

this descriptive  
time: methods  
next  
time: ↓

Season	soft drink	polio
F	M	M
W	L	L
Sp	M	M
Su	H	H

AMS7  
5 Apr  
14

①



H high  
L low  
M medium

strong  
positive  
association

but

doesn't prove  
causality

subjects (beer)  
↓

pop.

variables

↑  
N = population size  
↓

population  
all deer living  
at USC on 2/1/2017

1 row  
for  
each  
deer

sample  
the observed  
deer

idea: 2  
 $\hat{p} = \bar{y} = \bar{Y}$  is  
a good  
estimate of  
 $p = \theta$

$1 = Y$   
 $0 = N$   
  
 $N = 500$

disease?  
no/N/0  
no/N/0  
no/N/0  
yes/Y/1  
no/N/0  
:  
no/N/0

at  
random  
  
IID

disease?  
1  
2  
0  
5  
n = 85  
 ~~$\frac{3}{85} = 3.5\%$~~   
estimate

mean  $\bar{y} = \hat{p} = \hat{\theta}$  Sample size  
 $\bar{y}$ -hat  $\hat{p}$ -hat  $\hat{\theta}$ -hat

mean  $p = \theta = ?$   
proportion  
in pop.  
with  
disease  
unknown

goal of sampling:

representativeness:

mean of 1s & 0s  
keeps track  
automatically  
of % of 1s

without sample, unsample  
similar in all relevant  
ways  
population  
unsample } sample

at random with replacement: IID ← independent + identically distributed sampling

at random without replacement: SRS simple random sampling

(SRS) is more informative than (3)  
(IID), but (IID) has easier math

---

if  $n$  is a lot smaller than  $N$ ,

(SRS)  $\approx$  (IID)



is approximately  
equal to

---