

Adaptive Cluster Sampling

Math 255, St. Clair

1 / 28

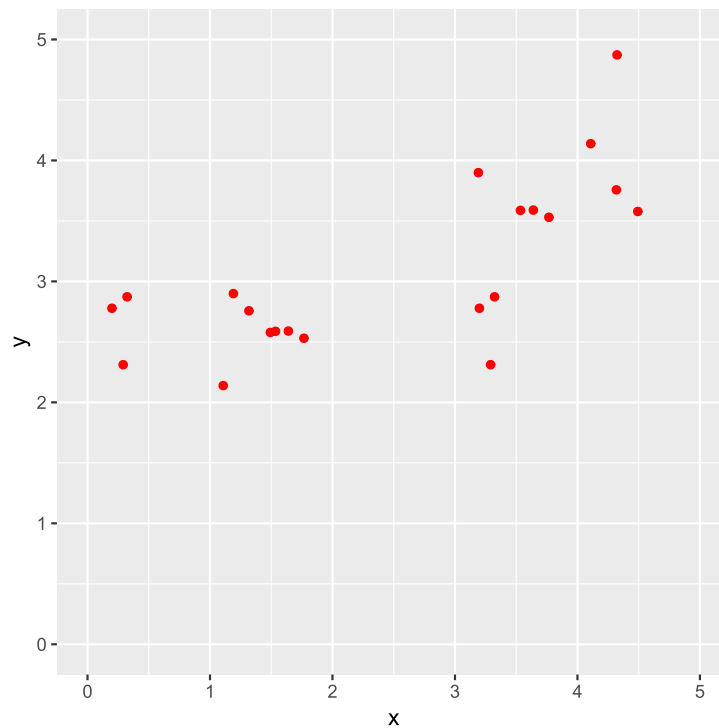
Adaptive designs

- *Adaptive* design: sampling units are determined "on the fly" based on observed characteristics of previously sampled units
- Adaptive *cluster* sampling (ACS): goal is to sample rare, clustered populations
- *cluster* denotes a spatial, social or genetic "closeness"
 - rare animal/plant species that are spatially clustered
 - rare disease/trait that are spatially *or* socially or genetically clustered
- Idea of ACS:
 - (1) get an initial SRS of units
 - (2) if sampled unit i meet a condition C and add i 's neighboring units to the sample
 - (3) repeat (2) until no more units can be added

2 / 28

Example

A simplified "strawberry field" example: how many plants (dots) in the field?



3 / 28

Example

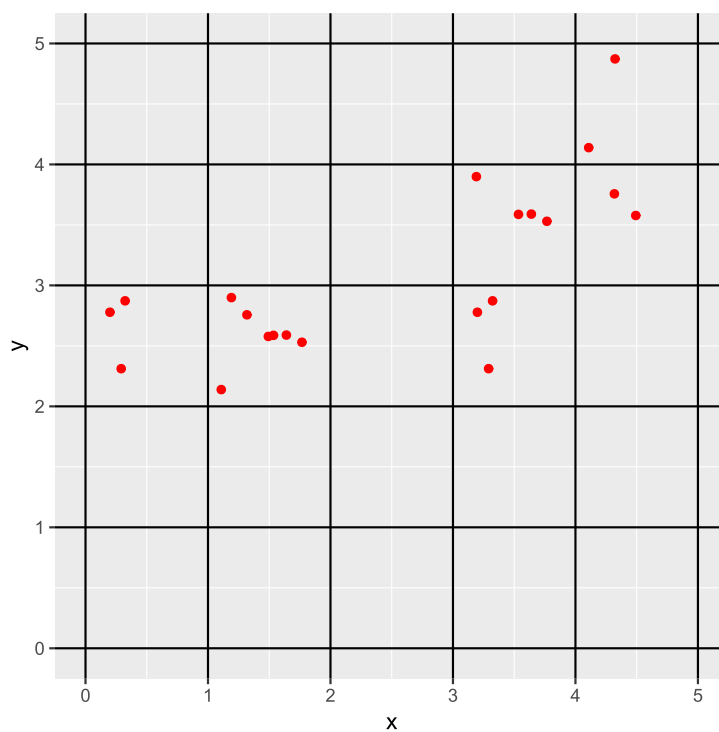
ACS design:

