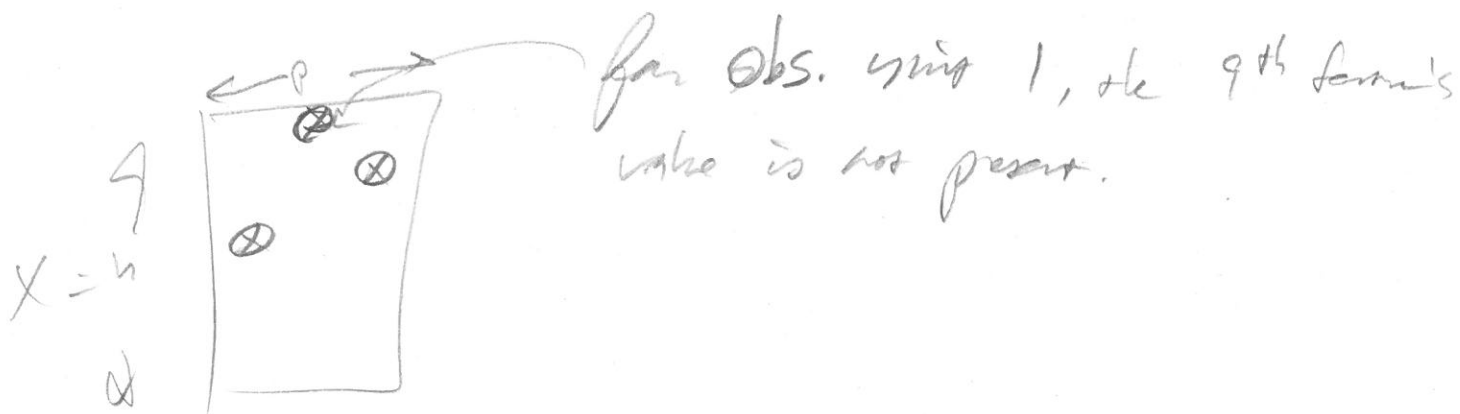


Math 370.8 Lec 25 5/10/18

1

# Missingness - source of data analysis



You cannot use any of the models in this class with this. Let  $M_j$  denote <sup>binary</sup> r.v. that  $j$ th female is missing

Types of missingness: missing data mechanisms, MDM

| MDM  | $P(M_j   X_{j,miss}, X_{j,miss}, X_{-j,obs}, V, \Sigma)$ | obs. vars<br>order pres |
|------|----------------------------------------------------------|-------------------------|
| MCAR | $P(M_j   X)$                                             |                         |
| MAR  | $P(M_j   X_{-j,miss}, X_{j,obs})$                        |                         |
| NMAR | does not simplify                                        |                         |

Ex of each?

MCAR: data corruption

MAR: old person surveys

NMAR: Chicago school starting GPA, age }  $X = \{$

How to handle?

(I) Listwise deletion: drop all units. Bad for  
(a) pred. performance

(b) generalization... any predictions will be ~~impossible~~

(II) Imputation: fill in holes with values. "X<sub>imp</sub>"  
(a) use  $\bar{X}_j$  for any value  $X_j$  missing. Not great!

(b) fit a model to  $X_j$  as outcome and predict, much better!  
Rec: MissForest: alg. which iteratively runs RF on  
all features filling in holes until "convergence".

$Y = f(x_1, \dots, x_p) + \epsilon$  but it also may be  $\epsilon$  being noisy for  
 $Y = f(x_1, \dots, x_p, m_1, \dots, m_p) + \epsilon$  where missingness itself matters.  
eg. Chicago public school.

Rec: create dummies  $m_1, \dots, m_p$ , then impute them from  
final design matrix is

$$X_{\text{final}} = [X_{\text{imp}}, m_1, m_2, \dots, m_p]$$

Even if  $X_j$  is known,  $2p$  dummies.  
you could just add a level for  $m_j$  and not impute...

