

Lecture 10.

$$\ln(\text{lagrate} \sim \underbrace{\text{lag1} + \text{lag2} + \text{lag3}}_{y} + \text{dens} + \text{vax}, \dots)$$

$y = (x_1, x_2, \dots, x_p)$

$$y \sim \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p.$$

Model: $E(y|x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

a collection of distns. $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon, E(\epsilon|x)=0$

* Survey $n=500$ students

$X = \# \text{ students in favor of policy A}$.

$$\rightarrow X \sim \text{Binom}(500, 0.6), P(X > 300) = ?$$

$\rightarrow \phi = \text{fraction of all students in favor of policy A.}$

$$X \sim \text{Binom}(500, \phi), \phi \in (0, 1) \leftarrow \text{statistical model}$$

Record $X=260$

It's not one single distribution explaining the data. Instead, we have an entire collection of distns fighting against each other to explain the same data.

① There are infinite $(\beta_0, \dots, \beta_p)$ we can take

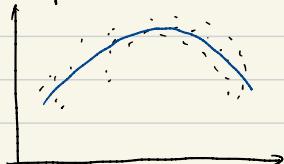
just by changing the coefficients

② Distrn of ϵ ? Candidates:

* $\epsilon \sim N(0, \sigma^2), \epsilon \perp \!\!\! \perp x$

* ϵ is cont., $E(\epsilon|x) = 0, \text{Var}(\epsilon|x) = \sigma^2$ (a huge number of distns in model space)

— Nonparametric



$E(Y|x) = \eta(x)$ restriction of models

$$Y = \eta(x) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

Each distribution is determined by $\eta(x)$ and ϵ .

$$Y_i \stackrel{iid}{\sim} N(\eta(x_i), \sigma^2),$$

$\eta \in$ Some collection of functions

$$\text{So } p(y) = \prod_{i=1}^n N(y_i | \eta(x_i), \sigma^2)$$

On the other hand, $E(Y|x) = \beta_0 + x^\top \beta, E(\epsilon|x) = 0$ is also non-parametric.

$$Y = \beta_0 + \beta^\top x + \epsilon, \epsilon \sim f(\epsilon|x) \text{ w/ } E(\epsilon|x) = 0, \text{Var}(\epsilon|x) < \infty.$$

Sample space is (β_0, β, f)

$\mathbb{R} \mathbb{R}^p \infty\text{-dim space}$ dimension is larger than any positive integer.

Concepts:

→ Data: $X \in \mathbb{R}^d$ Data can be a matrix $X = (x_1, \dots, x_n)^T$, $x_i \in \mathbb{R}^d$
Model: $X \sim P$, $P \in \mathcal{P}$

↓ model space a collection of distributions that are trying to
index by some simpler quantities explain the same data

$P(\cdot | \theta)$, $\theta \in \Theta \rightarrow$ parameter space

$$\text{Binom}(500, \phi), \phi \in (0, 1)$$

$$\bigotimes_{i=1}^n N(y(x_i), \sigma^2), y \in \text{function space}, \sigma > 0$$

$$y_i \stackrel{\text{iid}}{\sim} N(y(x_i), \sigma^2)$$

$\dim(\Theta) =$ finite or infinite

→ Three goals of Data analysis:

① Hypothesis testing:

$$H_0: P \in P_0 \subset \mathcal{P}$$

x_1, \dots, x_{22} (self-reported weekly expense on food)

Assume $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$

$$P = \bigotimes_{i=1}^n N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma > 0.$$

" $H_0: \mu \geq 175$ "

$$P_0 = \left\{ \bigotimes_{i=1}^n N(\mu, \sigma^2), \mu \geq 175, \sigma > 0 \right\}$$

* This is not exactly a formal hypothesis testing.

* This is the earlier version: want to prove it wrong

Formal Hypothesis, $P = P_0 \cup P_1$, $P_0 \cap P_1 = \emptyset$

$$H_0: P \in P_0, H_1: P \in P_1$$

② Prediction

Data: X $(x, x^*) \sim P$, $p \in \mathcal{P}$

Future obs: x^*

$$x^* \sim P(\cdot | x)$$

Typical

Clinical: Control $x_1, \dots, x_n, x^* \stackrel{\text{iid}}{\sim} N(\mu_1, \sigma^2)$] $H_0: \mu_1 = \mu_2$

Treatment $y_1, \dots, y_m, y^* \stackrel{\text{iid}}{\sim} N(\mu_2, \sigma^2)$] $H_1: \mu_1 \neq \mu_2$

$$y^* - x^*, P(Y^* > x^*)$$

real interests

③ Estimation

$$\left. \begin{array}{l} X_1, \dots, X_n, X_1^*, \dots, X_N^* \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \\ Y_1, \dots, Y_n, Y_1^*, \dots, Y_N^* \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \end{array} \right\} \quad \bar{Y}^* - \bar{X}^* \approx u_2 - u_1$$

Est = Large scale avg prediction.
↓
(| future obs)