- Divide the region into grid plots to create a sampling frame
- Sampling unit: grid plot (N=25)

4 / 28

Example

Sampling frame grid:



5 / 28

Example

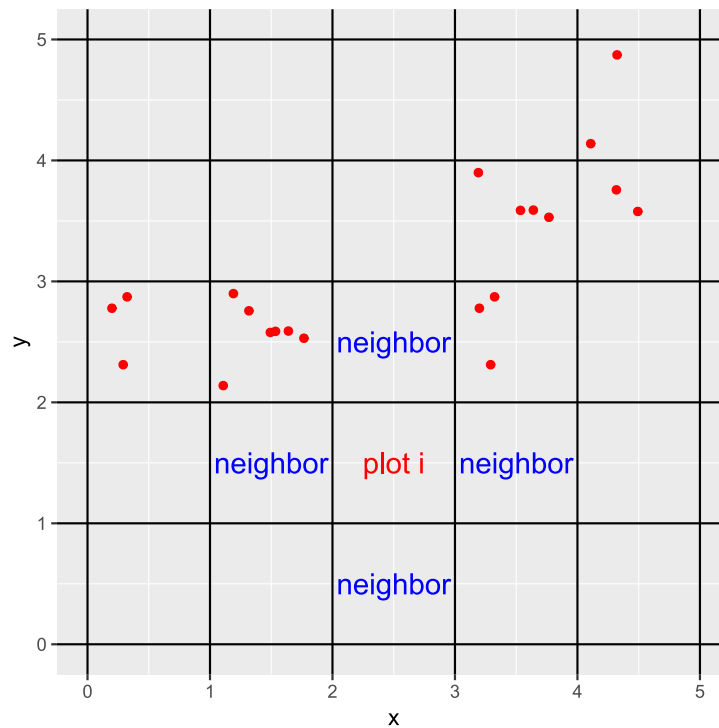
ACS design:

- define a neighborhood:
 - plot i 's neighbors are cells to the north/south/east/west

6 / 28

Example

Neighborhood of plot i



7 / 28

Example

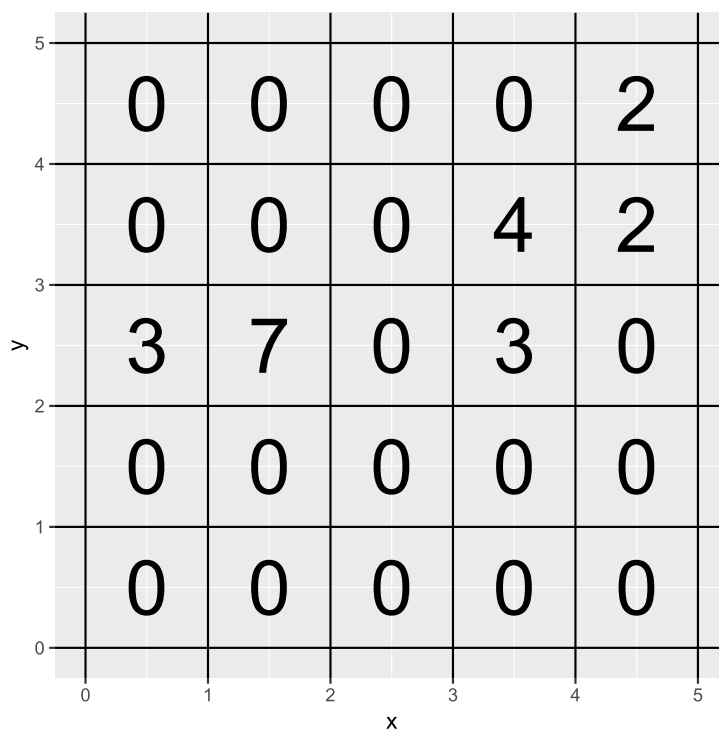
ACS design:

- Determine condition C : we want to find the plots with plants!
 - y_i = number of plants in plot i
 - $C : y_i > 0$, add neighbors if plot i contains plants
- All that matters in the population is units, neighborhoods and y_i values

8 / 28

Example

Let's just look at y_i s



9 / 28

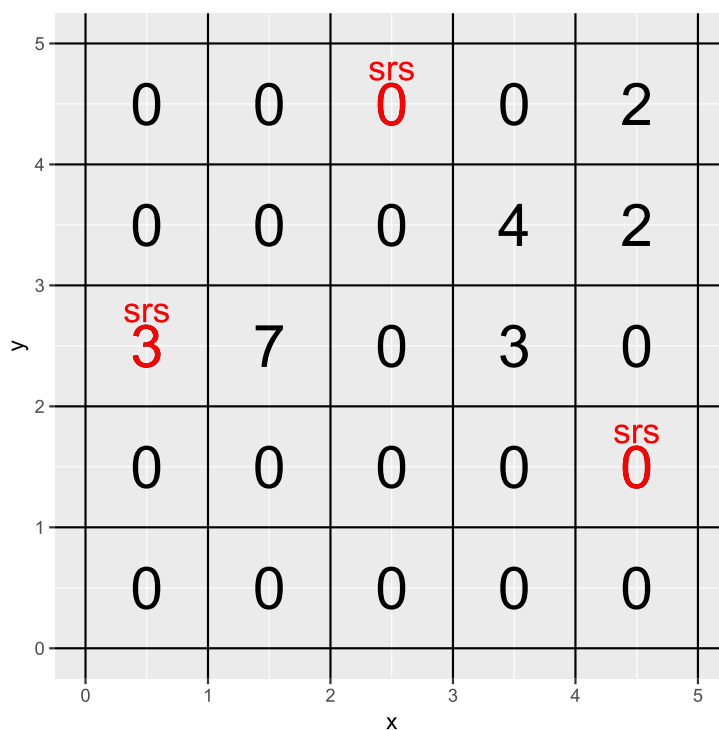
Example

ACS design:

- (1) Initial SRS of $n_1 = 3$ plots and count plants y_i
- (2) **adaptively add units:** If $y_i > 0$, add plot i 's neighbors
- (3) Repeat (2) until no more neighbors are adaptively added
 - all neighbors have $y_i = 0$

Example

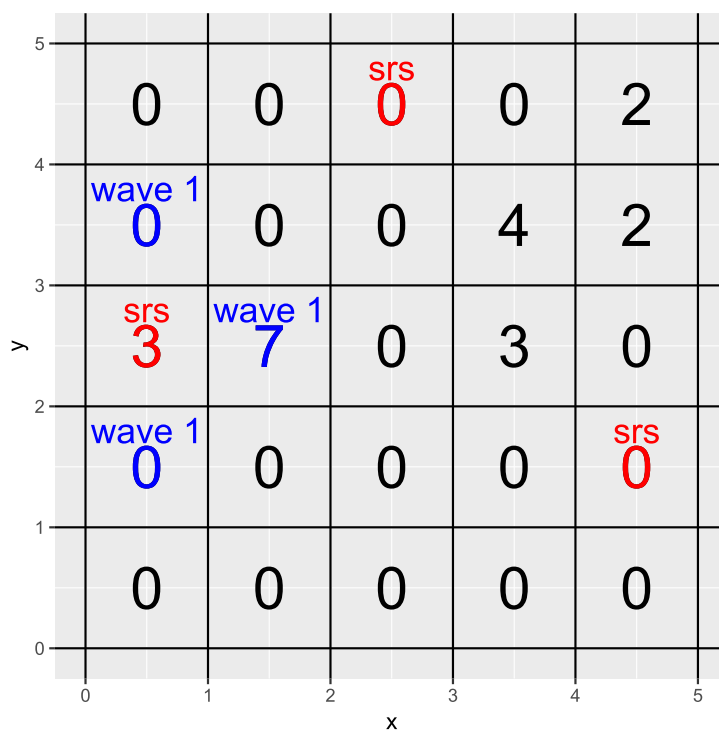
(1) Initial SRS of size 3 is highlighted