# Credit Card Decision

# Type of Supervised Learning

How much will individual pay back?  
 will the individual fully pay back or not?  
 what is the prob (↓) ?

Regression  
 Binary Classification  
 prob. estimation / prob. class.

primarily if  $y = \{0, 1\}$

$$y = t(x_1, \dots, x_p) = \underbrace{f(x_1, \dots, x_p)}_{\in \{0, 1\}} + \underbrace{\epsilon}_{\in \{-1, 0, 1\}} = \underbrace{h^*(x_1, \dots, x_p)}_{\text{ditto}} + \underbrace{\epsilon}_{\text{ditto}} = \underbrace{g(x_1, \dots, x_p)}_{\text{ditto}} + \underbrace{\epsilon}_{\text{ditto}}$$

$t, f, h^*, g$  are all outputs 0, 1 ; not probs!

Consider an alternate construction if  $f, h^*, g$  but not  $t$ !  
 Bernoullis with different  $\theta$ 's Let  $Y_1, \dots, Y_n$  be i.i.d.

$Y_i \sim \text{Bernoulli} \left( \underbrace{f(x_1, \dots, x_p)}_{\text{Prob}(Y_i=1 | \vec{x})} \right)$  the best you can do with the data at hand

$Y_i \sim \text{Bernoulli} (h_{pr}^*(x_1, \dots, x_p))$  is a noise model

$Y_i \sim \text{Bernoulli} (g_{pr}(x_1, \dots, x_p))$  is even noisier

Is there a tpr? Yes tpr = t since if all candidates known,

$\Rightarrow$  no randomness and  $\text{prob}(t | \vec{x}) = y$  itself

$t - f_{pr}$  : error due to ignorance

$f_{pr} - h_{pr}^*$  : error due to noise

$h_{pr}^* - g_{pr}$  : regression error

prob = 100%  $\Rightarrow$  prob = 90%

What is  $\text{spr}, 0$ ? The null model?  $\bar{y} = \frac{\sum y_i}{n} = \hat{p}$  (4)

As the model gets worse ... it moves towards  $\text{spr}, 0$ .

How to fit and get  $\text{spr}$ ?

Assume that each  $(\vec{x}_i, y_i) \in D$  are independent.

then 
$$P(D) = \underbrace{P(\vec{x}_1, y_1)}_{\text{from } f_{\text{pr}}(\vec{x}_1)} \underbrace{P(\vec{x}_2, y_2)}_{\text{from } f_{\text{pr}}(\vec{x}_2)} \dots \underbrace{P(\vec{x}_n, y_n)}_{\text{from } f_{\text{pr}}(\vec{x}_n)}$$
$$= f_{\text{pr}}(\vec{x}_1)^{y_1} (1 - f_{\text{pr}}(\vec{x}_1))^{1-y_1} \dots (1 - f_{\text{pr}}(\vec{x}_n))^{1-y_n}$$
$$= \prod_{i=1}^n f_{\text{pr}}(\vec{x}_i)^{y_i} (1 - f_{\text{pr}}(\vec{x}_i))^{1-y_i}$$

Now we want to find  $f_{\text{pr}}$  s.t.  $P(D)$  is maximal.  
Of course  $f_{\text{pr}}$  is very complicated. So... we make an assumption of  $\mathcal{H}_{\text{pr}}$ .  $\mathcal{H}_{\text{pr}}$  is a family? we need functions  $f_{\text{pr}}: \mathbb{R}^p \rightarrow [0, 1]$ . How about  $\{\vec{w} \cdot \vec{x} > 0, \vec{w} \in \mathbb{R}^{p+1}\}$

Not good! Only returns  $\{0, 1\}$ . Prob. of 50%, 59%, 93% not possible.  
What about  $\mathcal{H} = \{\vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{p+1}\}$ ?  $\vec{w} \cdot \vec{x} \in \mathbb{R}$ . We can return prob's of 107%! Or -30%....

