Session 2 — relationships between variables, regression

Review: random variables? $Y, X$    function/statistical model

$$Y = f(X)$$

"correlation does not equal causation"

Correlation : any statistical relationship between random variables

types of correlation: linear, nonlinear

Causal relationship : relationship between two variables where cause & effect can be established. Mathematically rigorous definition of probabilistic causality is complex: see Pearl, 1999

Intuition 1: cause proceeds effect

Intuition 2: causes are correlated with effect

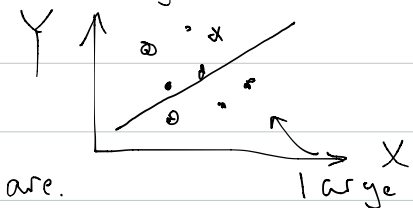Intuition 3: to most accurately evaluate cause & effect we have to experimentally manipulate cause.

Confounding : A relationship between two variables due to the presence of a third variable.

Examples of cause & effect: from paper/smoking

Discussion: is it possible to obtain data when cause & effects are not correlated even if the underlying causal relationship is REAL?

Linear Model : when the relationship at stakes can be expressed as $Y = \beta X$    $\beta$ is called "regression coefficient"

Linear regression : fits a line to a bunch of numbers. eg.



"$R^2$" is a measure of how scattered the points are. large $R^2$ also known as % variance small $R^2$ explained.
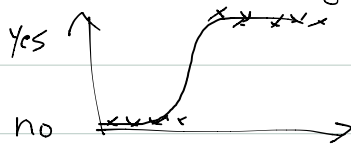
$\boxed{\text{multiple linear regression}}$ : fits a line to multiple variables

i.e. $Y = \beta_1 X_1 + \beta_2 X_2 + \dots$ etc   this is also called "adjusting"

Demo 1 : using R-statistical package to do regression.

$\boxed{\text{logistic regression}}$ : regression involving a binary variable

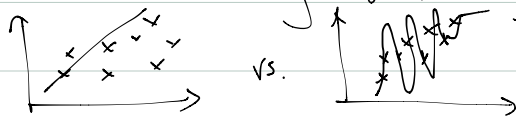odds of $(Y) = \text{logit}(\beta X) = \frac{1}{1+\exp(-\beta X)}$



p-values can be obtained for the rejecting hypothesis $\beta = 0$

can also do logistic regression on multiple variables.

$\boxed{\text{"overfitting":}}$ picking a complex model to reduce error, but with the resulting model having little validity. Example:



$\underline{R^2 \text{ increases}}$ as # of variables ↑

$\boxed{\text{Occam's razor:}}$ Simpler model is preferred.

Demo #2: adding terms to regression models.

Multiple regression addresses the problems of confounding to a certain extent, but not fully.

Final take home message: Complex, multivariate relationships often require large datasets & complex fitting procedures (i.e. Machine learning)

Exercise: in class demo of how regression works.

How to control for confounding?

(1) Stratification ⟶ divide the original data into subgroups ⟹ analyze

(2) regression based $\longrightarrow$ $Y = \beta X^{\swarrow causal} + \beta' X'^{\swarrow confounding}$

if $\beta$ is still significant $\longrightarrow$ $\beta$ is significant "adjusting for $X'$"

(3) more complex strategies : matching, etc.