11 / 28

Example

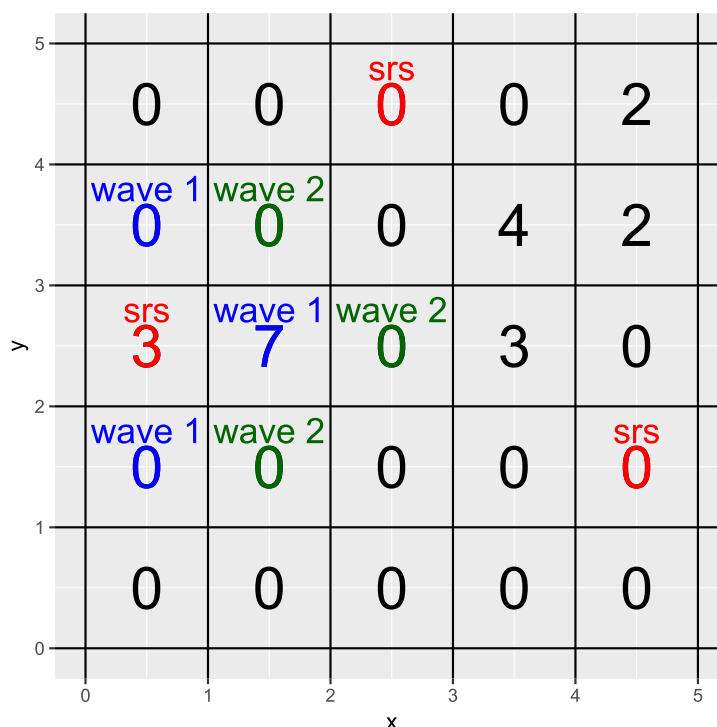
(2) add first round of neighbors:



12 / 28

Example

(2) add second (and final!) round of neighbors:



13 / 28

Estimation

- Units have unequal selection probabilities
 - need to use a Horvitz-Thompson estimate of total!
 - units with $y_i > 0$ have higher inclusion probabilities
- Inclusion probability for unit i looks like

$$\pi_i = P(\text{unit } i \text{ is in the SRS or adaptively added})$$

Problem: unless we see the *entire population*, we can't compute all π_i for observed units

- can't tell if unit i borders a cluster unless we've seen all units around it

14 / 28

Networks instead of plots

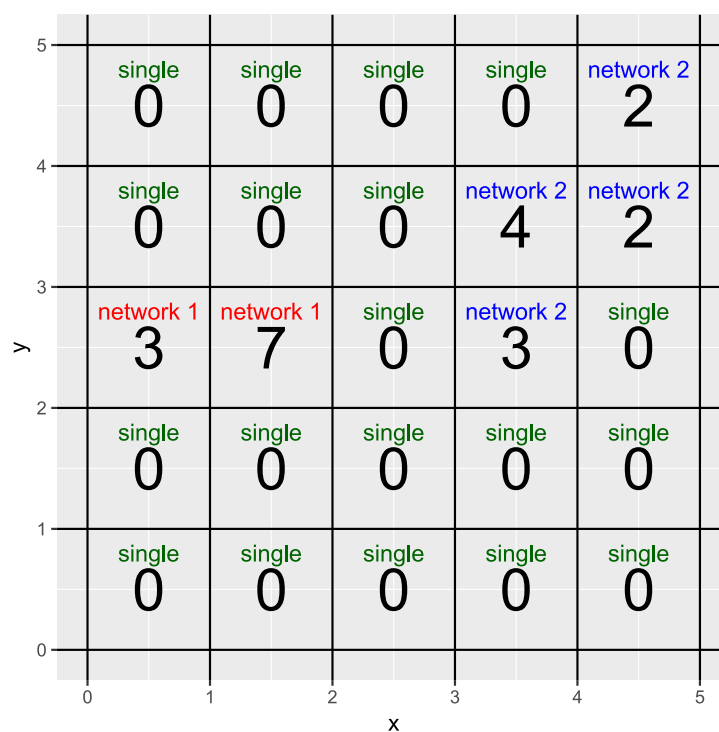
Solution: define observations in terms of *networks*

- **Network:** a cluster of units generated by selection of any of the units within the cluster
- networks either
 - contain units that all satisfy condition C
 - are a single unit where C is not satisfied

15 / 28

Example

Two networks satisfy $y_i > 0$, 19 other networks are single plot with $y_i = 0$.