What if we wanted to return the linear model but output is  $\in (0, 1)$ ?

Even the link function:  $\mathcal{H} = \{\phi(\vec{w} \cdot \vec{x}) : \vec{w} \in \mathbb{R}^p\}$  (5)

$\phi: \mathbb{R} \rightarrow (0,1)$  and monotonic and  $\phi' > 0$  always!  
if  $\vec{w} \cdot \vec{x} \uparrow \Rightarrow P(Y=1|\vec{x}) \uparrow$  Also  $P(Y=1|\vec{x}) \neq 1$   
 $\neq 0$

The link function we will study is the logistic link:  
why? Can even be completely same!

$$\phi(u) = \frac{e^u}{1+e^u} = \frac{1}{1+e^{-u}}$$

F.I. Another one is probit

Using this is logistic regression

$\phi(u) = F^{-1}(u)$  inverse CDF of standard normal

F.I. or complementary log-log:

$$\phi(u) = 1 - e^{-e^u}$$

$$\mathcal{H} = \left\{ \frac{e^{\vec{w} \cdot \vec{x}}}{1+e^{\vec{w} \cdot \vec{x}}} : \vec{w} \in \mathbb{R}^p \right\}$$

$$\Rightarrow \vec{b} = \underset{\vec{w} \in \mathbb{R}^p}{\text{argmax}} \left\{ \prod_{i=1}^n \left( \frac{e^{\vec{w} \cdot \vec{x}_i}}{1+e^{\vec{w} \cdot \vec{x}_i}} \right)^{y_i} \left( 1 - \frac{e^{\vec{w} \cdot \vec{x}_i}}{1+e^{\vec{w} \cdot \vec{x}_i}} \right)^{1-y_i} \right\}$$

$$\prod_{i=1}^n \left( \frac{1}{1+e^{-\vec{w} \cdot \vec{x}_i}} \right)^{y_i} \left( \frac{1}{1+e^{\vec{w} \cdot \vec{x}_i}} \right)^{1-y_i}$$

We cannot take  $\frac{d\mathcal{L}}{d\vec{w}} \stackrel{!}{=} 0$  and solve for  $\vec{w}$  we need to use opt. alg. e.g. gnt, descent, just the SVM.

Using numerical methods, we get  $\vec{b}$ .

$$\Rightarrow g_{pr}(\vec{x}) = \frac{e^{\vec{b} \cdot \vec{x}}}{1 + e^{\vec{b} \cdot \vec{x}}} = \text{prob cat's } \hat{P}(Y=1|\vec{x})$$

What does  $\vec{b} \cdot \vec{x}$  mean?

Let's see...

$$\hat{p} = \phi(\vec{b} \cdot \vec{x}) = \frac{1}{1 + e^{-\vec{b} \cdot \vec{x}}}$$

$$\Rightarrow \frac{1}{\hat{p}} = 1 + e^{-\vec{b} \cdot \vec{x}} \Rightarrow \frac{1}{\hat{p}} - 1 = e^{-\vec{b} \cdot \vec{x}}$$

$$\Rightarrow \frac{1-\hat{p}}{\hat{p}} = e^{-\vec{b} \cdot \vec{x}} \Rightarrow \frac{\hat{p}}{1-\hat{p}} = e^{\vec{b} \cdot \vec{x}} \Rightarrow \vec{b} \cdot \vec{x} = \ln\left(\frac{\hat{p}}{1-\hat{p}}\right)$$

Interpret  $b_j$ ! Same as last class except change "y"  
 increase by  $b_j$  to "log odds of  $P(Y=1)$  increase by  $b_j$ "

Validating  $g_{pr}$ . What is best prob function?  $g_{pr}$ .  
 So let's validate against other functions! Unknown!

The only thing you have is  $y_1, \dots, y_n \in Y = \{0, 1\}$