(\frac{d}{c+d}\right)} = \frac{ad}{bc}$

In a disease study, We select people disease
and non-disease and ask them the exposure and can't estimate
prevalance if These Probabilities. We can still conveniently use OR :-

We can estimate
prevalance of exposure $\left\{ \begin{array}{l} OR = \frac{\text{Odds}(E|D)}{\text{Odds}(E|\sim D)} = \frac{P(E|D)/P(\sim E|D)}{P(E|\sim D)/P(\sim E|\sim D)} \\ = \frac{(a/a+c)/(c/a+c)}{(b/b+d)/(d/b+d)} = \frac{ad}{bc} \end{array} \right.$

$OR \approx$ Rate Ratio (RR) — For RARE DISEASE!

(prevalence 5% or less)

## Non-Linearity In REGRESSION:

① Transform $Y$ eg. $\ln(Y)$
  - Can make linear relation
  - Can address non-constant variance

Can cause Problem in interpretability

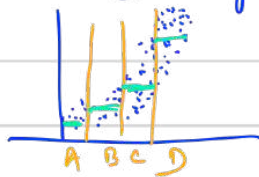② Transform $X \Rightarrow \sqrt{X}, \ln(X)$. Also
  ＊ Ladder of transformations

③ Polynomial / Quadratic fitting $\hat{y} = b_0 + b_1 X + b_2 X^2$

Again loose interpretability

(Adding powers, more inflection points)

$X^2$ |∩
$X^3$ |∿

④ Categorizing $X$



A B C D

Loss of Information

⑤ Non-linear Regression Model eg. SPLINE

Disadv. No Interpretability like coeff for $X$ feature imp

Assumptions in Linear Reg:
① Linearity
② Constant Variance (Homoscedasticity)

## $R^2 =$ coeff of Determination:

In a simple linear reg. $R^2$ is (Pearson Corr$^2$ coeff)$^2$
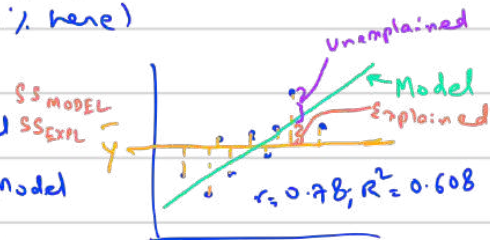
$= r$

$-1 \leq r \leq +1$    $0 \leq R^2 \leq 1$

$R^2 = \%$ age of variability in $Y$ explained by model.
(61% here)

like in ANOVA



$r = 0.78, R^2 = 0.608$

Total Var : $SS_{Total} = \sum_{all} (y_i - \bar{Y})^2$

$SS_{Total}$ — Explained by Model $SS_{EXPL}$
       — Unexplained by Model
         $\hookrightarrow SS_{UNEXPL}$ or SSERROR



Unexplained
←Model
Explained
$r = 0.78, R^2 = 0.608$
$SS_{MODEL}$

$$R^2 = \frac{SS_{ERROR}}{SS_{TOTAL}} = \frac{SS_{MODEL}}{SS_{MODEL} + SS_{ERROR}} = 1 - \frac{SS_{ERROR}}{SS_{TOTAL}}$$

Adjusted $R^2$ = $R^2$ - penalty for number of $x$'s in Model.

Used in multiple Linear Reg

to counteract some correl$^=$ b$^=$ $x$'s & $Y$.

$$Y \sim b_0 + b_1 x_1 + b_2 x_2$$

and Adj $R^2$ doesn't anymore explains % of variation as $R^2$.

NOTE: $R^2$ is generally calculated on the same data.

Use validation

## MEASURES OF VARIABILITY:

— RANGE (Max - Min.)

— IQR $Q_3 - Q_1$
Range middle 50% of data ordered
Insensitive to outliers ✓

— Sample variance ($s^2$)
$$\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}$$
If $x$ in kg
Unit = kg$^2$

※ sensitive to outliers ✓

— Sample S.D $(s) = \sqrt{s^2}$

## Monty hall PROBLEM:

Door

|  | 1 | 2. | 3 |
|---|---|---|---|
| Possibility 1 | $ | x | ↦ |
| 2 | * | $ | x |
| 3 | x | x | $ |

| Youchoose | Host | Outcome if you switch |
|---|---|---|
| 1 | 2 or 3 | Loose |
| 1 | has only 1 choice | Win |
| 1 | ''— | Win |

$\left. \right\} 2/3$

If you don't switch = Prob. of winning = 1/3