16 / 28

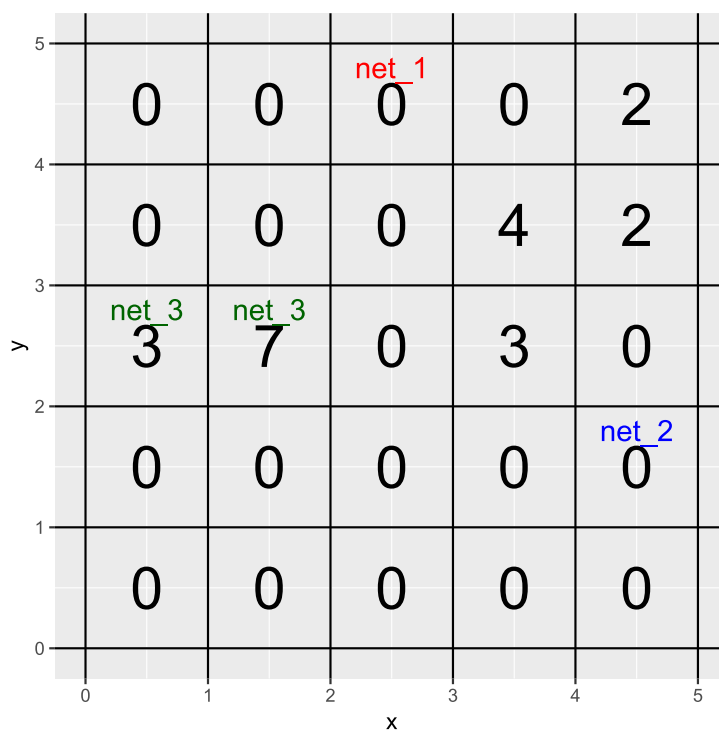
Networks

- Population has K distinct networks
 - Ex: $K = 19 + 2 = 21$
- Sample has κ distinct networks
 - Ex: $\kappa = 3$

17 / 28

Example

Three networks were sampled



18 / 28

Networks

- Let y_k^* be the total response of all units in network k

$$y_k^* = \sum_{i \in \text{net}_k} y_i$$

- We still have the same population total:

$$t = \sum_{i=1}^N y_i = \sum_{k=1}^K y_k^*$$

- Use HT estimator to weight the observed network totals y_k^*

19 / 28

Network inclusion probabilities

- α_k is the inclusion probability for network k :
 - network k is included in the ACS if **at least one** of its units is in the **initial SRS**
- Compute α_k using the complement rule:

$$\alpha_k = 1 - P(\text{no units in } k \text{ are in the initial SRS})$$

20 / 28

Network inclusion probabilities

- There are $\binom{N}{n_1}$ possible SRS of size n_1
- x_k = number of units in network k
- There are $\binom{N-x_k}{n_1}$ possible SRS of size n_1 that **don't contain any units** in network K

$$\alpha_k = 1 - P(\text{no units in } k \text{ are in the initial SRS}) = 1 - \frac{\binom{N-x_k}{n_1}}{\binom{N}{n_1}}$$

21 / 28

Example

- Three sampled networks:
- $net_1 = \{1\}, y_1^* = 0, x_1 = 1, \alpha_1 = 1 - \frac{\binom{25-1}{3}}{\binom{25}{3}} = 0.12$
- $net_2 = \{2\}, y_2^* = 0, x_2 = 1, \alpha_2 = 1 - \frac{\binom{25-1}{3}}{\binom{25}{3}} = 0.12$
- $net_3 = \{3, 4\}, y_3^* = 10, x_1 = 2, \alpha_3 = 1 - \frac{\binom{25-2}{3}}{\binom{25}{3}} = 0.23$
- Estimated total:

$$\hat{t}_{HT} = \sum_{k=1}^{\kappa} \frac{y_k^*}{\alpha_k} = \frac{0}{0.12} + \frac{0}{0.12} + \frac{10}{0.23} = 43.48$$

22 / 28

Network inclusion probabilities

- Joint inclusion probability that both networks j and k are in the ACS
 - Use the rule: $P(j \text{ or } k) = P(j) + P(k) - P(j \text{ and } k)$
- So the probability of j **and** k is

$$\begin{aligned}\alpha_{jk} &= \alpha_j + \alpha_k - P(j \text{ or } k \text{ in ACS}) \\ &= \alpha_j + \alpha_k - (1 - P(\text{neither } j, k \text{ in ACS})) \\ &= \alpha_j + \alpha_k - \left(1 - \frac{\binom{N-x_j-x_k}{n_1}}{\binom{N}{n_1}}\right)\end{aligned}$$

23 / 28

Example

- Joint prob for networks 1 and 2:

$$\alpha_{12} = 0.12 + 0.12 - \left(1 - \frac{\binom{25-1-1}{3}}{\binom{25}{3}}\right) = 0.01$$

- Joint prob for networks 1 and 3, and also 2 and 3:

$$\alpha_{13} = \alpha_{23} = 0.12 + 0.23 - \left(1 - \frac{\binom{25-1-2}{3}}{\binom{25}{3}}\right) = 0.01957$$

- SE for \hat{t}_{HT} : only need to sum over non-zero network responses (and all joint products are 0!)

$$SE_{HT}(\hat{t}_{HT}) = \sqrt{\frac{1 - 0.23}{0.23^2} 10^2 + 2(0)} = 38.15$$

24 / 28

Example

- To estimate in R, enter network level data: y_k^* and x_k

```
> acs_data <- data.frame(  
+   y_net = c(0,0,10),  
+   x_net = c(1,1,2) )
```

- Then get single network inclusion probabilities:

```
> n1 <- 3  
> N <- 25  
> acs_data$pi_single <- 1- choose(N - acs_data$x_net,n1)/choose(N,n1)  
> acs_data  
  y_net x_net pi_single  
1     0     1     0.12  
2     0     1     0.12  
3    10     2     0.23
```

25 / 28

Example

- Joint inclusion probabilities take more work
 - `jnt_fun` computes α_{jk} for all $k = 1, \dots, \kappa$

```
> jnt_fun <- function(xj,x=acs_data$x_net,N=25,n1=3)  
+ { 1- choose(N - xj,n1)/choose(N,n1) - choose(N - x,n1)/choose(N,n1) +  
+   choose(N - xj-x,n1)/choose(N,n1)}  
> jnt_fun(xj = 1)  
[1] 0.01000000 0.01000000 0.01956522  
> jnt_fun(xj = 2)  
[1] 0.01956522 0.01956522 0.03826087
```

26 / 28

Example

- Fill the rows of the inclusion matrix:

```
> jnt_mat <- matrix(
+   c(jnt_fun(acs_data$x_net[1]),
+     jnt_fun(acs_data$x_net[2]),
+     jnt_fun(acs_data$x_net[3])),
+   byrow=TRUE, nrow=3)
> diag(jnt_mat) <- acs_data$pi_single # fix diagonals
> jnt_mat
      [,1]      [,2]      [,3]
[1,] 0.12000000 0.01000000 0.01956522
[2,] 0.01000000 0.12000000 0.01956522
[3,] 0.01956522 0.01956522 0.23000000
```

27 / 28

Example

- Then use "pps" design:

```
> library(survey)
> acs_design <- svydesign(id = ~1, fpc= ~pi_single,
+                       pps=ppsmat(jnt_mat), data=acs_data)
> svytotal(~y_net, acs_design)
      total      SE
y_net 43.478 38.152
```

- Again, get $\hat{t}_{HT} = 43.48$ and SE of 38.15.
- Note: Unless n is very large and clusters not "too clustered", you can't trust conventional confidence intervals for ASC data!

28 / 28